

# Graph Representation Learning:

*Embedding, GNNs, and Pre-Training*

**Yuxiao Dong**

<https://ericdongyx.github.io/>

**Microsoft Research, Redmond**

# Joint Work with



**Jiezhong Qiu**  
Tsinghua  
(Jie Tang)



**Ziniu Hu**  
UCLA  
(Yizhou Sun)



**Hongxia Yang**  
Alibaba



**Jing Zhang**  
Renmin U. of China



**Jie Tang**  
Tsinghua



**Yizhou Sun**  
UCLA

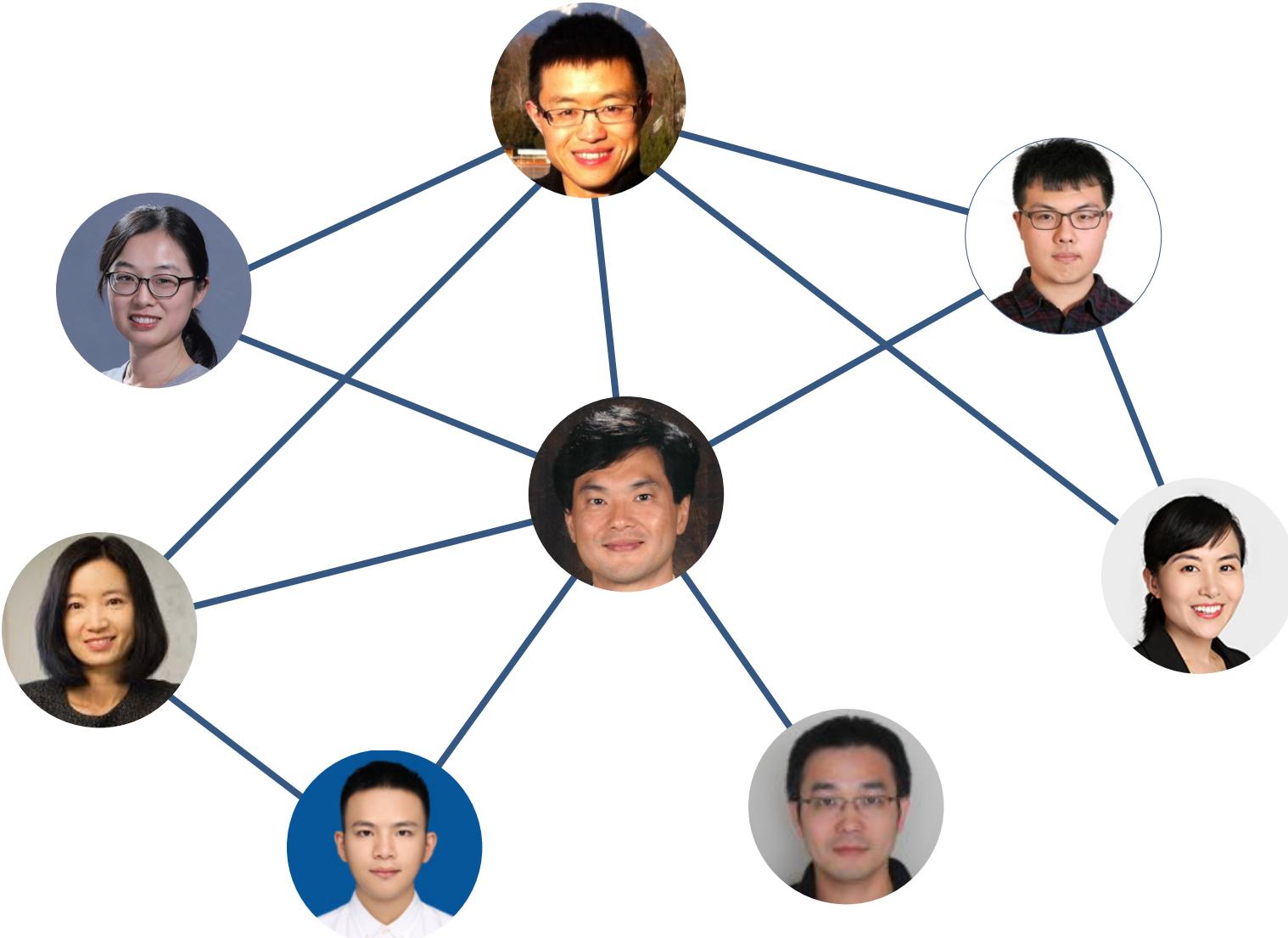


**Hao Ma**  
Facebook AI



**Kuansan Wang**  
Microsoft Research

# Why Graphs?



# Graphs



Academic Graph



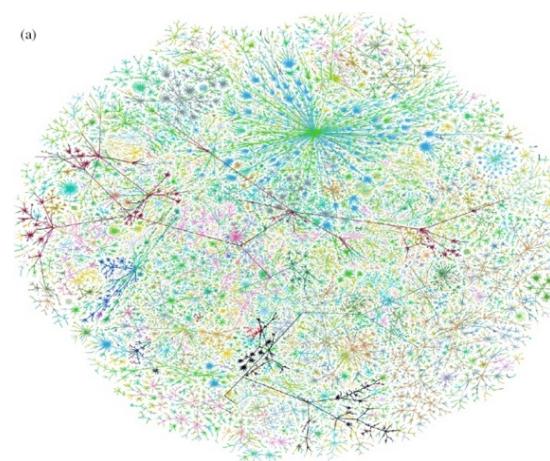
Office/Social Graph



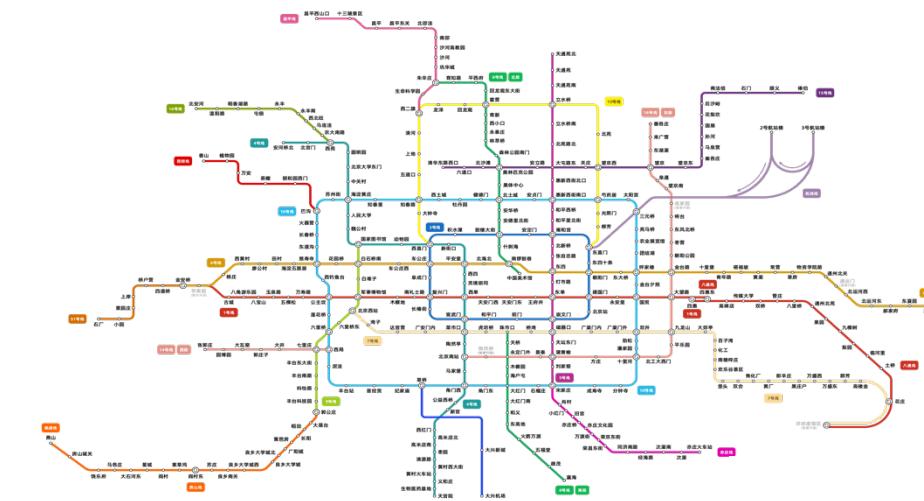
Biological Neural Networks



Knowledge Graph



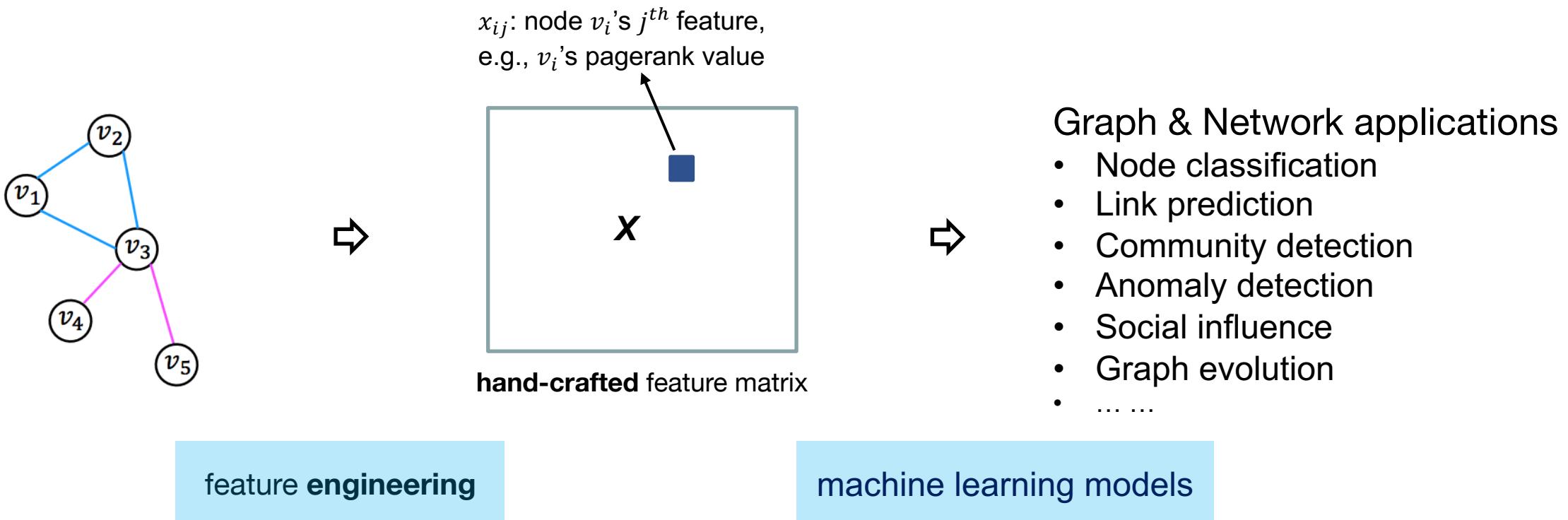
Internet



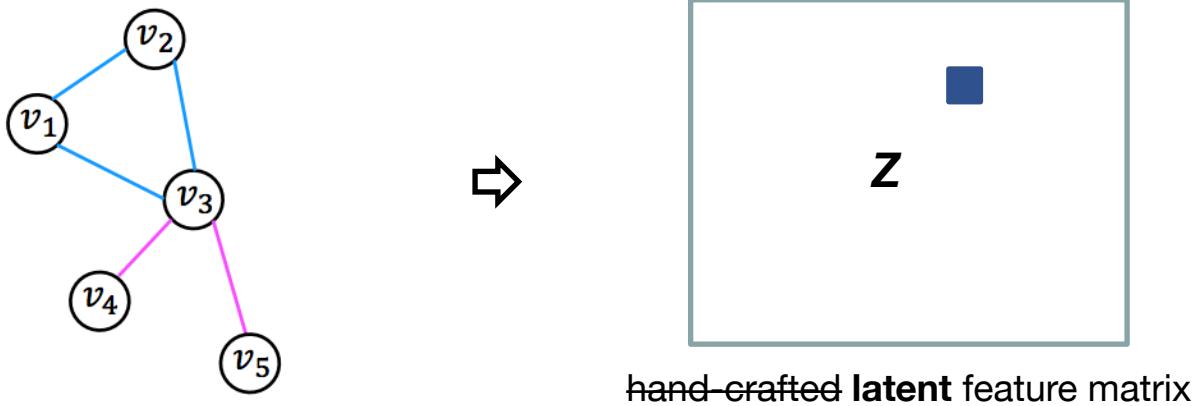
Transportation

figure credit: Web

# The Graph Mining Paradigm



# Graph Representation Learning



Graph & Network applications

- Node classification
- Link prediction
- Community detection
- Anomaly detection
- Social influence
- Graph evolution
- ... ...

Feature engineering learning

machine learning models

- Input: a network  $G = (V, E)$
- Output:  $\mathbf{Z} \in R^{|V| \times k}, k \ll |V|$ ,  $k$ -dim vector  $\mathbf{Z}_v$  for each node  $v$ .

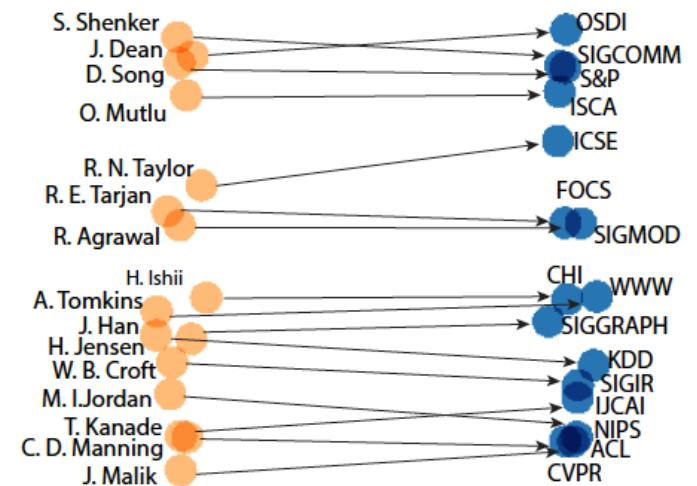
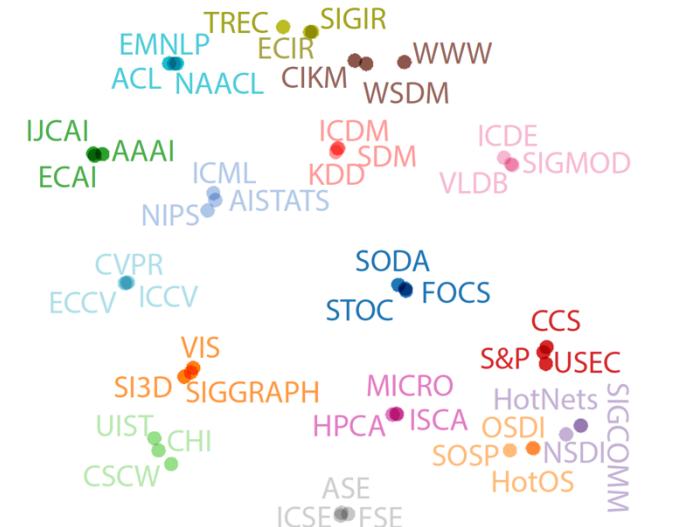
# Application: Embedding Heterogeneous Academic Graph

	225,572,477
	244,499,947
	664,891
	4,397
	48,758
	25,554

## Academic Graph



Graph Representation Learning



1. <https://academic.microsoft.com/>
2. Kuansan Wang et al. Microsoft Academic Graph: When experts are not enough. Quantitative Science Studies 1 (1), 396-413, 2020
3. Dong et al. metapath2vec: scalable representation learning for heterogeneous networks. In KDD 2017.
4. Code & data for metapath2vec: <https://ericdongyangx.github.io/metapath2vec/m2v.html>

# Application: Similarity Search & Recommendation

## Science

Science, also widely referred to as Science Magazine, is the peer-reviewed academic journal of the American Association for the Advancement of Science (AAAS) and one of the world's top academic journals. It was first published in 1880, is currently circulated weekly and has a subscriber base of around 130,000. Because institutional subscriptions and online access serve a larger audience, its estimated readership is 570,400 people.

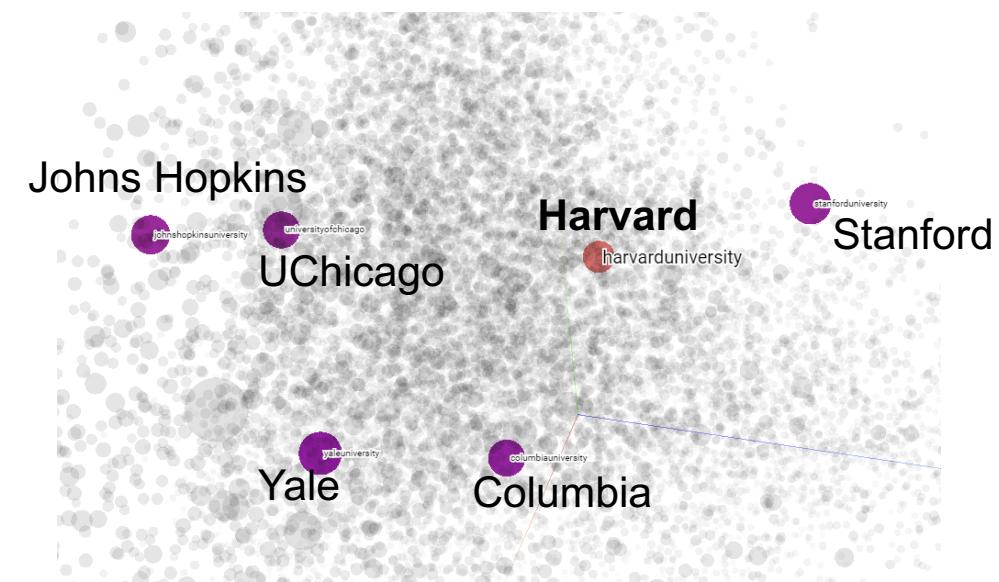
Website links: [scienmag.org](http://scienmag.org), [en.wikipedia.org](https://en.wikipedia.org)

### RELATED JOURNALS

[Nature](#)

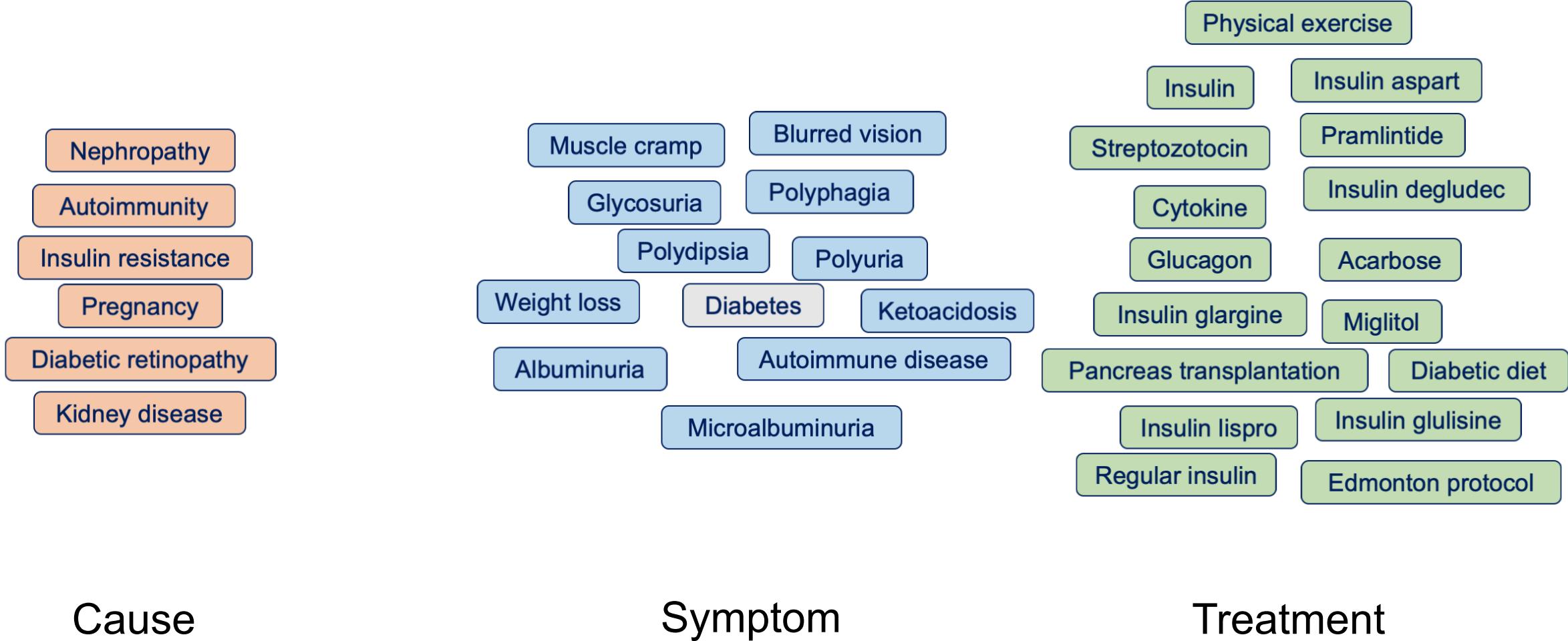
[Proceedings of the National Academy of Sciences of the United States of America](#)

[Nature Communications](#)

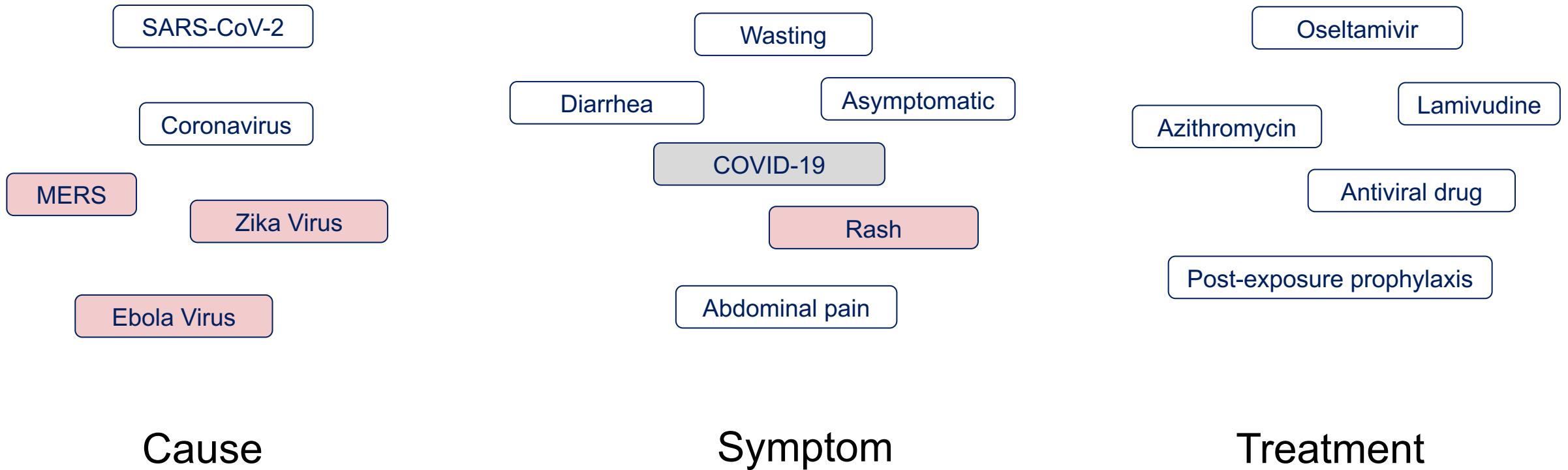


1. <https://academic.microsoft.com/>
2. Kuansan Wang et al. Microsoft Academic Graph: When experts are not enough. Quantitative Science Studies 1 (1), 396-413, 2020
3. Dong et al. metapath2vec: scalable representation learning for heterogeneous networks. In KDD 2017.
4. Code & data for metapath2vec: <https://ericdongyx.github.io/metapath2vec/m2v.html>

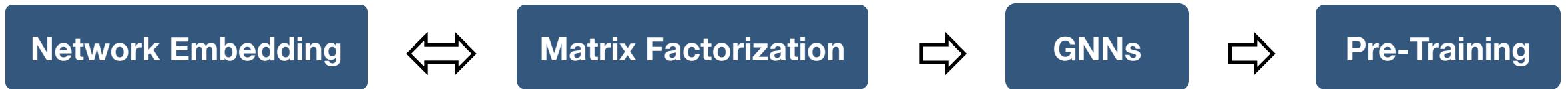
# Application: Reasoning about Diabetes from MAG



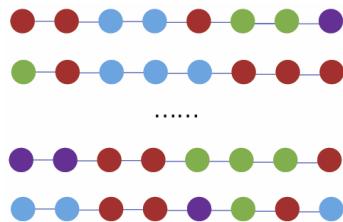
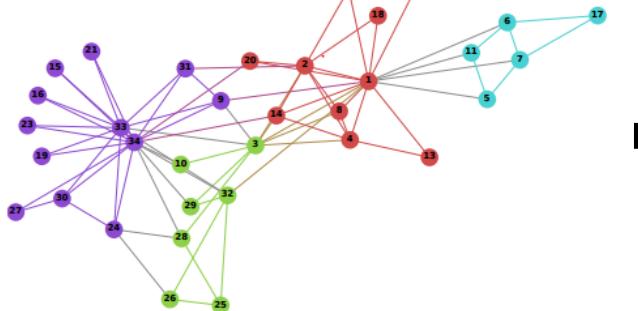
# Application: Reasoning about COVID-19 from MAG



# Graph Representation Learning



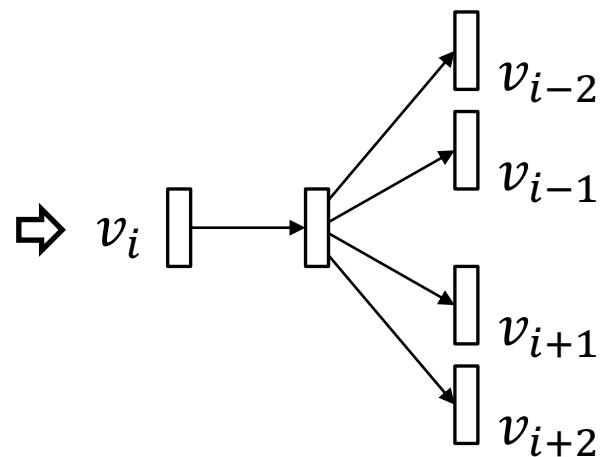
# Network Embedding



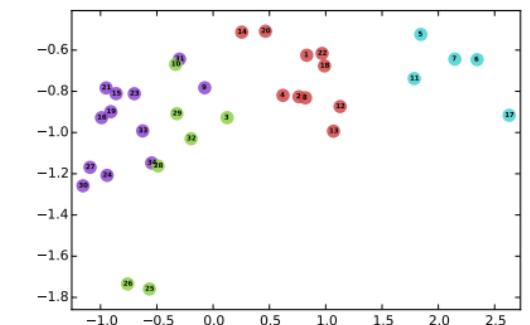
Sequences of objects

- Words in Text
- Nodes in graphs

Feature learning



Skip-Gram



1. Mikolov, et al. Efficient estimation of word representations in vector space. In *ICLR* 2013.
2. Perozzi et al. DeepWalk: Online learning of social representations. In *KDD' 14*, pp. 701–710.

# Distributional Hypothesis of Harris

- **Word embedding:** words in similar contexts have similar meanings (e.g., skip-gram in word embedding)
- **Node embedding:** nodes in similar structural contexts are similar
  - DeepWalk: structural contexts are defined by co-occurrence over random walk paths

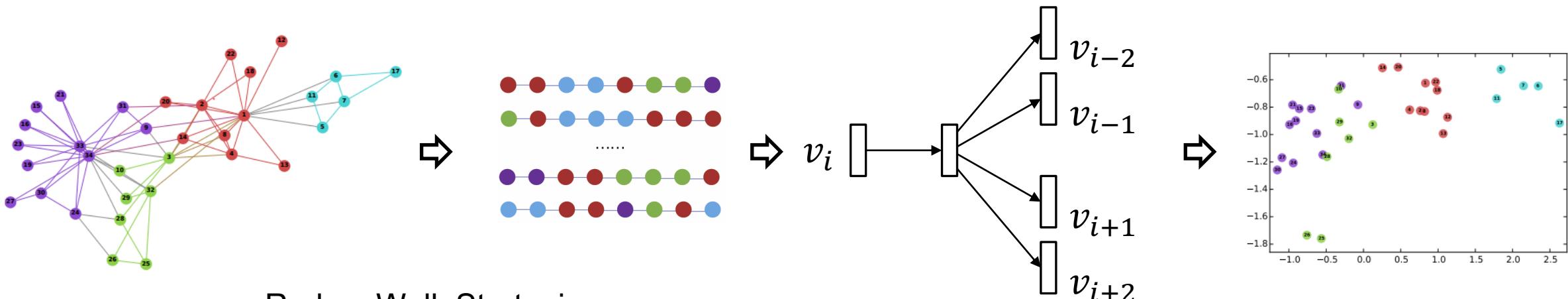
# The Objective

$$\mathcal{L} = \sum_{v \in V} \sum_{c \in N_{rw}(v)} -\log(P(c|v)) \quad \Leftrightarrow \quad p(c|v) = \frac{\exp(\mathbf{z}_v^\top \mathbf{z}_c)}{\sum_{u \in V} \exp(\mathbf{z}_v^\top \mathbf{z}_u)}$$

$\mathcal{L} \rightarrow$  to maximize the likelihood of node co-occurrence on a random walk path

$\mathbf{z}_v^\top \mathbf{z}_c \rightarrow$  the possibility that node  $v$  and context  $c$  appear on a random walk path

# Network Embedding: Random Walk + Skip-Gram

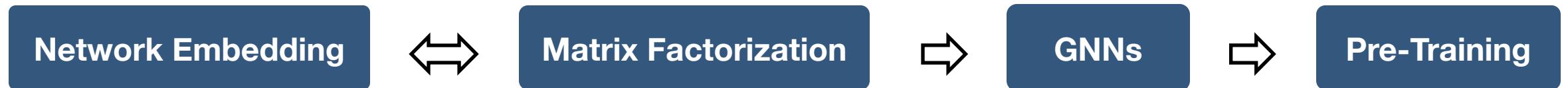


Random Walk Strategies:

- DeepWalk (walk length > 1)
- LINE (walk length = 1)
- PTE (walk length = 1)
- node2vec (biased random walk)
- metapath2vec (heterogeneous rw)

1. Perozzi et al. **DeepWalk**: Online learning of social representations. In *KDD'14*. **Most Cited Paper in KDD'14**.
2. Tang et al. **LINE**: Large scale information network embedding. In *WWW'15*. **Most Cited Paper in WWW'15**.
3. Grover and Leskovec. **node2vec**: Scalable feature learning for networks. In *KDD'16*. **2<sup>nd</sup> Most Cited Paper in KDD'16**.
4. Dong et al. **metapath2vec**: scalable representation learning for heterogeneous networks. In *KDD 2017*. **Most Cited Paper in KDD'17**.

# Graph Representation Learning



- DeepWalk
- LINE
- Node2vec
- PTE
- ...
- metapath2vec

# NetMF: Network Embedding as Matrix Factorization

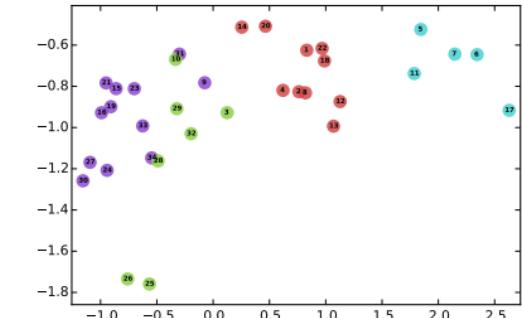
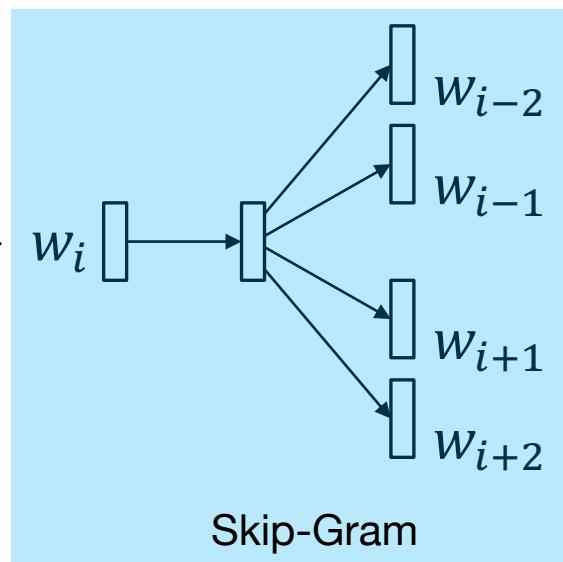
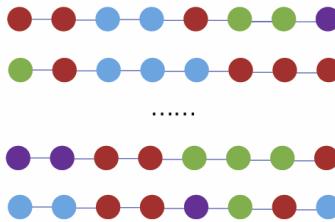
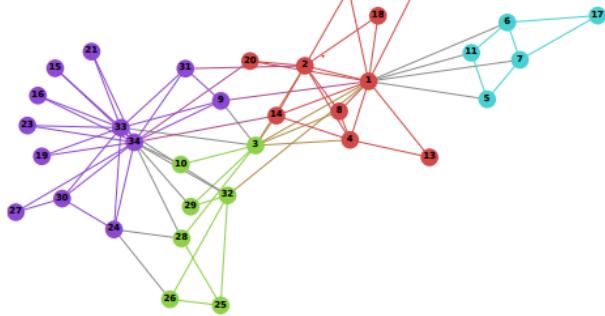
- DeepWalk  $\log \left( \frac{\text{vol}(G)}{b} \left( \frac{1}{T} \sum_{r=1}^T (\mathbf{D}^{-1} \mathbf{A})^r \right) \mathbf{D}^{-1} \right)$
- LINE  $\log \left( \frac{\text{vol}(G)}{b} \mathbf{D}^{-1} \mathbf{A} \mathbf{D}^{-1} \right)$
- PTE  $\log \left( \begin{bmatrix} \alpha \text{vol}(G_{\text{ww}})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

$\mathbf{A}$  Adjacency matrix  
 $\mathbf{D}$  Degree matrix

$b$ : #negative samples  
 $T$ : context window size

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

# Understanding Random Walk + Skip Gram



*Graph Language*

- $G$ : graph
- $A$ : adjacency matrix
- $D$ : degree matrix
- $\text{vol}(G)$ : volume of  $G$



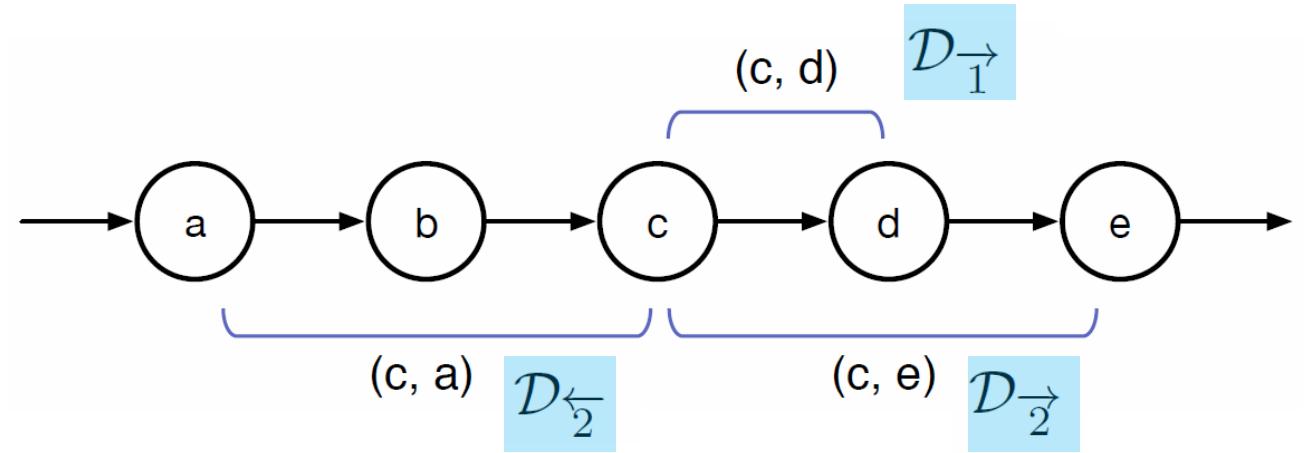
$$\log\left(\frac{\#(w, c)|\mathcal{D}|}{b\#(w)\#(c)}\right)$$

*NLP Language*

- $\#(w, c)$ : co-occurrence of  $w$  &  $c$
- $\#(w)$ : occurrence of word  $w$
- $\#(c)$ : occurrence of context  $c$
- $\mathcal{D}$ : word-context pair  $(w, c)$  multi-set
- $|\mathcal{D}|$ : number of word-context pairs

# Understanding Random Walk + Skip Gram

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



NLP Language

- $\#(w, c)$ : co-occurrence of w & c
- $\#(w)$ : occurrence of word w
- $\#(c)$ : occurrence of context c
- $\mathcal{D}$ : word–context pair  $(w, c)$  multi-set
- $|\mathcal{D}|$ : number of word-context pairs

Distinguish direction and distance

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

# Understanding Random Walk + Skip Gram

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

$$\frac{\#(w,c)}{|\mathcal{D}|} = \frac{1}{2T} \sum_{r=1}^T \left( \frac{\#(w,c)_{\vec{r}}}{|\mathcal{D}_{\vec{r}}|} + \frac{\#(w,c)_{\overleftarrow{r}}}{|\mathcal{D}_{\overleftarrow{r}}|} \right)$$

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

$$\frac{\#(w,c)_{\vec{r}}}{|\mathcal{D}_{\vec{r}}|} \xrightarrow{p} \frac{d_w}{\text{vol}(G)} (\mathbf{P}^r)_{w,c}$$

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

$$\frac{\#(w,c)_{\overleftarrow{r}}}{|\mathcal{D}_{\overleftarrow{r}}|} \xrightarrow{p} \frac{d_c}{\text{vol}(G)} (\mathbf{P}^r)_{c,w}$$

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

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

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

# Understanding Random Walk + Skip Gram

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

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

*Graph Language*

**A** Adjacency matrix

**D** Degree matrix

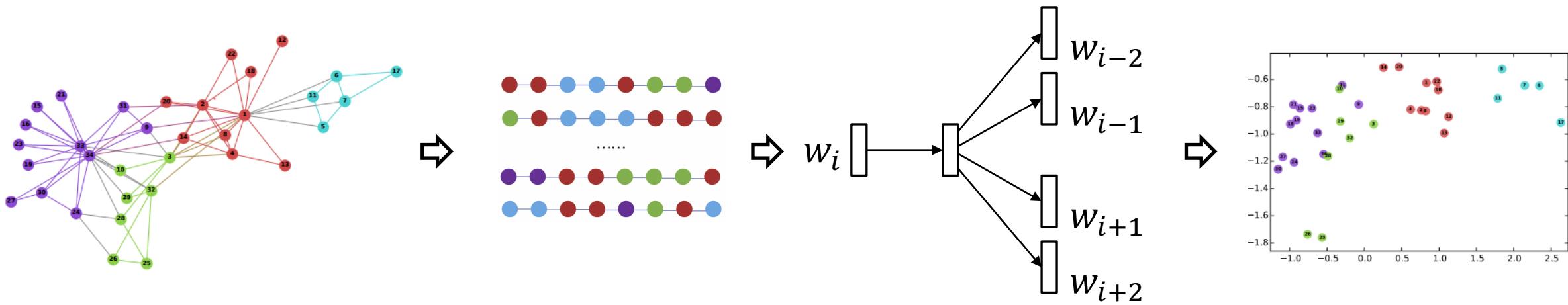
$\mathbf{P} = \mathbf{D}^{-1} \mathbf{A}$

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

b: #negative samples  
T: context window size

$$\log \left( \frac{\#(w, c) |\mathcal{D}|}{b \#(w) \cdot \#(c)} \right) = \log \left( \frac{\text{vol}(G)}{b} \left( \frac{1}{T} \sum_{r=1}^T (\mathbf{D}^{-1} \mathbf{A})^r \right) \mathbf{D}^{-1} \right)$$

# Understanding Random Walk + Skip Gram



DeepWalk is asymptotically and implicitly factorizing

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

$A$  Adjacency matrix

$D$  Degree matrix

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

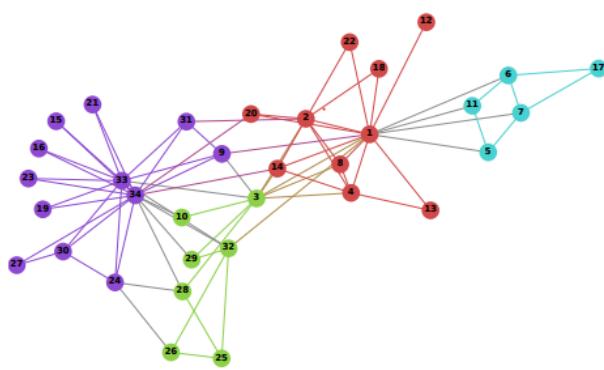
$b$ : #negative samples

$T$ : context window size

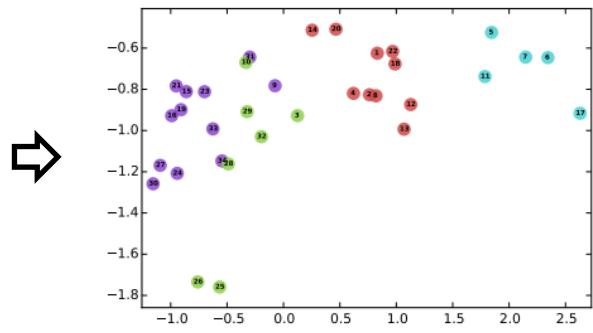
# Unifying DeepWalk, LINE, PTE, & node2vec as Matrix Factorization

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

# NetMF: Explicitly Factorizing the DeepWalk Matrix



Matrix  
Factorization



DeepWalk is asymptotically and implicitly factorizing

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

1. Qiu et al. Network embedding as matrix factorization: unifying deepwalk, line, pte, and node2vec. In *WSDM'18*.
2. Code & data for NetMF: <https://github.com/xptree/NetMF>

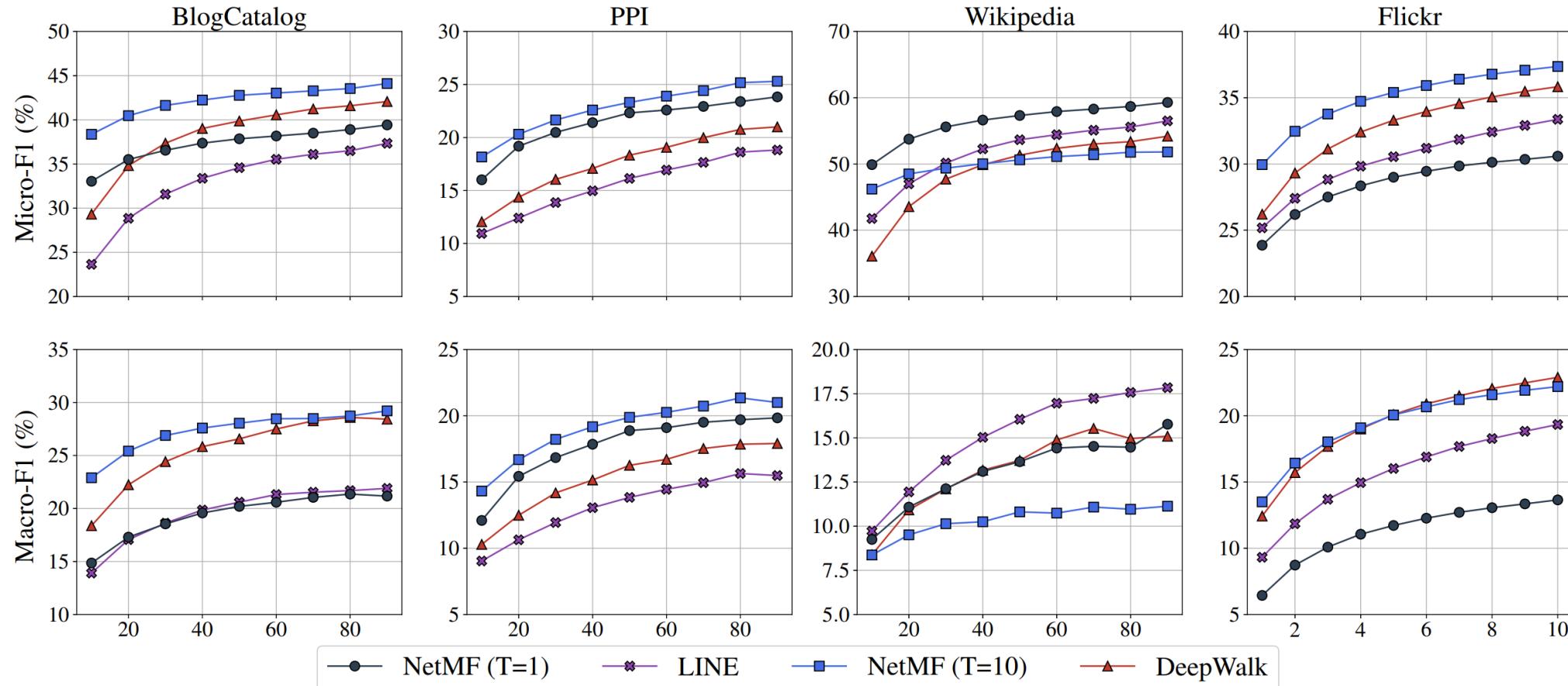
# NetMF

1. Construction
2. Factorization

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

1. Qiu et al. Network embedding as matrix factorization: unifying deepwalk, line, pte, and node2vec. In *WSDM'18*.
2. Code & data for NetMF: <https://github.com/xptree/NetMF>

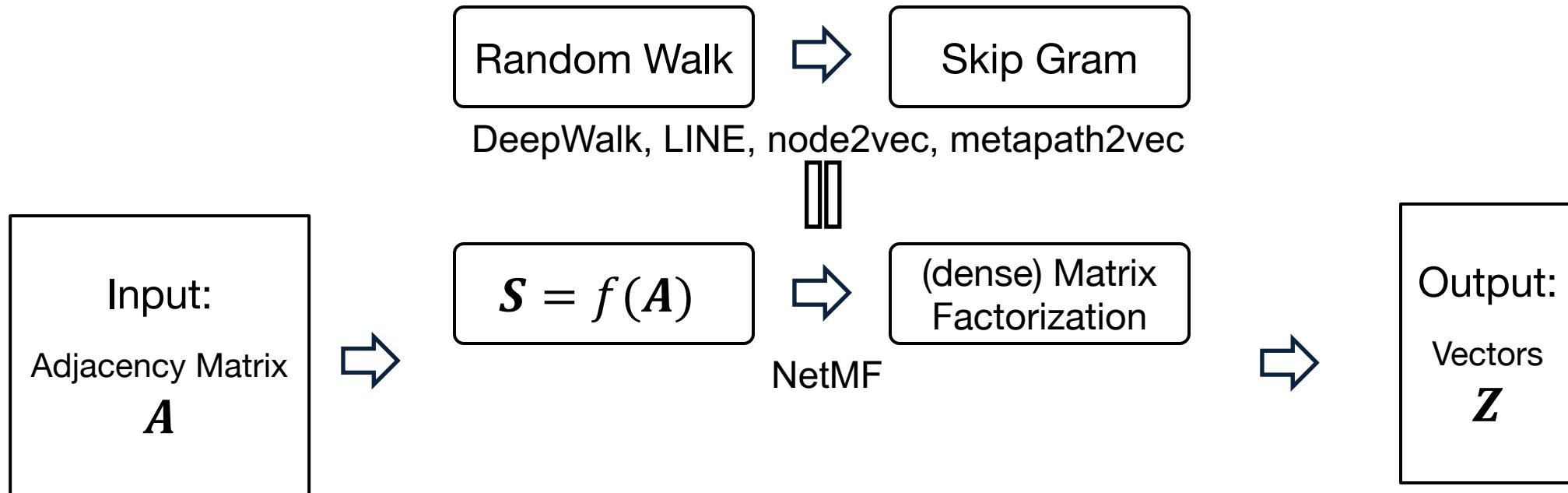
# Results



**Explicit matrix factorization (NetMF) offers performance gains over implicit matrix factorization (DeepWalk & LINE)**

1. Qiu et al. Network embedding as matrix factorization: unifying deepwalk, line, pte, and node2vec. In *WSDM'18*.
2. Code & data for NetMF: <https://github.com/xptree/NetMF>

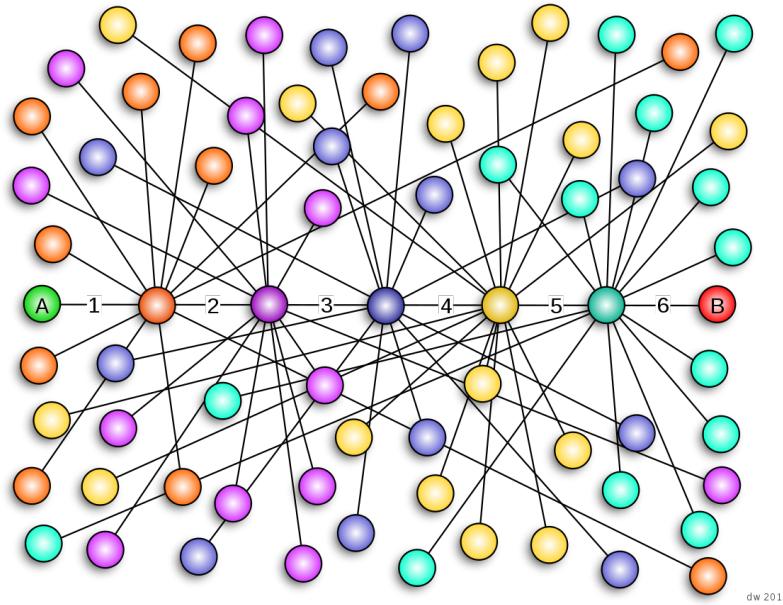
# Network Embedding



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

1. Qiu et al. Network embedding as matrix factorization: unifying deepwalk, line, pte, and node2vec. In *WSDM'18*.
2. Code & data for NetMF: <https://github.com/xptree/NetMF>

# Challenge?



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

$n^2$  non-zeros  
**Dense!!**

Time complexity  
 $O(n^3)$

# NetMF

How can we solve this issue?

1. Construction
2. Factorization

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

1. Qiu et al. NetSMF: Network embedding as sparse matrix factorization. In **WWW 2019**.
2. **Code & data** for NetSMF: <https://github.com/xptree/NetSMF>

# NetSMF--Sparse

How can we solve this issue?

1. **Sparse** Construction
2. **Sparse** Factorization

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

1. Qiu et al. NetSMF: Network embedding as sparse matrix factorization. In **WWW 2019**.
2. **Code & data** for NetSMF: <https://github.com/xptree/NetSMF>

# Sparsify $S$

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

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

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

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

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

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

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

$$\begin{aligned} S &= \log^\circ \left( \frac{\text{vol}(G)}{b} \left( \frac{1}{T} \sum_{r=1}^T (D^{-1} A)^r \right) D^{-1} \right) \\ \alpha_1 = \dots = \alpha_T = \frac{1}{T} \quad \Rightarrow \quad &= \log^\circ \left( \frac{\text{vol}(G)}{b} D^{-1} (D - L) D^{-1} \right) \\ &\approx \log^\circ \left( \frac{\text{vol}(G)}{b} D^{-1} (D - \tilde{L}) D^{-1} \right) \end{aligned}$$

1. Qiu et al. NetSMF: Network embedding as sparse matrix factorization. In **WWW 2019**.

2. Code & data for NetSMF: <https://github.com/xptree/NetSMF>

## NetSMF --- Sparse

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

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

- ▶ Factorize the constructed matrix

# NetSMF---Bounded Approximation Error

$$\begin{aligned} & \log^{\circ} \left( \frac{\text{vol}(G)}{b} \left( \frac{1}{T} \sum_{r=1}^T (\mathbf{D}^{-1} \mathbf{A})^r \right) \mathbf{D}^{-1} \right) \\ &= \log^{\circ} \left( \frac{\text{vol}(G)}{b} \mathbf{D}^{-1} (\mathbf{D} - \mathbf{L}) \mathbf{D}^{-1} \right) \longrightarrow \mathbf{M} \\ & \approx \log^{\circ} \left( \frac{\text{vol}(G)}{b} \mathbf{D}^{-1} (\mathbf{D} - \tilde{\mathbf{L}}) \mathbf{D}^{-1} \right) \longrightarrow \tilde{\mathbf{M}} \end{aligned}$$

## Theorem

The singular value of  $\tilde{\mathbf{M}} - \mathbf{M}$  satisfies

$$\sigma_i(\tilde{\mathbf{M}} - \mathbf{M}) \leq \frac{4\epsilon}{\sqrt{d_i d_{\min}}}, \forall i \in [n].$$

## Theorem

Let  $\|\cdot\|_F$  be the matrix Frobenius norm. Then

$$\left\| \text{trunc\_log}^{\circ} \left( \frac{\text{vol}(G)}{b} \tilde{\mathbf{M}} \right) - \text{trunc\_log}^{\circ} \left( \frac{\text{vol}(G)}{b} \mathbf{M} \right) \right\|_F \leq \frac{4\epsilon \text{vol}(G)}{b \sqrt{d_{\min}}} \sqrt{\sum_{i=1}^n \frac{1}{d_i}}.$$

1. Qiu et al. NetSMF: Network embedding as sparse matrix factorization. In **WWW 2019**.

2. Code & data for NetSMF: <https://github.com/xptree/NetSMF>

# Results

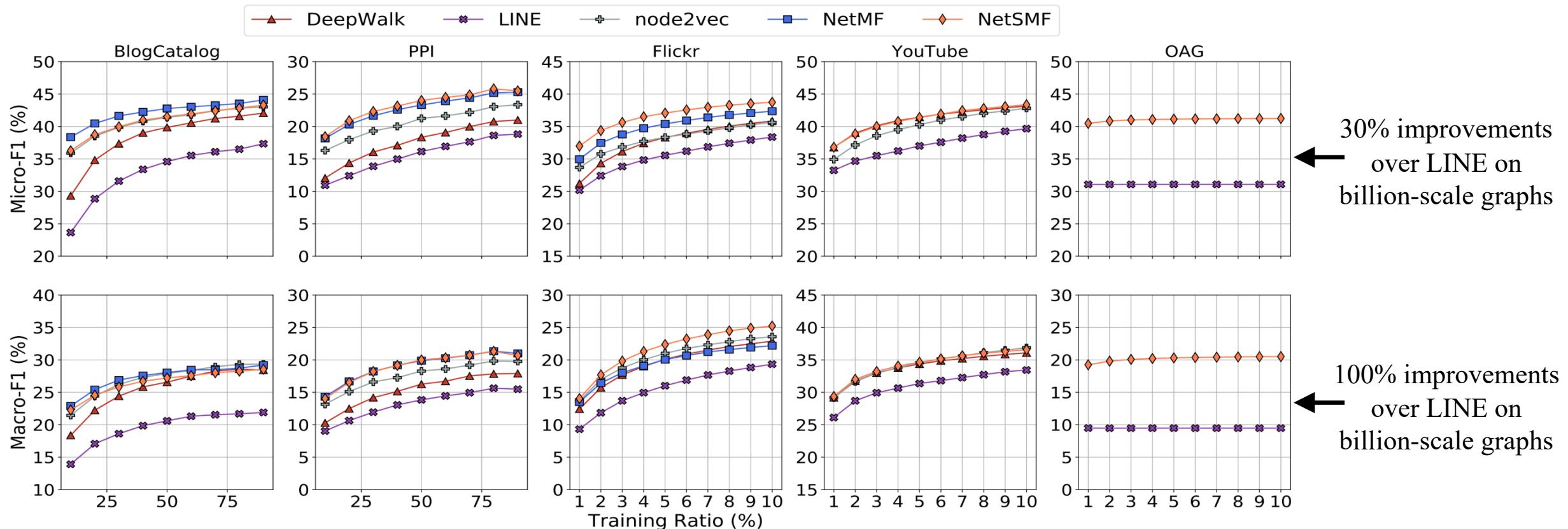
Dataset	BlogCatalog	PPI	Flickr	YouTube	OAG	
$ V $	10,312	3,890	80,513	1,138,499	67,768,244	•#non-zeros:
$ E $	333,983	76,584	5,899,882	2,990,443	895,368,962	•~4.5 Quadrillion → 45 Billion
#labels	39	50	195	47	19	

	LINE	DeepWalk	node2vec	NetMF	NetSMF
BlogCatalog	40 mins	12 mins	56 mins	19 mins	13 mins
PPI	41 mins	4 mins	4 mins	1 min	10 secs
Flickr	42 mins	2.2 hours	21 hours	5 days	48 mins
YouTube	46 mins	4.3 hours	4 days	×	4.1 hours
OAG	2.6 hours	–	–	×	24 hours

1. Qiu et al. NetSMF: Network embedding as sparse matrix factorization. In **WWW 2019**.

2. Code & data for NetSMF: <https://github.com/xptree/NetSMF>

# Results

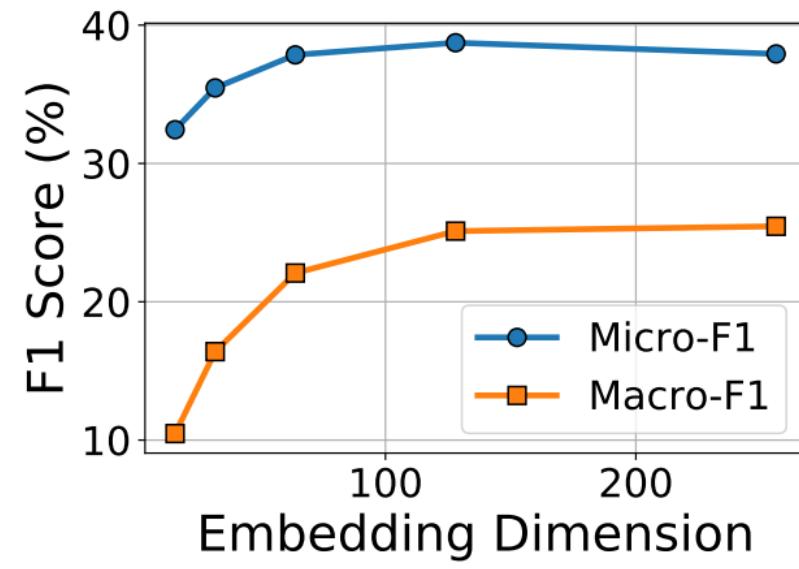


**Effectiveness:** NetSMF (sparse MF)  $\approx$  NetMF (explicit MF)  $>$  DeepWalk/LINE (implicit MF)

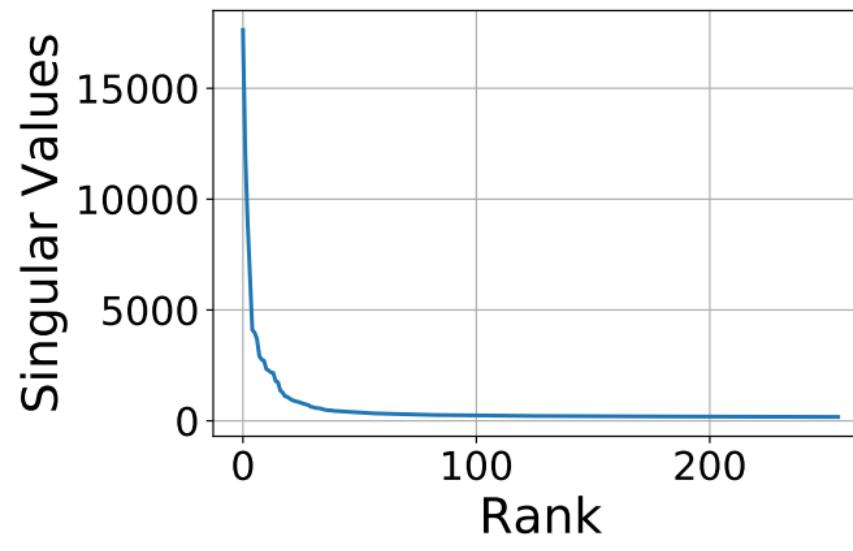
**Efficiency:** NetSMF (sparse MF) can handle billion-scale network embedding

1. Qiu et al. NetSMF: Network embedding as sparse matrix factorization. In **WWW 2019**.

2. **Code & data** for NetSMF: <https://github.com/xptree/NetSMF>



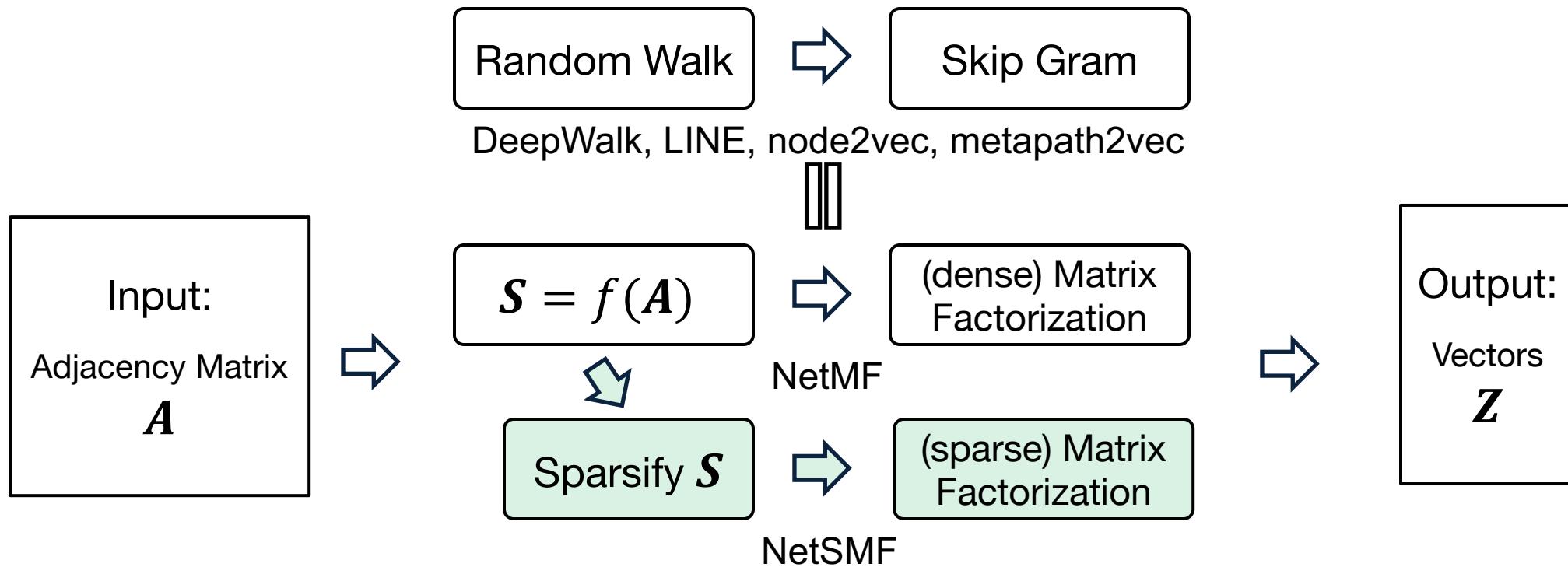
(a)



(b)

1. Qiu et al. NetSMF: Network embedding as sparse matrix factorization. In **WWW 2019**.
2. **Code & data** for NetSMF: <https://github.com/xptree/NetSMF>

# Network Embedding



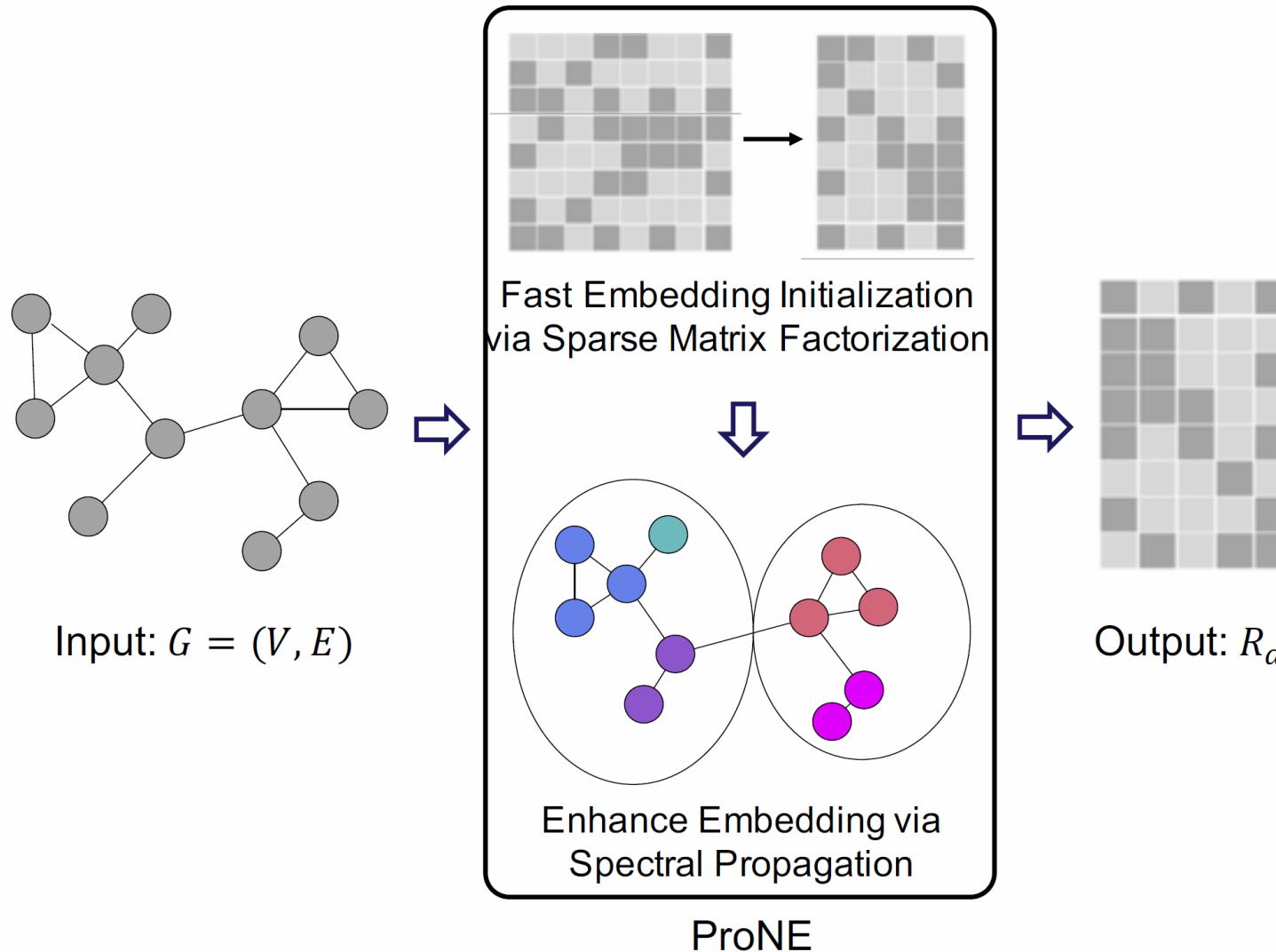
Incorporate network structures  $A$  into the similarity matrix  $S$ , and then factorize  $S$

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

1. Qiu et al. NetSMF: Network embedding as sparse matrix factorization. In **WWW 2019**.

2. Code & data for NetSMF: <https://github.com/xptree/NetSMF>

# ProNE: Propagation based Network Embedding



1. Zhang et al. ProNE: Fast and Scalable Network Representation Learning. In **IJCAI 2019**

2. **Code & data** for ProNE: <https://github.com/THUDM/ProNE>

# Spectral Propagation

$$R_d \leftarrow D^{-1}A(I_n - \tilde{L}) R_d$$

The idea of **Graph Neural Networks**

$D^{-1}A(I_n - \tilde{L})$  is  $D^{-1}A$  modulated by the filter in the spectrum

$\tilde{L} = Ug(\Lambda)U^T$  is the spectral filter of  $L = I_n - D^{-1}A$

1. Zhang et al. ProNE: Fast and Scalable Network Representation Learning. In **IJCAI 2019**

2. Code & data for ProNE: <https://github.com/THUDM/ProNE>

# Chebyshev expansion for efficiency

- To avoid explicit eigendecomposition and Fourier transform
  - Chebyshev expansion  $T_{i+1}(x) = 2xT_i(x) - T_{i-1}(x)$  with  $T_0(x) = 1, T_1(x) = x$

$$\begin{aligned}\widetilde{L} &= U \text{diag}([g(\lambda_1), \dots, g(\lambda_n)]) U^T &\rightarrow \widetilde{L} &\approx B_0(\theta)T_0(\bar{L}) + 2 \sum_{i=1}^{k-1} (-)^i B_i(\theta)T_i(\bar{L}) \\ &\approx U \sum_{i=0}^{k-1} c_i(\theta)T_i(\bar{\Lambda})U^T \\ &= \sum_{i=0}^{k-1} c_i(\theta)T_i(\bar{L})\end{aligned}$$

# Efficiency

Dataset	20 Threads			1 Thread
	DeepWalk	LINE	node2vec	ProNE
PPI	272	70	828	<b>3</b>
Wiki	494	87	939	<b>6</b>
BlogCatalog	1,231	185	3,533	<b>21</b>
DBLP	3,825	1,204	4,749	<b>24</b>
Youtube	68,272	5,890	>5days	<b>627</b>

1.1M nodes

19hours    98mins    10mins

Dataset	training ratio	0.01	0.03	0.05	0.07	0.09
		DeepWalk	LINE	node2vec	ProNE	GraRep
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
Youtube	ProNE ( $\pm\sigma$ )	48.8 ( $\pm 1.0$ )	<b>56.2</b> ( $\pm 0.5$ )	<b>58.0</b> ( $\pm 0.2$ )	<b>58.8</b> ( $\pm 0.2$ )	<b>59.2</b> ( $\pm 0.1$ )
	DeepWalk	38.0	40.1	41.3	42.1	42.8
	LINE	33.2	35.5	37.0	38.2	39.3
	ProNE (SMF)	36.5	40.2	41.2	41.7	42.1
	ProNE ( $\pm\sigma$ )	<b>38.2</b> ( $\pm 0.8$ )	<b>41.4</b> ( $\pm 0.3$ )	<b>42.3</b> ( $\pm 0.2$ )	<b>42.9</b> ( $\pm 0.2$ )	<b>43.3</b> ( $\pm 0.2$ )

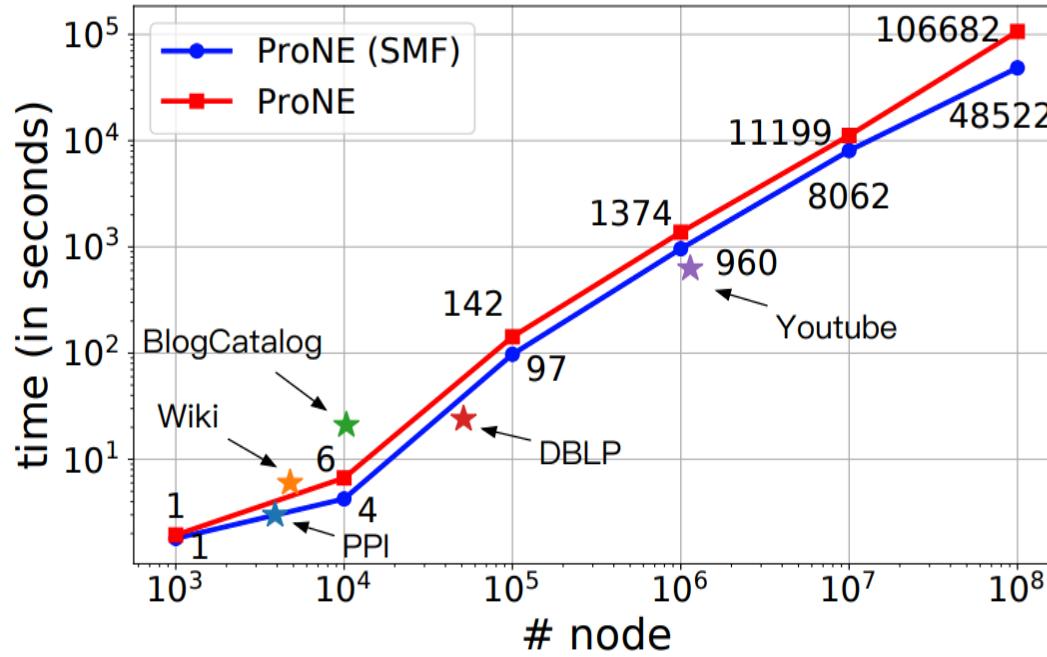
ProNE offers **10-400X** speedups (1 thread vs 20 threads)

Embed 100,000,000 nodes by 1 thread: 29 hours with performance superiority

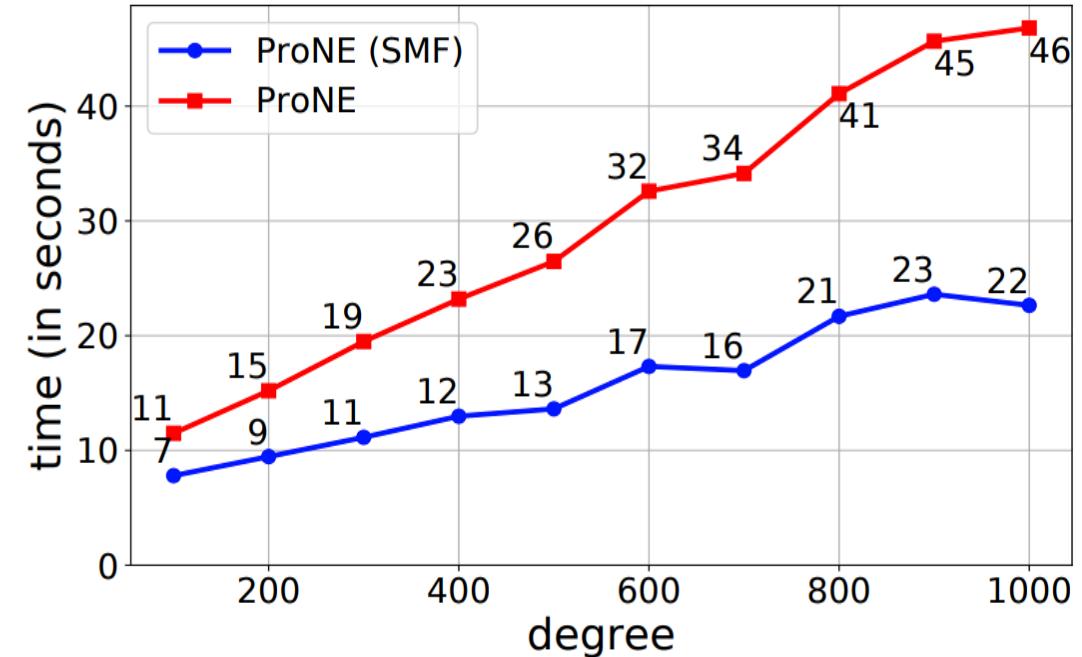
1. Zhang et al. ProNE: Fast and Scalable Network Representation Learning. In **IJCAI 2019**

2. Code & data for ProNE: <https://github.com/THUDM/ProNE>

# Scalability



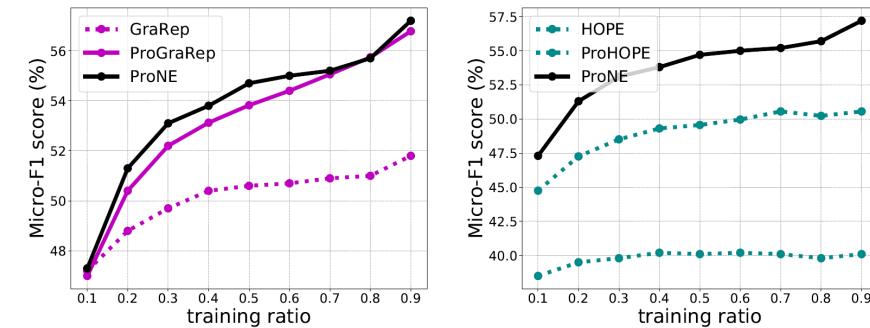
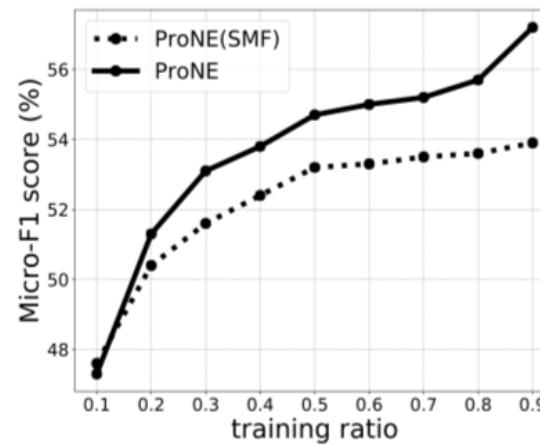
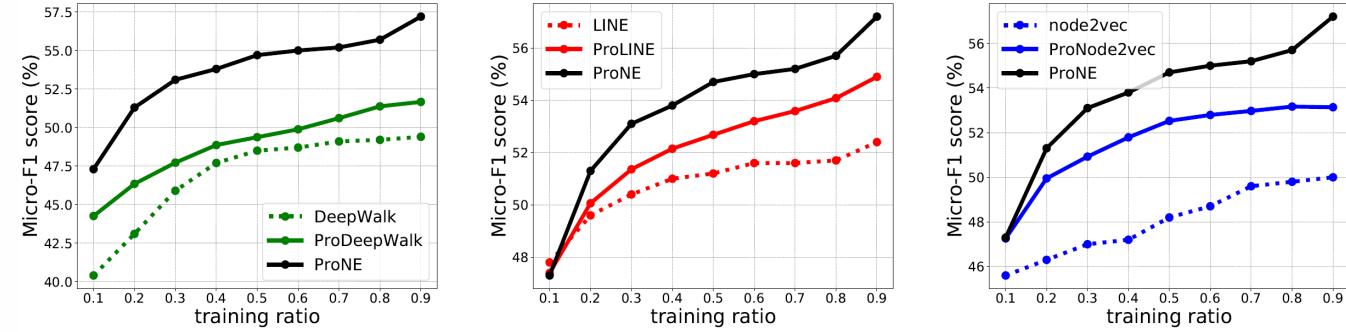
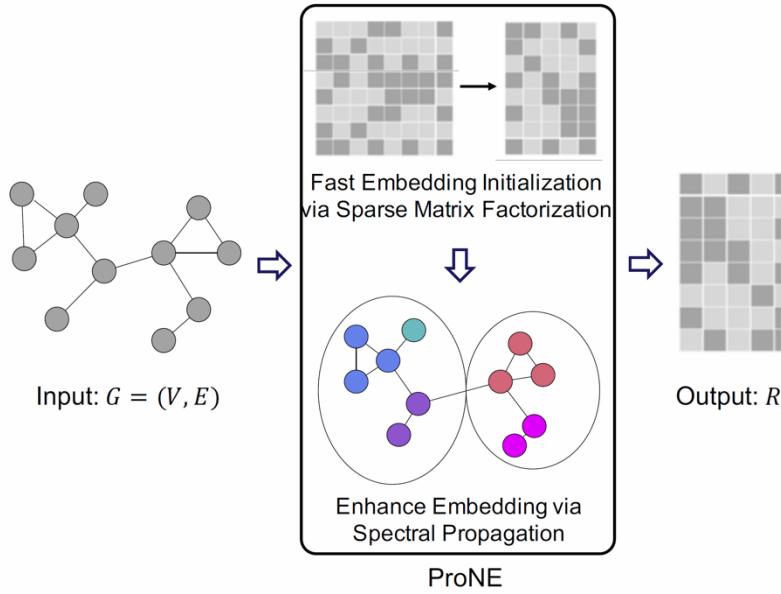
(a) The node degree is fixed to 10 and #nodes grows



(b) #nodes is fixed to 10,000 and the node degree grows

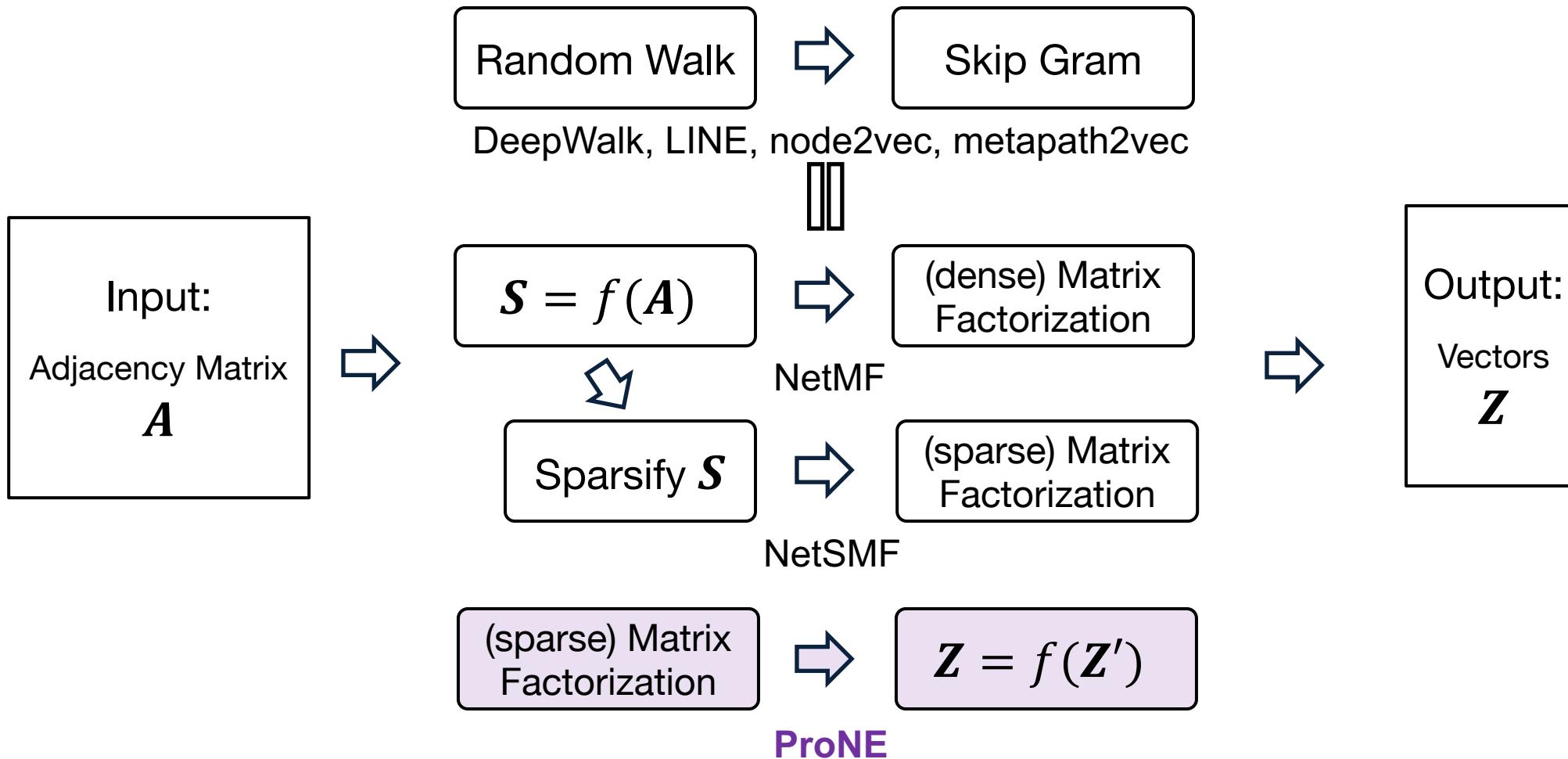
Embed 100,000,000 nodes by 1 thread: 29 hours with performance superiority

# ProNE: A General Propagation Framework



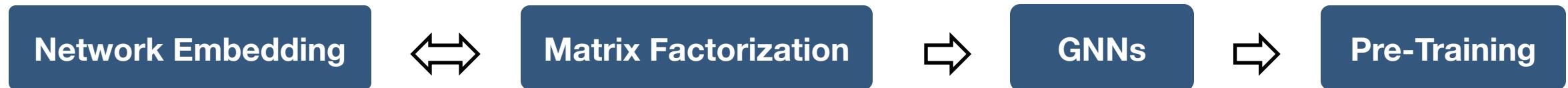
1. Zhang et al. ProNE: Fast and Scalable Network Representation Learning. In **IJCAI 2019**
2. **Code & data** for ProNE: <https://github.com/THUDM/ProNE>

# Network Embedding



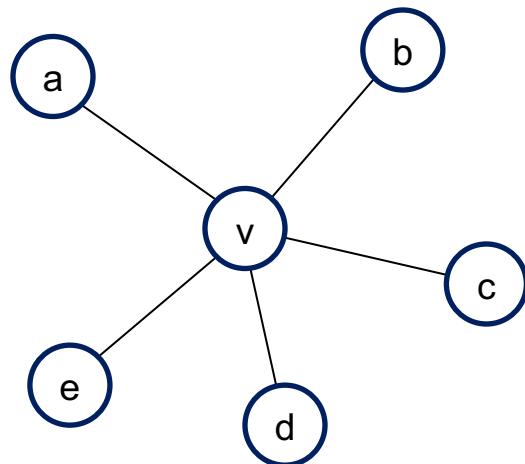
Factorize  $A$ , and then incorporate network structures via spectral propagation

# Graph Representation Learning



- DeepWalk
  - LINE
  - Node2vec
  - PTE
  - ...
  - metapath2vec
- NetMF
  - NetSMF
  - ...
  - ProNE (**Propagation**)

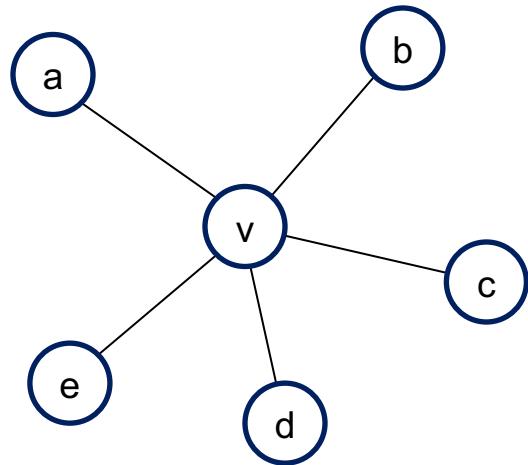
# Connecting NE with Graph Neural Networks



- **ProNE**
  - Propagation based network embedding
$$\mathbf{R}_d \leftarrow \mathbf{D}^{-1} \mathbf{A} (\mathbf{I}_n - \tilde{\mathbf{L}}) \mathbf{R}_d$$
- **GNN**
  - Neighborhood aggregation: aggregate neighbor information and pass into a neural network
$$\mathbf{h}_v = f(\mathbf{h}_v, \mathbf{h}_a, \mathbf{h}_b, \mathbf{h}_c, \mathbf{h}_d, \mathbf{h}_e)$$

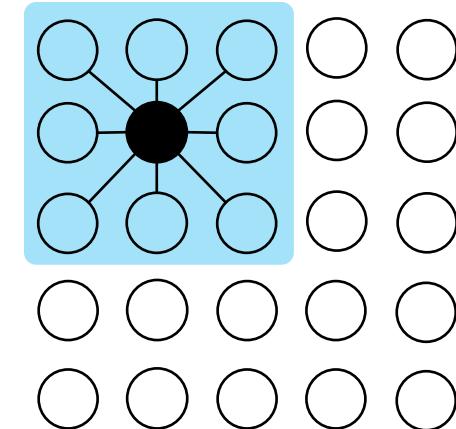
1. Justin Gilmer, et al. Neural message passing for quantum chemistry. In **ICML 2017**.
2. Zhang et al. ProNE: Fast and Scalable Network Representation Learning. In **IJCAI 2019**

# Graph Neural Networks



**Graph Convolution**

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



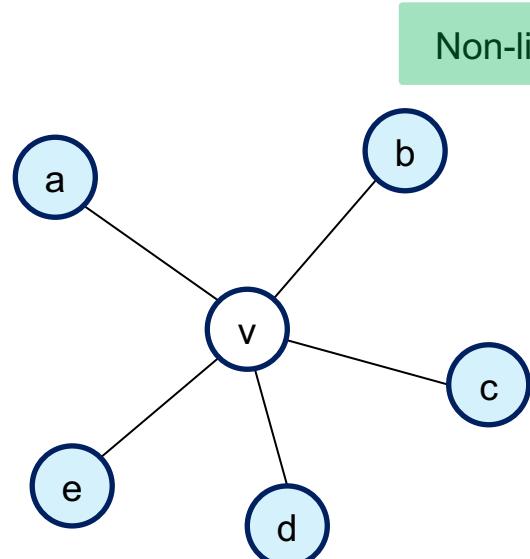
**CNN**

## 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**

# Graph Convolutional Networks



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

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

parameters in layer  $k$

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

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

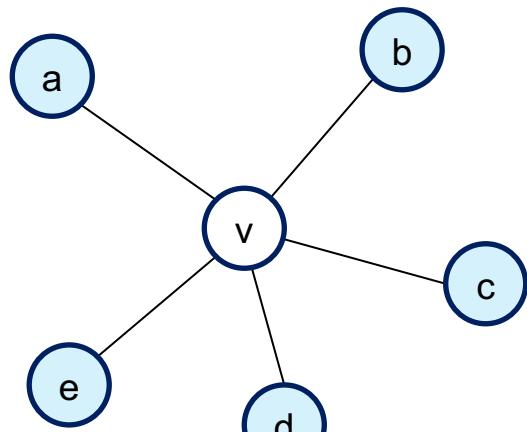
the neighbors of node  $v$

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

normalized Laplacian matrix

Aggregate info from neighborhood via the normalized Laplacian matrix

# Graph Convolutional Networks

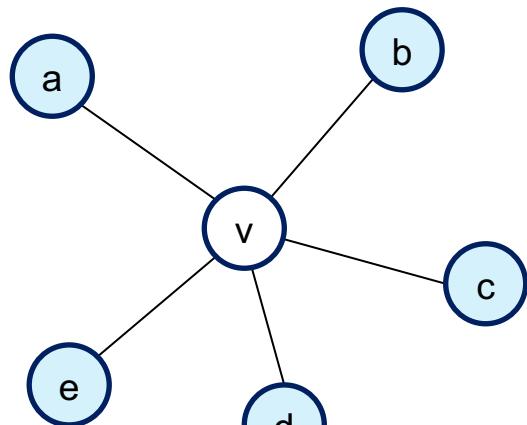


$$\mathbf{h}_v^k = \sigma(\mathbf{W}^k \sum_{u \in N(v)} \frac{\mathbf{h}_u^{k-1}}{\sqrt{|N(u)||N(v)|}} + \mathbf{W}^k \sum_v \frac{\mathbf{h}_v^{k-1}}{\sqrt{|N(v)||N(v)|}})$$

Aggregate from  $v$ 's neighbors

Aggregate from itself

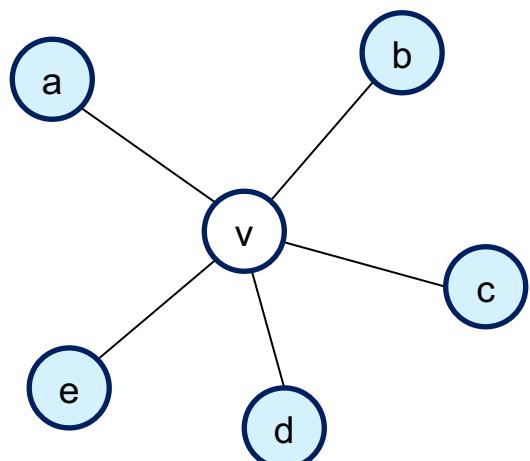
# Graph Convolutional Networks



The same parameters for both its neighbors & itself

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

# Graph Convolutional Networks

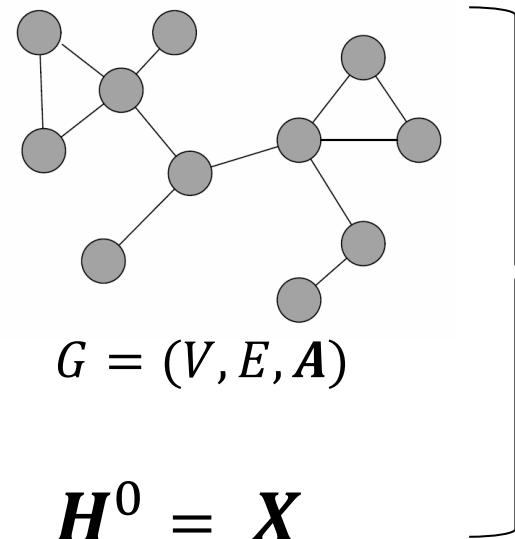


$$\mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \mathbf{H}^{(k-1)} \mathbf{W}^{(k)}$$

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

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

# Graph Convolutional Networks

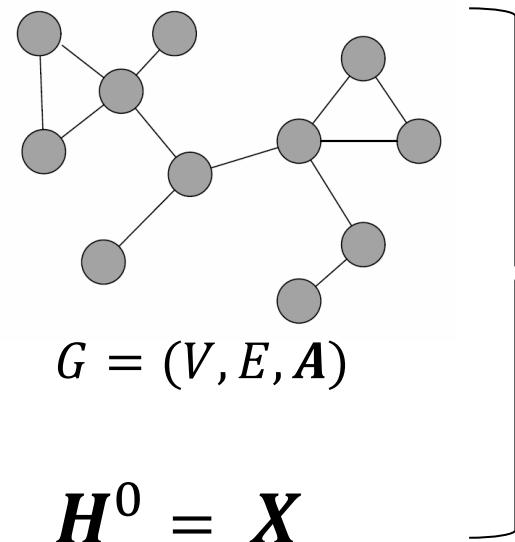


$$\Rightarrow H^k = \sigma \left( D^{-\frac{1}{2}}(A + I)D^{-\frac{1}{2}}H^{(k-1)}W^{(k)} \right) \Leftrightarrow Z = H^K$$

Input

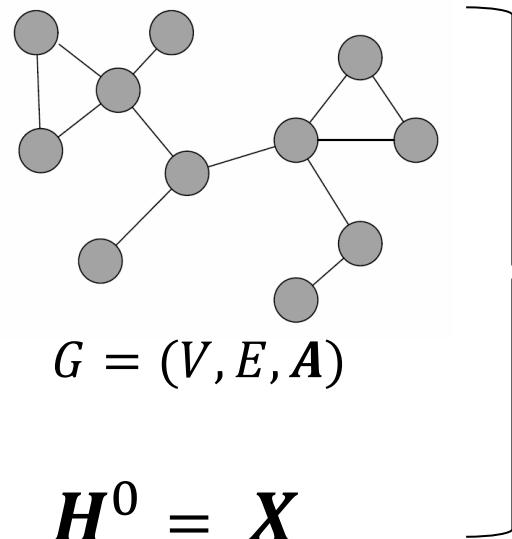
Output

# Graph Convolutional Networks



- Output
- $$\Leftrightarrow \mathbf{H}^k = \sigma \left( \mathbf{D}^{-\frac{1}{2}} (\mathbf{A} + \mathbf{I}) \mathbf{D}^{-\frac{1}{2}} \mathbf{H}^{(k-1)} \mathbf{W}^{(k)} \right) \Leftrightarrow \mathbf{Z} = \mathbf{H}^K$$
- **Model training**
    - The common setting is to have an end to end training framework with a supervised task
    - That is, define a loss function over  $\mathbf{Z}$

# Graph Convolutional Networks



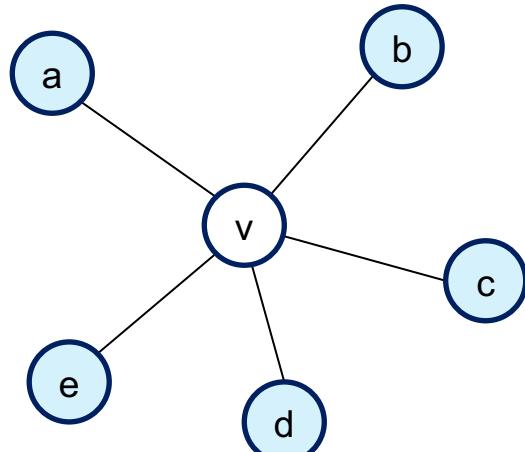
$$\Leftrightarrow \mathbf{H}^k = \sigma \left( \mathbf{D}^{-\frac{1}{2}} (\mathbf{A} + \mathbf{I}) \mathbf{D}^{-\frac{1}{2}} \mathbf{H}^{(k-1)} \mathbf{W}^{(k)} \right) \Leftrightarrow \mathbf{Z} = \mathbf{H}^K$$

Input

- Benefits: Parameter sharing for all nodes
  - #parameters is subline in  $|V|$
  - Enable inductive learning for new nodes

Output

# 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

# GraphSage

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

- Mean:

$$\text{AGG} = \sum_{u \in N(v)} \frac{\mathbf{h}_u^{k-1}}{|N(v)|}$$

- Pool

- Transform neighbor vectors and apply symmetric vector function.

element-wise mean/max

$$\text{AGG} = \gamma([\mathbf{Q} \mathbf{h}_u^{k-1}, \forall u \in N(v)])$$

- LSTM:

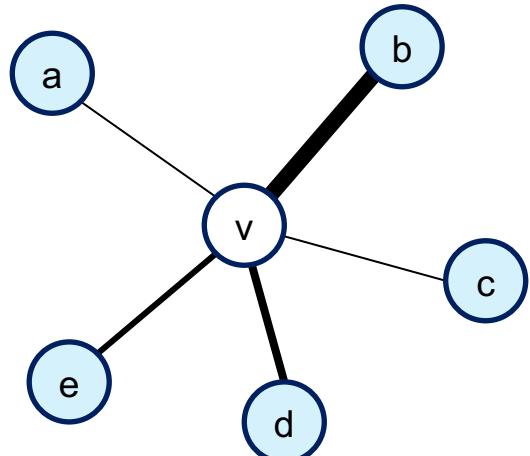
- Apply LSTM to random permutation of neighbors.

$$\text{AGG} = \text{LSTM}([\mathbf{h}_u^{k-1}, \forall u \in \pi(N(v))])$$

# Graph Neural Network

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

# Graph Attention



GCN

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

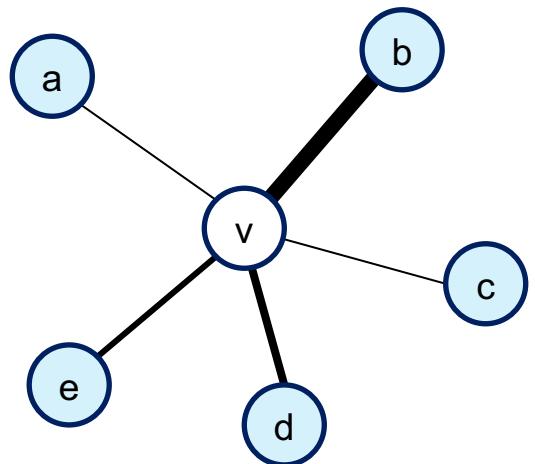
Aggregate info from neighborhood via the normalized Laplacian matrix

Graph Attention

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

Aggregate info from neighborhood via the learned attention

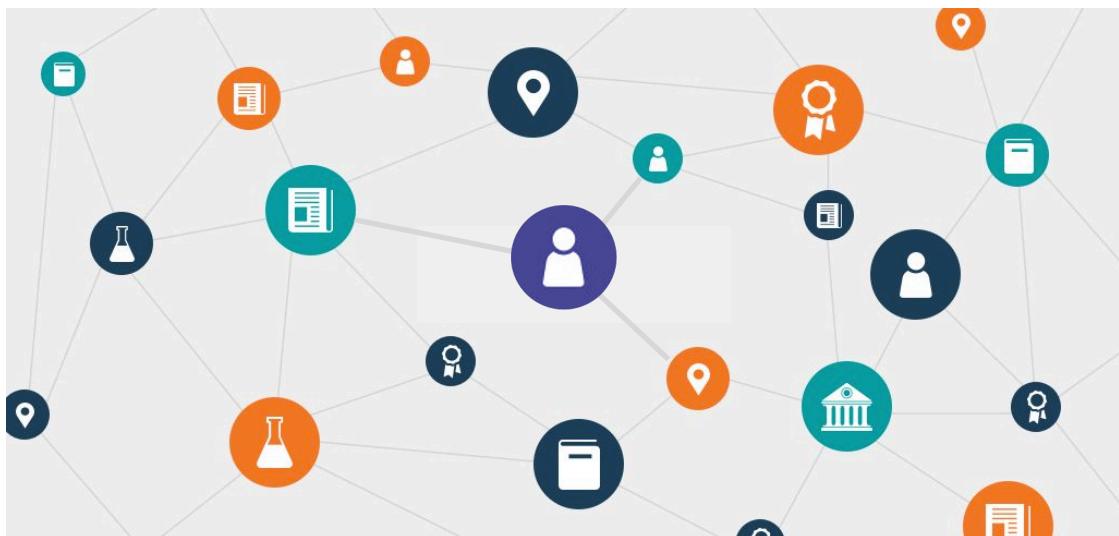
# Graph Attention



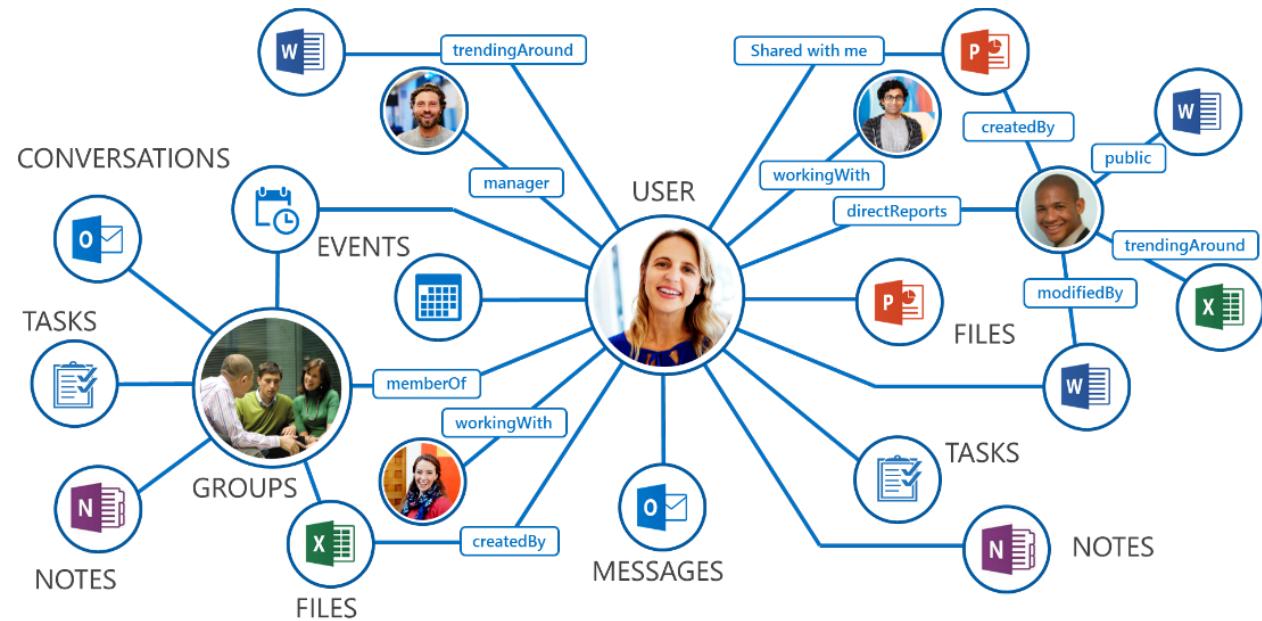
$$\alpha_{v,u} = \frac{\exp(\text{LeakyReLU}(\mathbf{a}^\top [\mathbf{Q}\mathbf{h}_v, \mathbf{Q}\mathbf{h}_u]))}{\sum_{u' \in N(v) \cup \{v\}} \exp(\text{LeakyReLU}(\mathbf{a}^\top [\mathbf{Q}\mathbf{h}_v, \mathbf{Q}\mathbf{h}_{u'}]))}$$

many ways to define attention!

# Attention over Heterogeneous Graphs?



heterogeneous academic graph



heterogeneous office graph

# Heterogeneous Graph Transformer (HGT)

- Current graph neural networks are not capable enough to capture graph **heterogeneity**
- **Heterogeneous Graph Transformer**
  - Unique parameters for each type of relationships



- meta relation of an edge  $e = (s, t)$   
 $\langle \tau(s), \phi(e), \tau(t) \rangle$

# Heterogeneous Graph Transformer (HGT)



- meta relation of an edge  $e = (s, t)$   
 $\langle \tau(s), \phi(e), \tau(t) \rangle$
- heterogeneous mutual attention

$$\text{Attention}_{HGT}(s, e, t) = \text{Softmax} \left( \parallel_{\forall s \in N(t)} \underset{i \in [1, h]}{\text{ATT-head}^i(s, e, t)} \right) \quad (3)$$

$$\text{ATT-head}^i(s, e, t) = \left( K^i(s) \ W_{\phi(e)}^{ATT} \ Q^i(t)^T \right) \cdot \frac{\mu \langle \tau(s), \phi(e), \tau(t) \rangle}{\sqrt{d}}$$

$$K^i(s) = \text{K-Linear}_{\tau(s)}^i \left( H^{(l-1)}[s] \right)$$

$$Q^i(t) = \text{Q-Linear}_{\tau(t)}^i \left( H^{(l-1)}[t] \right)$$

# Heterogeneous Graph Transformer



- meta relation of an edge  $e = (s, t)$   
 $\langle \tau(s), \phi(e), \tau(t) \rangle$
- heterogeneous message passing

$$\textbf{Message}_{HGT}(s, e, t) = \parallel_{i \in [1, h]} MSG\text{-}head^i(s, e, t)$$

$$MSG\text{-}head^i(s, e, t) = \text{M-Linear}_{\tau(s)}^i \left( H^{(l-1)}[s] \right) W_{\phi(e)}^{MSG}$$

# Heterogeneous Graph Transformer

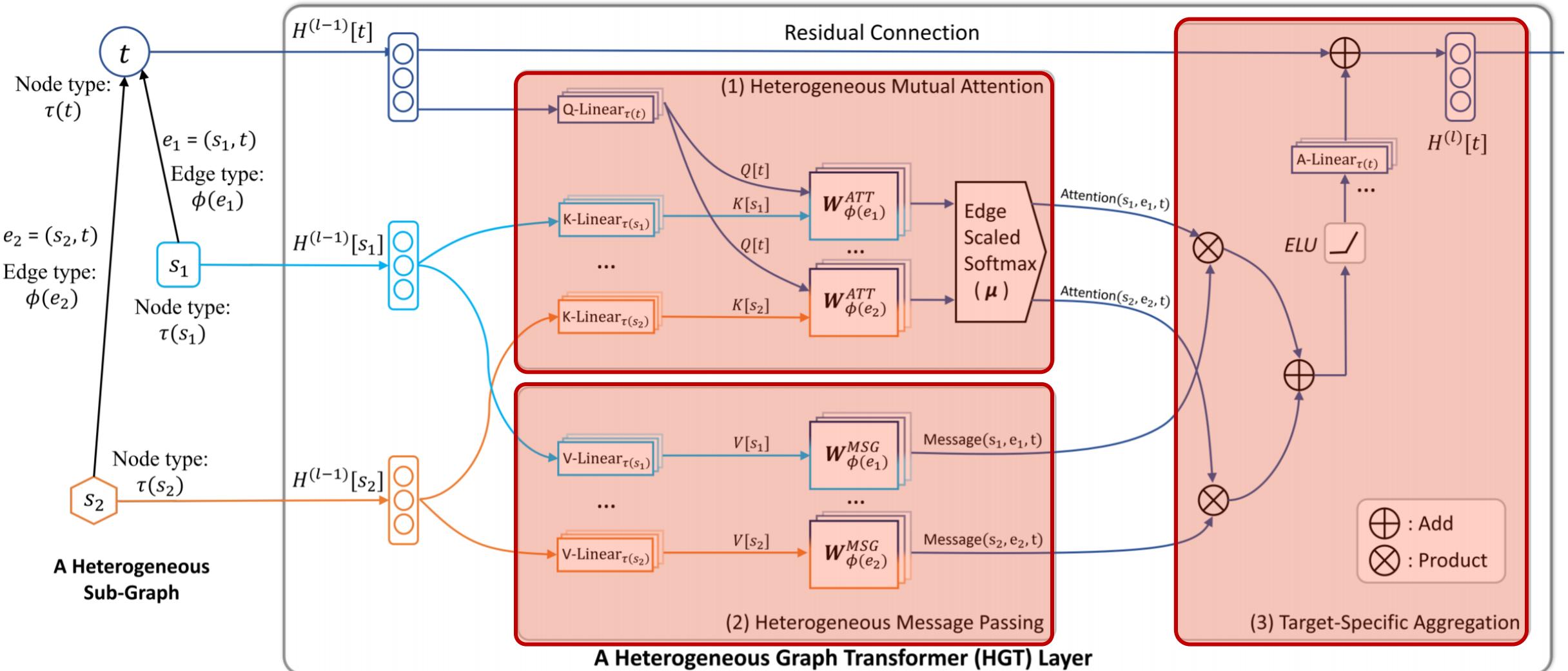


- meta relation of an edge  $e = (s, t)$   
 $\langle \tau(s), \phi(e), \tau(t) \rangle$
- target specific aggregation

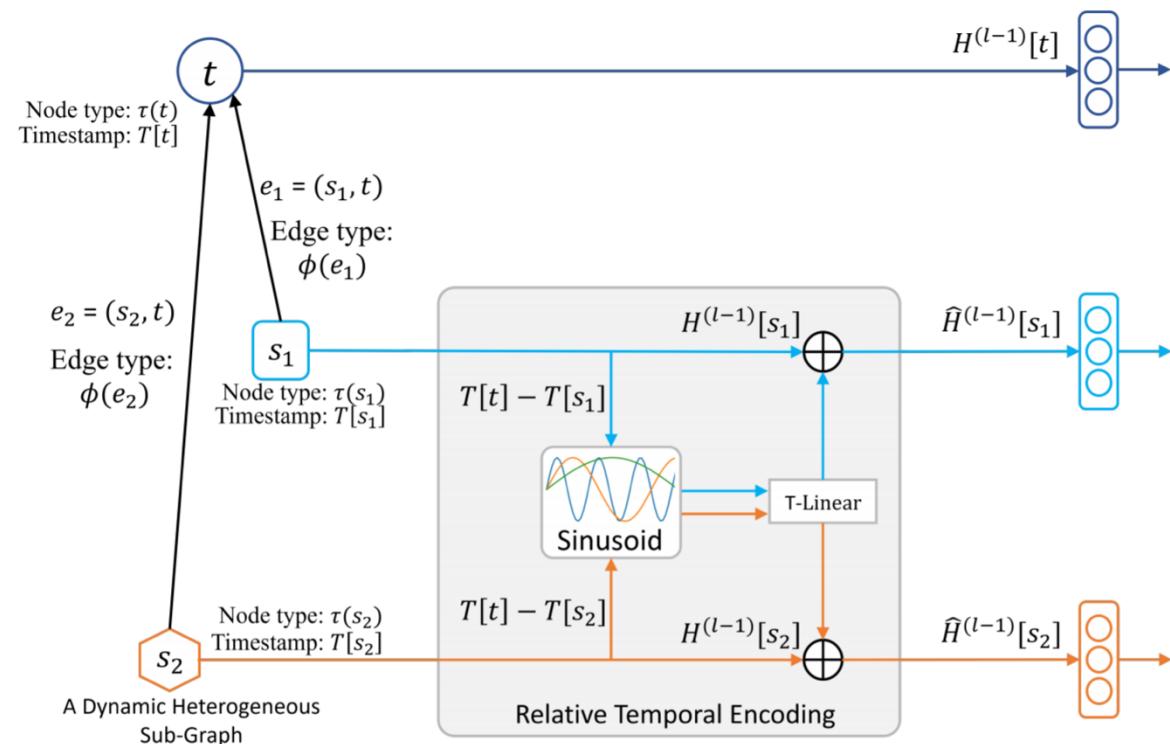
$$\tilde{H}^{(l)}[t] = \bigoplus_{\forall s \in N(t)} \left( \mathbf{Attention}_{HGT}(s, e, t) \cdot \mathbf{Message}_{HGT}(s, e, t) \right)$$

$$H^{(l)}[t] = \text{A-Linear}_{\tau(t)} \left( \sigma(\tilde{H}^{(l)}[t]) \right) + H^{(l-1)}[t]$$

# Heterogeneous Graph Transformer (HGT)



# Heterogeneous Graph Transformer (HGT)



- Relative Temporal Encoding

$$\hat{H}^{(l-1)}[s] = H^{(l-1)}[s] + RTE(\Delta T(t, s))$$

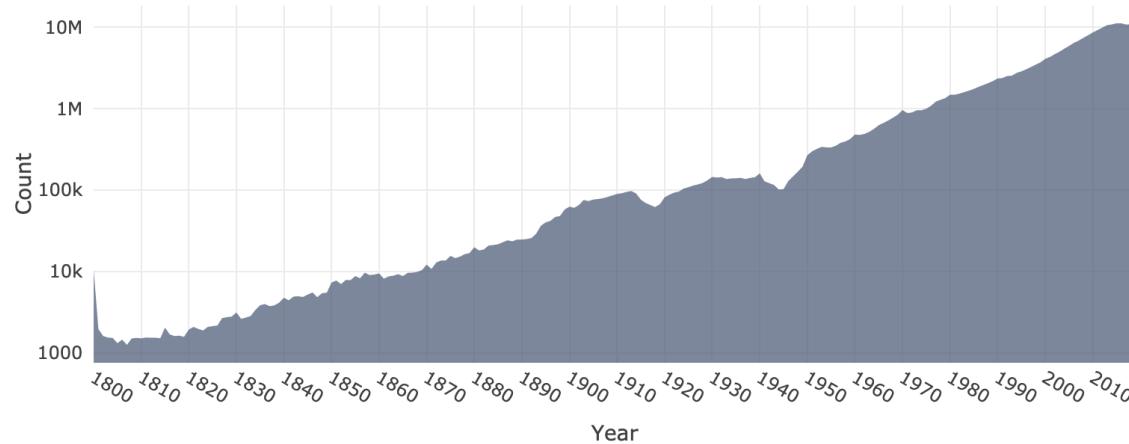
$$RTE(\Delta T(t, s)) = \text{T-Linear}\left(Base(\Delta T_{t,s})\right)$$

$$Base(\Delta T(t, s), 2i) = \sin\left(\Delta T_{t,s}/10000^{\frac{2i}{d}}\right)$$

$$Base(\Delta T(t, s), 2i + 1) = \cos\left(\Delta T_{t,s}/10000^{\frac{2i+1}{d}}\right)$$

# Experiments

	<b>236,963,398</b>
	Publications
	<b>240,667,697</b>
	Authors
	<b>740,048</b>
	Topics
	<b>4,467</b>
	Conferences
	<b>48,876</b>
	Journals
	<b>25,759</b>
	Institutions



- Sampling subgraphs from large-scale graphs
  - From homogeneous graphs → LADIES algorithm
  - From heterogeneous graphs → HGSampling algo

1. Ziniu Hu, et al. Heterogeneous Graph Transformer. **WWW 2020**.

2. Difan Zou, et al. Layer-Dependent Importance Sampling for Training Deep and Large Graph Convolutional Networks. In NeurIPS'19.

# Results

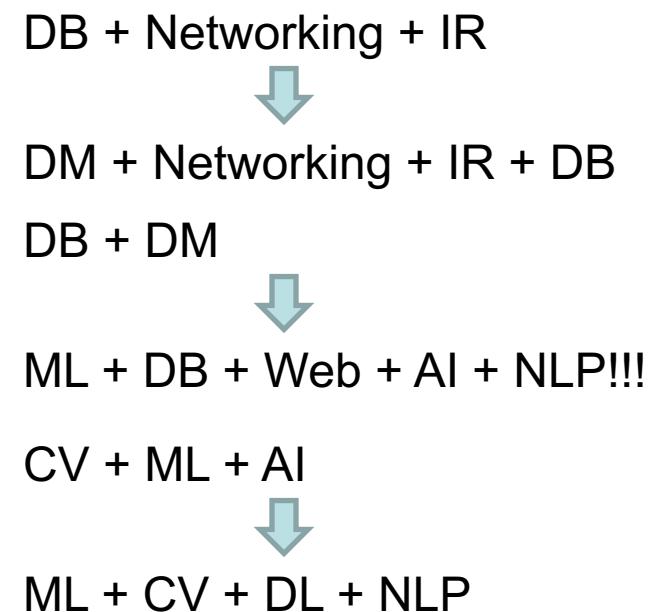
GNN Models		GCN [7]	RGCN [12]	GAT [21]	HetGNN [25]	HAN [22]	HGT <sub>noHeter</sub>	HGT <sub>noTime</sub>	HGT
# of Parameters		1.69M	8.80M	1.69M	8.41M	9.45M	3.12M	7.44M	8.20M
Paper-Field (L1)	NDCG	.558±.141	.563±.128	.601±.103	.615±.084	.617±.096	.674±.086	.702±.089	<b>.735±.084</b>
	MRR	.513±.136	.526±.105	.587±.096	.595±.076	.604±.092	.652±.078	.676±.082	<b>.713±.081</b>
Paper-Field (L2)	NDCG	.241±.074	.258±.046	.276±.049	.271±.062	.281±.051	.301±.046	.307±.052	<b>.332±.048</b>
	MRR	.192±.067	.206±.052	.228±.045	.231±.053	.242±.049	.257±.058	.260±.064	<b>.276±.071</b>
Paper-Venue	NDCG	.303±.066	.354±.051	.461±.057	.447±.071	.478±.062	.515±.059	.538±.064	<b>.551±.062</b>
	MRR	.114±.070	.198±.047	.244±.052	.226±.059	.269±.067	.295±.060	.322±.048	<b>.334±.061</b>
Author Disambiguation	NDCG	.730±.064	.742±.057	.785±.063	.792±.052	.810±.049	.834±.058	.849±.066	<b>.857±.054</b>
	MRR	.762±.042	.786±.048	.843±.044	.852±.058	.876±.056	.903±.041	.911±.043	<b>.918±.048</b>

**HGT offers ~9–21% improvements over existing (heterogeneous) GNNs**

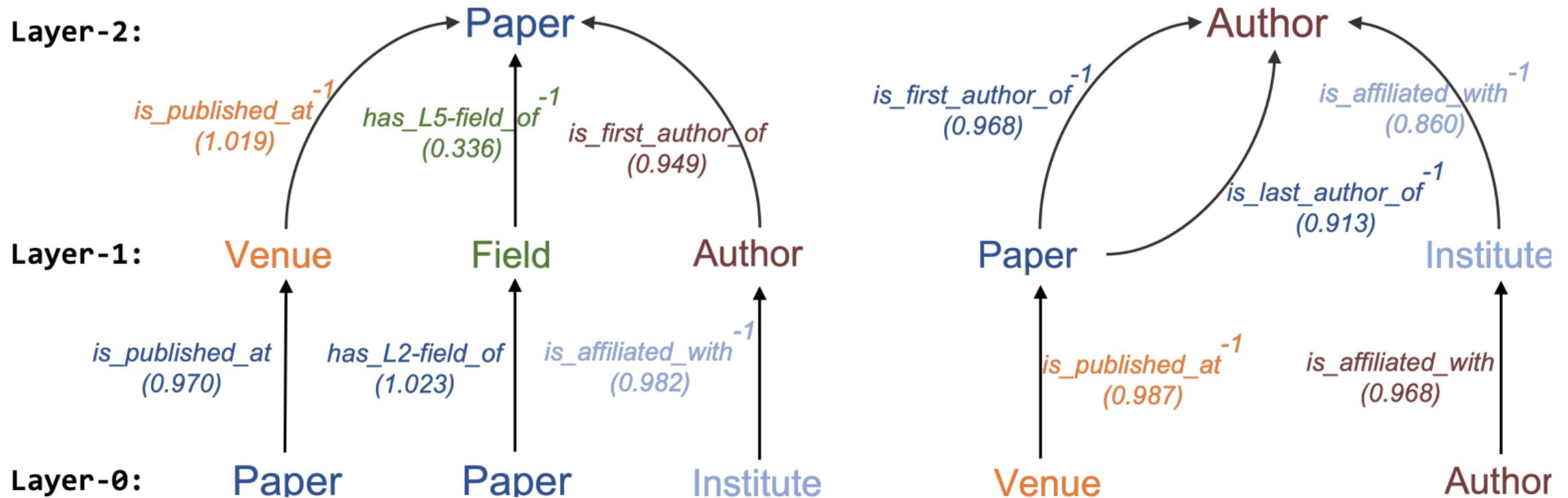
# Case Study

Experiments done in **2019!**

Venue	Time	Top-5 Most Similar Venues
WWW	2000	SIGMOD, VLDB, NSDI, GLOBECOM, SIGIR
	2010	GLOBECOM, KDD, CIKM, SIGIR, SIGMOD
	2020	KDD, GLOBECOM, SIGIR, WSDM, SIGMOD
KDD	2000	SIGMOD, ICDE, ICDM, CIKM, VLDB
	2010	ICDE, WWW, NeurIPS, SIGMOD, ICML
	2020	NeurIPS, SIGMOD, WWW, AAAI, EMNLP
NeurIPS	2000	ICCV, ICML, ECCV, AAAI, CVPR
	2010	ICML, CVPR, ACL, KDD, AAAI
	2020	ICML, CVPR, ICLR, ICCV, ACL

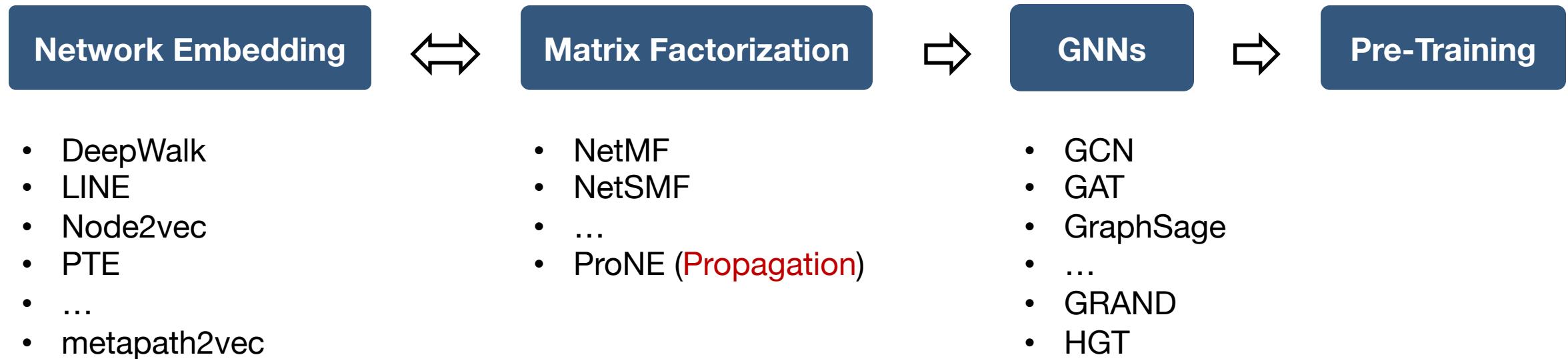


# What is the Best Part of HGT?



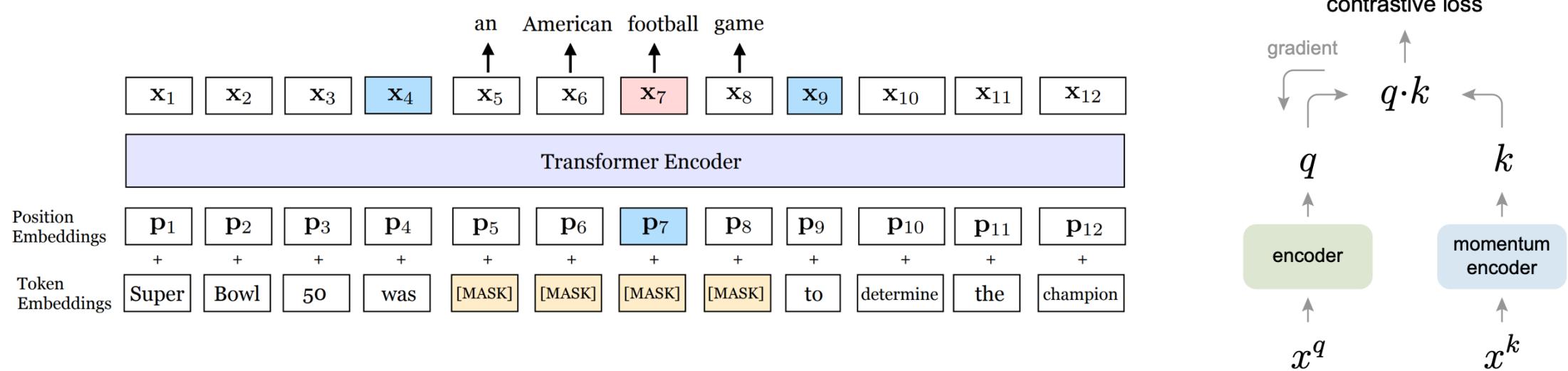
Learn meta-paths & their weights implicitly!

# Graph Representation Learning



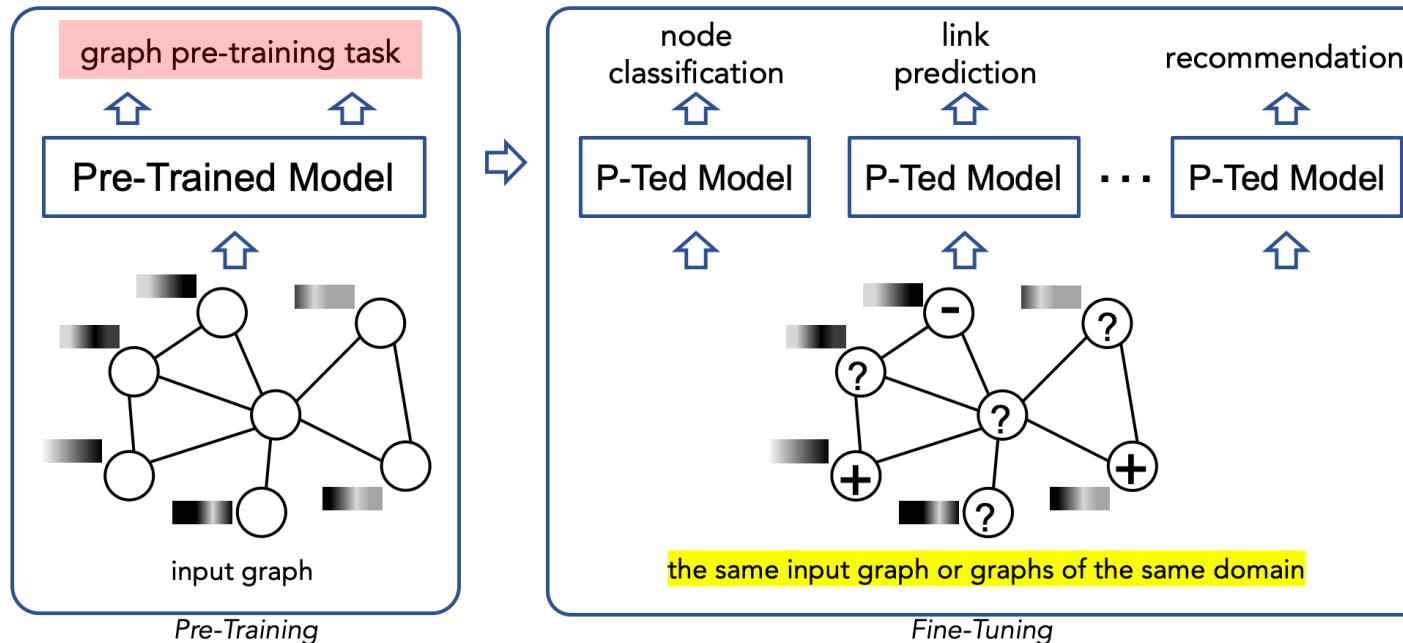
# Language and Image Pre-Training, Graphs?

- Recent progress of pre-training models in NLP & CV
  - ELMO, BERT, XLNet, MoCo, etc.
    - Model level: Transformer
    - **Pre-training Task:** masked language modeling & next sentence prediction



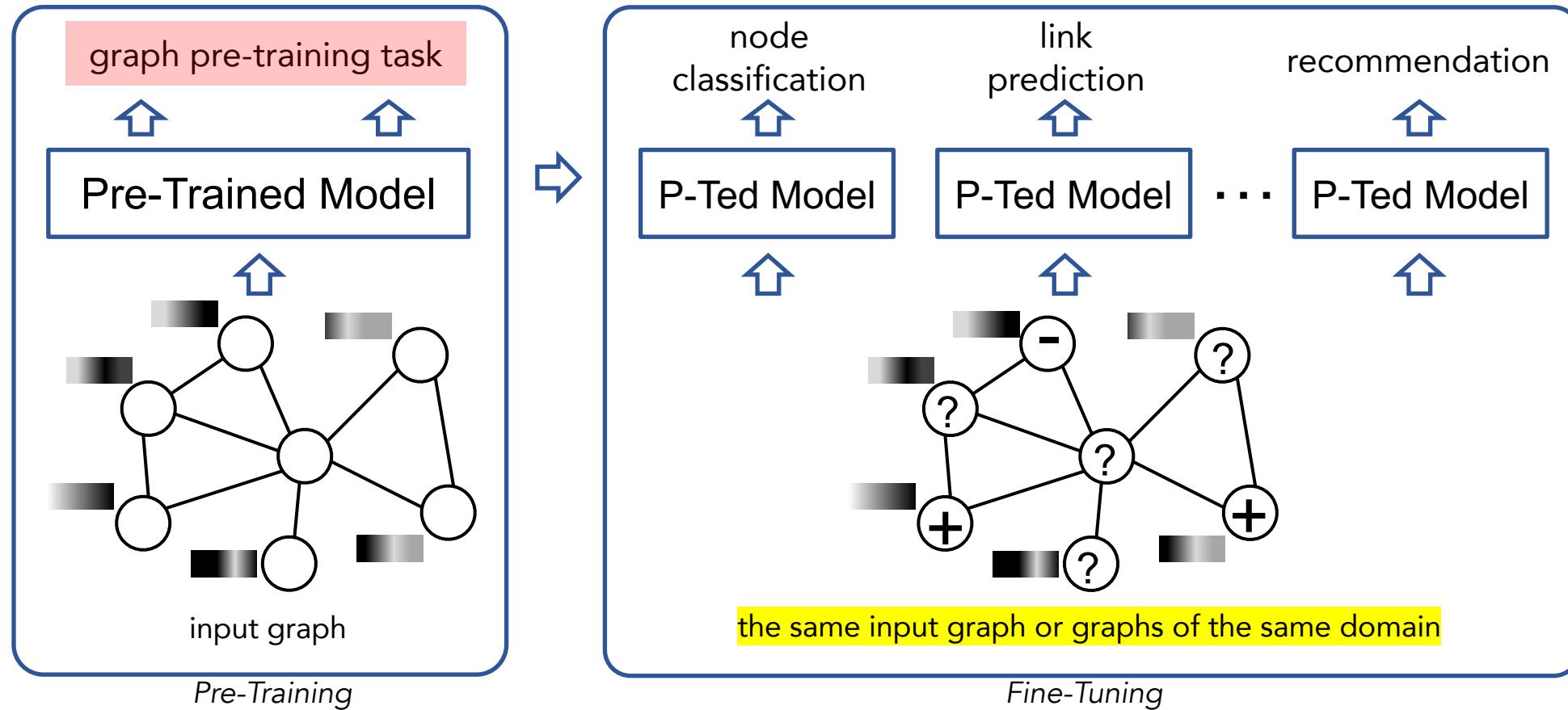
# GNN Pre-Training

- The **FIRST** graph pre-training setting:
  - To pre-train from **one graph**
  - To fine-tune for unseen tasks on **the same graph or graphs of the same domain.**



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

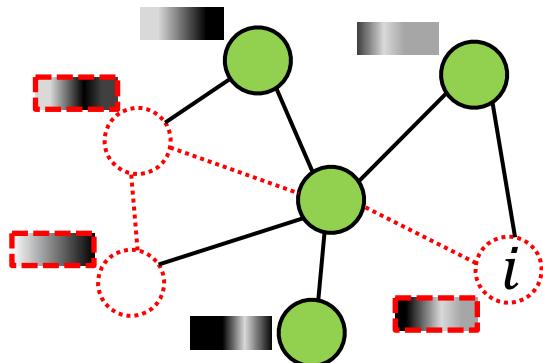
# GNN Pre-Training



# GPT-GNN: Generative Pre-Training of GNNs

- Model the graph distribution  $p(G; \theta)$  by learning to reconstruct the input graph.
  - Factorize the graph likelihood into two terms:
    - Attribute Generation
    - Edge Generation

$$\log p_\theta(X, E) = \sum_{i=1}^{|V|} \log p_\theta(X_i, E_i | X_{<i}, E_{<i}).$$



attribute and edge **masked**  
input graph

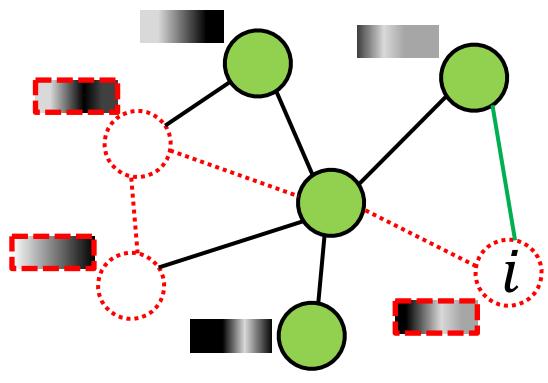
$$\begin{aligned} & p_\theta(X_i, E_i | X_{<i}, E_{<i}) \\ &= p_\theta(X_i | X_{<i}, E_{<i}) \cdot p_\theta(E_i | X_{<i}, E_{<i}) \end{aligned}$$

?

# GPT-GNN: Generative Pre-Training of GNNs

- Model the graph distribution  $p(G; \theta)$  by learning to reconstruct the input graph.
  - Factorize the graph likelihood into two terms:
    - Attribute Generation
    - Edge Generation

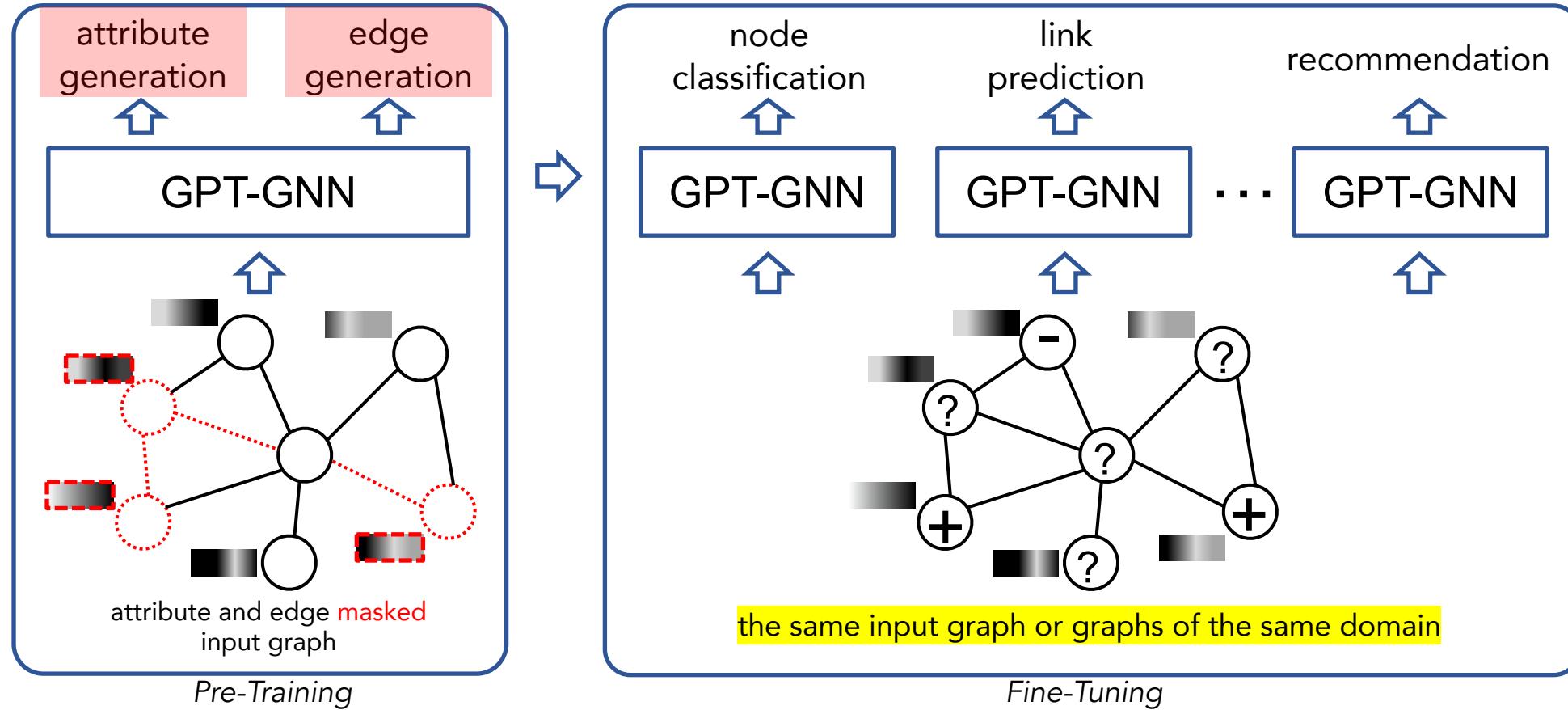
$$\log p_\theta(X, E) = \sum_{i=1}^{|V|} \log p_\theta(X_i, E_i | X_{<i}, E_{<i}).$$



attribute and edge **masked**  
input graph

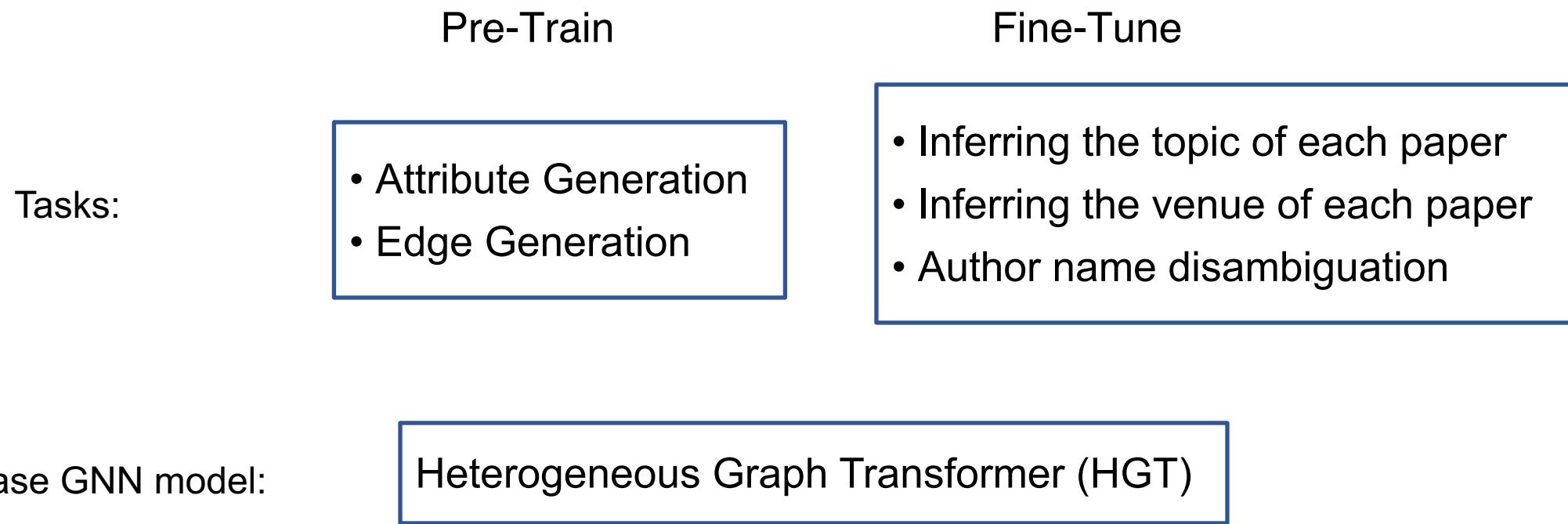
$$\begin{aligned} & p_\theta(X_i, E_i | X_{<i}, E_{<i}) \\ &= \sum_o p_\theta(X_i, E_{i,\neg o} | E_{i,o}, X_{<i}, E_{<i}) \cdot p_\theta(E_{i,o} | X_{<i}, E_{<i}) \\ &= \mathbb{E}_o \left[ p_\theta(X_i, E_{i,\neg o} | E_{i,o}, X_{<i}, E_{<i}) \right] \\ &= \mathbb{E}_o \left[ \underbrace{p_\theta(X_i | E_{i,o}, X_{<i}, E_{<i})}_{\text{1) generate attributes}} \cdot \underbrace{p_\theta(E_{i,\neg o} | E_{i,o}, X_{\leq i}, E_{<i})}_{\text{2) generate edges}} \right]. \end{aligned}$$

# GPT-GNN: Generative Pre-Training of GNNs



# GPT-GNN: Generative Pre-Training of GNNs

- Data 1: Open Academic Graph

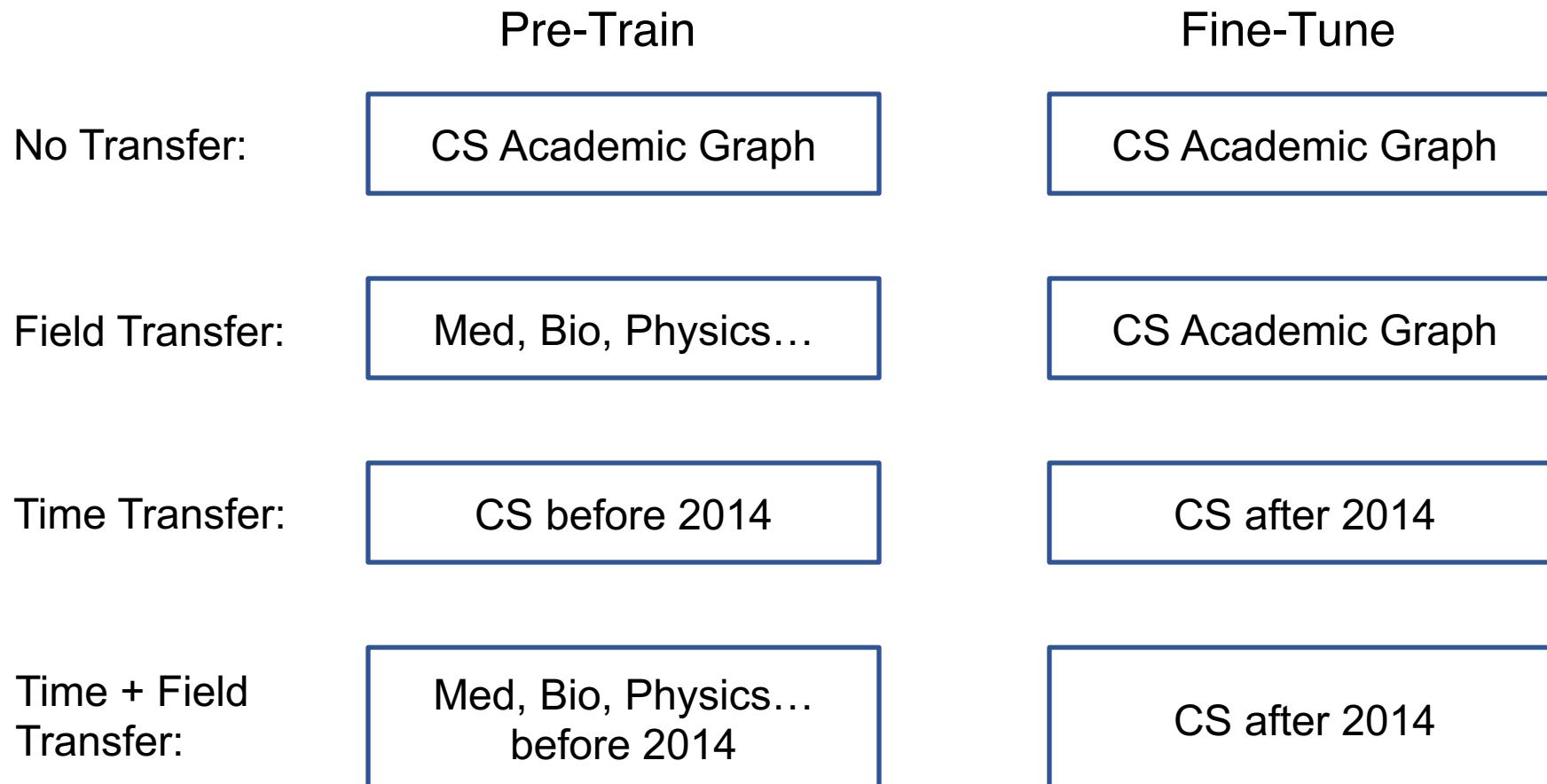


1.Ziniu Hu et al. GPT-GNN: Generative Pre-Training of Graph Neural Networks. **KDD** 2020.

2.Code & data for GPT-GNN: <https://github.com/acbull/GPT-GNN>

# GPT-GNN: Generative Pre-Training of GNNs

- Data 1: Open Academic Graph



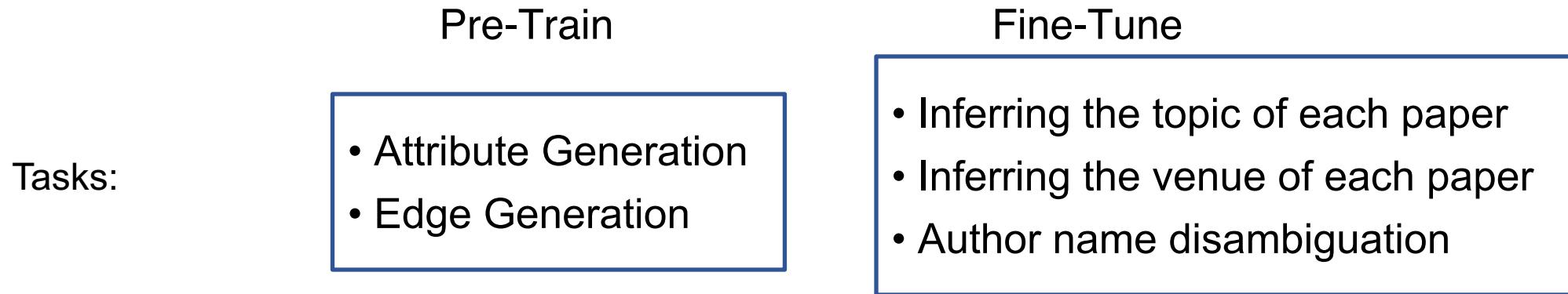
# GPT-GNN: Generative Pre-Training of GNNs

Downstream Dataset	OAG			
Evaluation Task	Paper-Field	Paper-Venue	Author ND	
Field Transfer	No Pre-train	.346±.149	.598±.122	.813±.105
	GAE	.403±.114	.626±.093	.836±.084
	GraphSAGE (unsp.)	.368±.125	.609±.096	.818±.092
	Graph Infomax	.387±.112	.612±.097	.827±.084
	GPT-GNN (Attr)	.396±.118	.623±.105	.834±.086
	GPT-GNN (Edge)	.413±.109	.635±.096	.842±.093
	GPT-GNN	<b>.420±.107</b>	<b>.641±.098</b>	<b>.848±.102</b>
Time Transfer	GAE	.384±.117	.619±.101	.828±.095
	GraphSAGE (unsp.)	.352±.121	.601±.105	.815±.093
	Graph Infomax	.369±.116	.606±.102	.821±.089
	GPT-GNN (Attr)	.374±.114	.614±.098	.826±.089
	GPT-GNN (Edge)	.397±.105	.629±.102	.836±.088
	GPT-GNN	<b>.405±.108</b>	<b>.635±.101</b>	<b>.840±.093</b>
	GAE	.371±.124	.611±.108	.821±.102
Time + Field Transfer	GraphSAGE (unsp.)	.349±.130	.602±.118	.812±.097
	Graph Infomax	.360±.121	.600±.102	.815±.093
	GPT-GNN (Attr)	.364±.115	.609±.103	.824±.094
	– (w/o node separation)	.347±.128	.601±.102	.813±.108
	GPT-GNN (Edge)	.390±.116	.622±.104	.830±.105
	– (w/o adaptive queue)	.376±.121	.617±.115	.828±.104
	GPT-GNN	<b>.397±.112</b>	<b>.628±.108</b>	<b>.833±.102</b>

- All pre-training frameworks help the performance of GNNs
  - GAE, GraphSage, Graph Infomax
  - GPT-GNN
- GPT-GNN helps the most by achieving a relative performance gain of 9.1% over the base model without pre-training
- Both self-supervised tasks in GPT-GNN help the pre-training framework
  - Attribute generation
  - Edge generation

# GPT-GNN: Generative Pre-Training of GNNs

- Data 1: Open Academic Graph



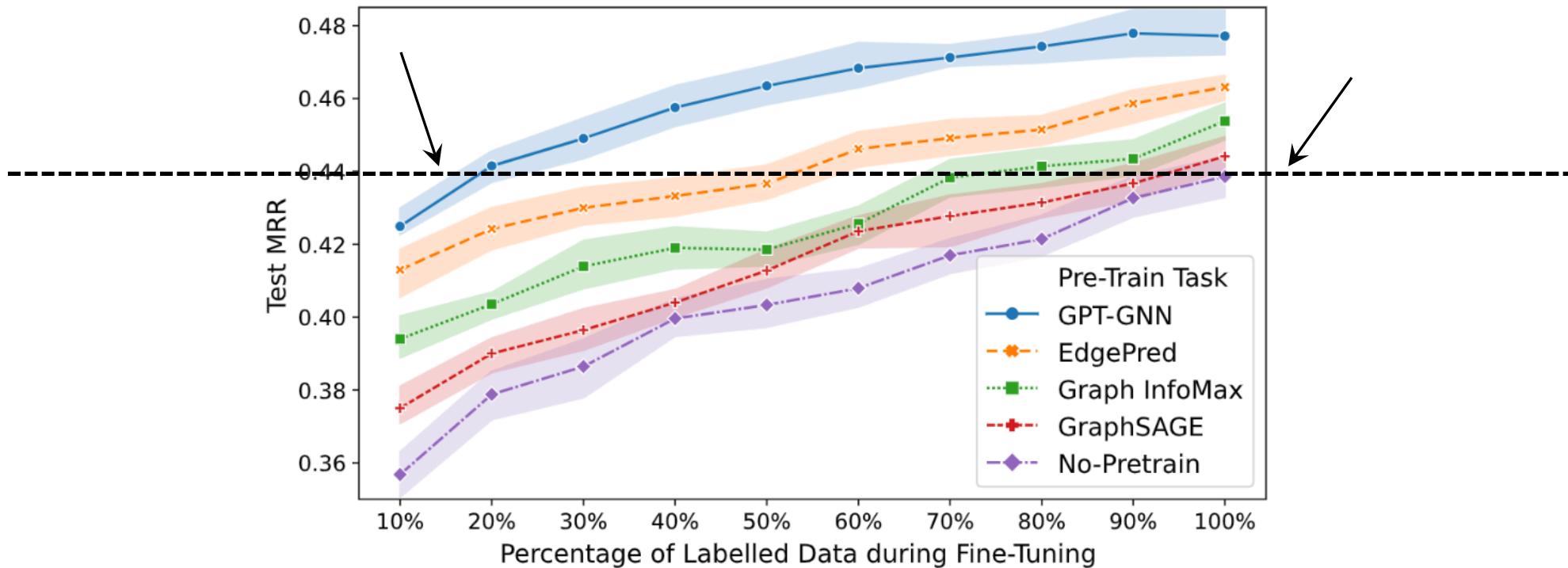
Base GNN model:

Heterogeneous Graph Transformer (HGT)

Model	HGT	GCN	GAT	RGCN	HAN
No Pre-train	.346	.327	.318	.296	.332
GPT-GNN	.420	.359	.382	.351	.406
Relative Gain	21.4%	9.8%	20.1%	18.9%	22.3%

# The Promise of Graph Pre-Training!

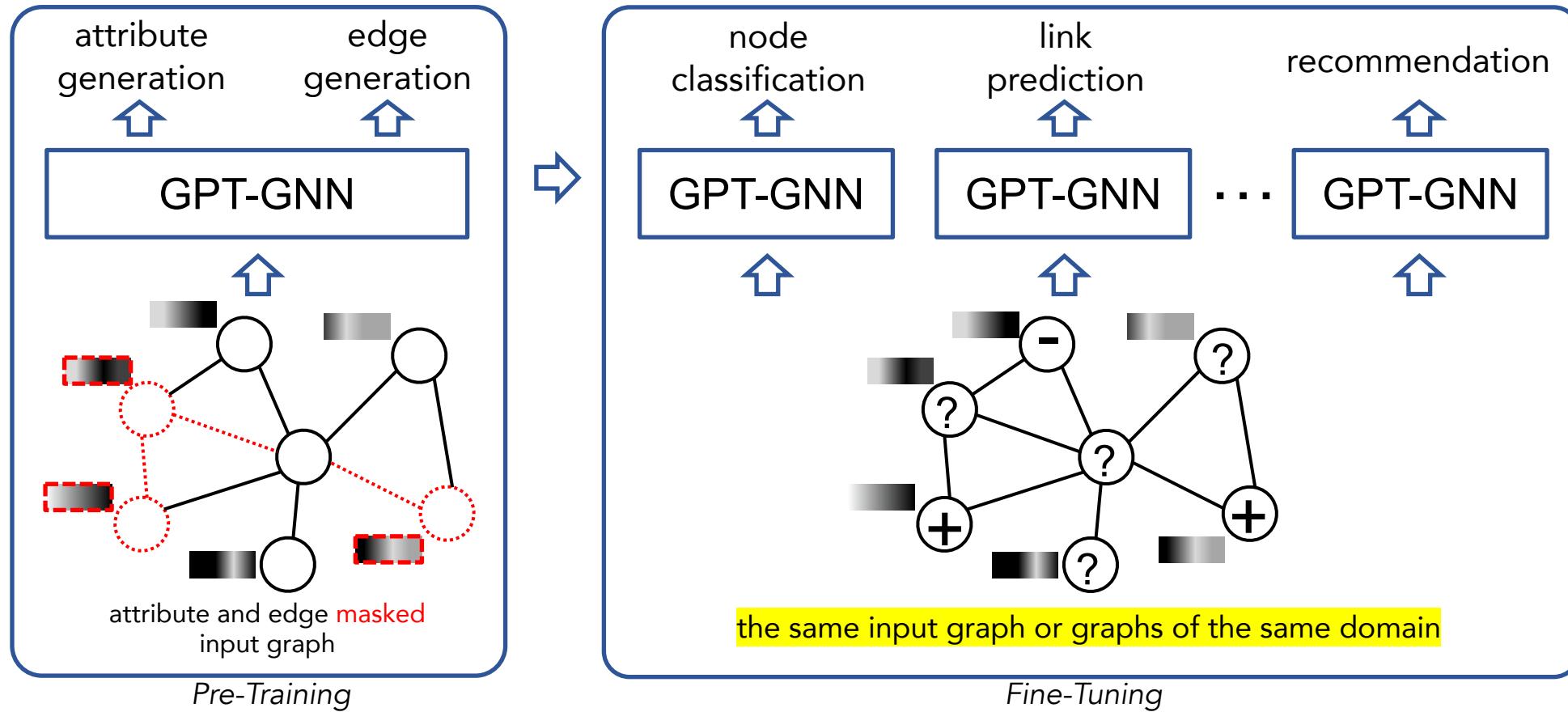
- During fine-turning



The GNN model **w/o** pre-training by using **100% training data**  
VS

The GNN model **with** pre-training by using **10-20% training data**

# GNN Pre-Training on the “Same” Networks



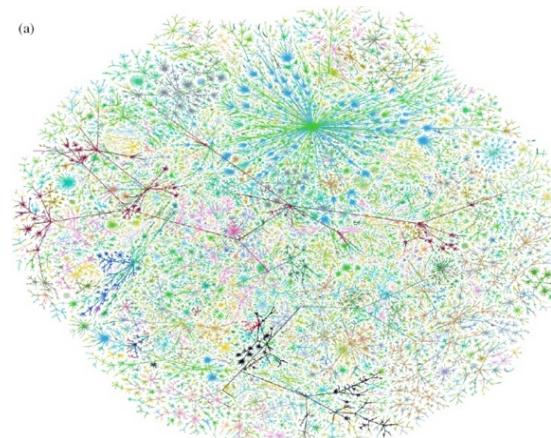
# Graphs



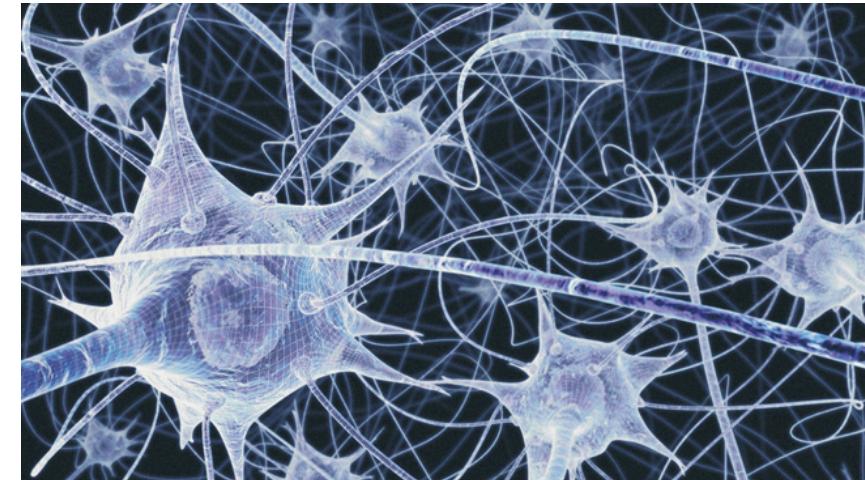
## Office/Social Graph



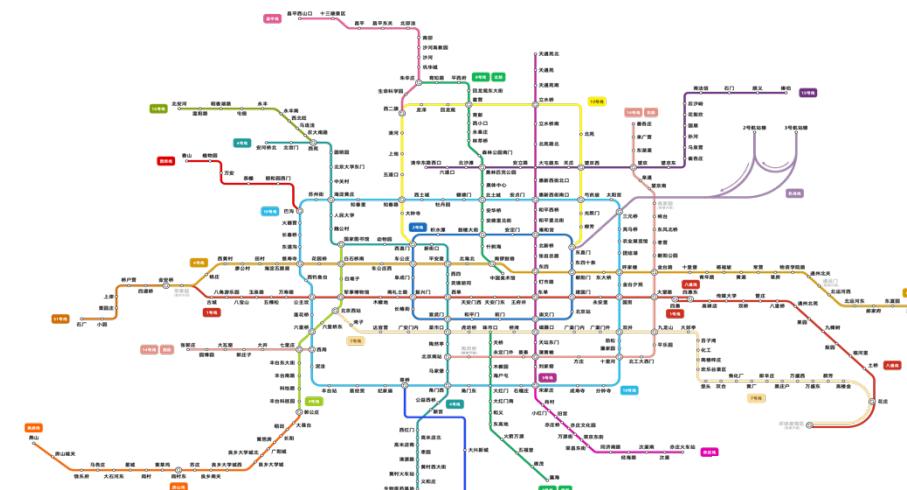
# Knowledge Graph



## Internet



# Biological Neural Networks

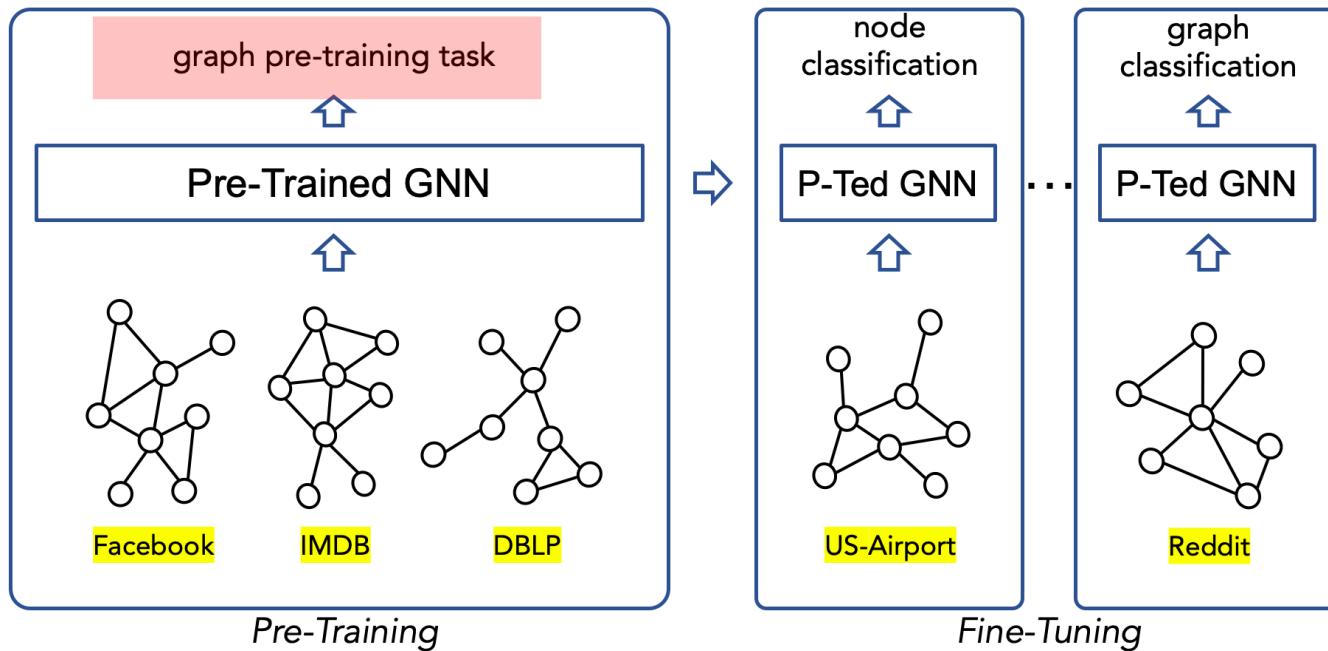


## Transportation

figure credit: Web

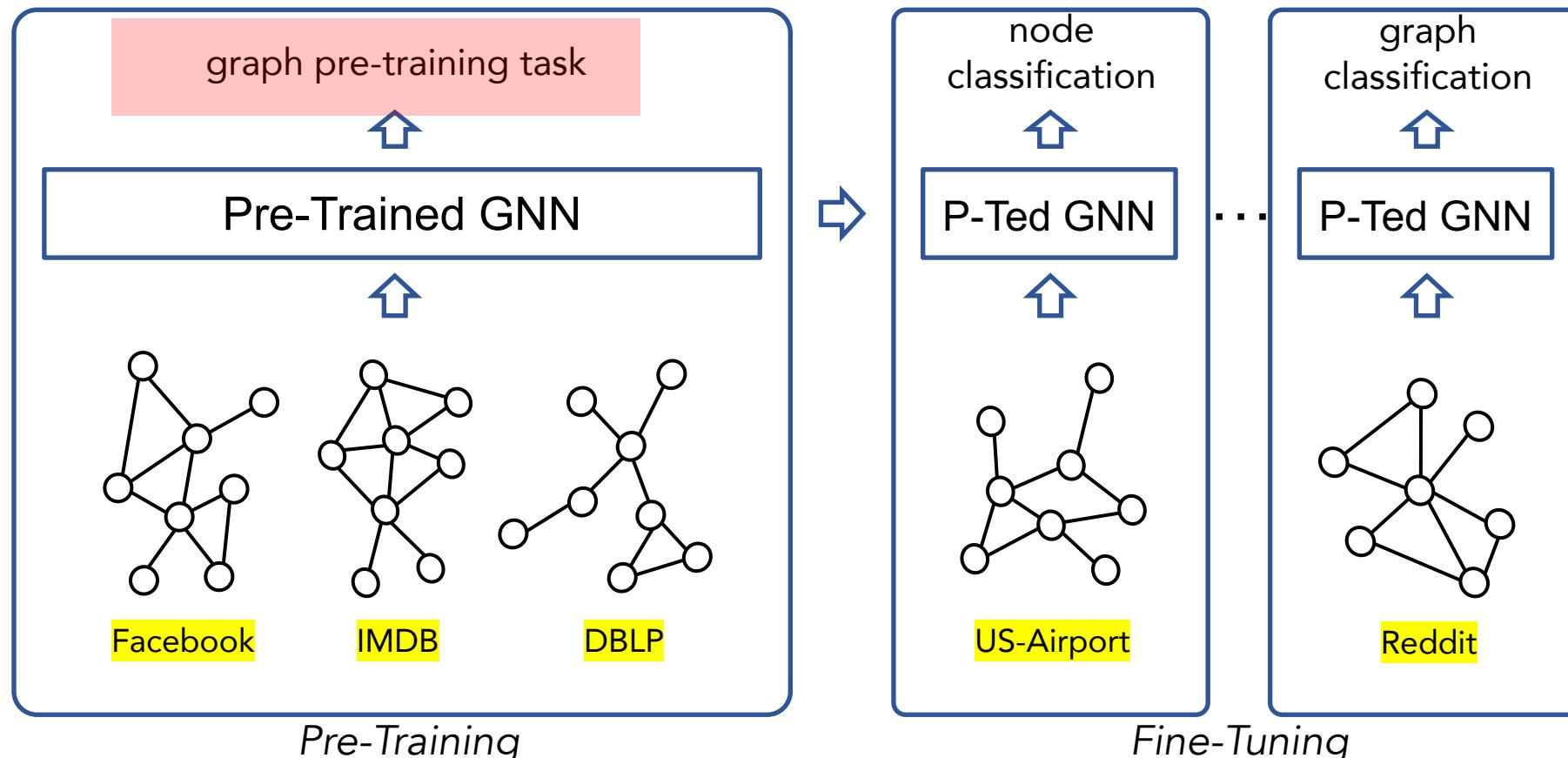
# GNN Pre-Training

- The **SECOND** 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



# 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?

1. Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In CVPR '18.

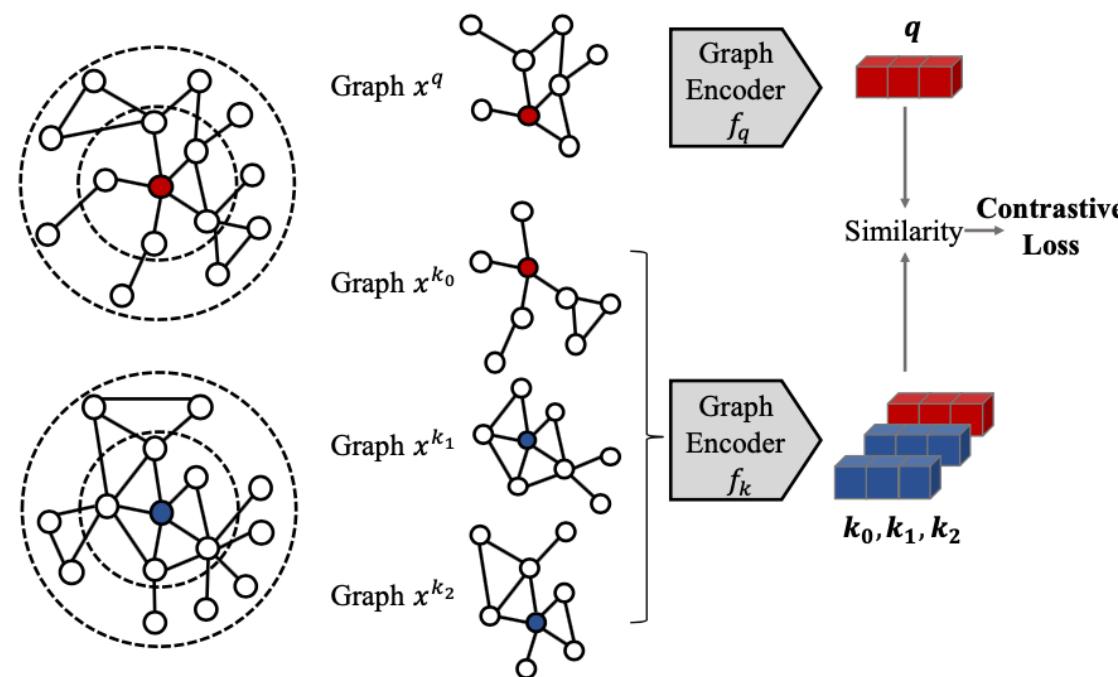
2. Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. 2020. Momentum contrast for unsupervised visual representation learning. In CVPR '20.

# Graph Contrastive Coding (GCC)

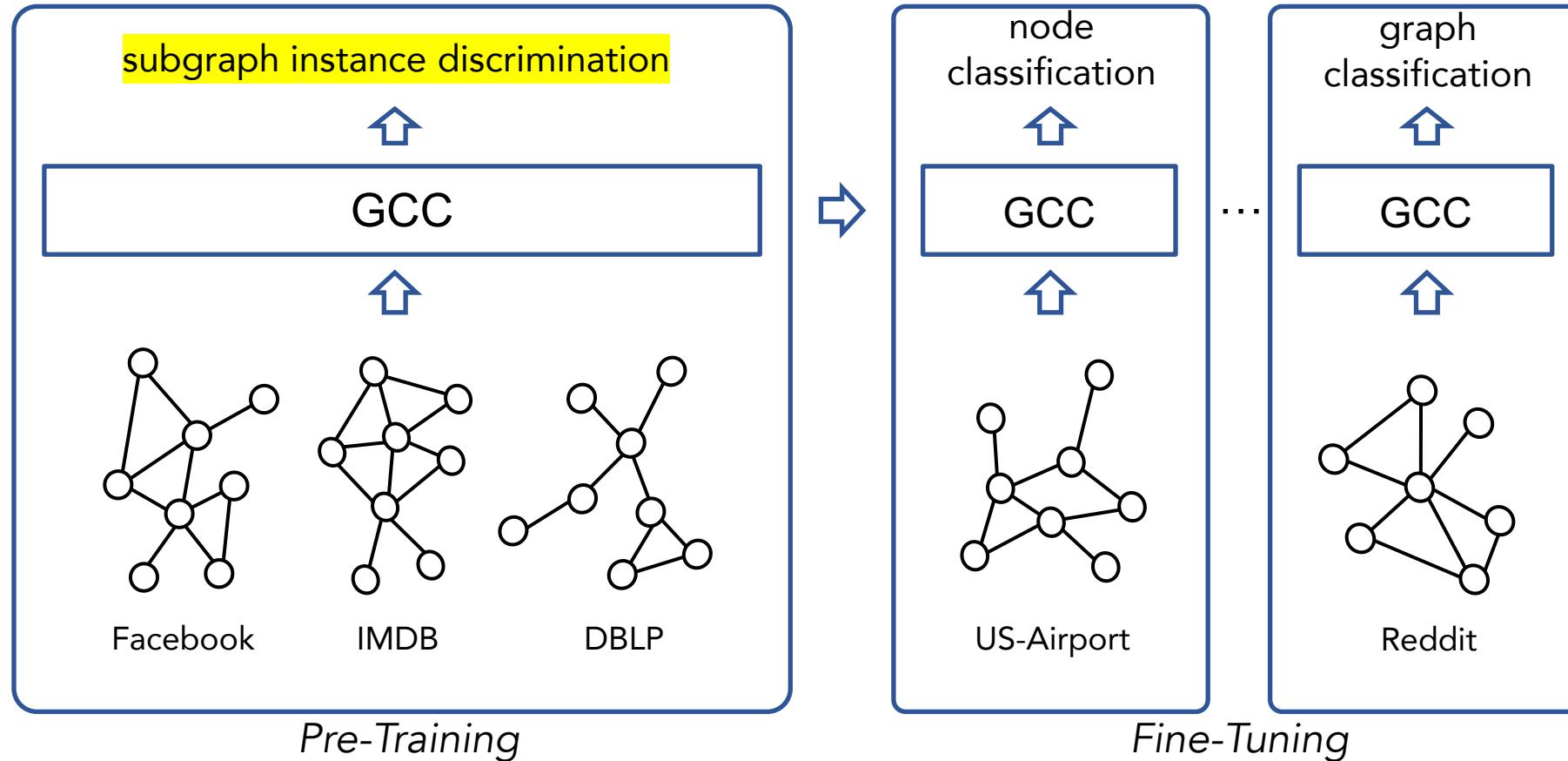
- Contrastive learning for graphs
  - Q1: How to define instances in graphs?
  - Q2: How to define (dis) similar instance pairs in and across graphs?
  - Q3: What are the proper graph encoders?

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

Subgraph instance discrimination



# Graph Contrastive Coding (GCC)

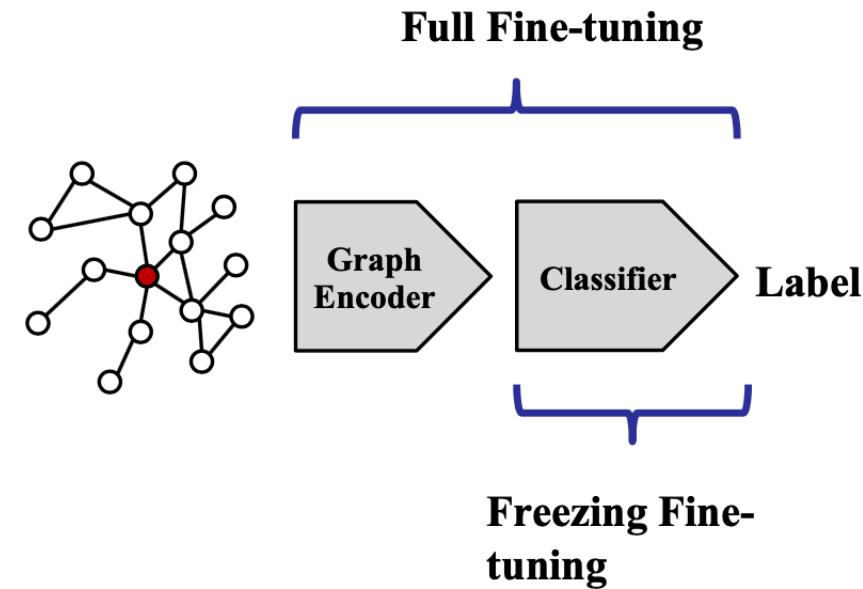


# GCC Pre-Training / Fine-Tuning

- pre-train on six graphs

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, IMDB-M
    - Graph classification
  - AMiner academic graph
    - Similarity search
- The base GNN
  - Graph Isomorphism Network (GIN)



1.Jiezhang 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>

# Results

## Node Classification

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>

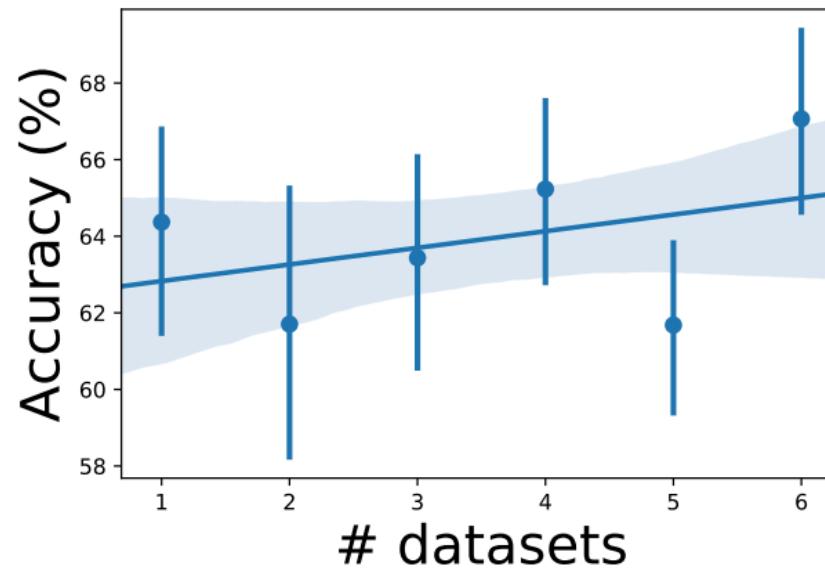
## Similarity Search

	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
# ground 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

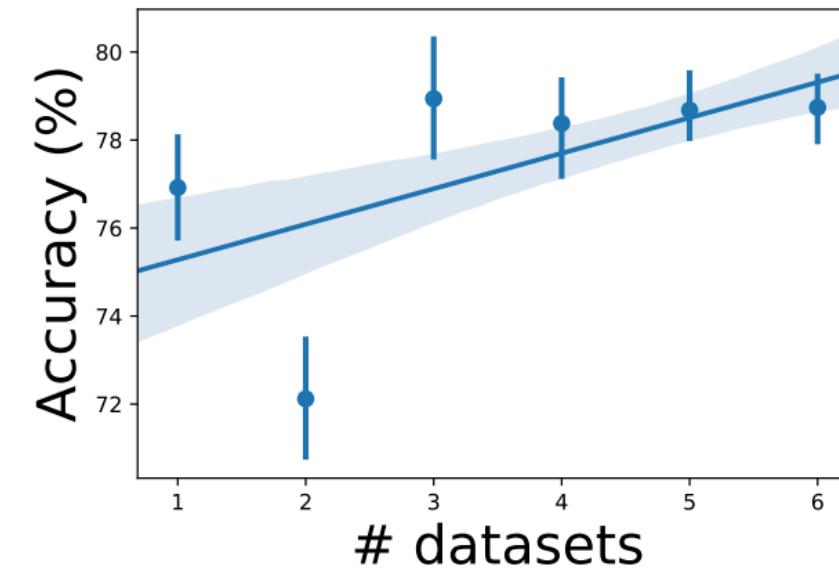
## Graph Classification

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

# Results

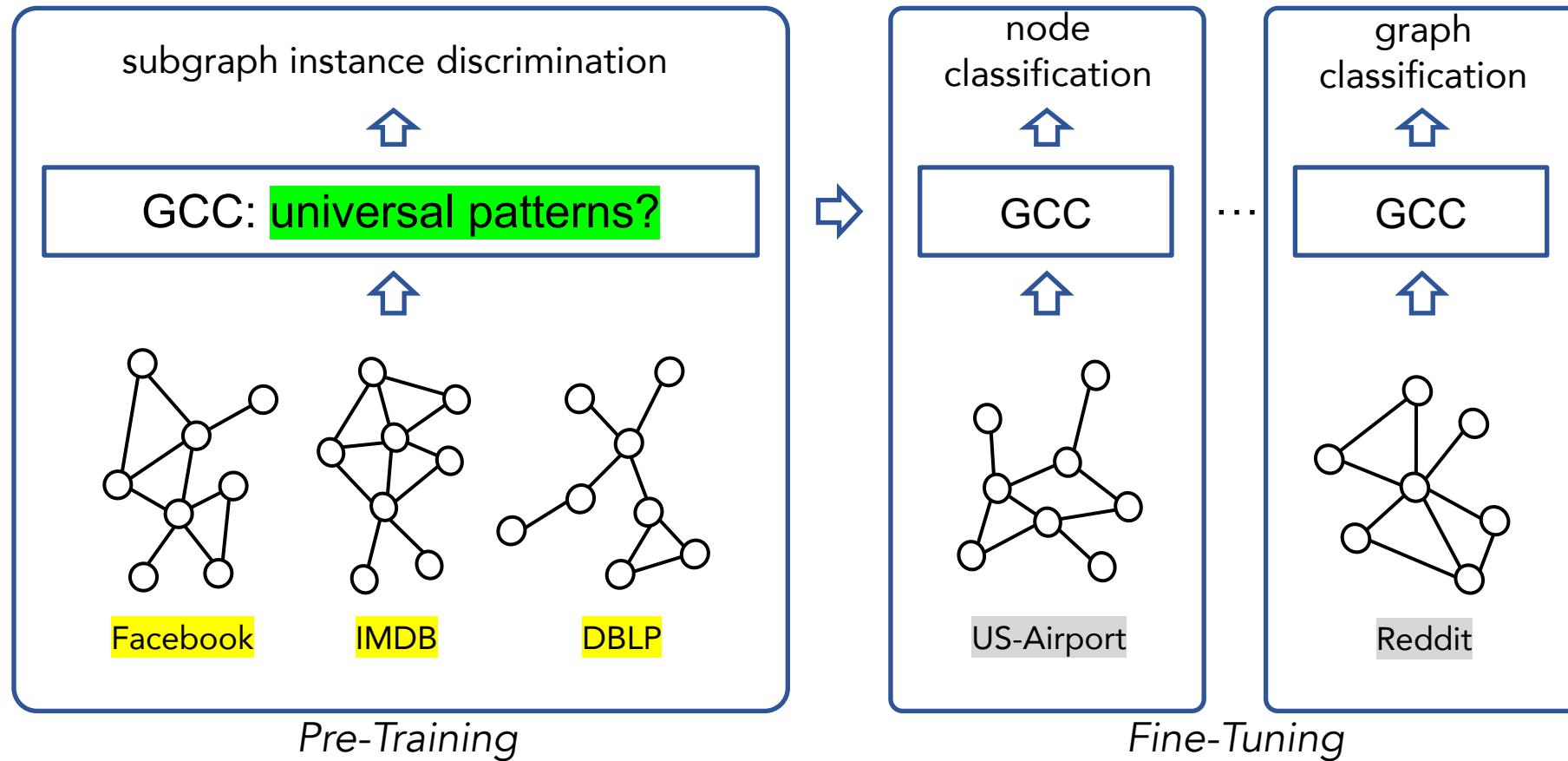


(a) **US-Airport**:  $y = 0.4344x + 62.394$



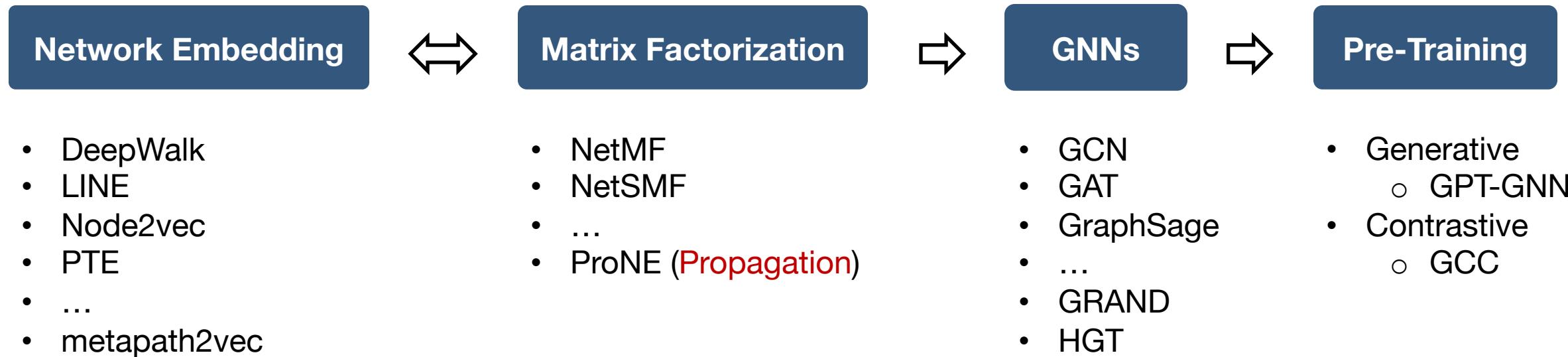
(b) **COLLAB**:  $y = 0.8065x + 74.4737$

# Results



Does the pre-training of GNNs learn the **universal structural patterns** across networks?

# Graph Representation Learning



What graph data to use?



Realistic datasets



Flexible data loaders



Unified evaluation

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

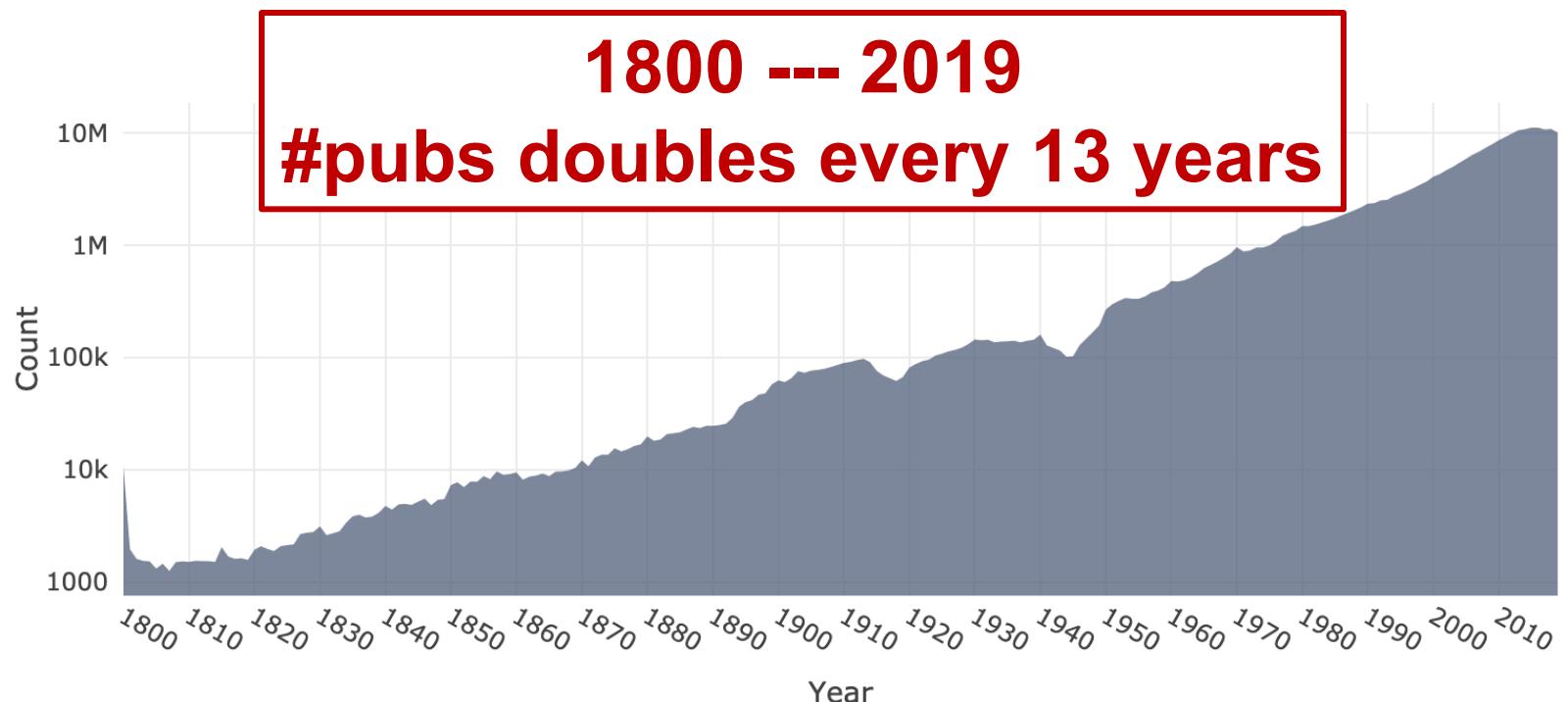
## Leaderboard for ogbn-arxiv

The classification accuracy on the test set. The higher, the better.

Rank	Method	Accuracy	Contact	References	Date
1	GCN	$0.7174 \pm 0.0029$	Matthias Fey – OGB team	Paper, Code	May 1, 2020
2	GraphSAGE	$0.7149 \pm 0.0027$	Matthias Fey – OGB team	Paper, Code	May 1, 2020
3	Node2vec	$0.7007 \pm 0.0013$	Matthias Fey – OGB team	Paper, Code	May 1, 2020
4	MLP	$0.5550 \pm 0.0023$	Matthias Fey – OGB team	Paper, Code	May 1, 2020

# Microsoft Academic Graph & AMiner & OAG

	236,963,398
	Publications
	240,667,697
	Authors
	740,048
	Topics
	4,467
	Conferences
	48,876
	Journals
	25,759
	Institutions



1. Kuansan Wang et al. Microsoft Academic Graph: When experts are not enough. Quantitative Science Studies 1 (1), 2020
2. Fanjin Zhang et al. OAG: Toward Linking Large-scale Heterogeneous Entity Graphs. KDD 2019.
3. Jie Tang et al. Arnetminer: extraction and mining of academic social networks. In KDD 2008.

# References

1. Ziniu Hu et al. GPT-GNN: Generative Pre-Training of Graph Neural Networks. **KDD** 2020.
2. Jiezhong Qiu et al. GCC: Graph Contrastive Coding for Graph Neural Network Pre-Training. **KDD** 2020.
3. Kuansan Wang et al. Microsoft Academic Graph: When experts are not enough. *Quantitative Science Studies* 1 (1), 396-413, 2020.
4. Weihua Hu et al. Open Graph Benchmark: Datasets for Machine Learning on Graphs. arXiv 2020.
5. Feng et al. Graph Random Neural Networks. arXiv 2020.
6. Ziniu Hu et al. Heterogeneous Graph Transformer. **WWW** 2020.
7. Yuxiao Dong et al. Heterogeneous Network Representation Learning. **IJCAI** 2020.
8. Jiezhong Qiu et al. *NetSMF: Large-Scale Network Embedding as Sparse Matrix Factorization*. **WWW** 2019.
9. Jie Zhang et al. *ProNE: Fast and Scalable Network Representation Learning*. **IJCAI** 2019.
10. Fanjing Zhang et al. OAG: Toward Linking Large-scale Heterogeneous Entity Graphs. ACM **KDD** 2019.
11. Xian Wu et al. Neural Tensor Decomposition. **WSDM** 2019.
12. Jiezhong Qiu et al. *Network Embedding as Matrix Factorization: Unifying DeepWalk, LINE, PTE, and node2vec*. **WSDM** 2018.
13. Yuxiao Dong et al. *metapath2vec: Scalable Representation Learning for Heterogeneous Networks*. **KDD** 2017.
14. Perozzi et al. DeepWalk: Online learning of social representations. In **KDD'14**.
15. Tang et al. LINE: Large scale information network embedding. In **WWW'15**.
16. Grover and Leskovec. node2vec: Scalable feature learning for networks. In **KDD'16**.
17. Harris, Z. (1954). Distributional structure. *Word*, 10(23): 146-162.
18. Kipf et al. Semisupervised Classification with Graph Convolutional Networks. **ICLR** 2017
19. Velickovic et al. Graph Attention Networks. **ICLR** 2018
20. Hamilton et al. Inductive Representation Learning on Large Graphs. **NeurIPS** 2017
21. Defferrard et al. Convolutional Neural Networks on Graphs with Fast Localized Spectral Filtering. In **NeurIPS** 2016
22. Jacob Devlin et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. **NAACL-HLT** 2019.
23. Justin Gilmer, et al. Neural message passing for quantum chemistry. arXiv: 2017.
24. Kaiming He, et al. Momentum contrast for unsupervised visual representation learning. arXiv: 2019
25. Tomas Mikolov et al. Distributed representations of words and phrases and their compositionality. **NeurIPS** 2013.
26. Petar Velickovic et al. Deep Graph Infomax. In **ICLR** 19.
27. Zhen Yang et al. Understanding Negative Sampling in Graph Representation Learning. **KDD** 2020.

# Thank you!



**Jiezhong Qiu**  
Tsinghua  
(Jie Tang)



**Ziniu Hu**  
UCLA  
(Yizhou Sun)



**Hongxia Yang**  
Alibaba



**Jing Zhang**  
Renmin U. of China



**Jie Tang**  
Tsinghua



**Yizhou Sun**  
UCLA



**Hao Ma**  
Facebook AI



**Kuansan Wang**  
Microsoft Research

Papers & data & code available at <https://ericdongyangx.github.io/>  
[ericdongyangx@gmail.com](mailto:ericdongyangx@gmail.com)