# CS 190I
# Deep Learning
# Graph Neural Networks

Lei Li (leili@cs)

UCSB

# Recap

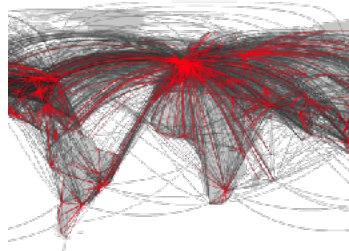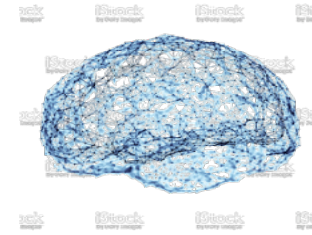| | training objective | backbone | size(#params) | training data (#tokens) |
|---|---|---|---|---|
| ELMo | next token prediction | two separate LSTM | 94M | 5.5 billion |
| BERT | masked token prediction + next sentence prediction | Transformer Encoder | 110M 340M | 3.3 billion |
| GPT-3 | next token prediction | Transformer Decoder | 175B | 500 billion |

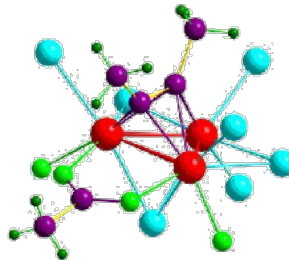# Graph Data is everywhere
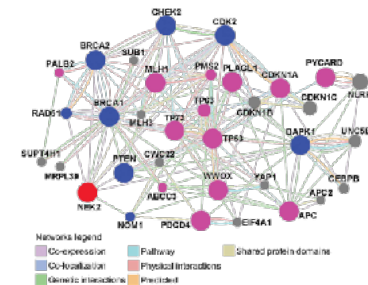
Social Graphs

Transportation Graphs
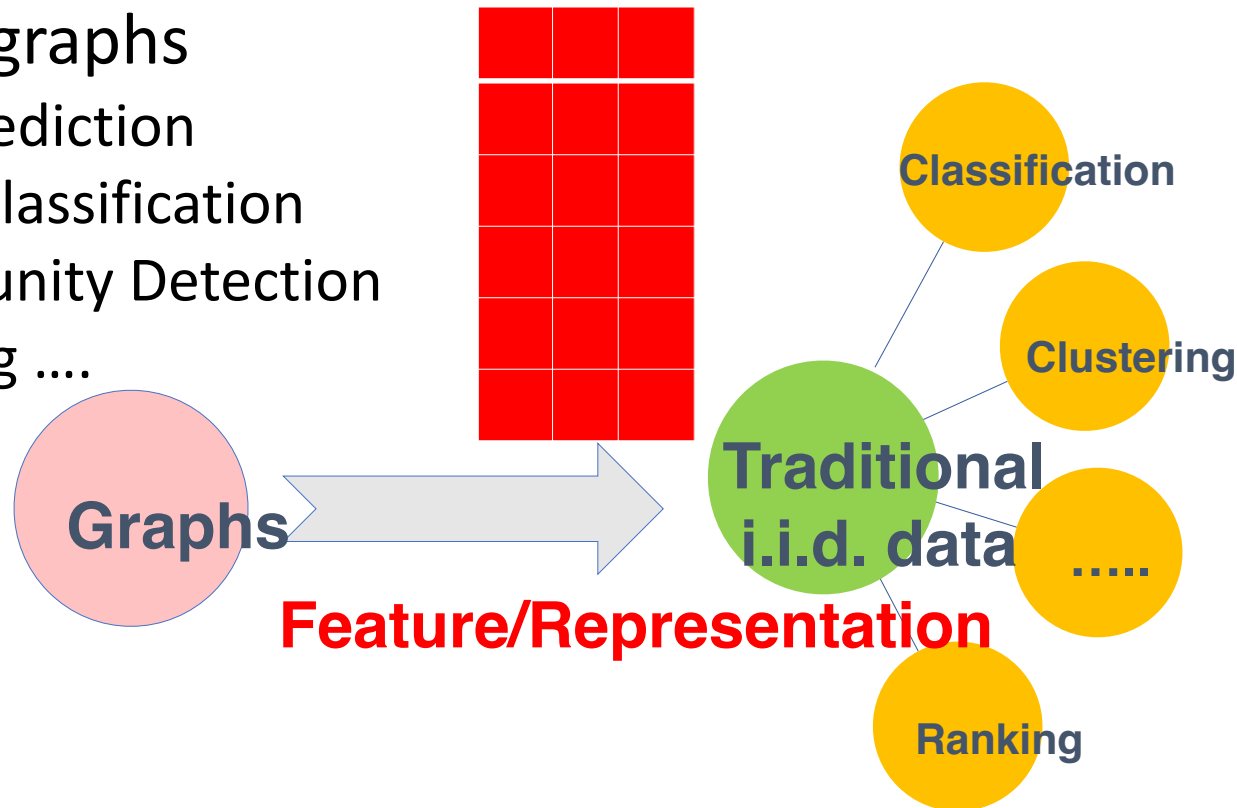
Brain Graphs

Web Graphs

Molecular Graphs

Gene Graphs

# ML on Graphs

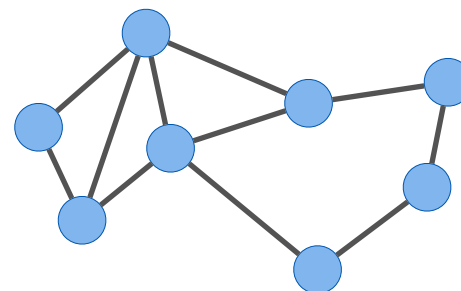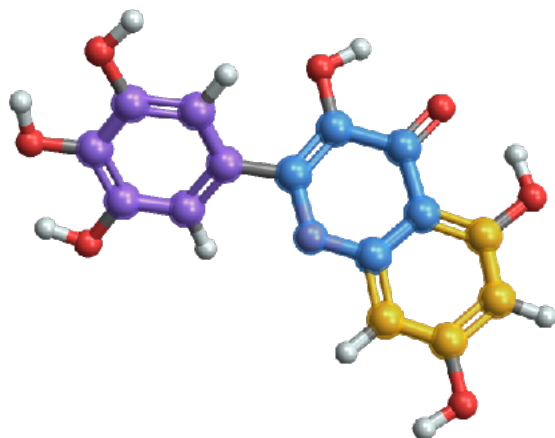Numerous real-world problems can be summarized as a set of tasks on graphs

- Link prediction
- Node Classification
- Community Detection
- Ranking ....

**Graphs**

**Feature/Representation**

**Traditional i.i.d. data**

**Classification**

**Clustering**

**.....**

**Ranking**

Slides adapted from Yao Ma & Jiliang Tang@MSU
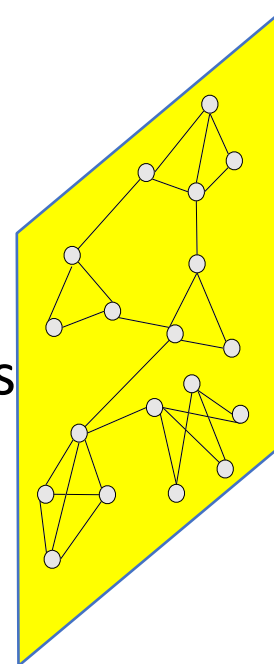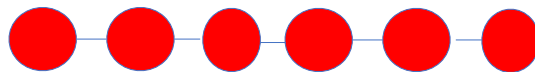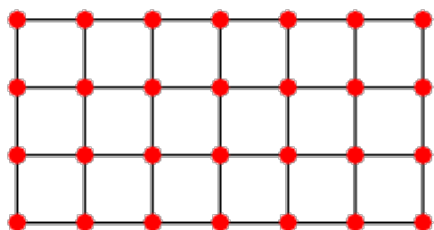
# Example: predict toxicity of a drug

Toxic?

# Deep Learning Meets Graphs: Challenges

Traditional DL is designed for simple grids or sequences

- CNNs for fixed-size images/grids
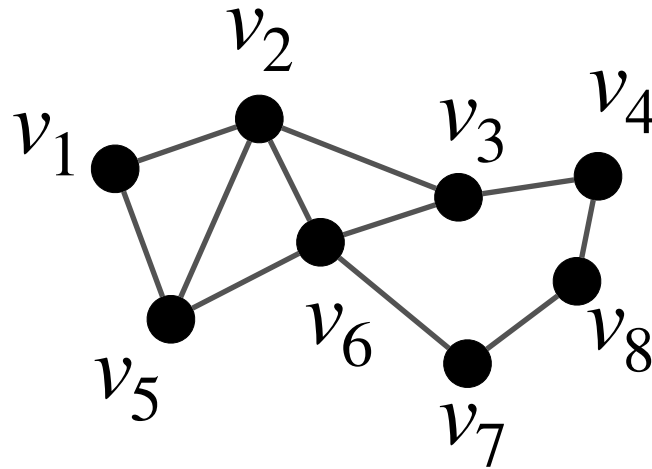- RNNs for text/sequences



But nodes on graphs have different connections

- Arbitrary neighbor size
- Complex topological structure
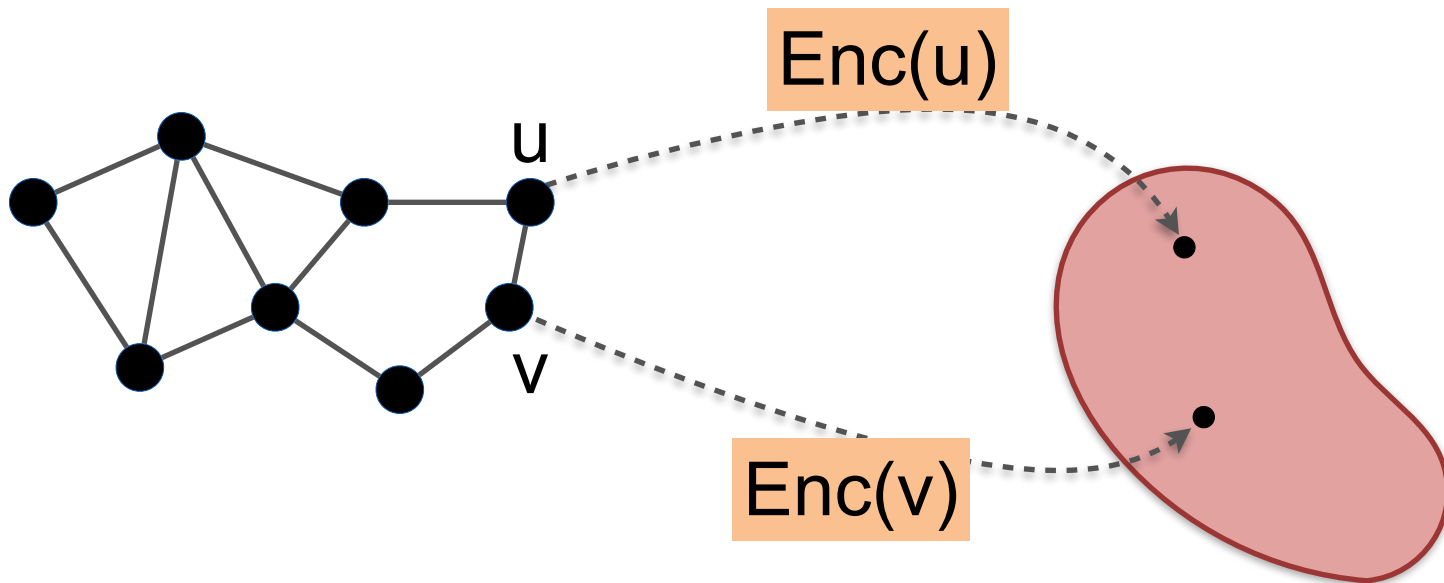- No fixed node ordering

# Graph Representation



Graph: $G = \{V, E\}$

Nodes: $V = \{v_1, v_2, \ldots, v_N\}$

Edges: $E = \{e_1, e_2, \ldots, e_M\} \subset V \times V$

# Node Embedding
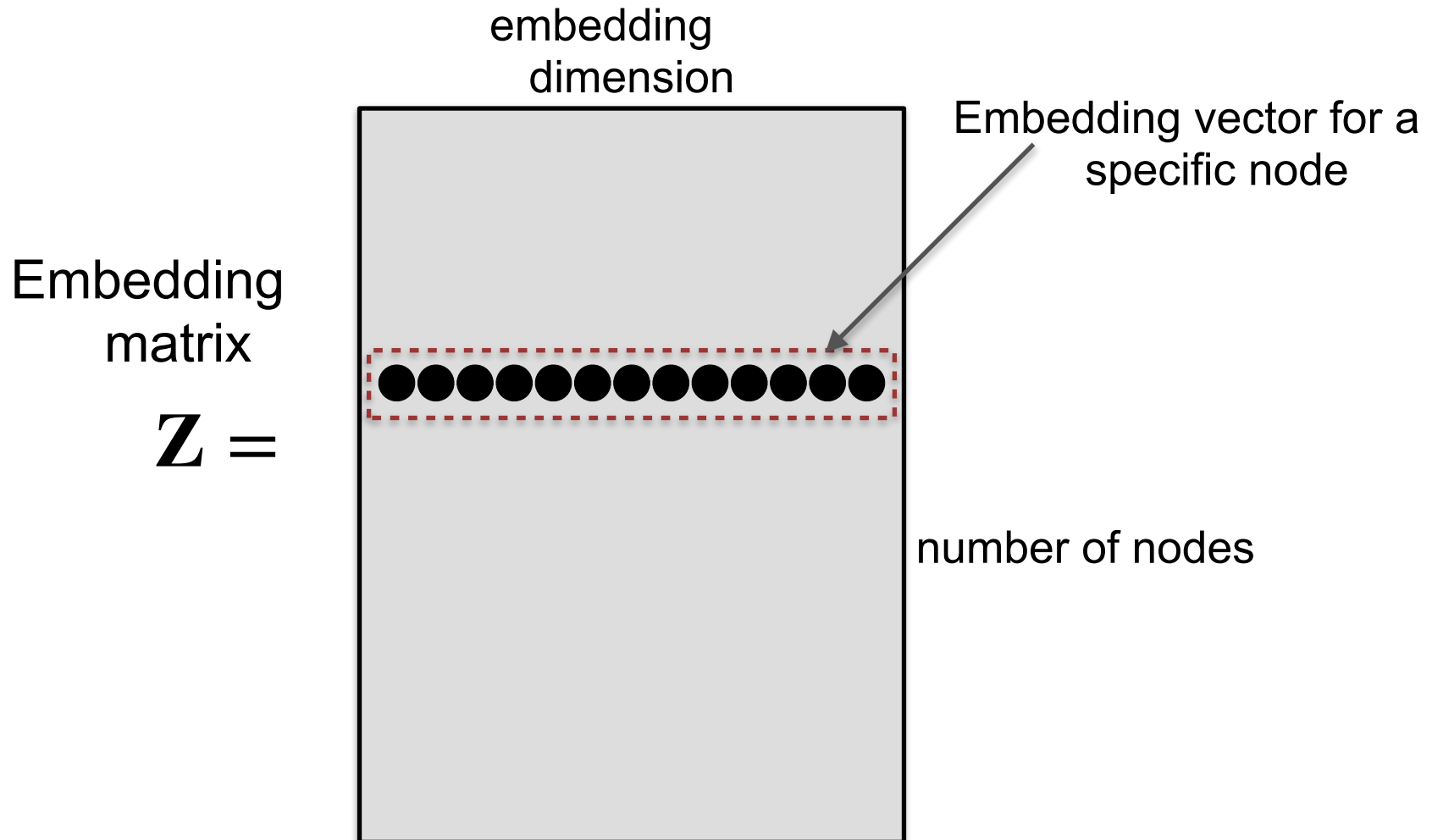
$$Enc(\,\cdot\,) : V \rightarrow \mathbb{R}^d$$

# "Shallow" Node Embedding

• is just a lookup-table

embedding
dimension

Embedding vector for a
specific node

Embedding
matrix

$$\mathbf{Z} =$$

number of nodes

# Deep Graph Neural Network

Graph
Convolution

Graph
Convolution

Activation
function

· · ·

Output is embedding matrix for nodes
for further downstream tasks: e.g. node classification

# Graph Neural Network

Every node's neighbor defines a convolutional kernel



aggregate information
from its neighbors
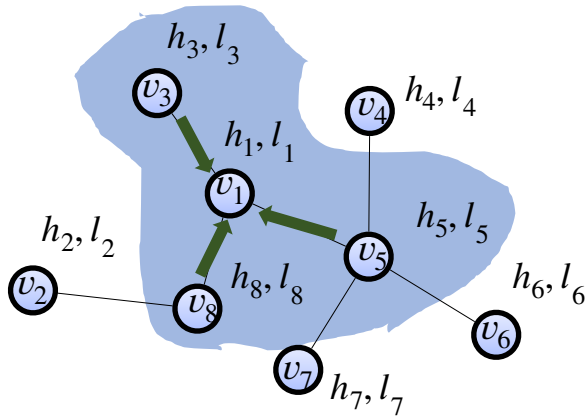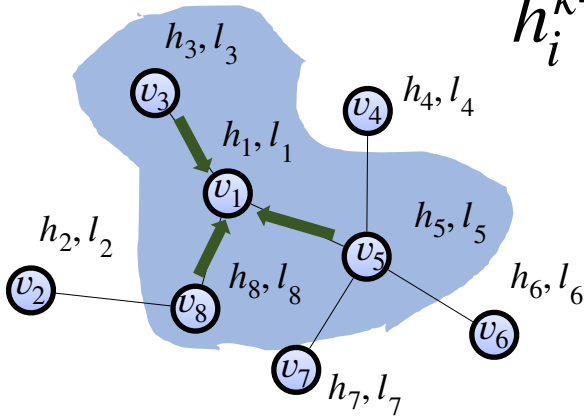
# Aggregate Neighbors

$h_i$: node (hidden) embedding vector
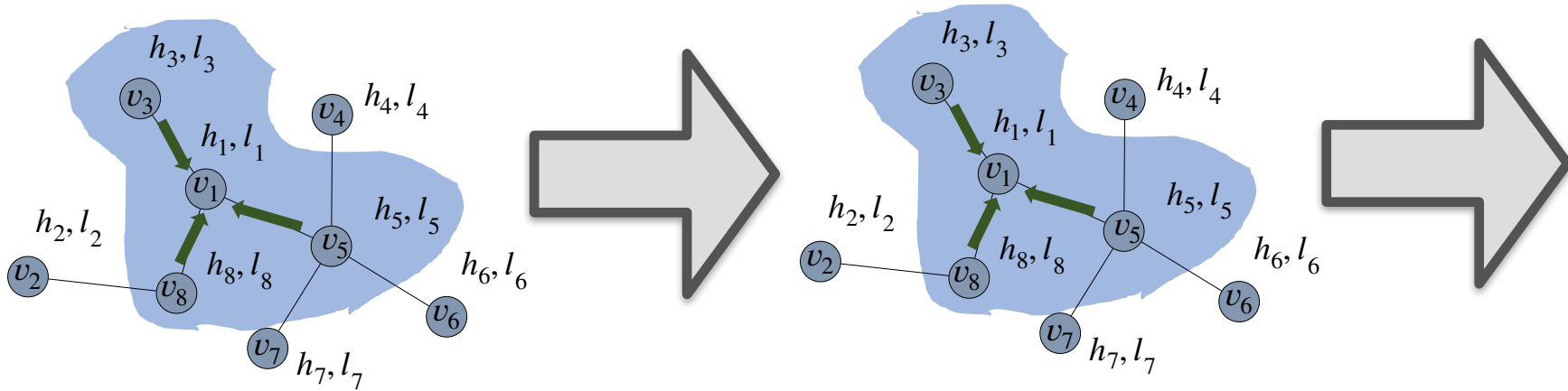
aggregate information from its neighbors

$$h_i^{k+1} = \text{Aggregate}_{v_j \in N(v_i)} f(h_i^k, h_j^k), \forall v_i \in V$$



$N(v_i)$: Neighbors of the node $v_i$.

$f(\cdot)$: Feedforward network.

# Multiple Computation Layers
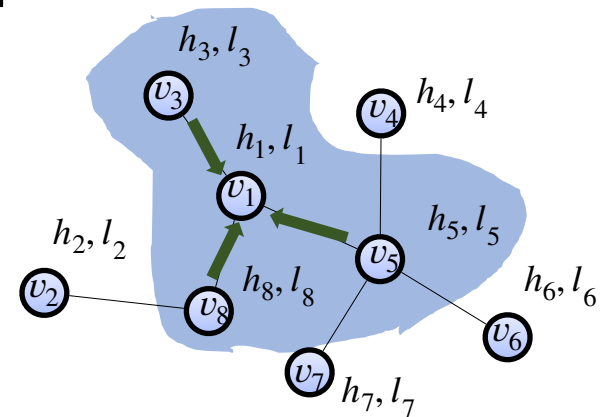
# A Simple Graph Convolution Layer

- Simple approach: averaging neighbor's message and apply nonlinear transformation

initial embedding: $\quad h_i^0 = x_i$

computing next layer:

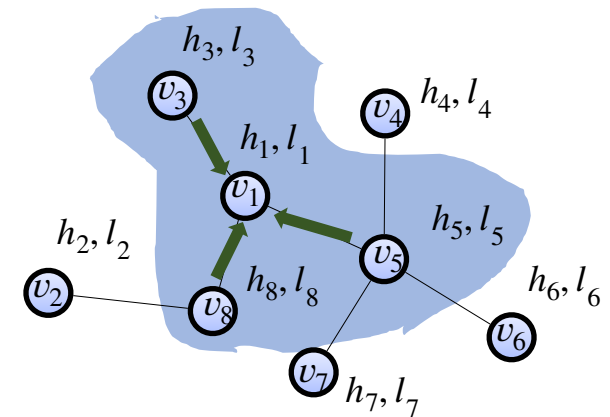$$h_i^{k+1} = \sigma(W_k \frac{1}{|N(v_i)|} \sum_{v_j \in N(v_i)} h_j^k + B_k h_i^k)$$

$$h_1^2 = tanh\left(W_1 \cdot \frac{1}{3}(h_3^1 + h_5^1 + h_8^1) + B_1 h_1^1\right)$$
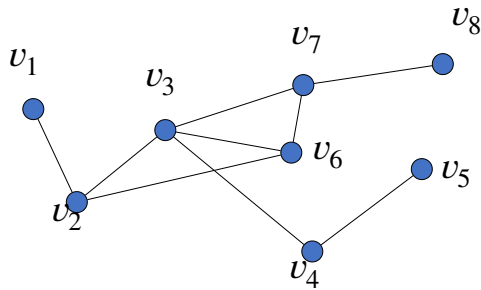
14

# A Simple Graph Convolution Layer

- More layers:



$$h_1^{(3)} = tanh\left(W_2 \cdot \frac{1}{3}(h_3^{(2)} + h_5^{(2)} + h_8^{(2)}) + B_2 h_1^{(2)}\right)$$

# Matrix Representations of Graphs

Adjacency Matrix: $A[i,j] = 1$ if $v_i$ is adjacent to $v_j$

$$A[i,j] = 0, \text{ otherwise}$$



Adjacency Matrix $A$

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

**Spectral graph theory.** American Mathematical Soc.; 1997.

# Matrix Representation of GCN

- Neighbor Aggregation can be performed efficiently using matrix operations
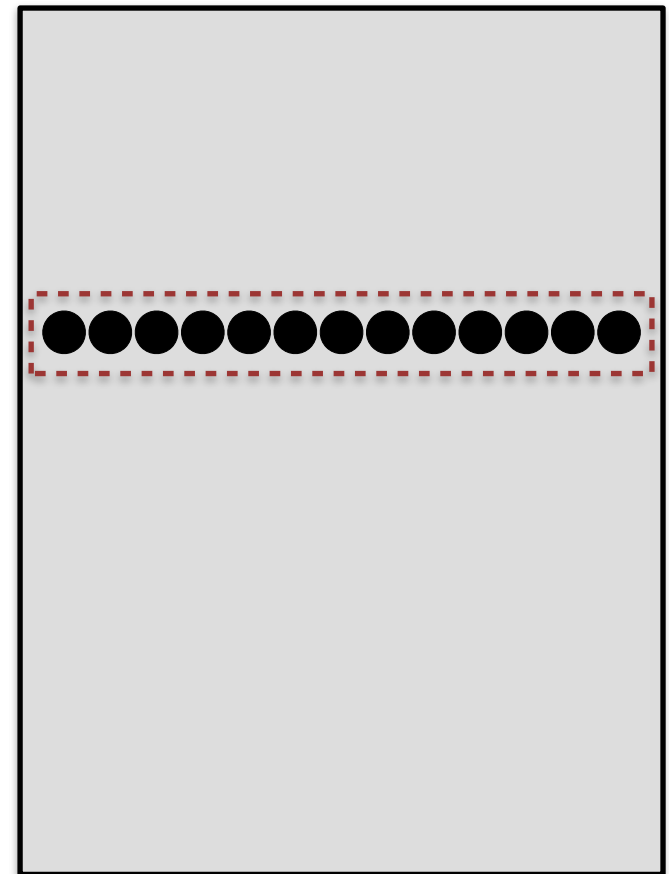
$$H^k = [h_1^k, \ldots, h_{|V|}^k]^T$$

Then $\displaystyle\sum_{v_j \in N(v_i)} h_j^k = A_{i,:} H^k$

Let D be diagonal matrix

$$D_{i,i} = \text{Degree}(v_i) = \sum_j A_{i,j}$$

Then $\displaystyle\frac{1}{|N(v_i)|} \sum_{v_j \in N(v_i)} h_j^k = D^{-1} A H^k$

# Matrix Representation of GCN

- Neighbor Aggregation can be performed efficiently using matrix operations

$$H^k = [h_1^k, \ldots, h_{|V|}^k]^T$$

$$\tilde{A} = D^{-1}A$$

$$H^{k+1} = \sigma(\tilde{A}H^k \cdot W_k^T + H^k B_k^T)$$

# Graph Convolution Network

- Neighbor Aggregation can be performed efficiently using matrix operations

- To make $\tilde{A}$ symmetric

$$H^k = [h_1^k, \ldots, h_{|V|}^k]^T$$

$$\tilde{A} = D^{-\frac{1}{2}} A D^{-\frac{1}{2}}$$
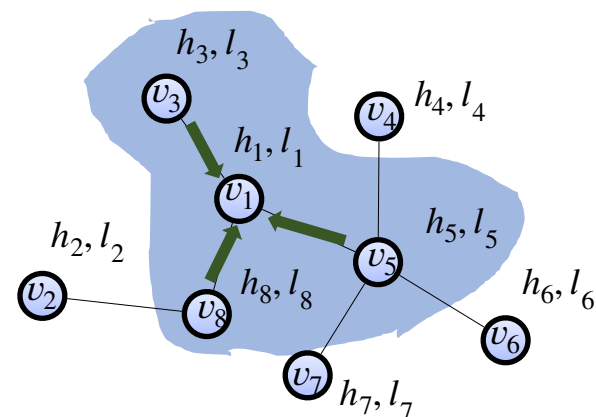
$$H^{k+1} = \sigma(\tilde{A} H^k \cdot W_k^T + H^k B_k^T)$$

# **Prediction Layer**

- For node classification:

$$o_i = \text{Softmax}(h_i^{(m)})$$

- For graph classification:

$$o = \text{Softmax}(\frac{1}{N}\sum_i h_i^{(m)})$$

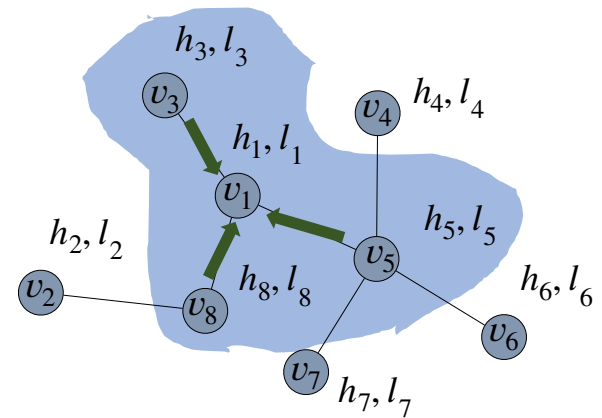# Property: Equivariant

- the embeddings computed from graph convolution layers is invariant to node permutation

$$h_i^0 = x_i$$

$$h_i^{k+1} = \sigma(W_k \frac{1}{|N(v_i)|} \sum_{v_j \in N(v_i)} h_j^k + B_k h_i^k)$$

# Model Training

- Parameters: weight matrix for each layer

$$h_i^{k+1} = \sigma(W_k \frac{1}{|N(v_i)|} \sum_{v_j \in N(v_i)} h_j^k + B_k h_i^k)$$

- Supervised training: e.g. Node classification
  - Linked nodes have similar embedding

$$L = \sum_i CE(y_i, f(h_i^K)) \qquad f_i = \text{Softmax}(h_i^{(K)})$$

  - $y_i$ is node label

# Model Training

- Parameters: weight matrix for each layer

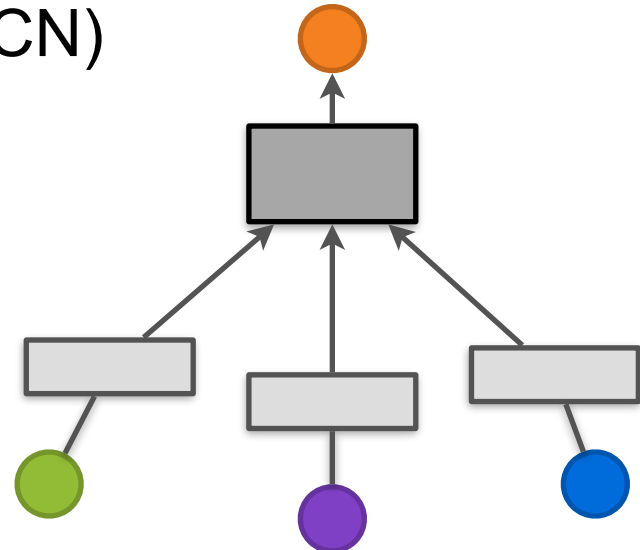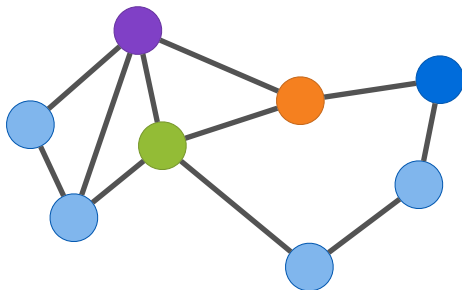$$h_i^{k+1} = \sigma(W_k \frac{1}{|N(v_i)|} \sum_{v_j \in N(v_i)} h_j^k + B_k h_i^k)$$

- Unsupervised training:

  – Linked nodes have similar embedding

$$L = \sum_{i,j} CE(y_{i,j}, Sim(h_i^K, h_j^K))$$

  – $y_{i,j} = 1$ if there is edge from v_i to v_j

  – Similarity can be defined in many ways: e.g. inner product $h_i \cdot h_j$

# Generic GNN framework

- GNN layer = message passing + Aggregation
  - different design choices under this framework
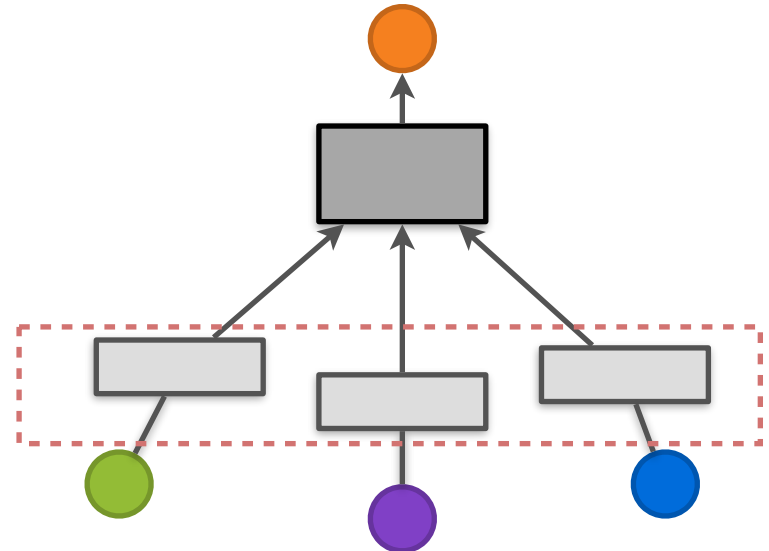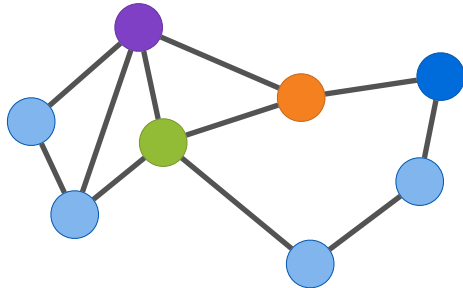  - Graph convolutional network (GCN)
  - GraphSAGE
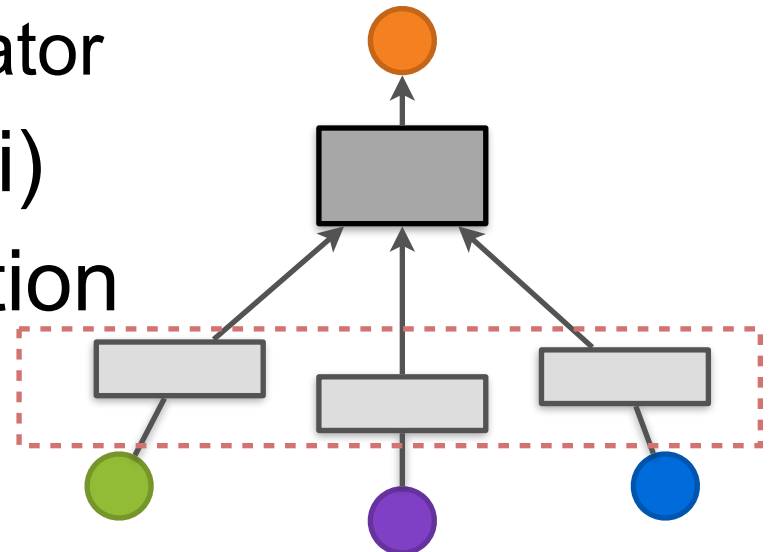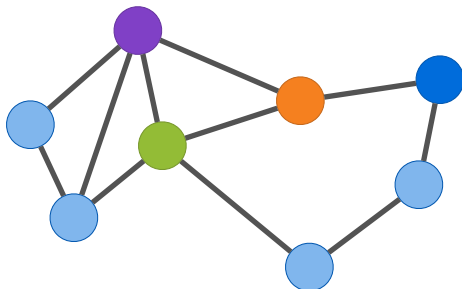  - GAT

# **Message Computation**

- Each node will create a message
- e.g. Linear projection

$$m_i^k = W_k \cdot h_i^{(k)}$$

# Aggregation/Pooling
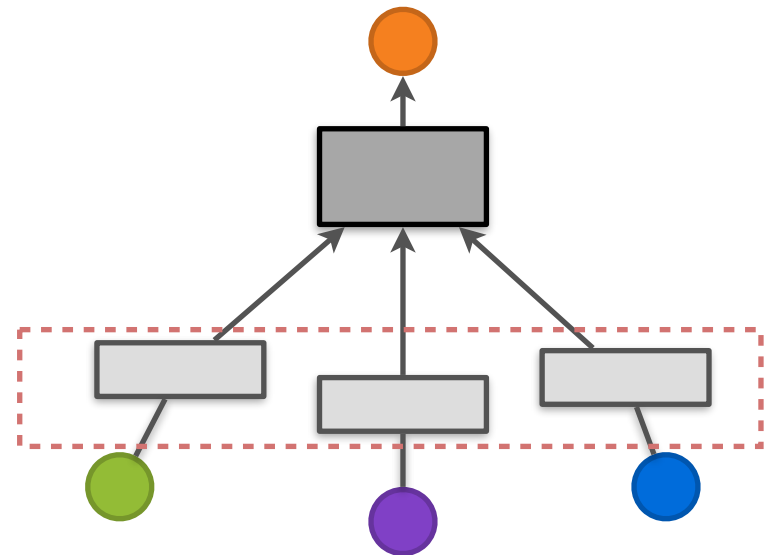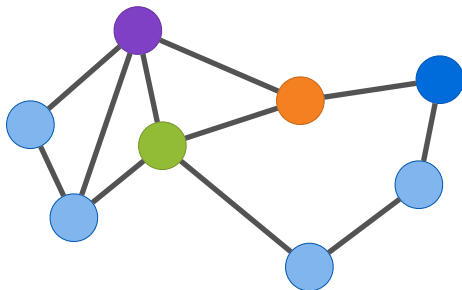
- Each node will aggregate messages from its neighbors
- e.g.
  - Sum, Mean, Max operator
- Concat(AGG{m_j}, m_i)
- Apply nonlinear activation

# GraphSAGE

$$h_i^{k+1} = \sigma\left(W_k \cdot \text{CONCAT}\left(h_i^k, \text{AGG}(\{h_j^k, \forall v_j \in N(v_i)\})\right)\right)$$

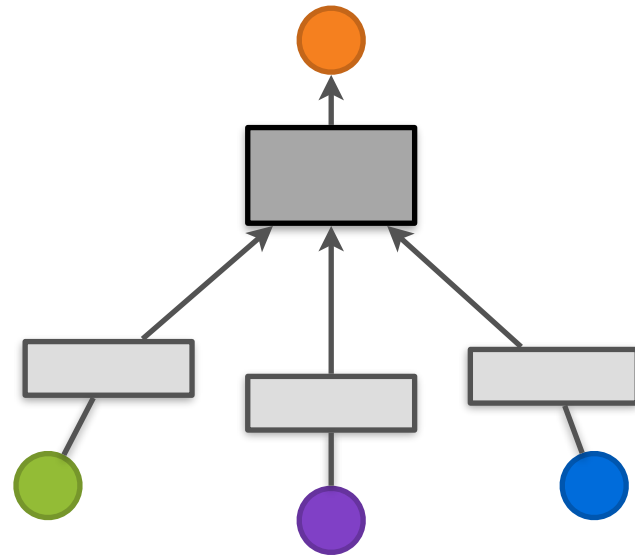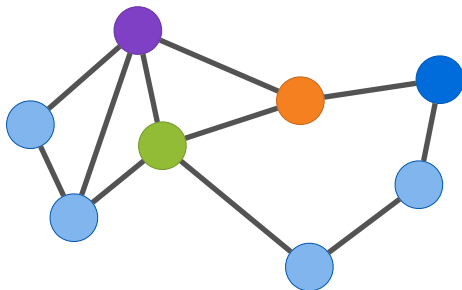AGG can be designed in multiple ways, like pooling (sum, avg, max)

# Graph Attention Network (GAT)

$$h_i^{k+1} = \sigma(\sum_{v_j \in N(v_i)} \alpha_{ij} W_k h_{v_j}^k)$$
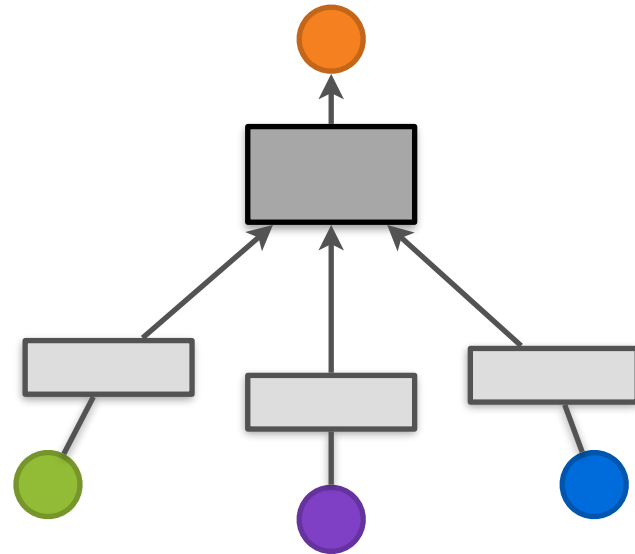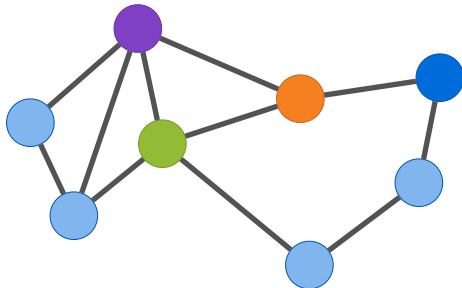
attention weight

$$\alpha_{ij} = \text{Attention}(W_k h_i, W_k h_j) = \frac{\exp(W_k h_i)^T W_k h_j}{\sum_{j'} \exp(W_k h_i)^T W_k h_{j'}}$$

# Multi-head Attention for GAT? Yes

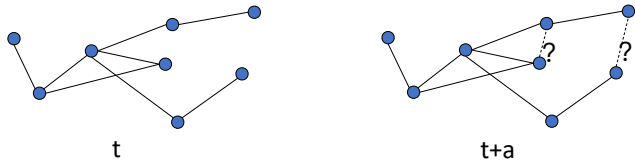$$h_i^{k+1} = \sigma\left( \sum_{v_j \in N(v_i)} \alpha_{ij} W_k h_{v_j}^k \right)$$

$$\alpha_{ij} = \text{Attention}(W_k h_i, W_k h_j) = \frac{\exp(W_k h_i)^T W_k h_j}{\sum_{j'} \exp(W_k h_i)^T W_k h_{j'}}$$
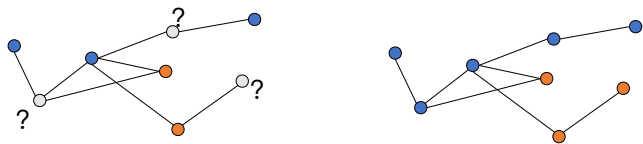
# Tasks on Graph-Structured Data

## Node-level

### Link Prediction



t                    t+a

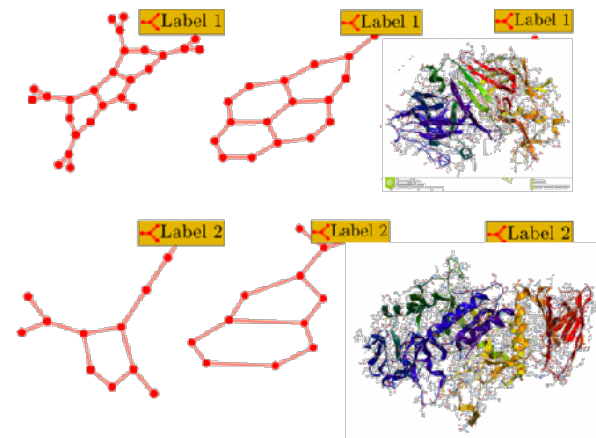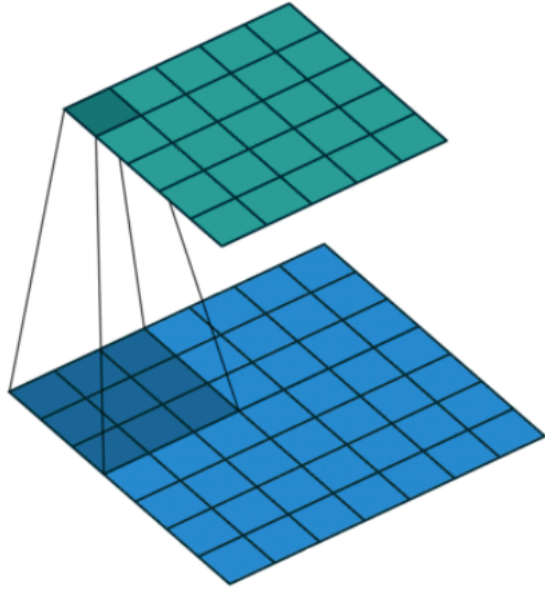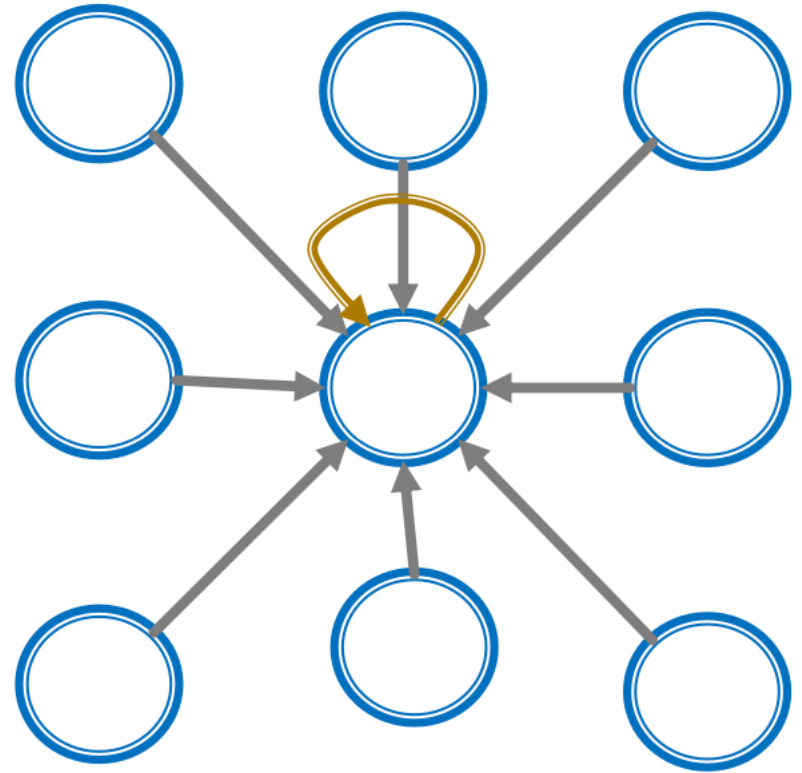### Node Classification



## Graph-level

### Graph Classification

# Relation between GNN and CNN



Image

Graph

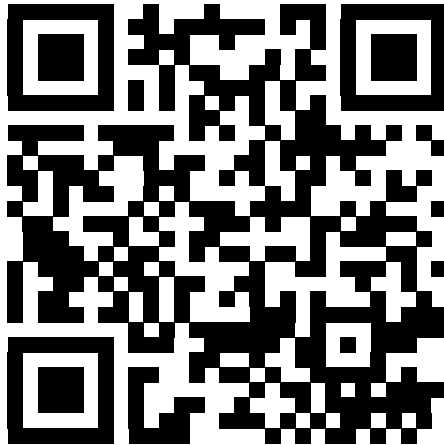CNN can be viewed as a special GNN on grid graph[31]

# GNN vs. Transformer
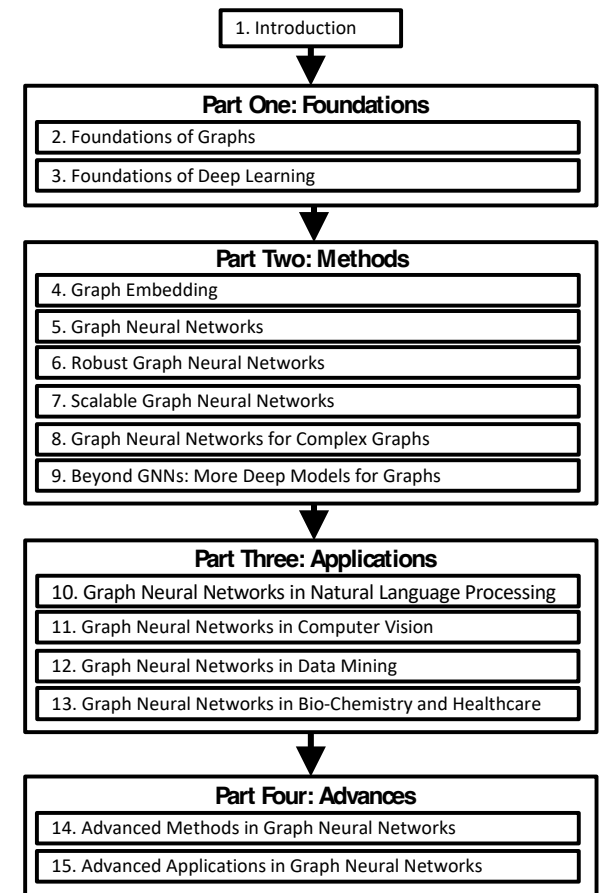
- Transformer is special GNN on a full-connected graph

# Book: Deep Learning on Graphs

https://cse.msu.edu/~mayao4/dlg_book/

1. Introduction

**Part One: Foundations**
2. Foundations of Graphs
3. Foundations of Deep Learning

**Part Two: Methods**
4. Graph Embedding
5. Graph Neural Networks
6. Robust Graph Neural Networks
7. Scalable Graph Neural Networks
8. Graph Neural Networks for Complex Graphs
9. Beyond GNNs: More Deep Models for Graphs

**Part Three: Applications**
10. Graph Neural Networks in Natural Language Processing
11. Graph Neural Networks in Computer Vision
12. Graph Neural Networks in Data Mining
13. Graph Neural Networks in Bio-Chemistry and Healthcare

**Part Four: Advances**
14. Advanced Methods in Graph Neural Networks
15. Advanced Applications in Graph Neural Networks

# Summary

- Graph neural network
  - message passed along graph edges
  - aggregate message/embedding by FFN
  - many variants

# **Next Up**

- Variational Auto-Encoder