

DEEP LEARNING

Lecture 12: Deep Learning on Graphs

Dr. Yang Lu

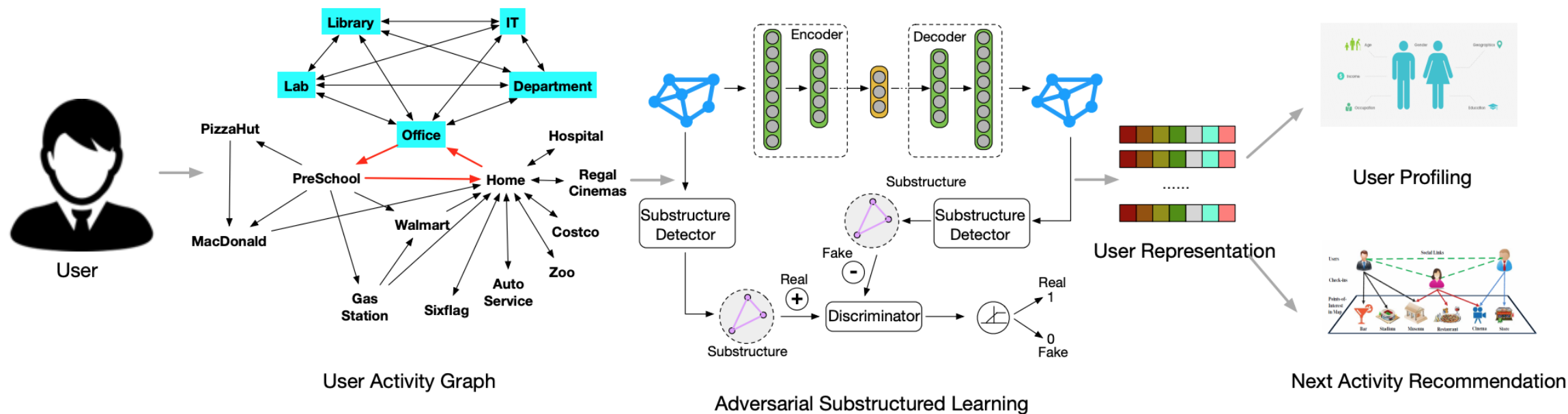
Department of Computer Science and Technology

luyang@xmu.edu.cn



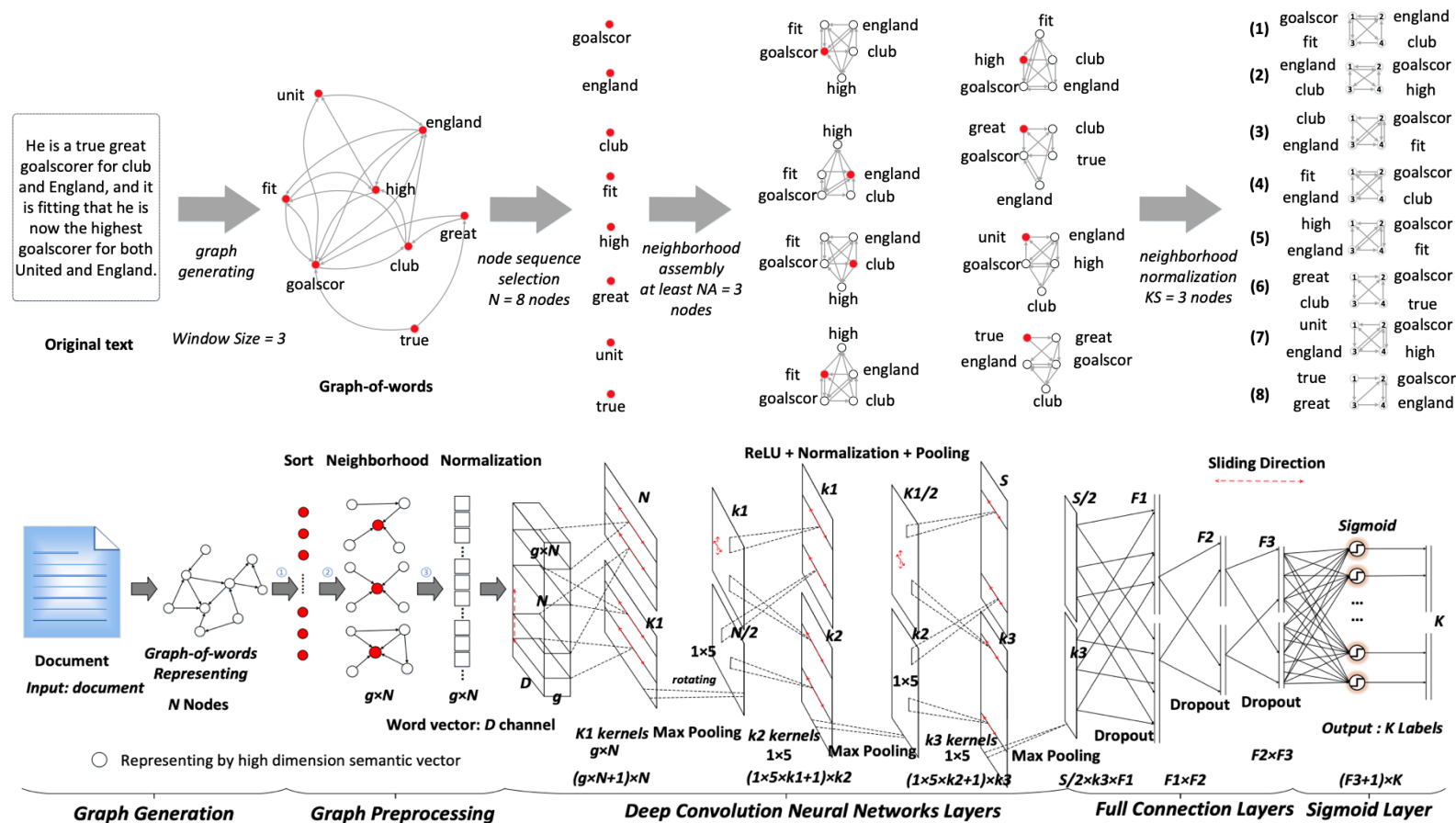
GNN Applications

■ User profiling



GNN Applications

Text classification

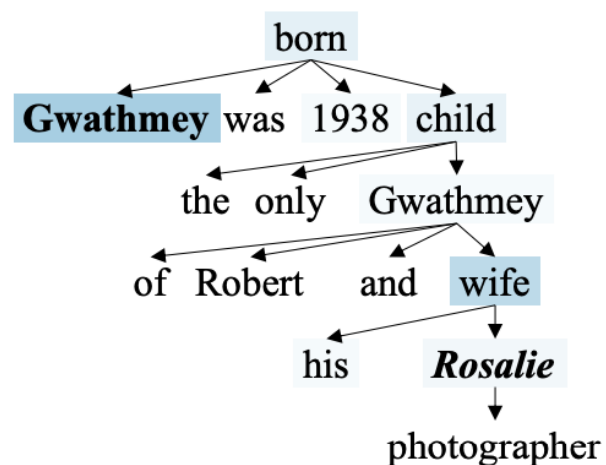


GNN Applications

■ Relation extraction

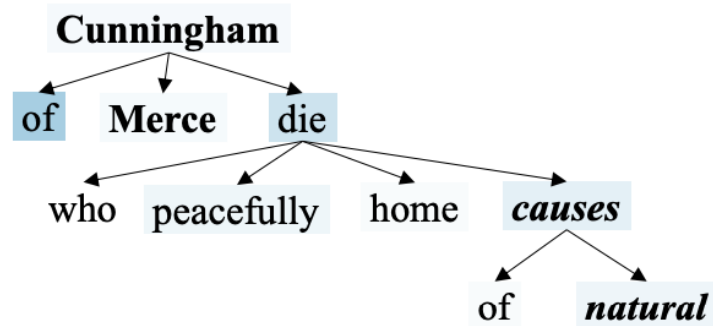
Relation: *per:parents*

Gwathmey was born in 1938, the only child of painter Robert Gwathmey and his wife, **Rosalie**, a photographer.



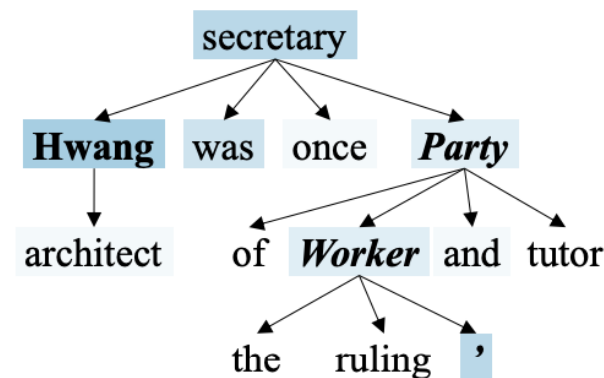
Relation: *per:cause_of_death*

"It is with great sorrow that we note the passing of **Merce Cunningham**, who died peacefully in his home last night of **natural causes**", the Cunningham Dance Foundation and the Merce Cunningham Dance Company said in a statement.



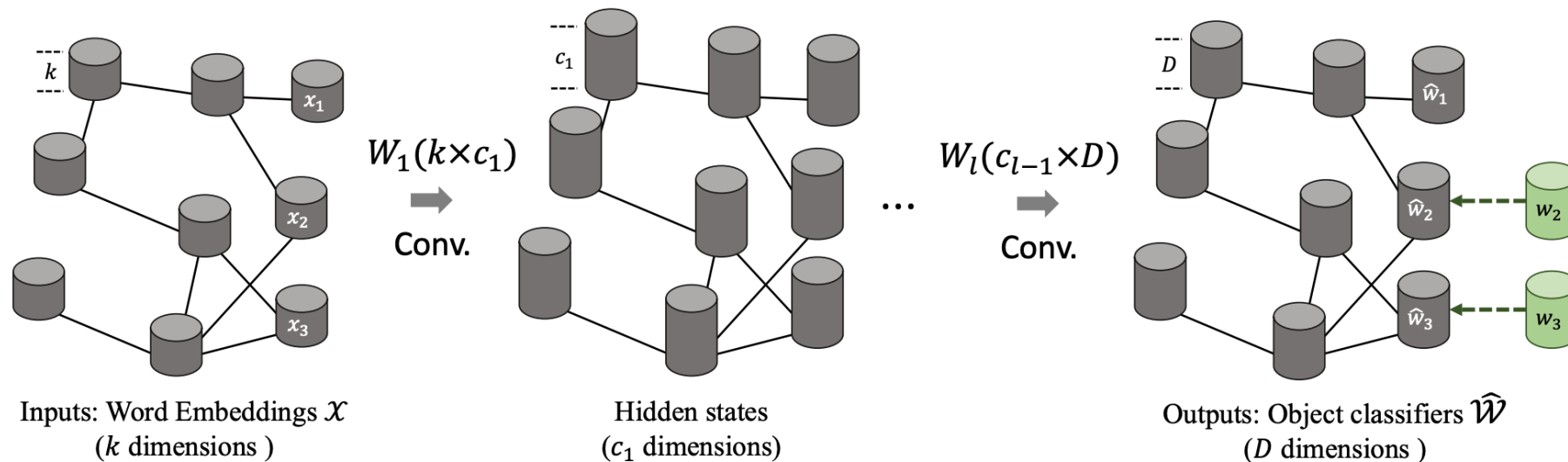
Relation: *per:employee_of*

Hwang, architect of the Pyongyang regime's ideology of "juche" or self-reliance, was once secretary of the ruling **Workers' Party** and a tutor to current leader Kim Jong-Il.



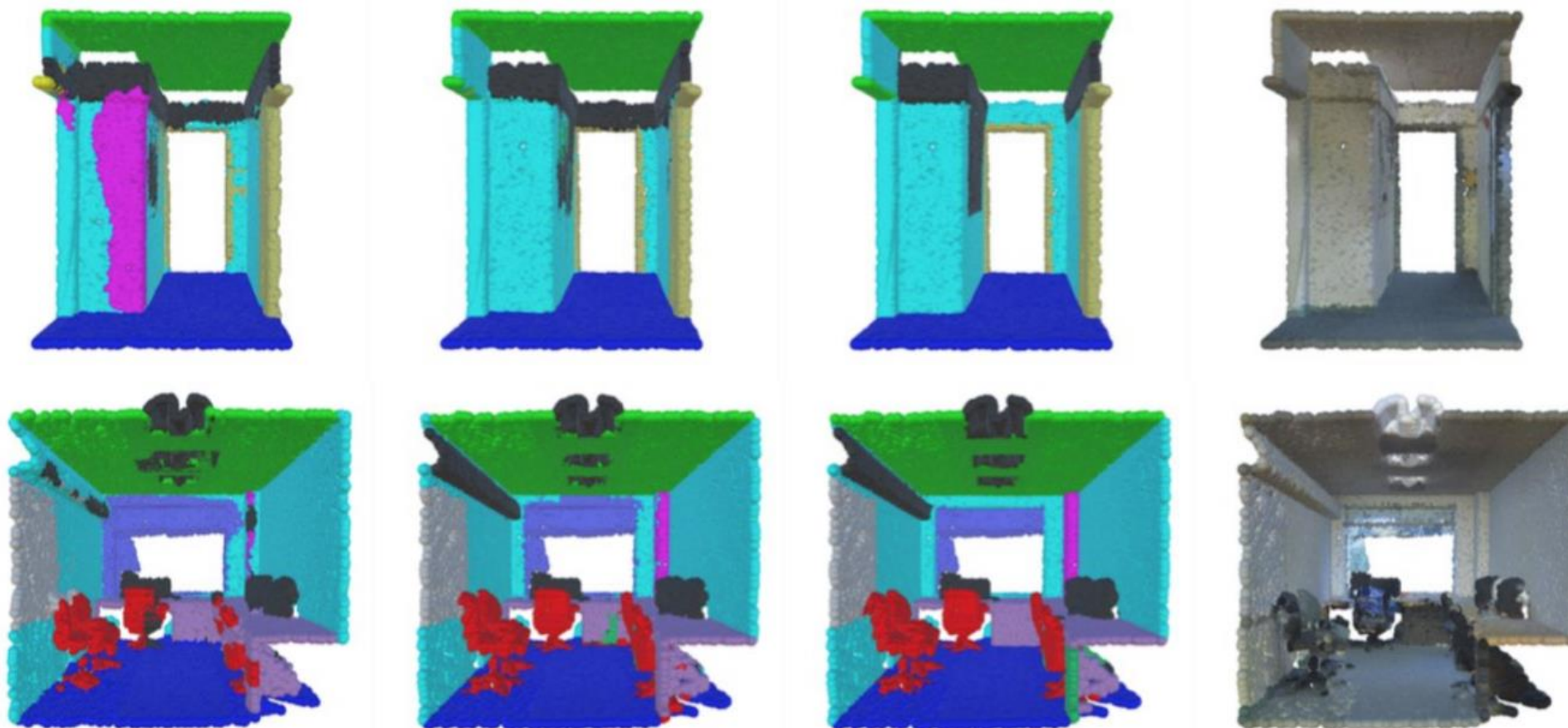
GNN Applications

■ Zero-shot image classification



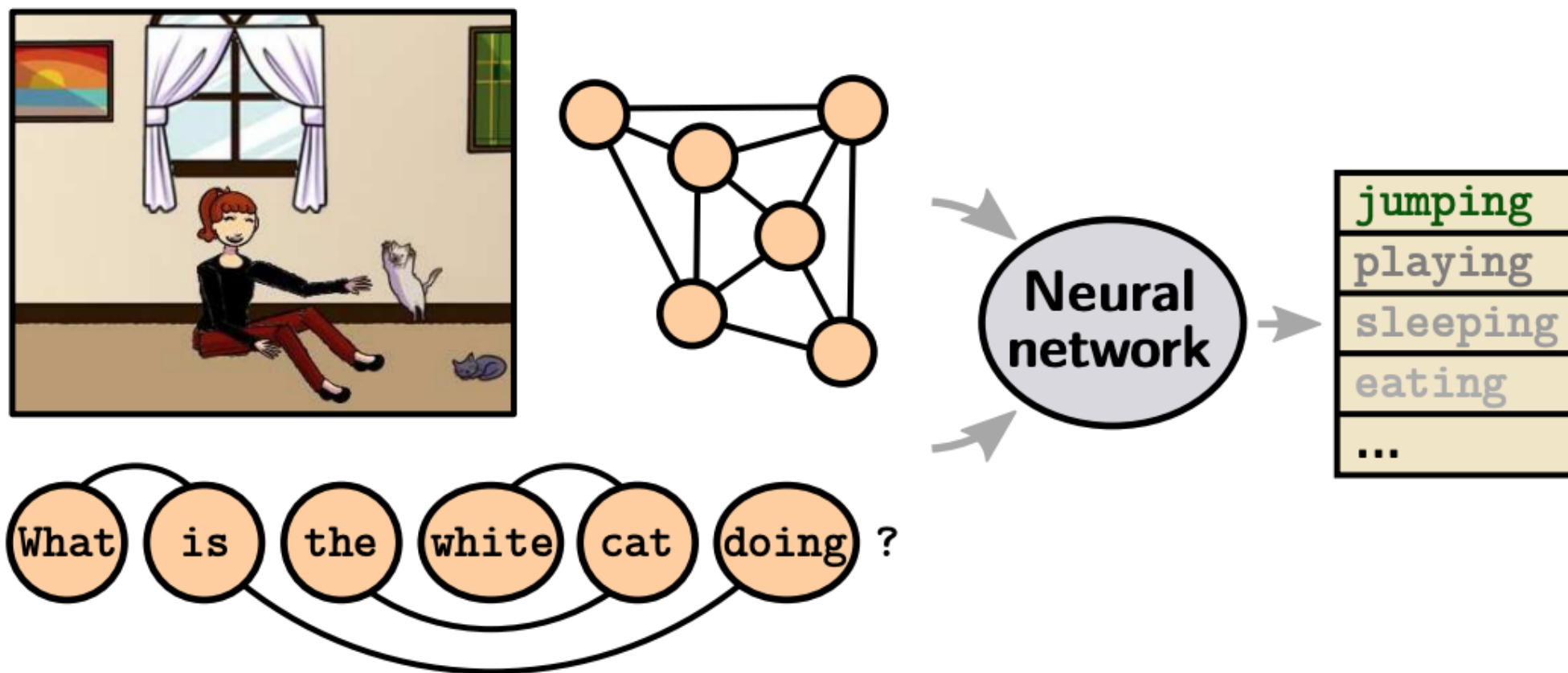
GNN Applications

■ Point cloud semantic segmentation



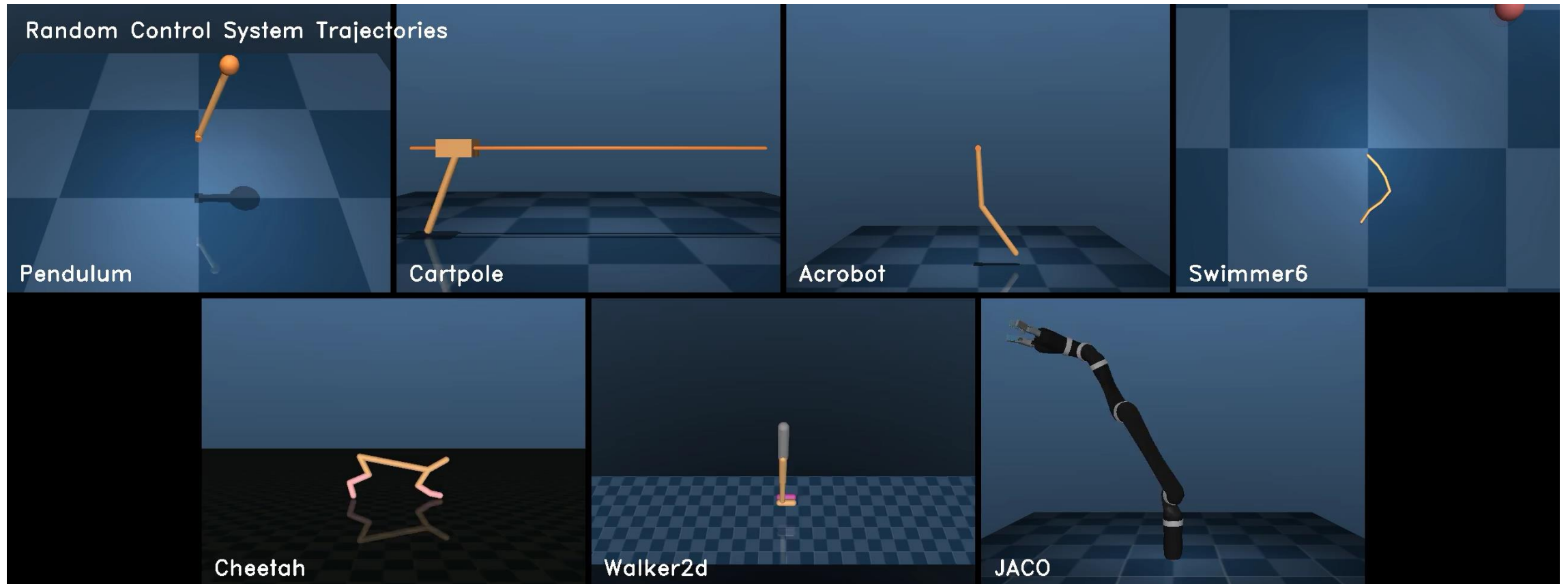
GNN Applications

■ Visual question answering



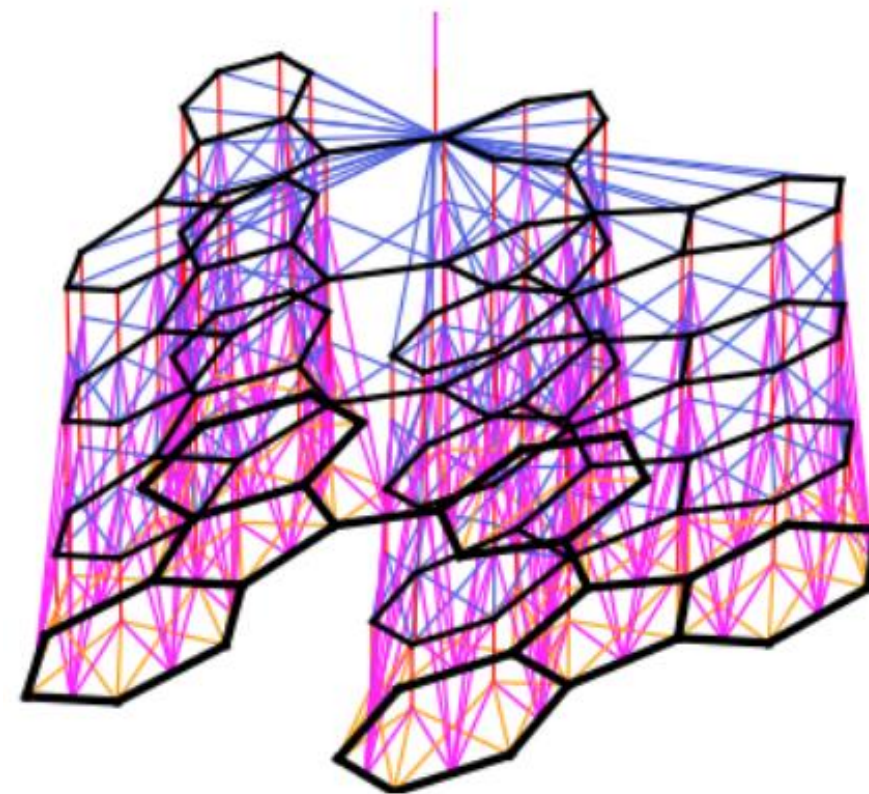
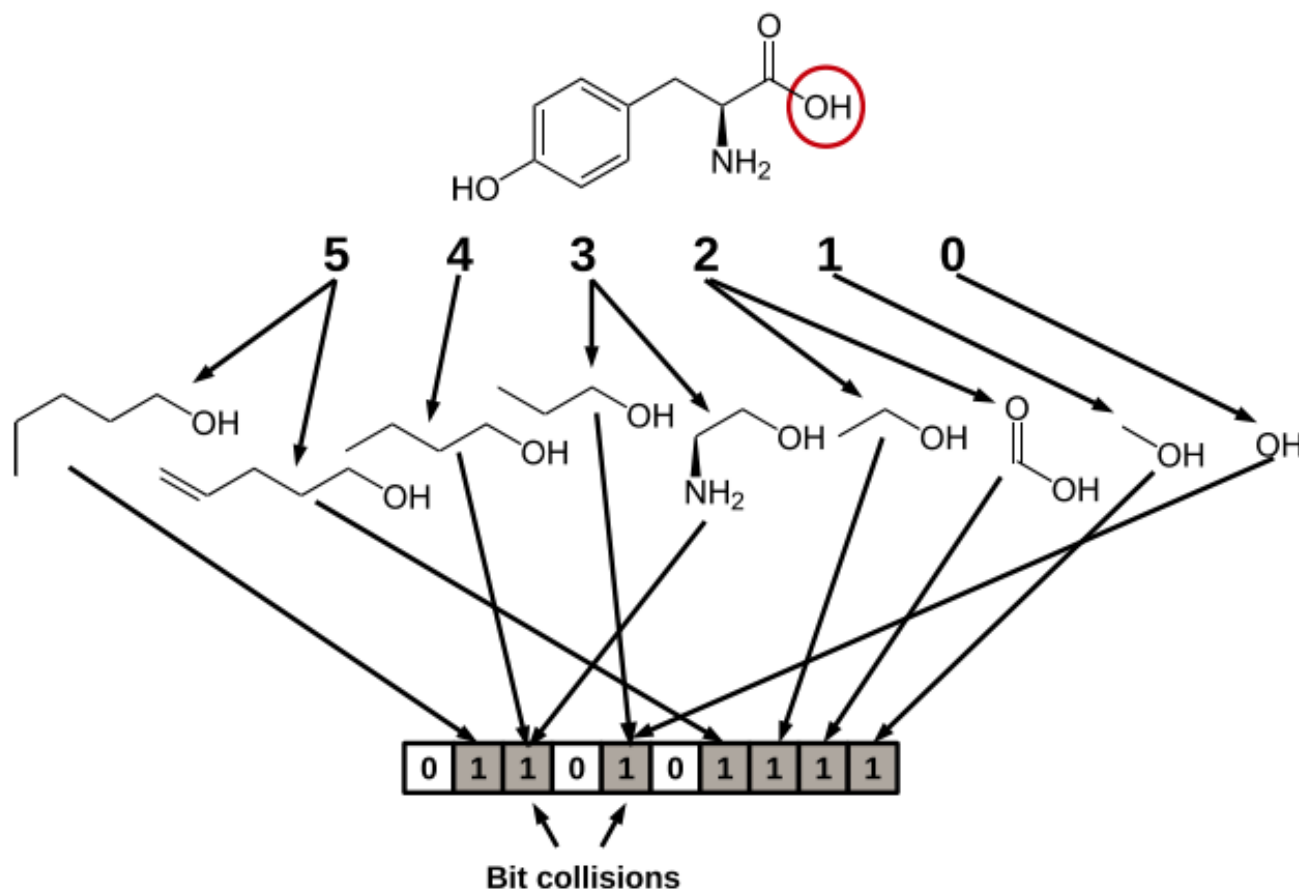
GNN Applications

■ Physics systems



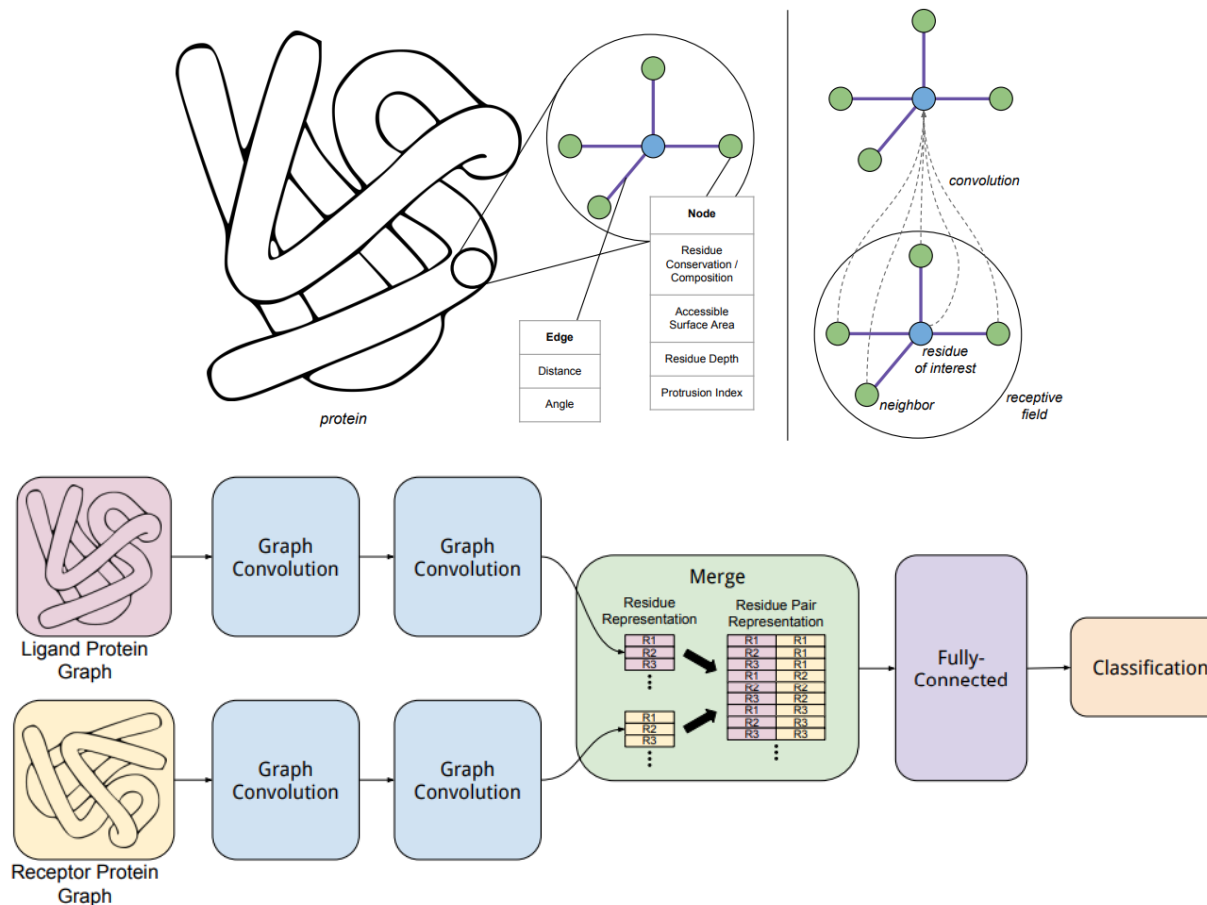
GNN Applications

■ Molecular fingerprints



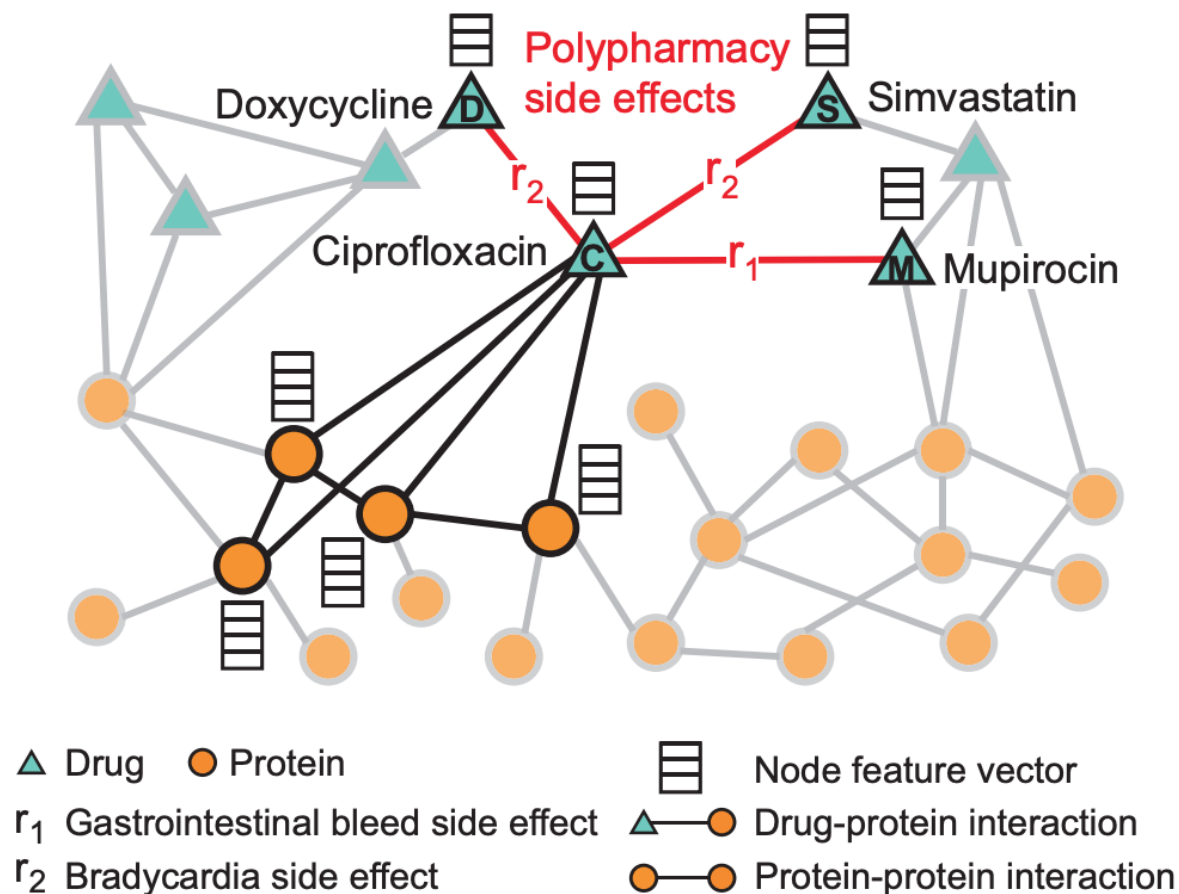
GNN Applications

■ Protein interface prediction



GNN Applications

■ Polypharmacy side effects



Outlines

- Graph Representation Learning
 - Deepwalk
 - LINE
 - Node2vec
- Graph Neural Networks
 - GCN
 - GraphSAGE
 - GAT
- Application to Recommender System
- Recent Advances

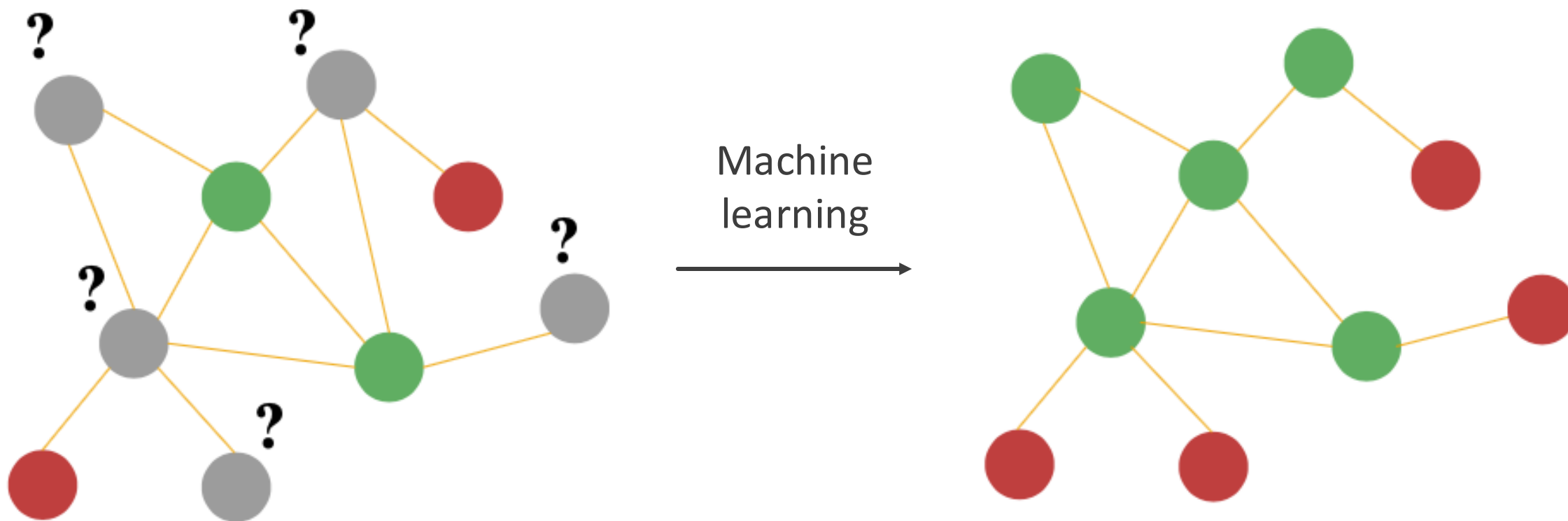




GRAPH REPRESENTATION LEARNING



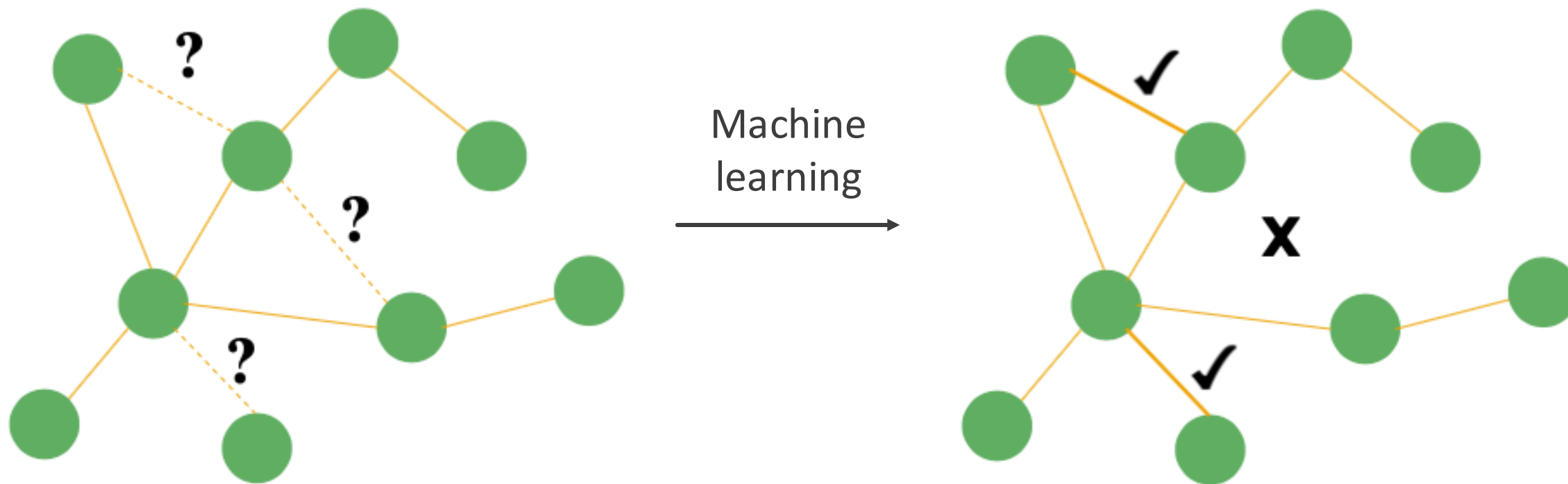
Graph Representation Learning



Node classification



Graph Representation Learning



Link prediction



Graph Representation Learning

- Given the graph, the only information we have is

$$G = \langle V, E \rangle$$

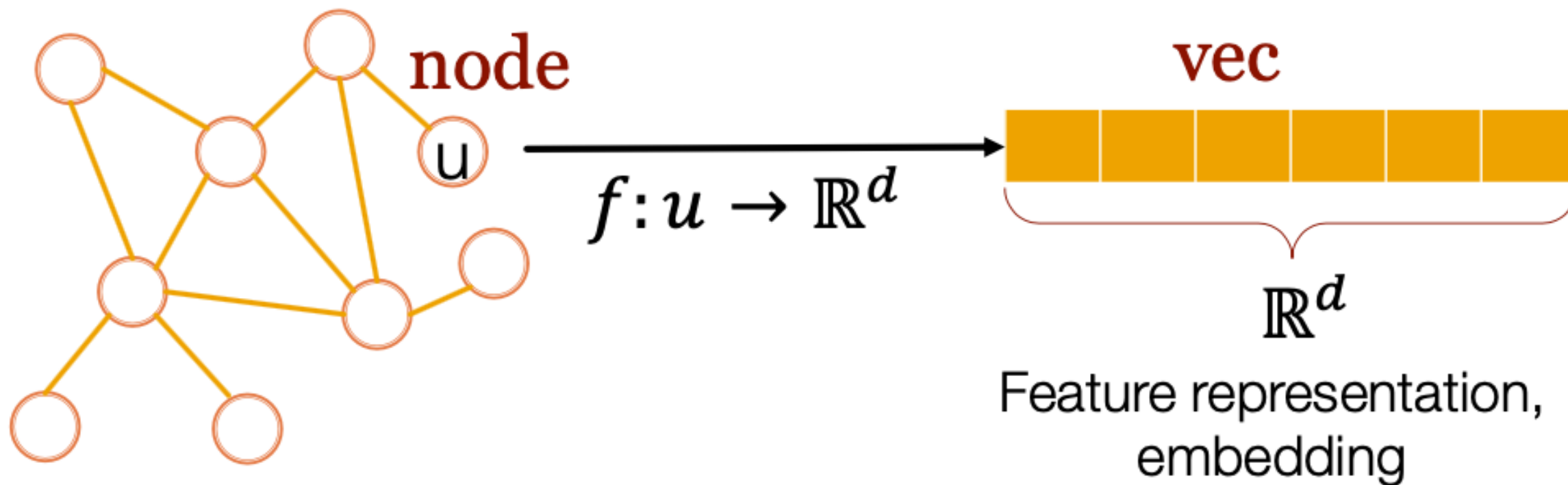
and their corresponding labels.

- What are the features?
- We can do feature engineering:
 - degree (count of adjacent nodes);
 - mean of degrees of neighbor nodes;
 - number of triangles a node forms with other nodes;
 - ...



Graph Representation Learning

- Goal: Efficiently learn task-independent features (embeddings) from graphs.



Graph Representation Learning

- Can we directly apply CNN or RNN on graphs?
- Probably no, because images and texts are structured.
 - Images are 2d matrices.
 - Texts are sequences.
- Graphs are far more complex.



Graph Representation Learning

Before the age of deep learning, we have some traditional machine learning methods:

- **Locally Linear Embedding**: low dimensional representations of each node can be represented by the linear combination of its neighbors.

$$\min \frac{1}{2} \sum_i \left| \mathbf{x}_i - \sum_j W_{ij} \mathbf{x}_j \right|^2$$

- **Laplacian Eigenmaps**: low dimensional representations of connected nodes are similar.

$$\min \frac{1}{2} \sum_{i,j} |\mathbf{x}_i - \mathbf{x}_j|^2 W_{ij}$$

- **Graph Factorization**: matrix factorization.



Graph Representation Learning

- Problem of these methods: can't scale!
- Key idea: If we assume that the connected nodes share similar properties (e.g. labels) in a graph, we should make their representations similar.
- Recall something?
- **Word2vec and xxx2vec!**
- But how to generate training pairs?





DEEPWALK

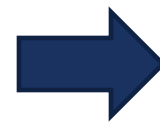
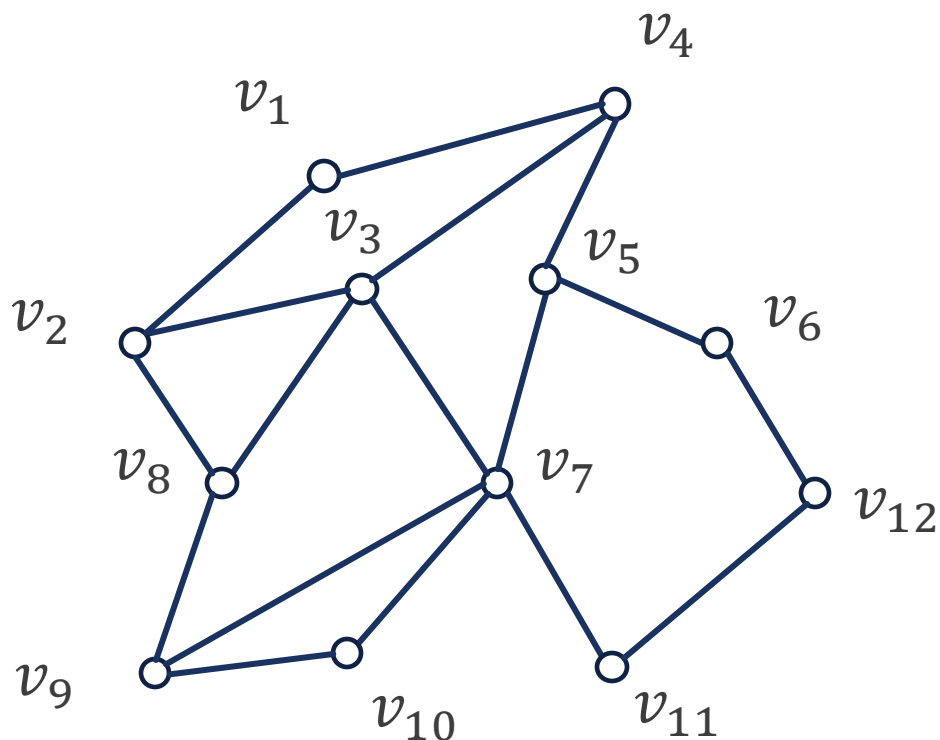
Deepwalk

Deepwalk: Online learning of social representations

[B Perozzi](#), [R Al-Rfou](#), [S Skiena](#) - Proceedings of the 20th ACM SIGKDD ..., 2014 - dl.acm.org

... **DeepWalk**, a ... **DeepWalk** generalizes recent advancements in language modeling and unsupervised feature learning (or deep learning) from sequences of words to graphs. **DeepWalk** ...

☆ Save ↗ Cite Cited by 10106 Related articles All 22 versions



(v_3, v_4, v_5, v_6)



(v_7, v_{10}, v_9, v_8)

Then what?

Skipgram!



Algorithm 1 DEEPWALK(G, w, d, γ, t)

Input: graph $G(V, E)$

 window size w

 embedding size d

 walks per vertex γ

 walk length t

Output: matrix of vertex representations $\Phi \in \mathbb{R}^{|V| \times d}$

1: Initialization: Sample Φ from $\mathcal{U}^{|V| \times d}$

2: Build a binary Tree T from V

3: **for** $i = 0$ to γ **do**

4: $\mathcal{O} = \text{Shuffle}(V)$

5: **for each** $v_i \in \mathcal{O}$ **do**

6: $\mathcal{W}_{v_i} = \text{RandomWalk}(G, v_i, t)$

7: SkipGram($\Phi, \mathcal{W}_{v_i}, w$)

8: **end for**

9: **end for**

Algorithm 2 SkipGram($\Phi, \mathcal{W}_{v_i}, w$)

1: **for each** $v_j \in \mathcal{W}_{v_i}$ **do**

2: **for each** $u_k \in \mathcal{W}_{v_i}[j - w : j + w]$ **do**

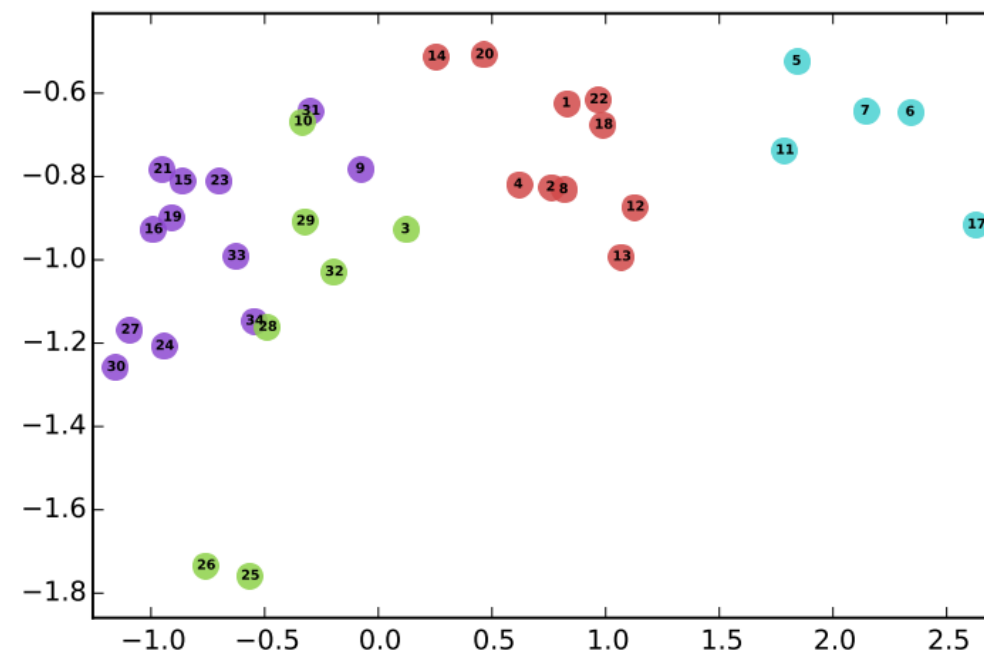
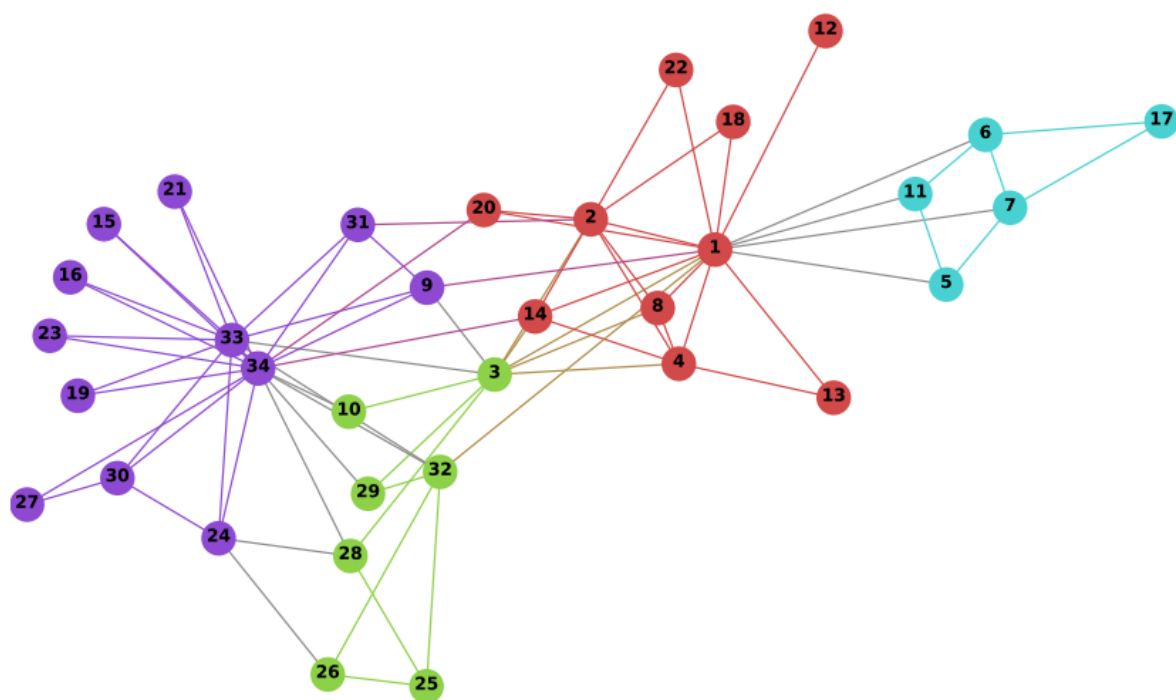
3: $J(\Phi) = -\log \Pr(u_k \mid \Phi(v_j))$

4: $\Phi = \Phi - \alpha * \frac{\partial J}{\partial \Phi}$

5: **end for**

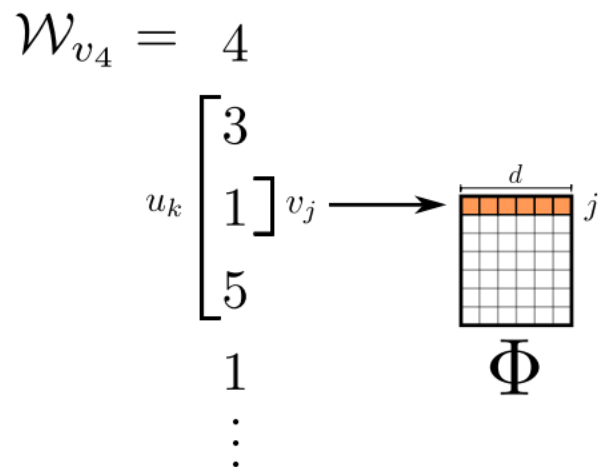
6: **end for**

Deepwalk

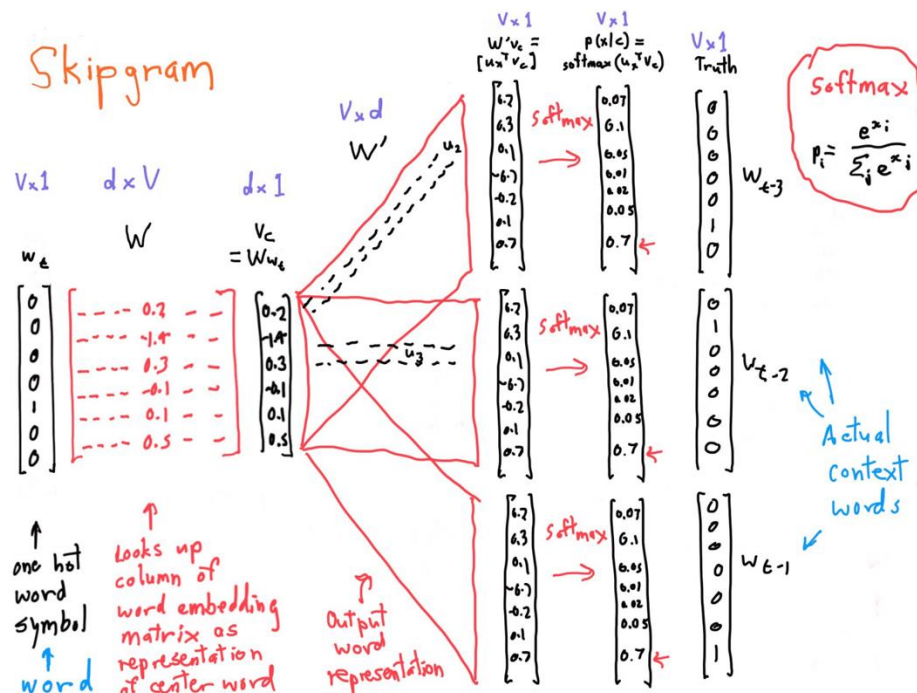


Difference

- Notice any difference between Deepwalk and Word2vec?



Deepwalk



Word2vec



Datasets

- Task: multi-label classification.
- Take Flickr as an example:
 - Nodes: users.
 - Links: following between users.
 - Categories: subscribe to different interest groups (e.g. black and white photos, or animals).

Data	BlogCatalog	Flickr	YouTube
Categories	39	195	47
Nodes (n)	10, 312	80, 513	1, 138, 499
Links (m)	333, 983	5, 899, 882	2, 990, 443
Network Density	6.3×10^{-3}	1.8×10^{-3}	4.6×10^{-6}
Maximum Degree	3, 992	5, 706	28, 754
Average Degree	65	146	5



Problems of Deepwalk

- Deepwalk is a pioneer work that builds a bridge between graph representation and word2vec.
- However, it is not specifically designed for graphs.
- How about directed graph? Weighted graph?





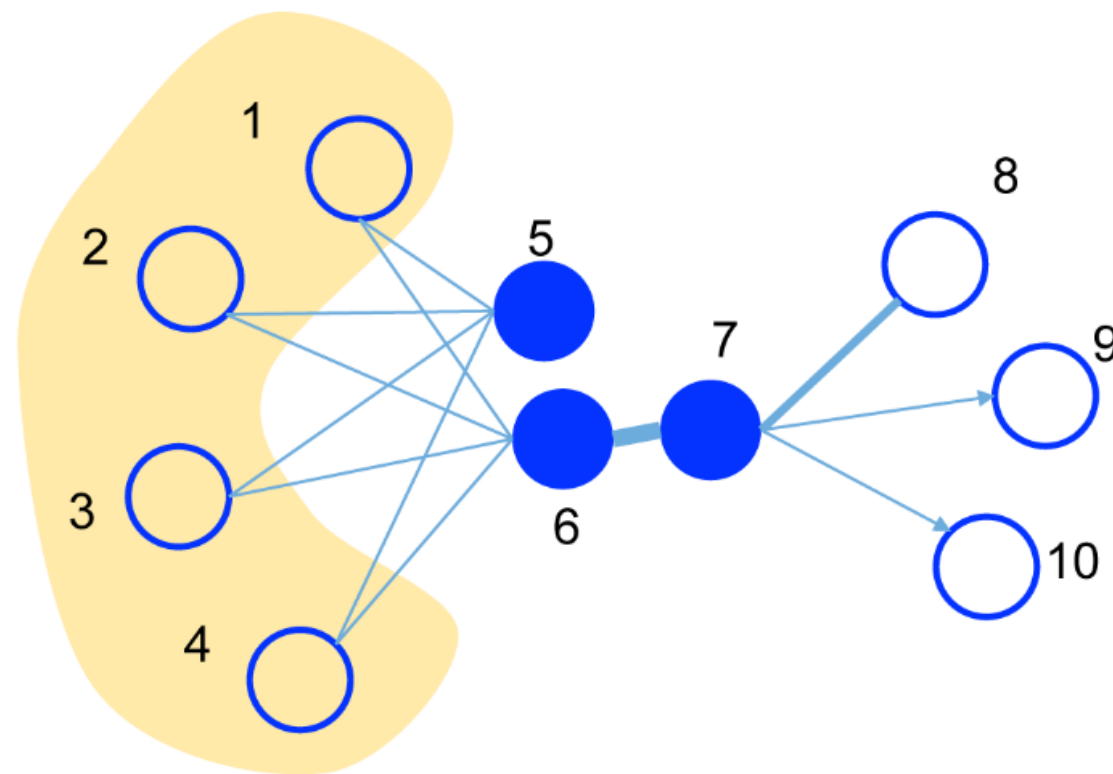
LINE

- **First-order proximity** in the real world data is not sufficient for preserving the global network structures.
- **Second-order proximity** is also very important.
 - It can be interpreted as nodes with shared neighbors being likely to be similar.
- The degree of overlap of two people's friendship networks correlates with the strength of ties between them.



LINE

- Vertex 6 and 7 should be placed closely in the low-dimensional space as they are connected through a strong tie.
- Vertex 5 and 6 should also be placed closely as they share similar neighbors.



LINE with First-Order Proximity

- For each **undirected** edge (i, j) , the **joint** probability between vertex v_i and v_j as follows:

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

- Use their edge weight as the label, W is total weight in the graph.

$$\hat{p}_1(v_i, v_j) = \frac{w_{ij}}{W}.$$

- Minimize the KL-divergence between p_1 and \hat{p}_1 :

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

The constant W can be omitted in minimization



LINE with First-Order Proximity

- For each **directed** edge (i, j) , we first define the **conditional** probability of “context” v_j generated by vertex v_i as:

$$p_1(v_j|v_i) = \frac{1}{1 + \exp(-\mathbf{u}_j'^T \mathbf{u}_i)}$$

- Here, we use difference representations for center and context just like Word2vec, why?



LINE with Second-Order Proximity

- If we consider the second-order proximity, v_j can be the neighbor of v_i 's neighbor.

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

where $|V|$ is the number of vertices or “contexts.”



LINE with Second-Order Proximity

- Similarly, minimize the KL-divergence:

$$- \sum_{j \in \mathcal{N}(i)} w_{ij} \log p_2(v_j | v_i).$$

where $\mathcal{N}(i)$ is the neighborhood of v_i , **including first and second order.**

- w_{ij} depends on two situation:

- v_j is the neighbor of v_i : w_{ij} is simply the weight.

- v_j is the neighbor of neighbor of v_i : $w_{ij} = \sum_{k \in \mathcal{N}(i)} w_{ik} \frac{w_{kj}}{d_k}$, d_k is the out-degree of v_k .



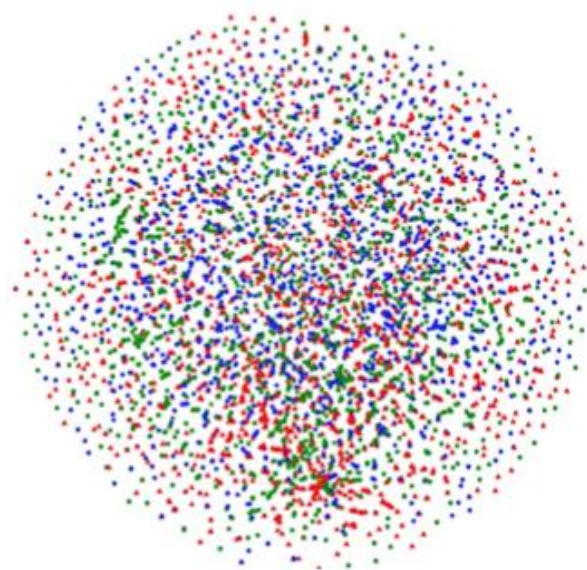
Combining First-Order and Second-Order Proximities

Two ways:

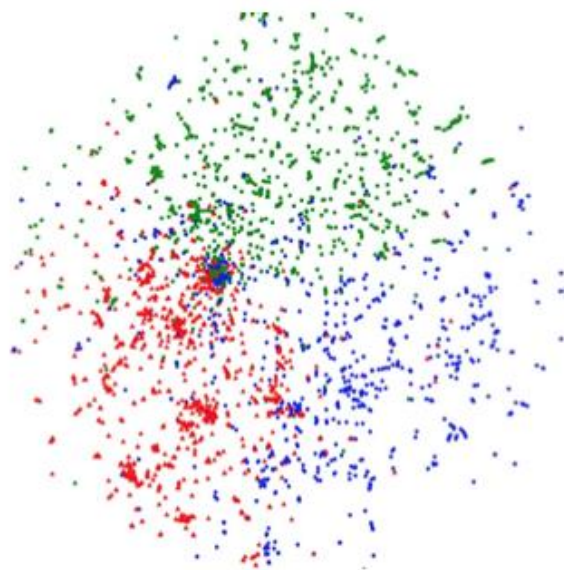
- Train separately and then concatenate.
- Jointly train the objective function.



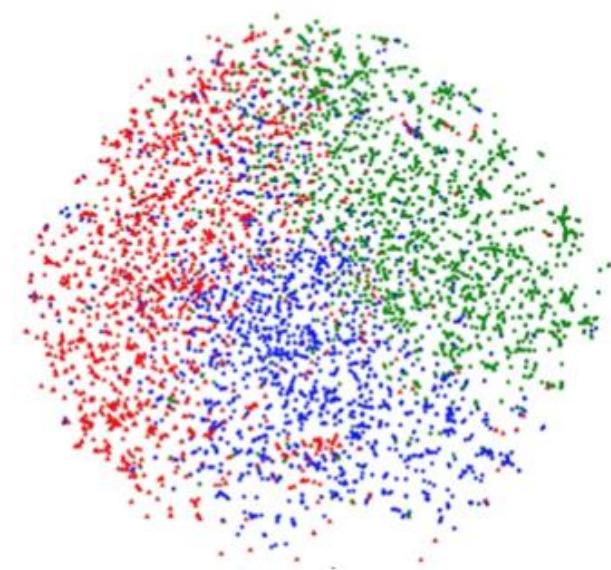
LINE



(a) GF



(b) DeepWalk



(c) LINE(2nd)

Visualization of the co-author network. The authors are mapped to the 2-D space using the t-SNE package with learned embeddings as input. Color of a node indicates the community of the author. Red: “data Mining,” blue: “machine learning,” green: “computer vision.”



Deepwalk vs LINE

- Deepwalk is actually a returnable DFS.
- LINE is a 2-level BFS.

Can we combine DFS and BFS?





NODE2VEC

- Motivation: It is now either DFS (Deepwalk) or BFS (LINE). It is too rigid to explore the network neighborhood.
- Can we make it flexible?

Node2vec



Jure Leskovec

Professor of Computer Science, [Stanford University](#)
在 [cs.stanford.edu](#) 的电子邮件经过验证 - [首页](#)



[Data mining](#) [Machine Learning](#) [Graph Neural Networks](#) [Knowledge Graphs](#) [Complex Networks](#)

标题

引用次数

年份

Inductive representation learning on large graphs

W Hamilton, Z Ying, J Leskovec

Advances in neural information processing systems 30

17429

2017

node2vec: Scalable feature learning for networks

A Grover, J Leskovec

Proceedings of the 22nd ACM SIGKDD international conference on Knowledge ...

12961

2016

How powerful are graph neural networks?

K Xu, W Hu, J Leskovec, S Jegelka

arXiv preprint arXiv:1810.00826

8972

2018

SNAP Datasets: Stanford large network dataset collection

J Leskovec, A Krevl

4670

2014

On the opportunities and risks of foundation models

R Bommasani, DA Hudson, E Adeli, R Altman, S Arora, S von Arx, ...

arXiv preprint arXiv:2108.07258

4128

2021

Graph convolutional neural networks for web-scale recommender systems

R Ying, R He, K Chen, P Eksombatchai, WL Hamilton, J Leskovec

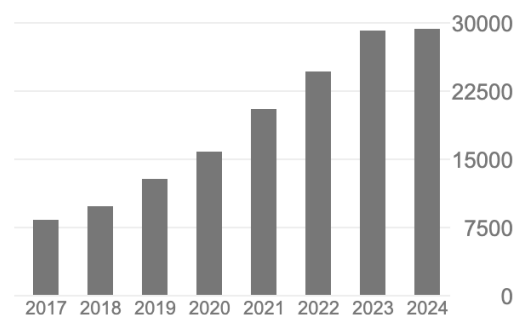
4045

2018

引用次数

查看全部

	总计	2019 年至今
引用	182897	132383
h 指数	155	136
i10 指数	372	347



开放获取的出版物数量

查看全部



无法查看的文章

可查看的文章

根据资助方的强制性开放获取政策



厦门大学信息学院 (特色化示范性软件学院)

School of Informatics Xiamen University (National Characteristic Demonstration Software School)

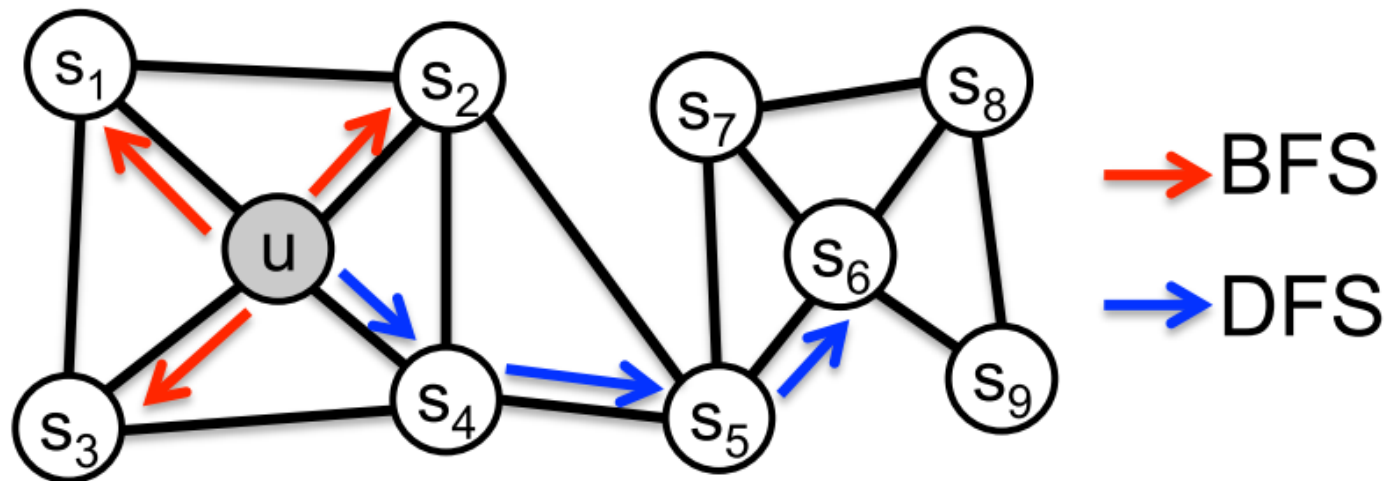


厦门大学 计算机科学与技术系

Department of Computer Science and Technology, Xiamen University

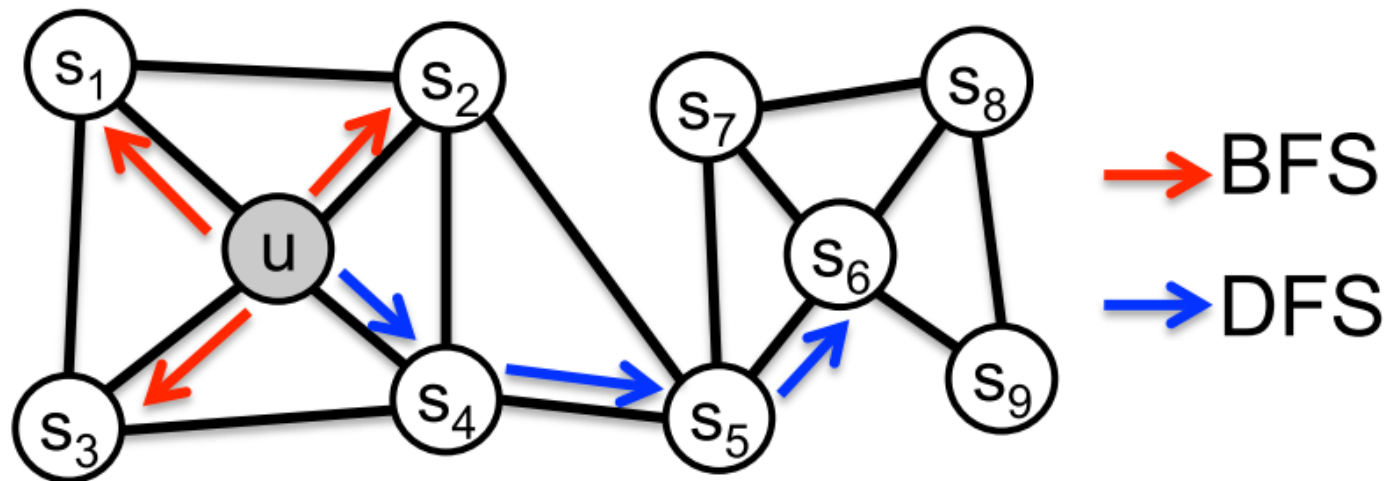
Node2vec

- Nodes u and s_1 belonging to the same tightly knit community.
- Nodes u and s_6 in the two distinct communities share the same structural role of a hub node.
- Should u be similar to s_1 or s_6 ?
- Both, but in different perspective.



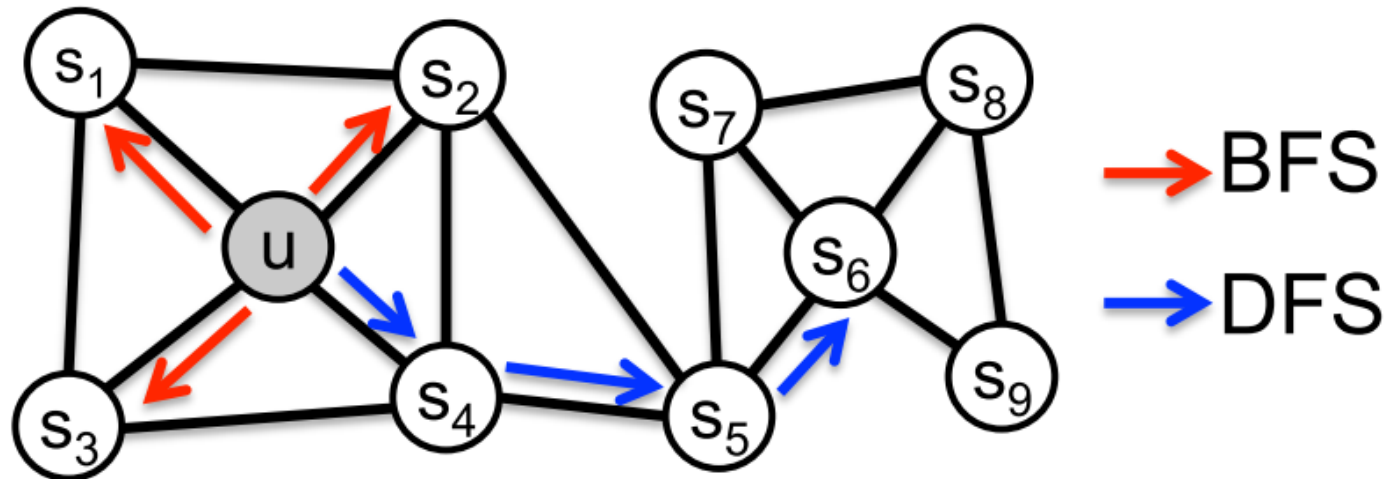
Node2vec

- Real-world networks commonly exhibit a mixture of such equivalences.
- The representations should be flexible to have similar embeddings for:
 - nodes from the same network community;
 - nodes that share similar roles.



Node2vec

- Idea: use flexible, biased random walks that can trade off between local and global views of the network.
- Walk of length 3 ($N_R(u)$ of size 3):
 - $N_{BFS}(u) = \{s_1, s_2, s_3\}$, local microscopic view.
 - $N_{DFS}(u) = \{s_4, s_5, s_6\}$, global macroscopic view.



- The probability from c_{i-1} to c_i is:

$$P(c_i = x | c_{i-1} = v) = \begin{cases} \frac{\pi_{vx}}{Z} & \text{if } (v, x) \in E \\ 0 & \text{otherwise} \end{cases}$$

- The first step is same for both DFS and BFS, by simply setting:

$$\pi_{vx} = w_{vx}$$



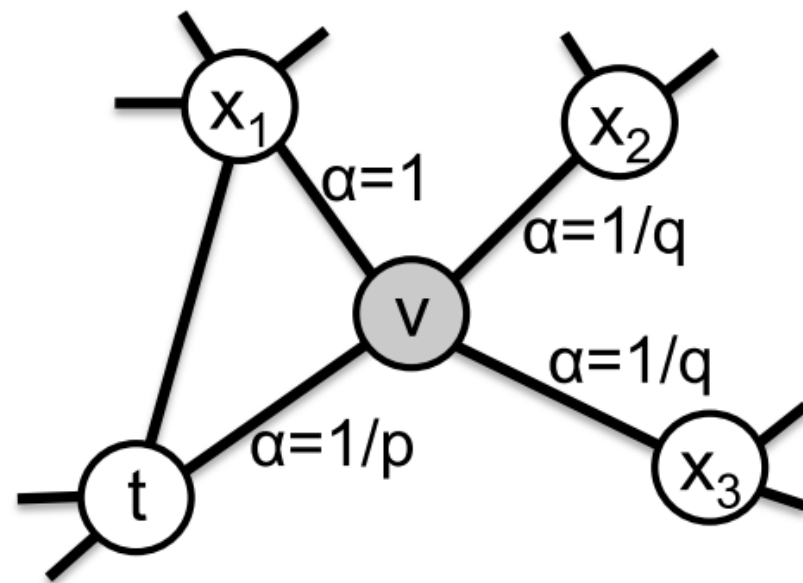
Node2vec

- Consider a random walk that just traversed edge (t, v) and now resides at node v .
- For the steps after the second step, we set

$$\pi_{vx} = \alpha_{pq}(t, x) w_{vx}$$

where

$$\alpha_{pq}(t, x) = \begin{cases} \frac{1}{p} & \text{if } d_{tx} = 0 \\ 1 & \text{if } d_{tx} = 1 \\ \frac{1}{q} & \text{if } d_{tx} = 2 \end{cases}$$

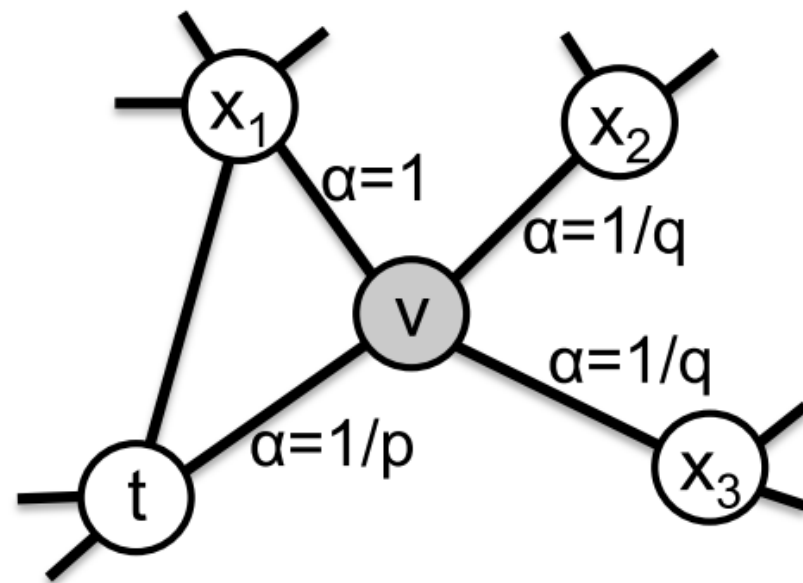


and $d_{tx} \in \{0, 1, 2\}$ denotes the shortest path distance between nodes t and x .



Node2vec

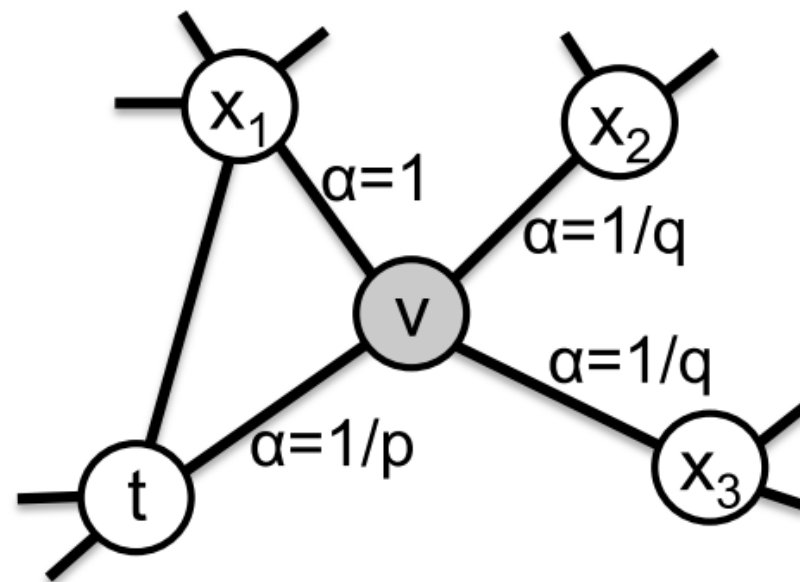
- p and q are hyperparameters, to control how we move from the second step.
- Return parameter p :
 - Return back to the previous node.
- In-out parameter q :
 - Moving outwards (DFS) vs. inwards (BFS)
 - Intuitively, q is the “ratio” of BFS vs. DFS



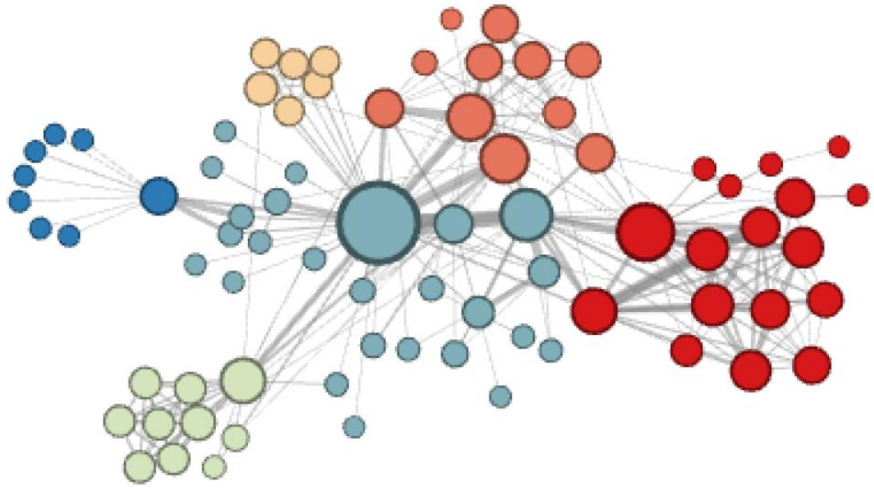
Node2vec

Cases:

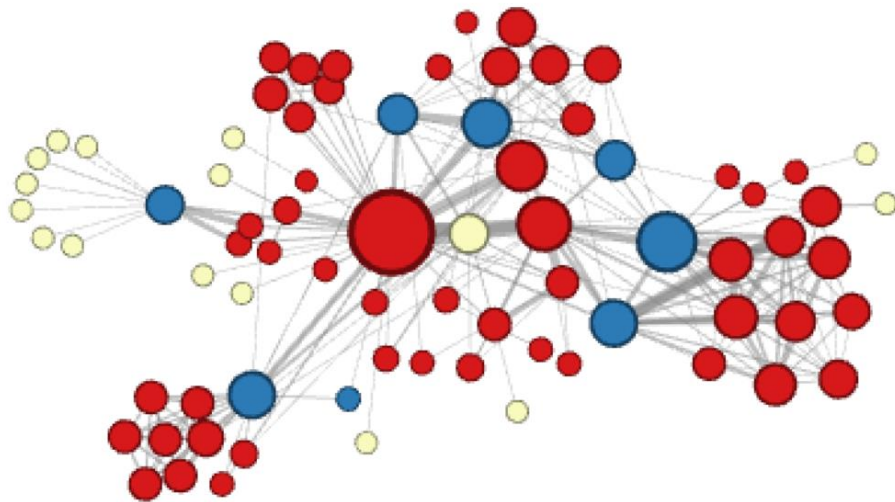
- p large, q large: non-returnable BFS.
- p large, q small: non-returnable DFS.
- p small, q large: returnable BFS.
- p small, q small: returnable DFS.
- $p = q = 1$: random walk.



Node2vec



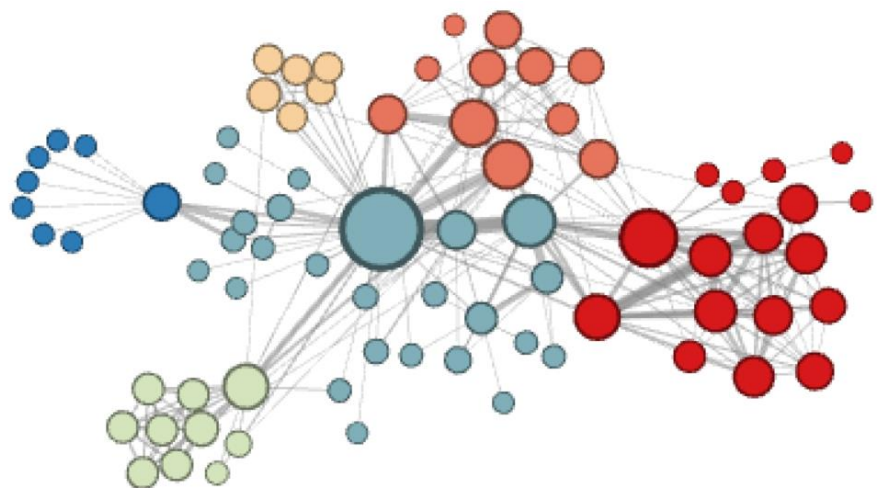
DFS
 $p = 1, q = 0.5$



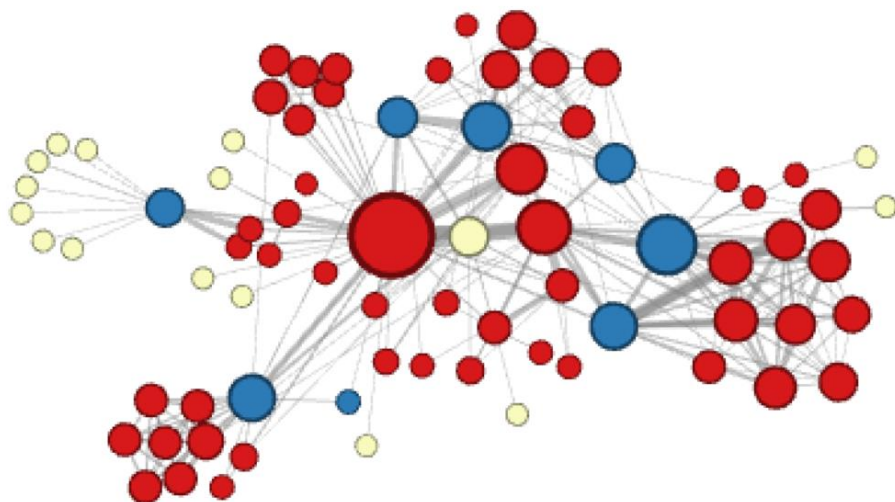
BFS
 $p = 1, q = 2$



Node2vec



- DFS travels to the world, therefore know the difference.



- BFS only sees the neighborhood, therefore only know the difference between itself and its neighborhood.



How to Use Node Embeddings

After we obtain the embedding \mathbf{z}_i for node i , how to use?

- Clustering/community detection: Clustering on nodes \mathbf{z}_i .
- Node classification: Predict label $f(\mathbf{z}_i)$ of node i based on \mathbf{z}_i .
- Link prediction: Predict edge (i, j) based on $f(\mathbf{z}_i, \mathbf{z}_j)$ by concatenate, avg, product, or take a difference between the embeddings:
 - Concatenate: $f(\mathbf{z}_i, \mathbf{z}_j) = g([\mathbf{z}_i, \mathbf{z}_j])$
 - Hadamard: $f(\mathbf{z}_i, \mathbf{z}_j) = g(\mathbf{z}_i \otimes \mathbf{z}_j)$
 - Sum/Avg: $f(\mathbf{z}_i, \mathbf{z}_j) = g(\mathbf{z}_i + \mathbf{z}_j)$
 - Distance: $f(\mathbf{z}_i, \mathbf{z}_j) = g(\|\mathbf{z}_i - \mathbf{z}_j\|_2)$

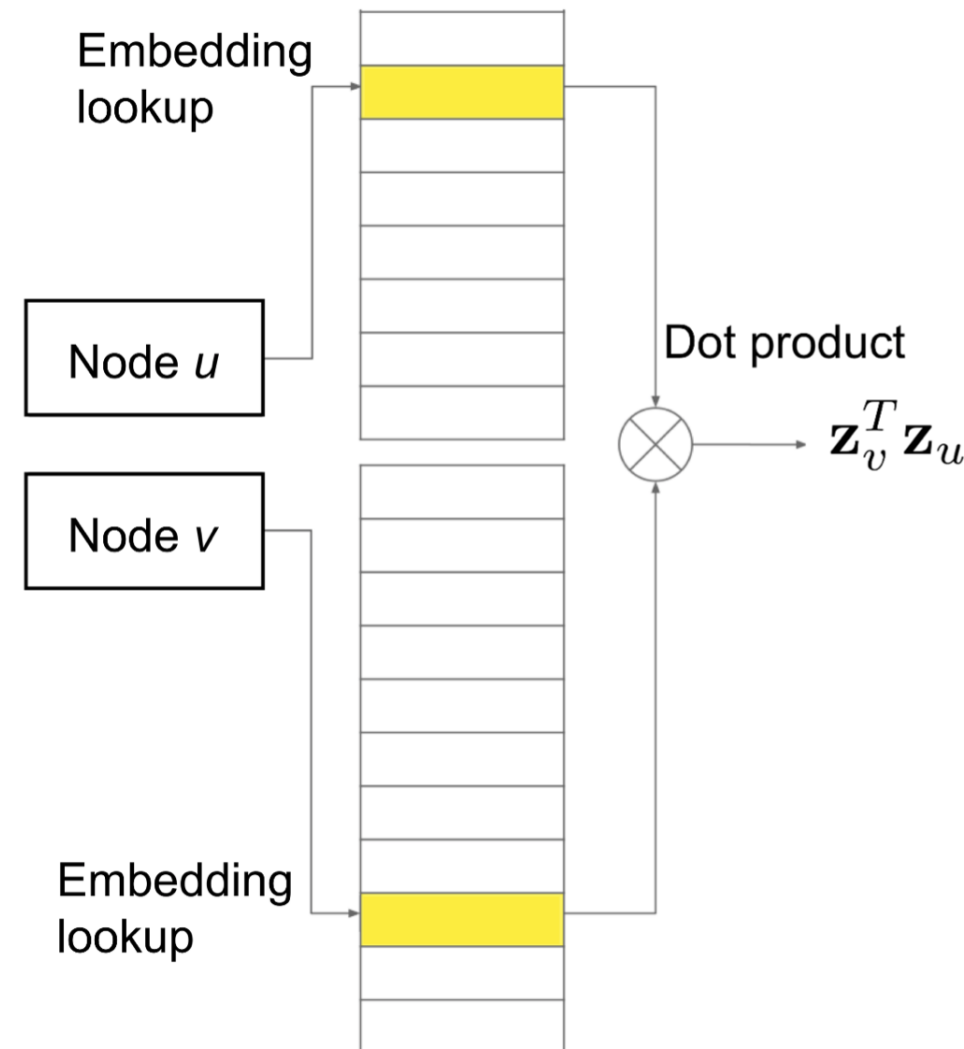


Shallow Encoders

Shallow encoders:

- One-layer of data transformation.
- A single hidden layer maps node u to embedding \mathbf{z}_u by

$$\mathbf{z}_u = f(\mathbf{z}_v, v \in N_R(u)).$$



Shallow Encoders

Limitations of shallow embedding methods:

- No parameter sharing:
 - Every node has its own unique embedding.
- Transductive, not inductive:
 - Cannot generate embeddings for nodes that are not seen during training.
- Do not incorporate node features:
 - Many graphs have features that we can and should leverage.
- Separated from downstream tasks.
 - Training is not end-to-end.

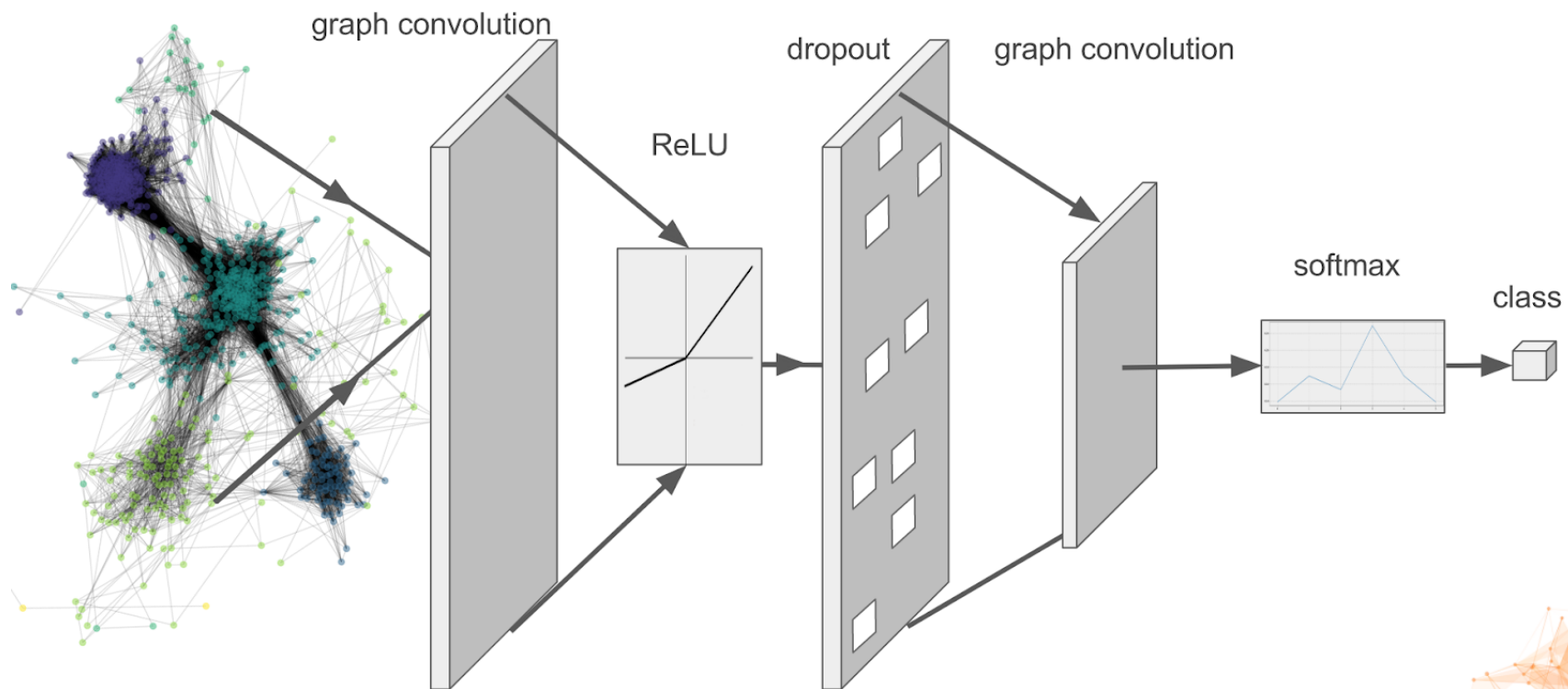




GRAPH NEURAL NETWORKS

Deep Graph Encoder

- Instead of directly learning embedding, can we learn mapping to generate embedding?



© 2016 Expero, Inc. All Rights Reserved



廈門大學信息學院 (特色化示范性软件学院)

School of Informatics Xiamen University (National Characteristic Demonstration Software School)



廈門大學 计算机科学与技术系

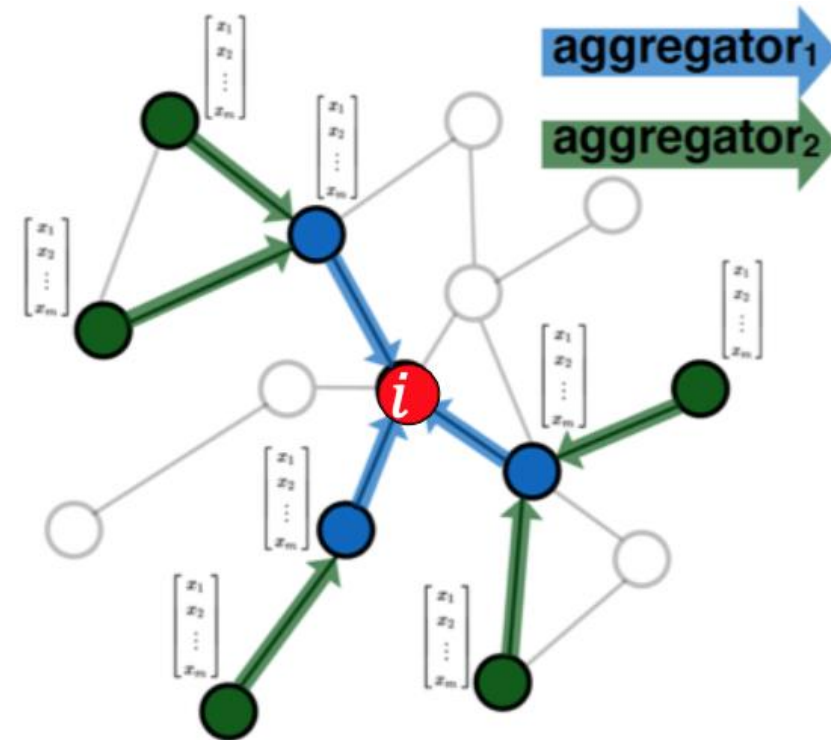
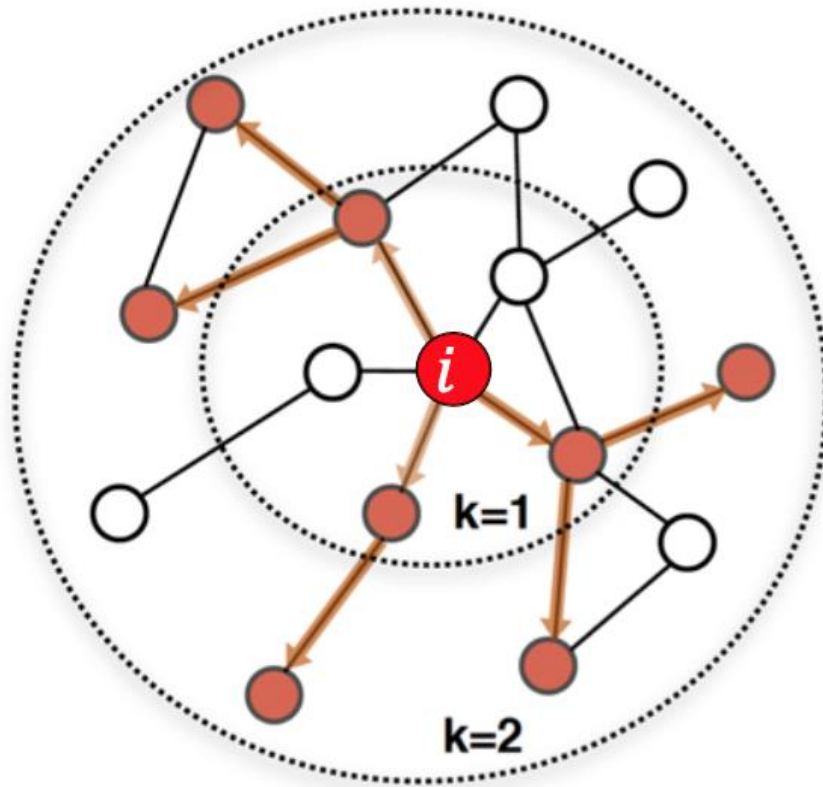
Department of Computer Science and Technology, Xiamen University

Image source: <https://www.expero.com/post/node-classification-by-graph-convolutional-network>



GCN

- Idea: Node's neighborhood defines a computation graph.



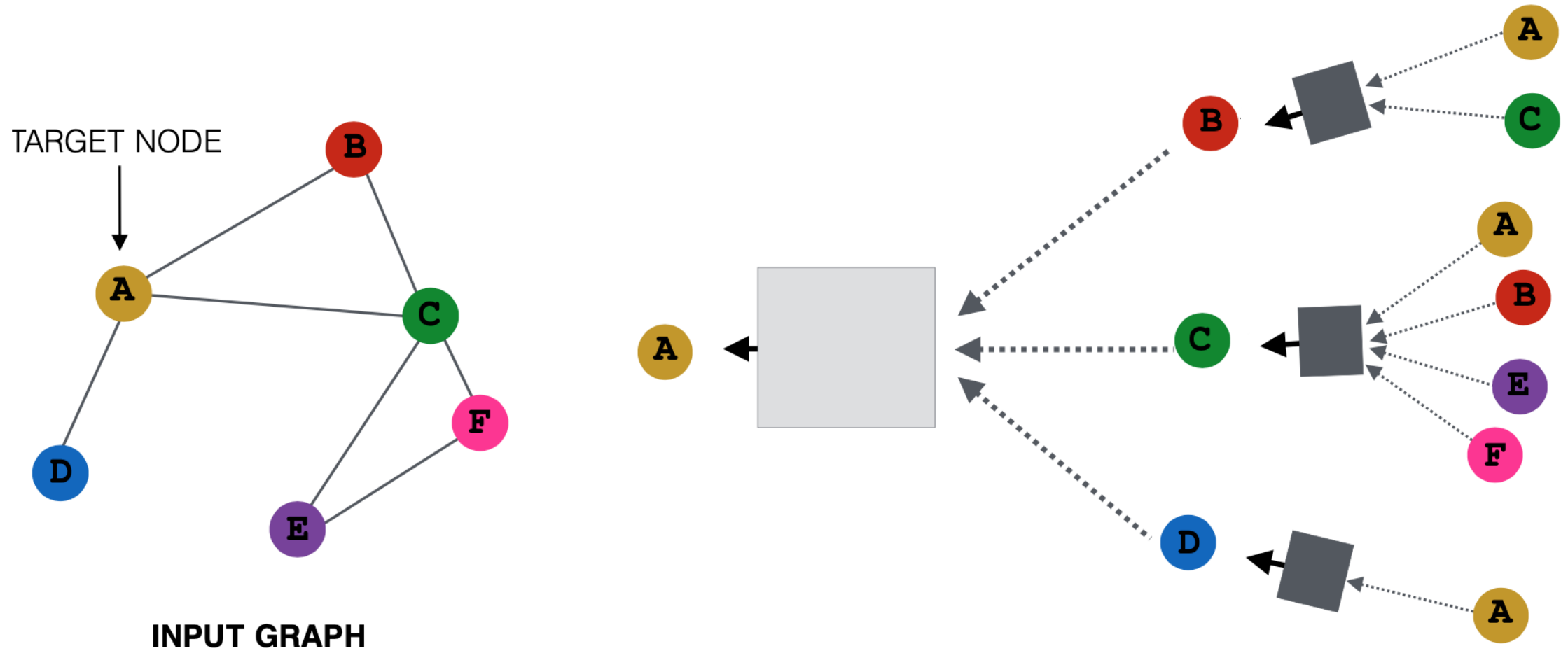
GCN: Basic Setting

Assume we have a graph G :

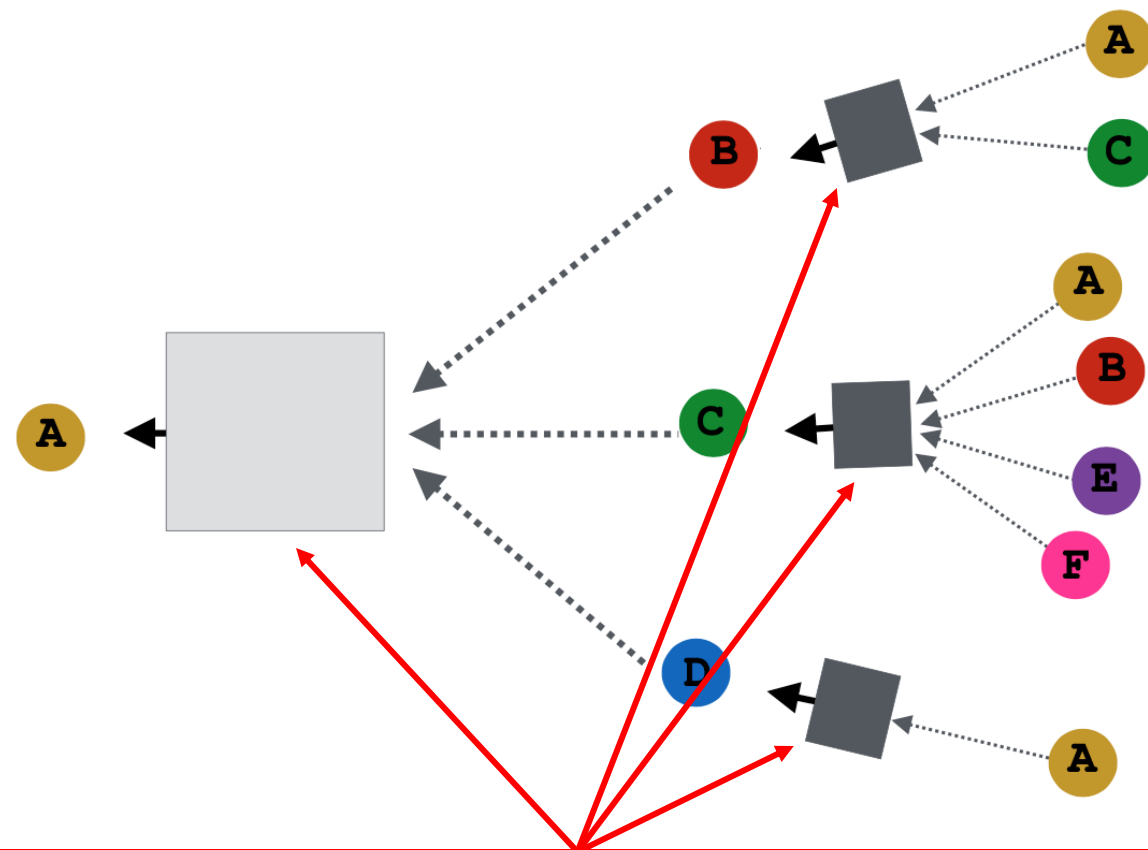
- V is the vertex set.
- A is the adjacency matrix (assume binary).
- $X \in \mathbb{R}^{m \times |V|}$ is a matrix of node initial features.
- Node initial features:
 - Social networks: user profile, user image.
 - Biological networks: gene expression profiles, gene functional information.
 - No features: one-hot vector or constant vector.



GCN: Architecture



GCN: Architecture

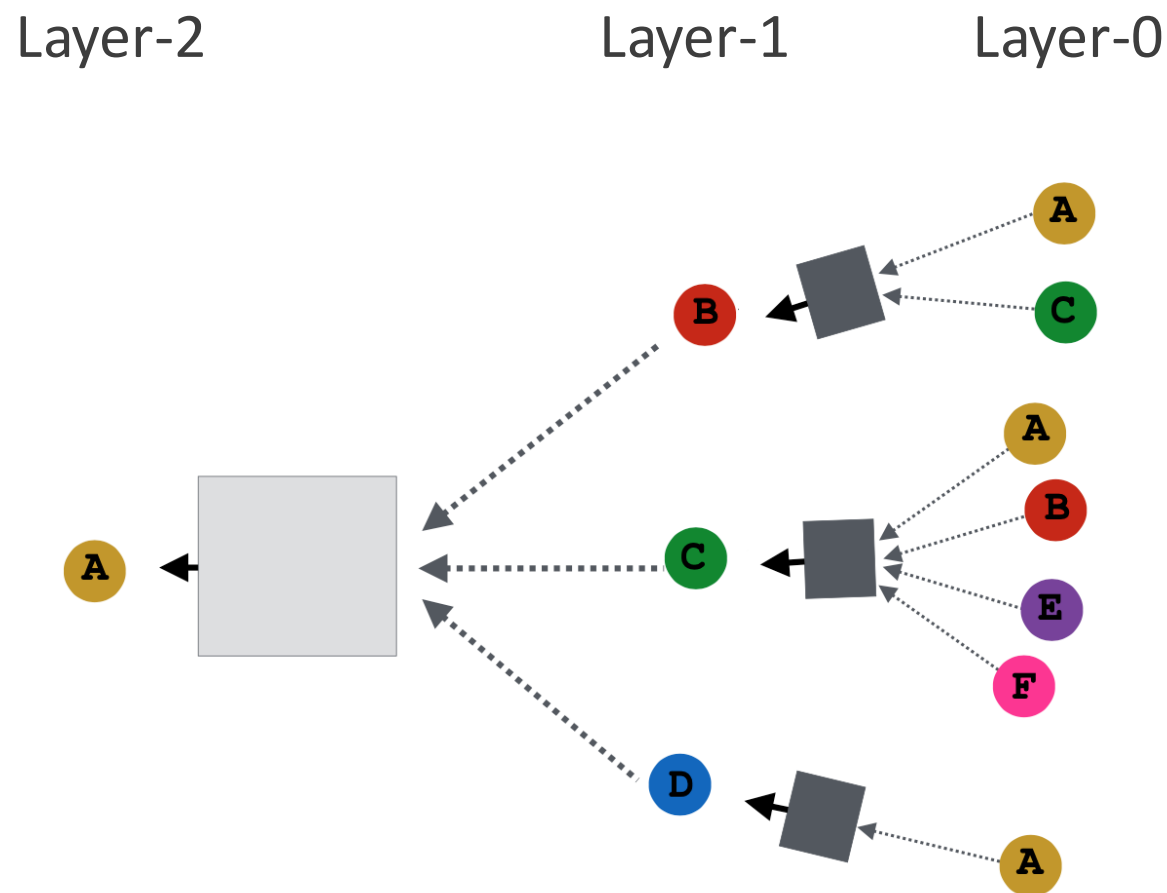


Average the information from the previous layer and apply neural network



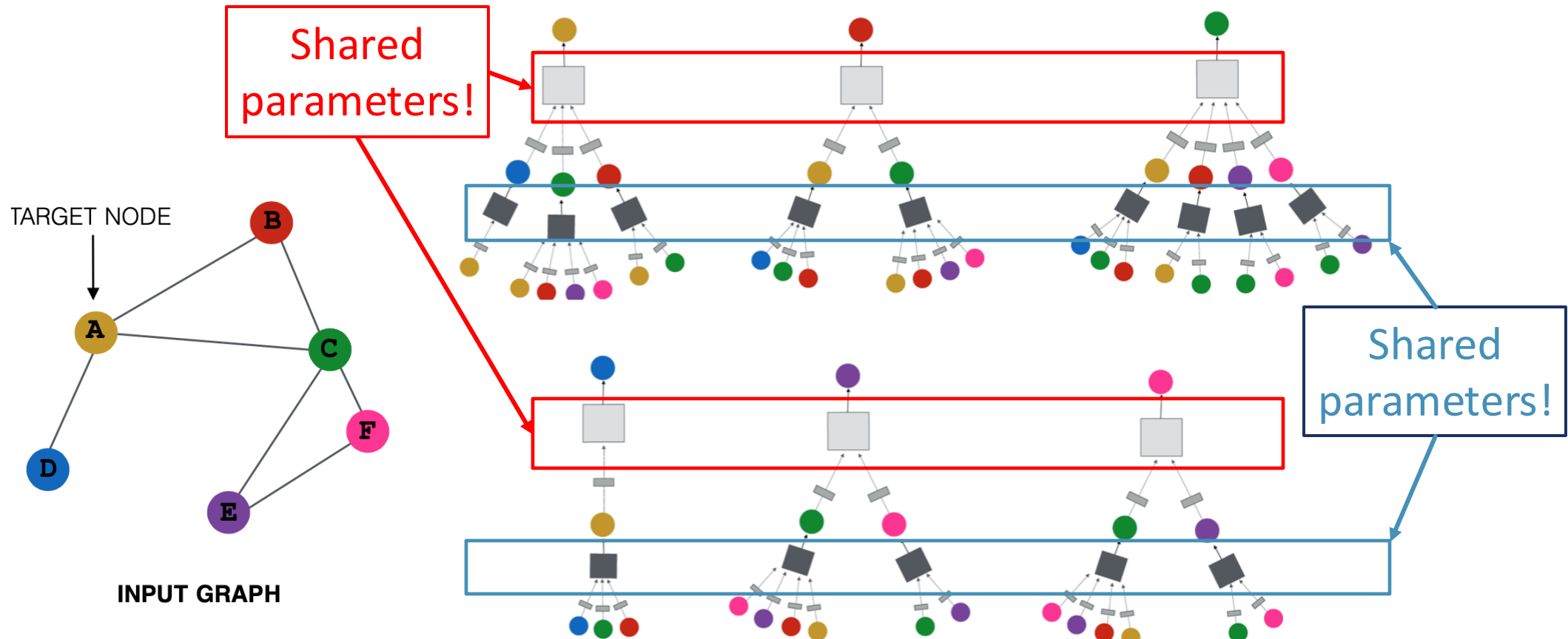
GCN: Multiple Layers

- Model can be of arbitrary depth:
- Nodes have embeddings at each layer.
- Layer-0 embedding of node u is its input feature, x_u .
- Layer-K embedding gets information from nodes that are K hops away.



GCN: Parameter Sharing

- Every node defines a computation graph based on its neighborhood!



GCN: Deep Encoder

- For each node v , its embedding at Layer- k is \mathbf{h}_v^k :

$$\begin{aligned}\mathbf{h}_v^0 &= \mathbf{x}_v \\ \mathbf{h}_v^k &= \sigma \left(\mathbf{W}_k \sum_{u \in N(v)} \frac{\mathbf{h}_u^{k-1}}{|N(v)|} + \mathbf{B}_k \mathbf{h}_v^{k-1} \right), k = 1, \dots, K \\ \mathbf{z}_v &= \mathbf{h}_v^K\end{aligned}$$

- \mathbf{W}_k is the parameter at Layer- k for the averaged neighborhood of node v ;
- \mathbf{B}_k is the parameter at Layer- k for node v itself.



GCN: Deep Encoder

- In original GCN paper, the neural network is represented by (sparse) matrix operations.

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

can be formulated as

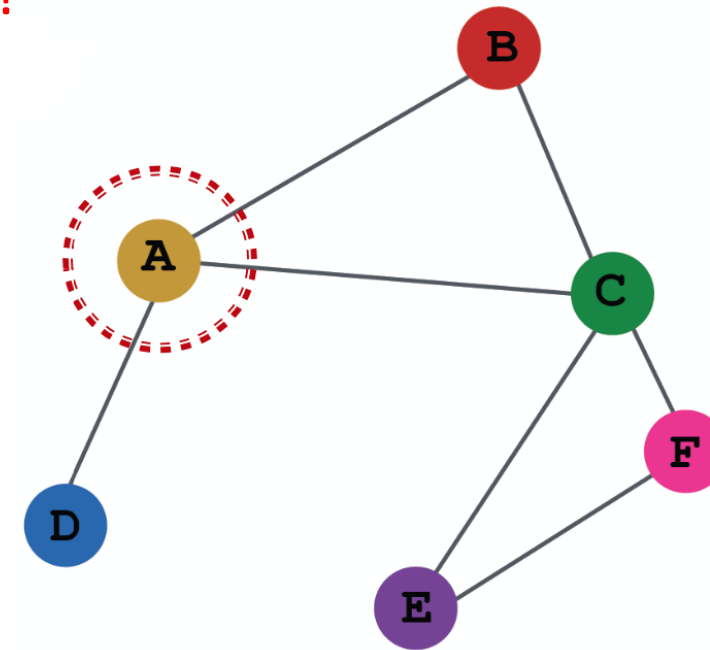
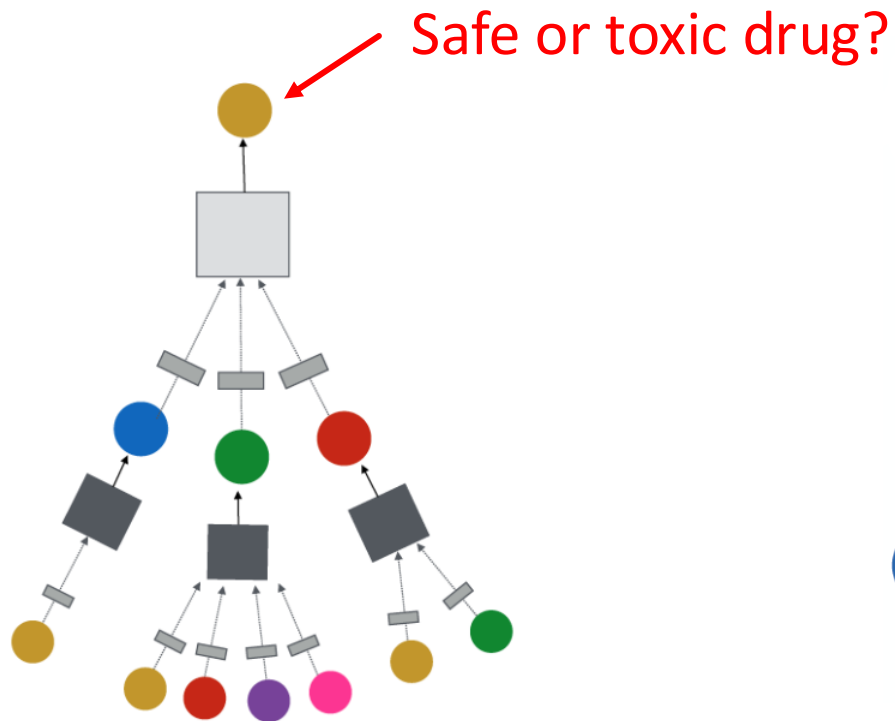
$$\mathbf{H}^k = \sigma(\tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{H}^{k-1} \mathbf{W}_k)$$

where $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}_N$ is the adjacency matrix with added self-connections, $\tilde{\mathbf{D}}_{ii} = \sum_j \tilde{\mathbf{A}}_{ij}$ is the degree matrix.



GCN: Training

- Train in a supervised manner: Directly train the end-to-end model for a supervised task (e.g., node classification).



drug-drug interaction network



- Train in an unsupervised manner:
 - Use only the graph structure.
 - “Similar” nodes have similar embedding.
- How to find similar nodes?
 - Deepwalk, node2vec, ...





GRAPHSAGE

- So far we have aggregated the neighbor messages by taking their (weighted) average.

$$\mathbf{h}_v^k = \sigma \left(\mathbf{W}_k \sum_{u \in N(v)} \frac{\mathbf{h}_u^{k-1}}{|N(v)|} + \mathbf{B}_k \mathbf{h}_v^{k-1} \right).$$

It is very straightforward and simple.

- Can we make it more sophisticated to learn more latent information from a graph?

GraphSAGE

Algorithm 1: GraphSAGE embedding generation (i.e., forward propagation) algorithm

Input : Graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$; input features $\{\mathbf{x}_v, \forall v \in \mathcal{V}\}$; depth K ; weight matrices $\mathbf{W}^k, \forall k \in \{1, \dots, K\}$; non-linearity σ ; differentiable aggregator functions $\text{AGGREGATE}_k, \forall k \in \{1, \dots, K\}$; neighborhood function $\mathcal{N} : v \rightarrow 2^{\mathcal{V}}$

Output : Vector representations \mathbf{z}_v for all $v \in \mathcal{V}$

```
1  $\mathbf{h}_v^0 \leftarrow \mathbf{x}_v, \forall v \in \mathcal{V}$  ;
2 for  $k = 1 \dots K$  do
3   for  $v \in \mathcal{V}$  do
4      $\mathbf{h}_{\mathcal{N}(v)}^k \leftarrow \text{AGGREGATE}_k(\{\mathbf{h}_u^{k-1}, \forall u \in \mathcal{N}(v)\})$ ;
5      $\mathbf{h}_v^k \leftarrow \sigma \left( \mathbf{W}^k \cdot \text{CONCAT}(\mathbf{h}_v^{k-1}, \mathbf{h}_{\mathcal{N}(v)}^k) \right)$ 
6   end
7    $\mathbf{h}_v^k \leftarrow \mathbf{h}_v^k / \|\mathbf{h}_v^k\|_2, \forall v \in \mathcal{V}$ 
8 end
9  $\mathbf{z}_v \leftarrow \mathbf{h}_v^K, \forall v \in \mathcal{V}$ 
```

Generalizes the aggregation function

Replay sum by concat



Mean aggregator

$$\text{AGGREGATE}_k = \sum_{u \in N(v)} \frac{h_u^{k-1}}{|N(v)|}$$

- Nearly equivalent to the convolutional propagation rule used in GCN.
- This concatenation can be viewed as a simple form of a “skip connection” between the different layers.



LSTM aggregator

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

- LSTMs have the advantage of larger expressive capability.
- Apply LSTM to random permutation of the node's neighbors $\pi(N(v))$.



Pooling aggregator

$$\text{AGGREGATE}_k = \max(\{\sigma(\mathbf{W}_{pool} \mathbf{h}_{u_i}^k + \mathbf{b}), \forall u_i \in N(v)\})$$

where max is taken element-wise.

- By applying the max-pooling operator to each of the computed features, the model effectively captures different aspects of the neighborhood set.



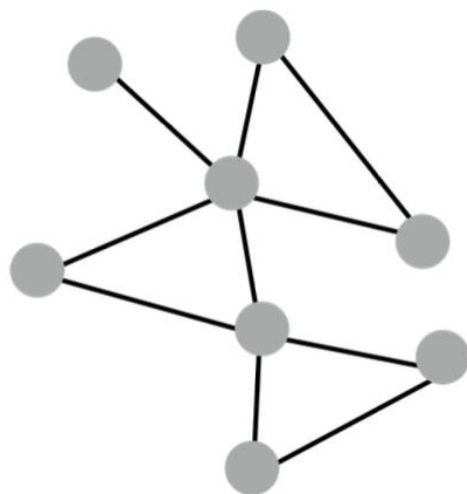
Transductive vs. Inductive

- GNN is usually in a semi-supervised learning manner.
 - The unlabelled node is involved during training.
- Semi-supervised learning can be grouped into two categories:
 - **Transductive**: The testing data is from the unlabelled data.
 - **Inductive**: The testing data is unseen in training.

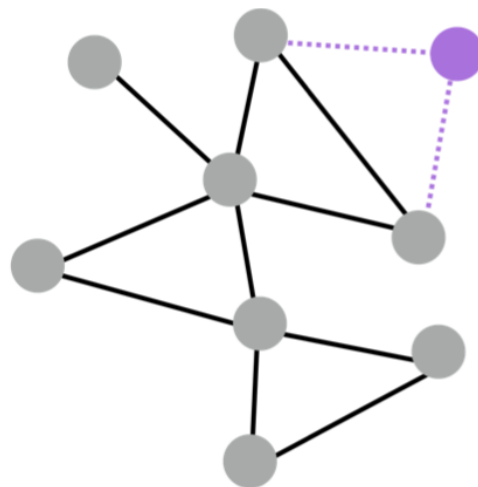


Inductive Capacity for New Nodes

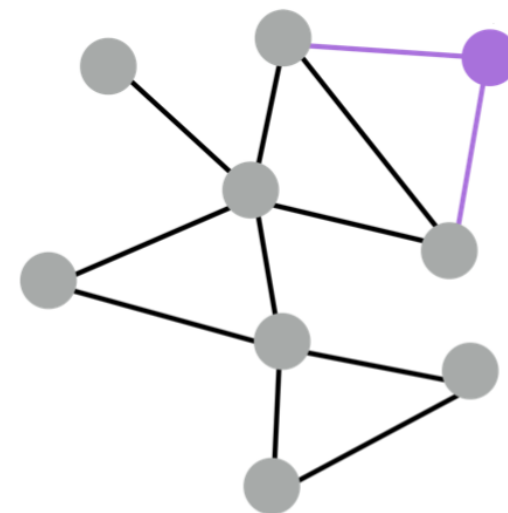
- Many application settings constantly encounter previously unseen nodes.
- E.g. new user and new item in a recommendation system.



Train on known graph



New node arrives

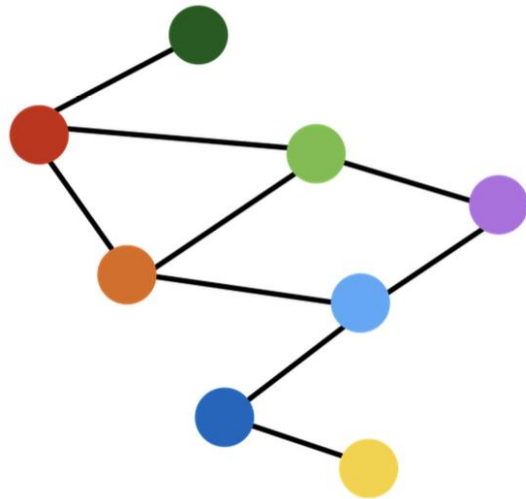


Generate embedding
for new node

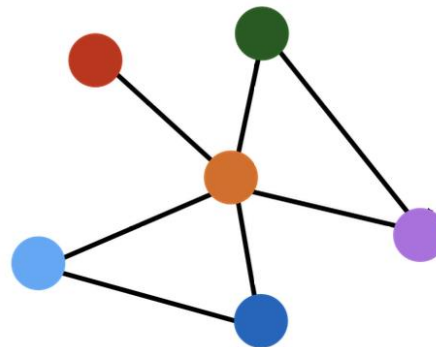


Inductive Capacity for New Graphs

- The trained GCN parameters can also be used to generalize to entirely unseen graphs.
- E.g. train on protein interaction graph from model organism A and generate embeddings on newly collected data about organism B.



Train on one graph



Generalize to new graph





GAT

- Check the neighborhood aggregation of GCN again:

$$\mathbf{h}_v^k = \sigma \left(\mathbf{W}_k \sum_{u \in N(v)} \frac{\mathbf{h}_u^{k-1}}{|N(v)|} + \mathbf{B}_k \mathbf{h}_v^{k-1} \right).$$

- What is the weight of each neighbor $u \in N(v)$ that contributes to node v ?

$$\frac{1}{|N(v)|}$$

- It simply assumes that all neighbors are equally important to node v .

- Can we simply learn a weight for each node in the graph?
 - Important node (e.g. with large degree) deserves large weight.
- Probably not.
- The importance of each node to each neighbor should be different.
- Goal: Specify **arbitrary importance** to different neighbors of each node in the graph.
- Idea: Compute embedding \mathbf{h}_v^k of each node in the graph following an **attention network**.



- First compute attention coefficients of e_{vu} across node v , and its neighbor u based on their representation at layer $k - 1$:

$$e_{vu} = a(\mathbf{W}_k \mathbf{h}_u^{k-1}, \mathbf{W}_k \mathbf{h}_v^{k-1})$$

- e_{vu} indicates the importance of node u message to node v .
- The attention network a can just be a simple single-layer neural network:

$$a(\mathbf{p}, \mathbf{q}) = \mathbf{A}^T [\mathbf{p}, \mathbf{q}]$$

where \mathbf{A} is a learnable parameter.



- Then normalize over all neighbors to get the weight α_{vu} :

$$\alpha_{vu} = \frac{\exp e_{vu}}{\sum_{k \in N(v)} \exp e_{vk}}$$

- The final attention-weighted aggregation is:

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



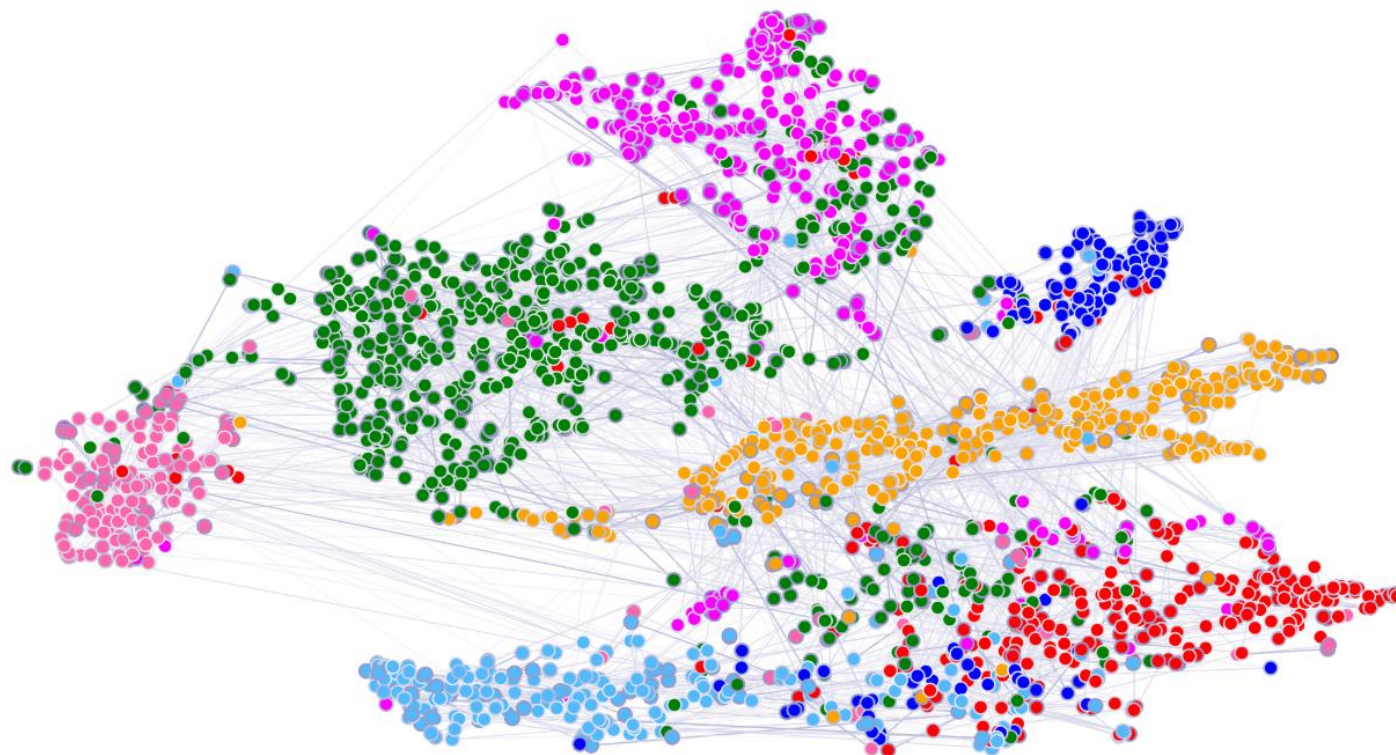
GAT: Multi-Head Attention

- Borrow the idea of multi-head attention from Transformer:

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

- We got T head and each head t has its own weights.





A t-SNE plot of the computed feature representations of a pre-trained GAT model's first hidden layer on the Cora dataset. Node colors denote classes. Edge thickness attention coefficient.

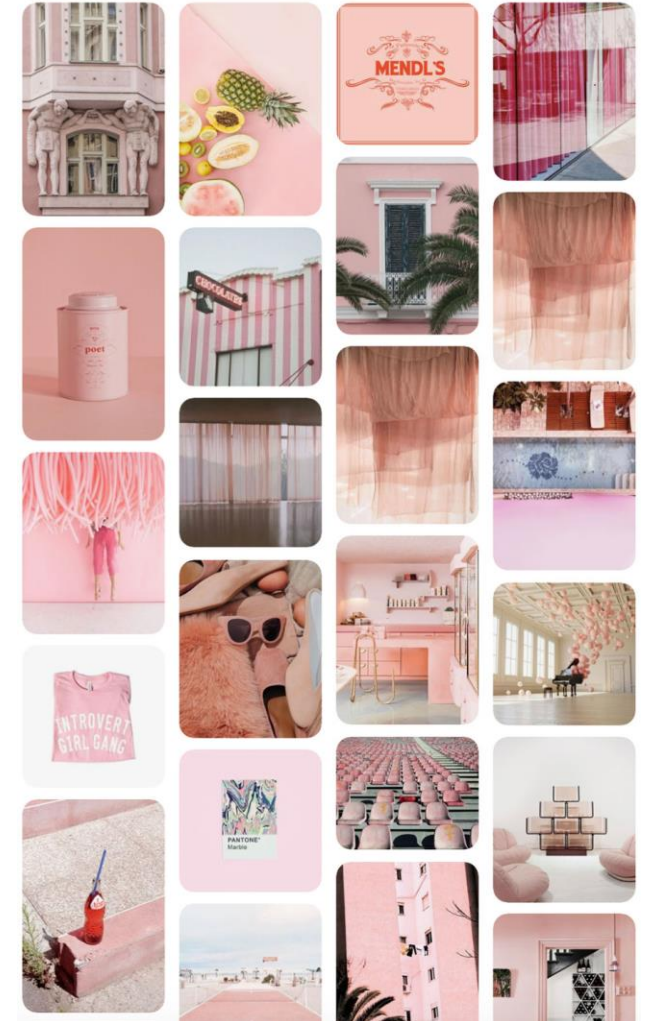




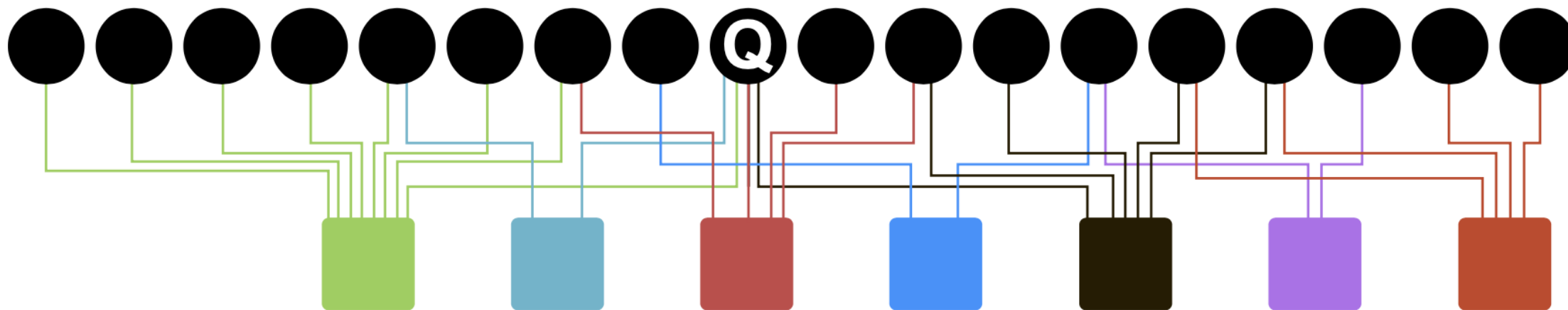
APPLICATION TO RECOMMENDER SYSTEM



- Pinterest is an American image sharing and social media service.
- Users can save and discover images, GIFs and videos in the form of pinboards.
- 300M users, 4+B pins, 2+B pinboards.



Pinterest Graph

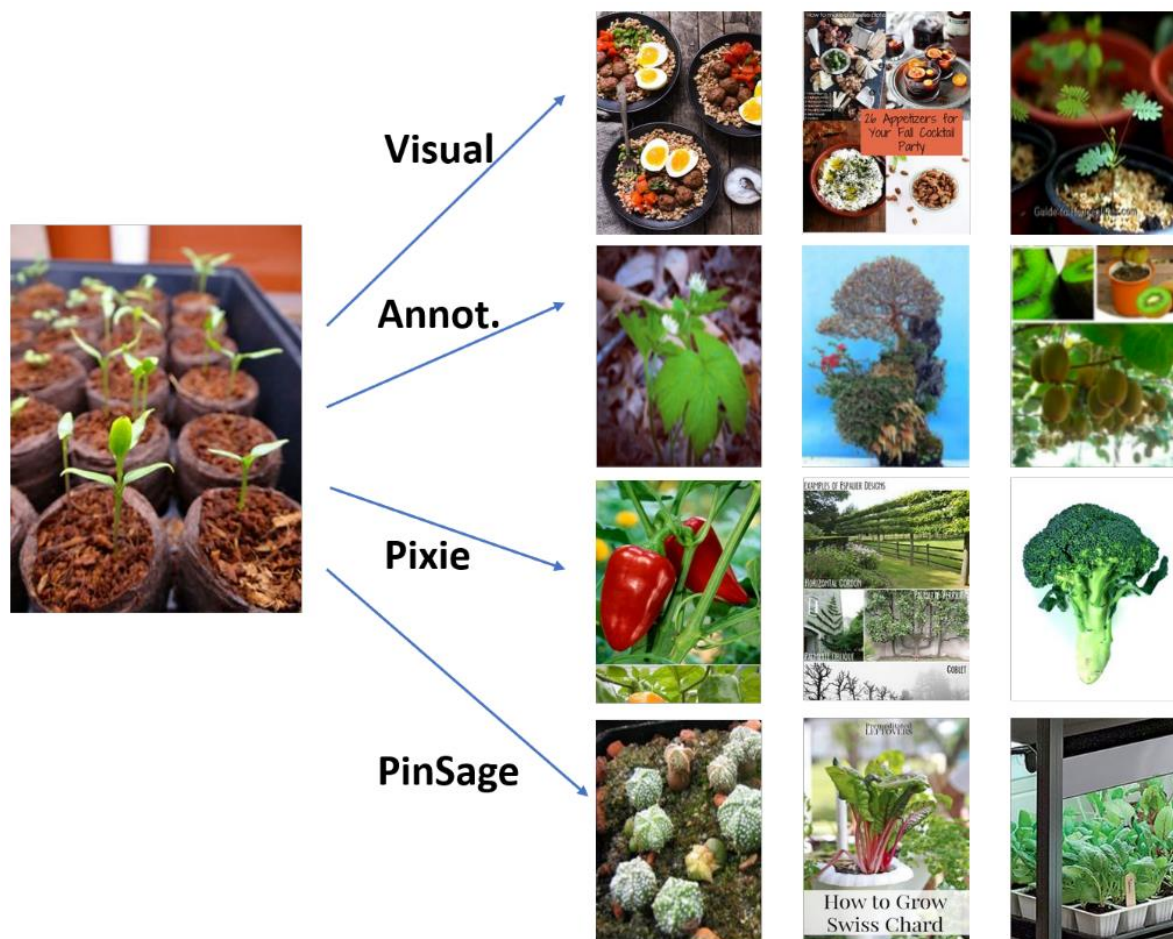


- Graph: 2B pins, 1B boards, 20B edges.
 - Graph is dynamic: Need to apply to new pins and new boards without model retraining.
- Rich node features: images, text with pins.

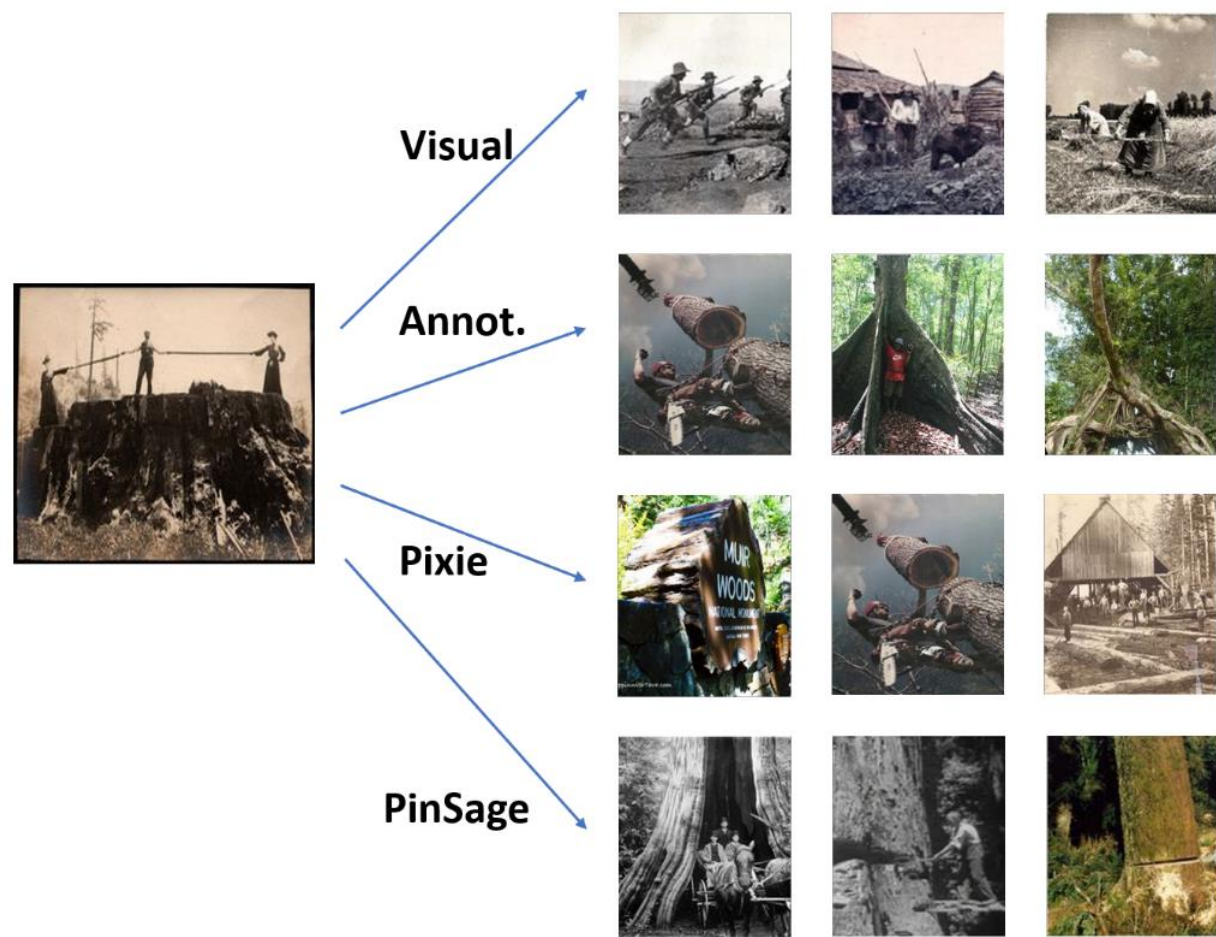


- Goal: Generate pin embeddings in a web-scale Pinterest graph containing billions of objects.
- Pin embeddings are essential to various tasks like pin recommendation, classification, clustering, ranking.
- Services like “Related Pins”, “Search”, “Shopping”, “Ads”.

PinSage: Result



PinSage: Result





ADVANCES



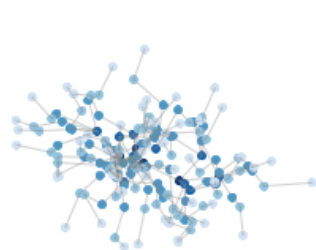
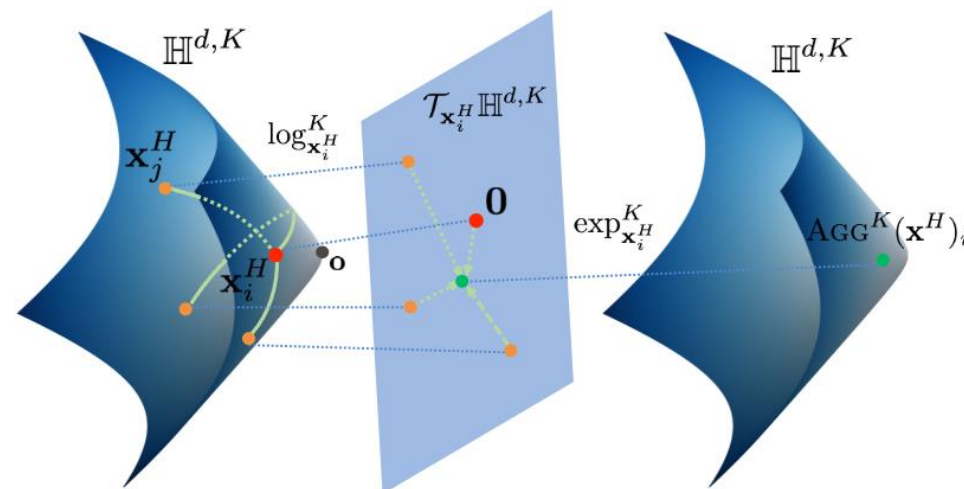
HGCN

Hyperbolic graph convolutional neural networks

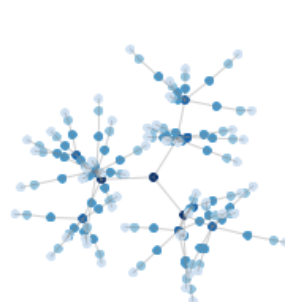
[I Chami, Z Ying, C Ré...](#) - Advances in neural ..., 2019 - proceedings.neurips.cc

... and scale-free **graphs** in inductive settings: (1) We ... **hyperbolic** space to transform input features which lie in Euclidean space into **hyperbolic** embeddings; (2) We introduce a **hyperbolic** ...

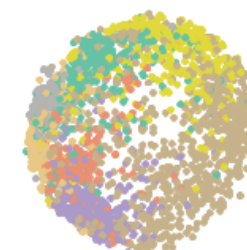
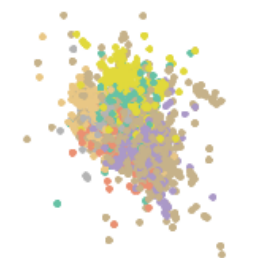
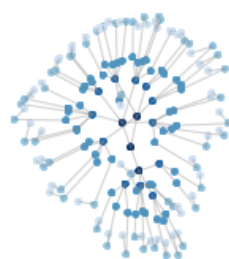
☆ Save ㉞ Cite Cited by 541 Related articles All 17 versions ㉞



(a) GCN layers.



(b) HGCN layers.



(c) GCN (left), HGCN (right).



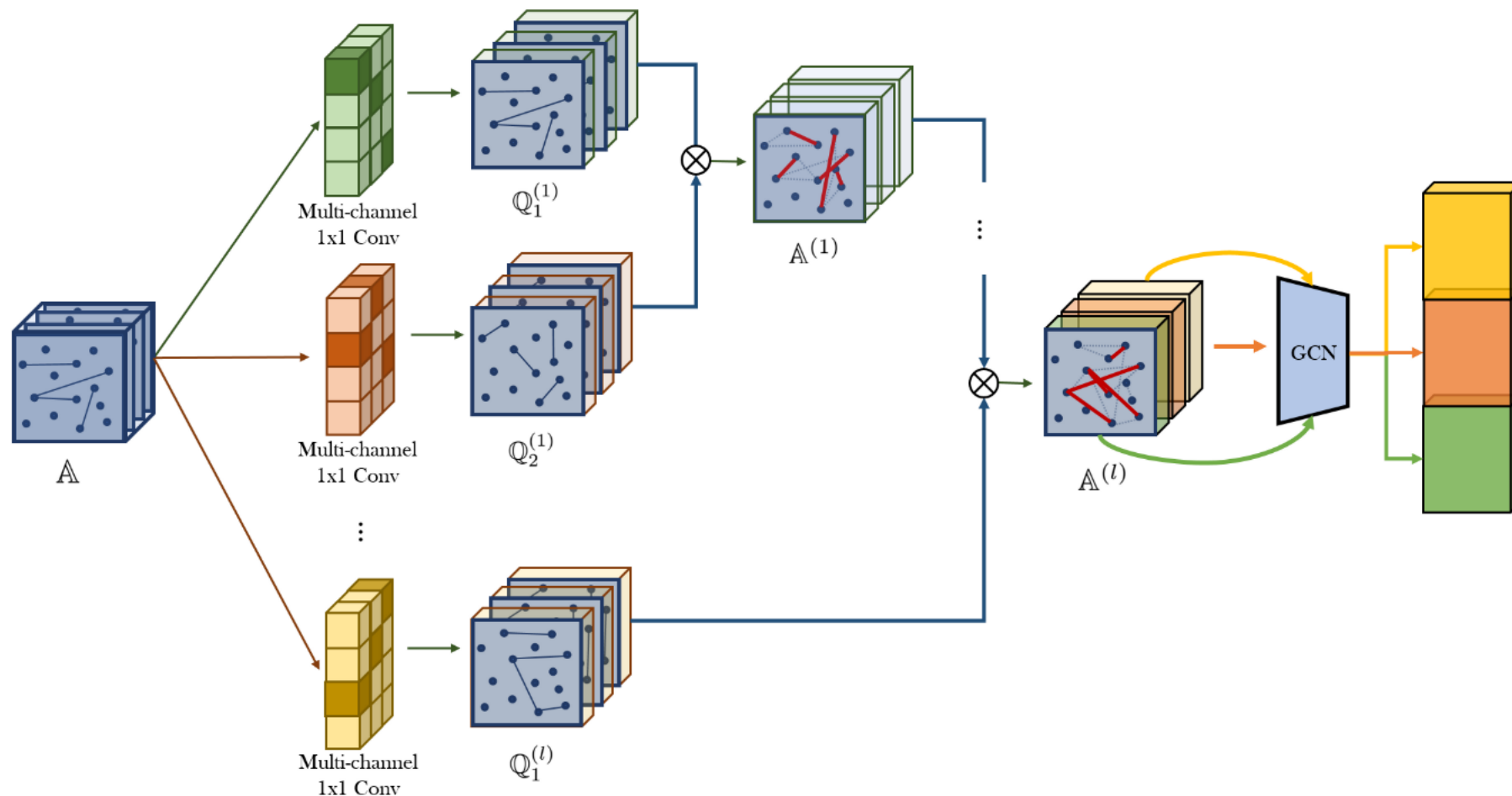
GTN

Graph transformer networks

S Yun, M Jeong, R Kim, J Kang... - Advances in neural ..., 2019 - proceedings.neurips.cc

... **Graph Transformer Network (GTN)** that learns to transform a heterogeneous input **graph** into useful meta-path **graphs** for each task and learn node representation on the **graphs** in an ...

☆ Save Cite Cited by 733 Related articles All 11 versions



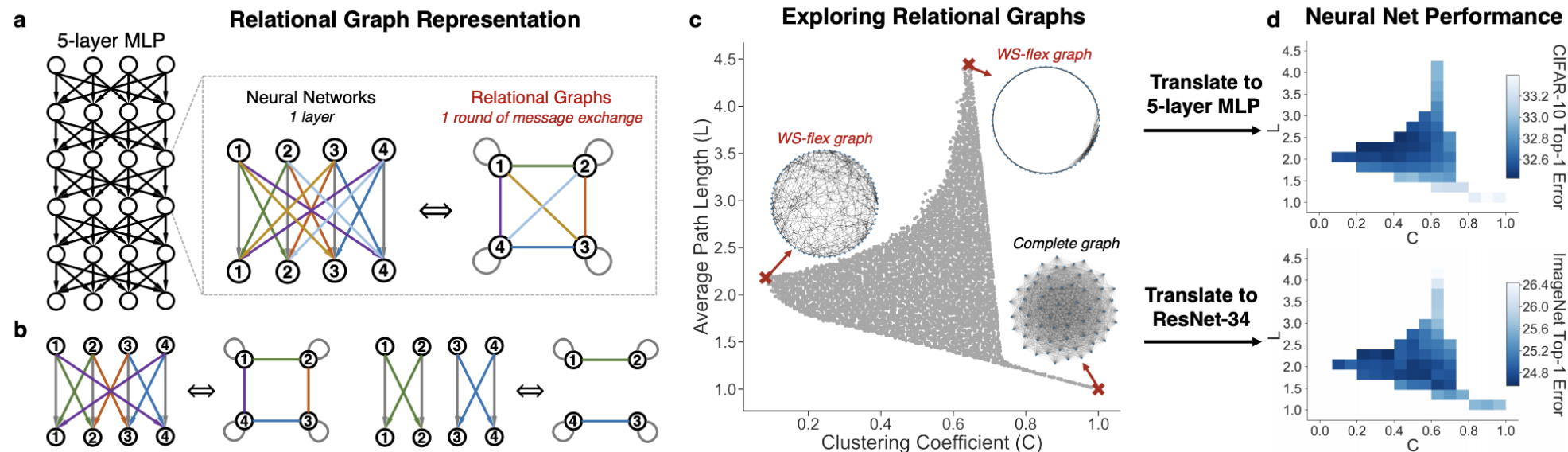
Relational Graph

Graph structure of neural networks

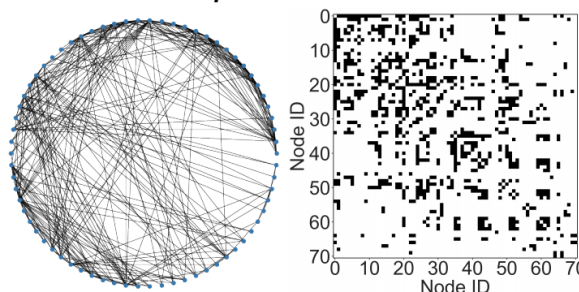
J You, J Leskovec, K He, S Xie - International Conference on ..., 2020 - proceedings.mlr.press

... **graphs**. Here we systematically study the relationship between the **graph structure** of a **neural network** ... of representing a **neural network** as a **graph**, which we call **relational graph**. Our ...

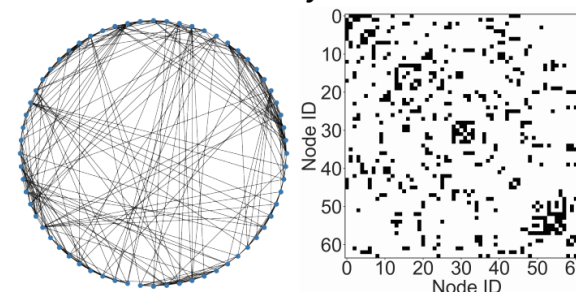
☆ Save 📄 Cite Cited by 132 Related articles All 11 versions 🔗



Biological neural network:
Macaque whole cortex



Artificial neural network:
Best 5-layer MLP



Conclusion

After this lecture, you should know:

- What is a graph representation?
- How does random walk help generate graph representation?
- What kind of role do BFS and DFS play in node2vec?
- What is the basic architecture of GNN?
- How is attention applied to GNN?



Suggested Reading

- 深度学习中不得不学的Graph Embedding方法
- 关于Node2vec算法中Graph Embedding同质性和结构性的进一步探讨



- Tutorial at WWW 2019 on Representation Learning on Networks
- CS224W Machine Learning with Graphs



Thank you!

- Any question?
- Don't hesitate to send email to me for asking questions and discussion. 😊

