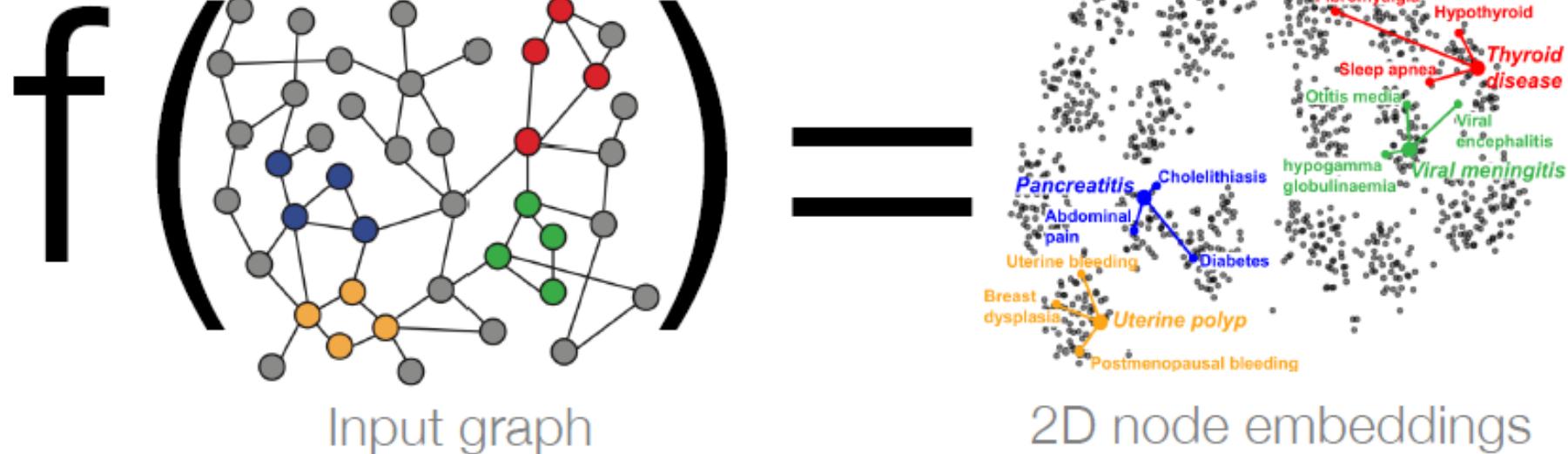


# Graph Neural Networks

Slides adapted from Jure Leskovec's CS224W lecture on Machine Learning on Graphs

# Recap – Node Embeddings

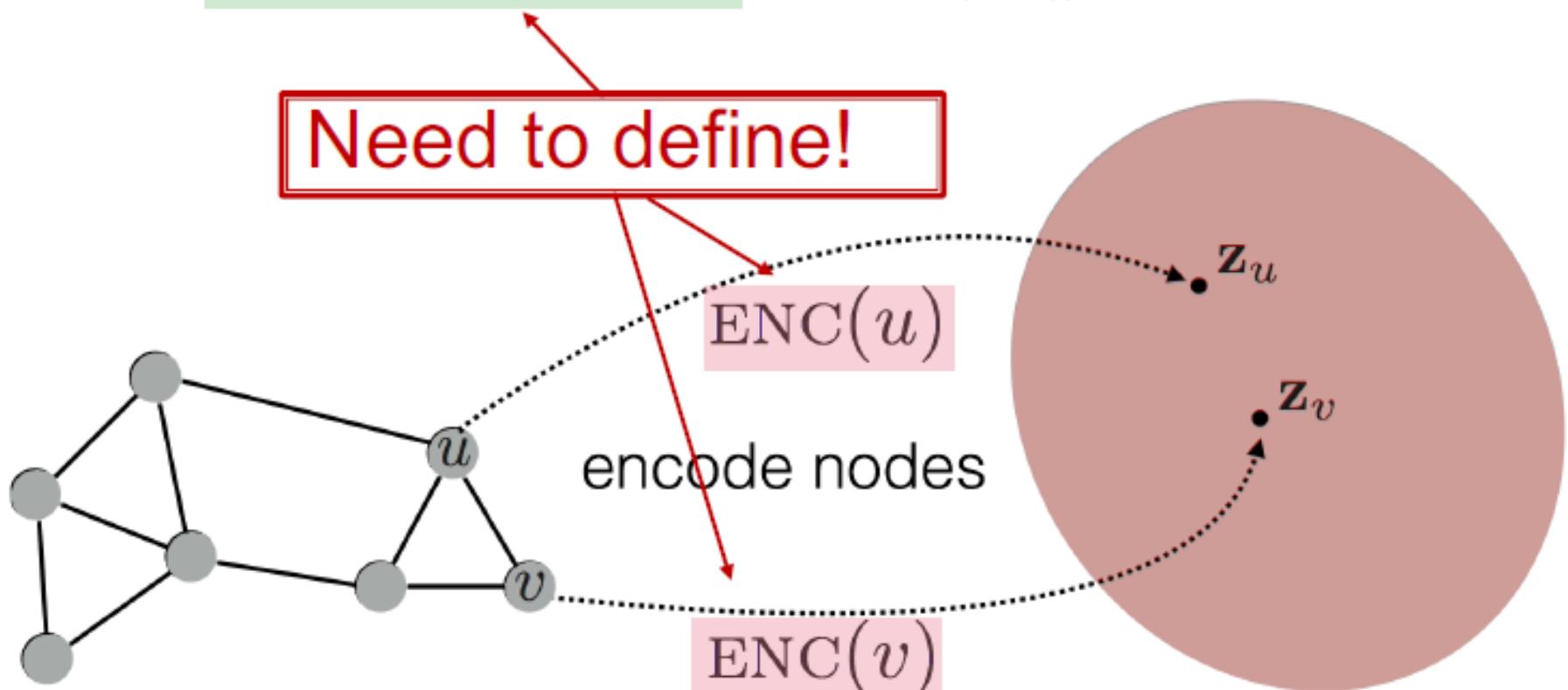
- **Intuition:** Map nodes to  $d$ -dimensional embeddings such that similar nodes in the graph are embedded close together



How to learn mapping function  $f$ ?

# Recap – Node Embeddings

Goal:  $\text{similarity}(u, v) \approx \mathbf{z}_v^T \mathbf{z}_u$



Input network

d-dimensional  
embedding space

# Recap – Node Embeddings

- **Encoder:** maps each node to a low-dimensional vector

$\text{ENC}(v) = \mathbf{z}_v$

*d*-dimensional embedding  
node in the input graph

- **Similarity function:** specifies how the relationships in vector space map to the relationships in the original network

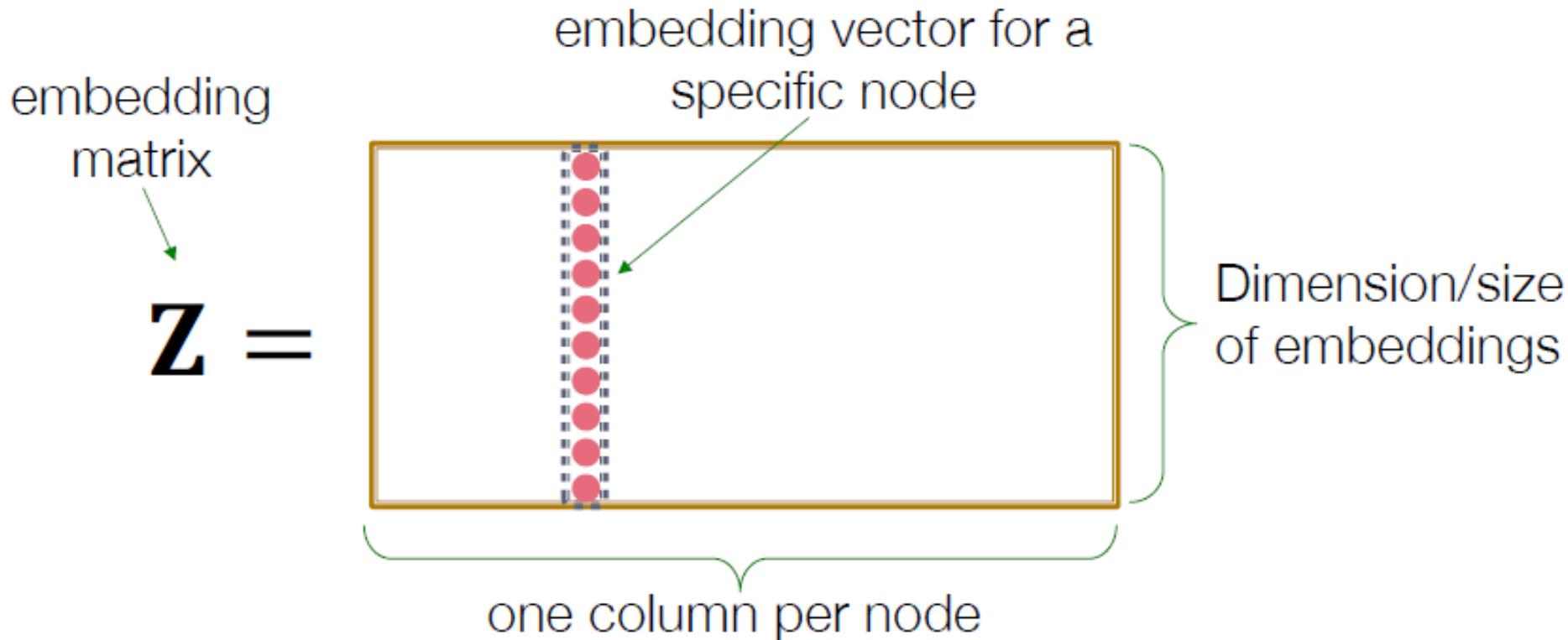
$$\text{similarity}(u, v) \approx \mathbf{z}_v^T \mathbf{z}_u$$

Similarity of  $u$  and  $v$  in  
the original network

**Decoder**  
dot product between node  
embeddings

# Recap – Shallow Encoding

Simplest encoding approach: **encoder is just an embedding-lookup**



# Recap – Shallow Encoders

## ■ Limitations of shallow embedding methods:

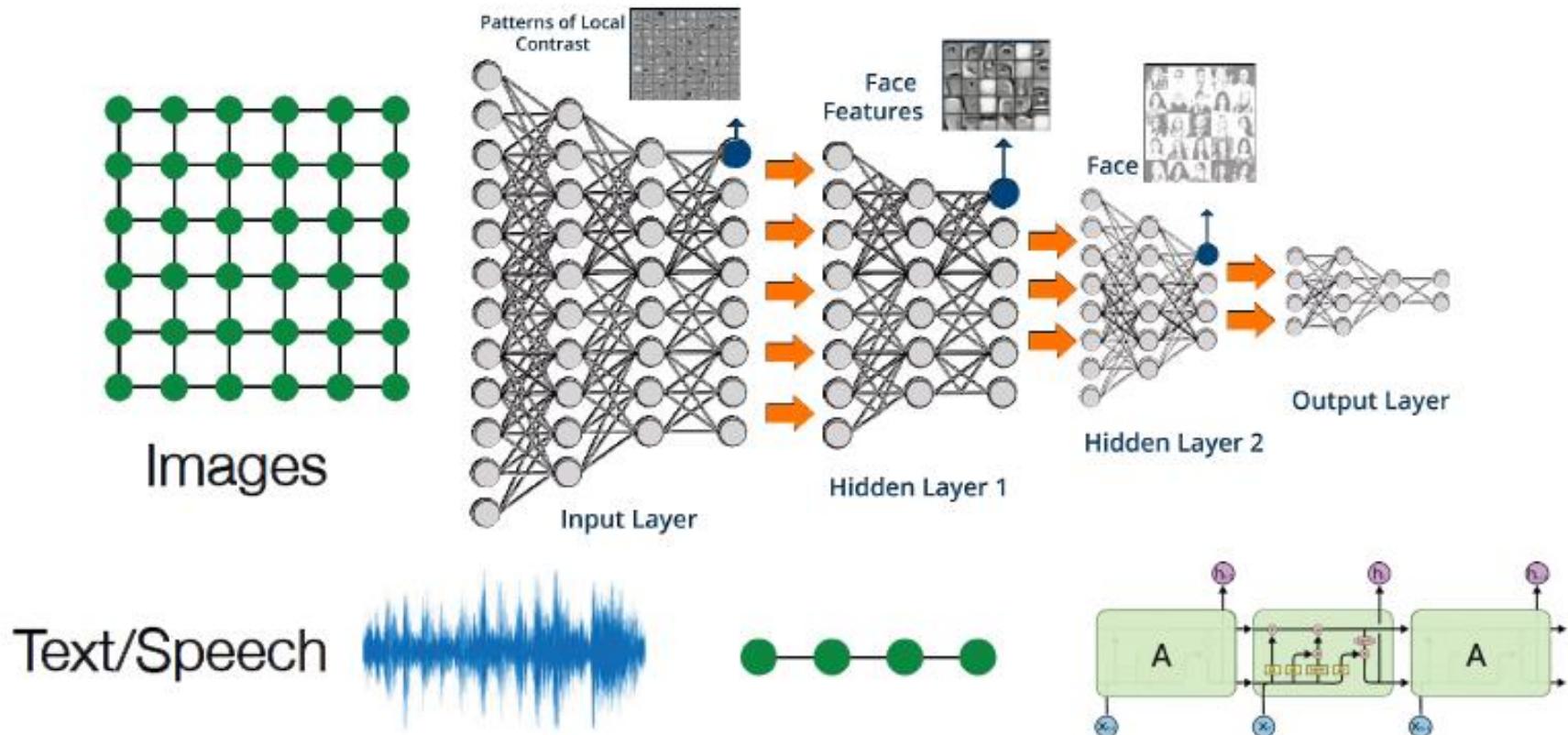
- **$O(|V|)$  parameters are needed:**
  - No sharing of parameters between nodes
  - Every node has its own unique embedding
- **Inherently “transductive”:**
  - Cannot generate embeddings for nodes that are not seen during training
- **Do not incorporate node features:**
  - Many graphs have features that we can and should leverage

# Tasks on Networks

## Tasks we will be able to solve:

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

# Modern ML Toolbox

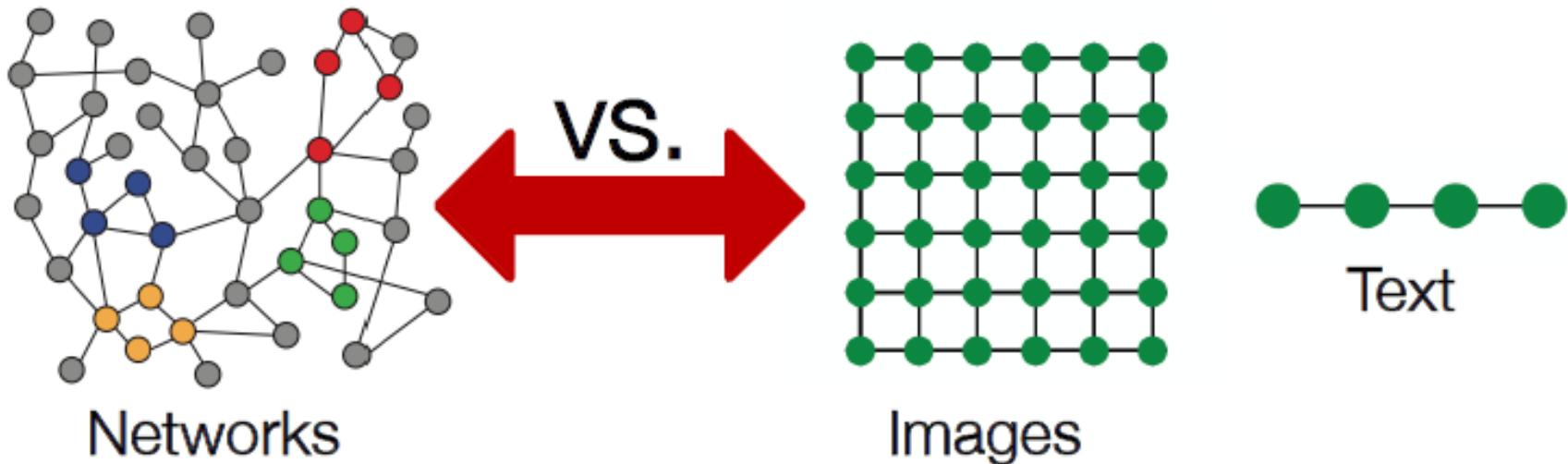


Modern deep learning toolbox is designed  
for simple sequences & grids

# Why deep learning on Graphs is hard?

**But networks are far more complex!**

- Arbitrary size and complex topological structure (i.e., no spatial locality like grids)



- No fixed node ordering or reference point
- Often dynamic and have multimodal features

# Machine Learning as Optimization

- **Supervised learning:** we are given input  $\mathbf{x}$ , and the goal is to predict label  $\mathbf{y}$
- **Input  $\mathbf{x}$  can be:**
  - Vectors of real numbers
  - Sequences (natural language)
  - Matrices (images)
  - Graphs (potentially with node and edge features)
- **We formulate the task as an optimization problem**

# Machine Learning as Optimization

- **Formulate the task as an optimization problem:**

$$\min_{\Theta} \boxed{\mathcal{L}(y, f(x))}$$

Objective function



- $\Theta$ : a set of **parameters** we optimize
  - Could contain one or more scalars, vectors, matrices ...
  - E.g.  $\Theta = \{Z\}$  in the shallow encoder (the embedding lookup)
- $\mathcal{L}$ : **loss function**. Example: L2 loss

$$\mathcal{L}(y, f(x)) = \|y - f(x)\|_2$$

- Other common loss functions:
  - L1 loss, huber loss, max margin (hinge loss), cross entropy ...
  - See <https://pytorch.org/docs/stable/nn.html#loss-functions>

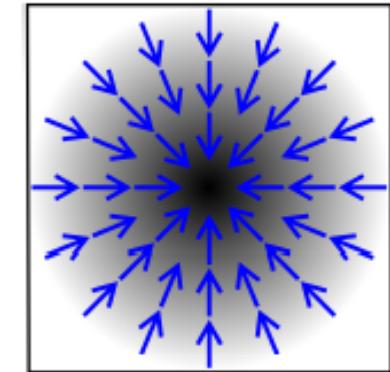
# Loss function example

- One common loss for classification: **cross entropy (CE)**
- Label  $\mathbf{y}$  is a categorical vector (**one-hot encoding**)
  - e.g.  $\mathbf{y} = \begin{array}{|c|c|c|c|c|} \hline 0 & 0 & 1 & 0 & 0 \\ \hline \end{array}$   $\mathbf{y}$  is of class "3"
- $f(\mathbf{x}) = \text{Softmax}(g(\mathbf{x}))$ 
  - Recall from lecture 3:  $f(\mathbf{x})_i = \frac{e^{g(\mathbf{x})_i}}{\sum_{j=1}^C e^{g(\mathbf{x})_j}}$ ,  
where  $C$  is the number of classes.  
 $g(\mathbf{x})_i$  denotes  $i$ -th coordinate of the vector output of func.  $g(\mathbf{x})$
  - e.g.  $f(\mathbf{x}) = \begin{array}{|c|c|c|c|c|} \hline 0.1 & 0.3 & 0.4 & 0.1 & 0.1 \\ \hline \end{array}$
- $\text{CE}(\mathbf{y}, f(\mathbf{x})) = -\sum_{i=1}^C (y_i \log f(\mathbf{x})_i)$ 
  - $y_i, f(\mathbf{x})_i$  are the **actual** and **predicted** value of the  $i$ -th class.
  - **Intuition:** the lower the loss, the closer the prediction is to one-hot
- **Total loss over all training examples:**
  - $\mathcal{L} = \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{T}} \text{CE}(\mathbf{y}, f(\mathbf{x}))$
  - $\mathcal{T}$ : training set containing all pairs of data and labels  $(\mathbf{x}, \mathbf{y})$

# Machine Learning as Optimization

- How to optimize the **objective function**?
- **Gradient vector:** Direction and rate of fastest increase  
**Partial derivative**

$$\nabla_{\Theta} \mathcal{L} = \left( \frac{\partial \mathcal{L}}{\partial \Theta_1}, \frac{\partial \mathcal{L}}{\partial \Theta_2}, \dots \right)$$



<https://en.wikipedia.org/wiki/Gradient>

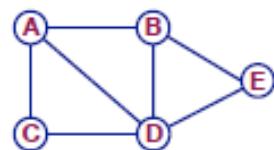
- $\Theta_1, \Theta_2 \dots$  : components of  $\Theta$
- Recall **directional derivative** of a multi-variable function (e.g.  $\mathcal{L}$ ) along a given vector represents the instantaneous rate of change of the function along the vector.
- Gradient is the directional derivative in the **direction of largest increase**

# Setup

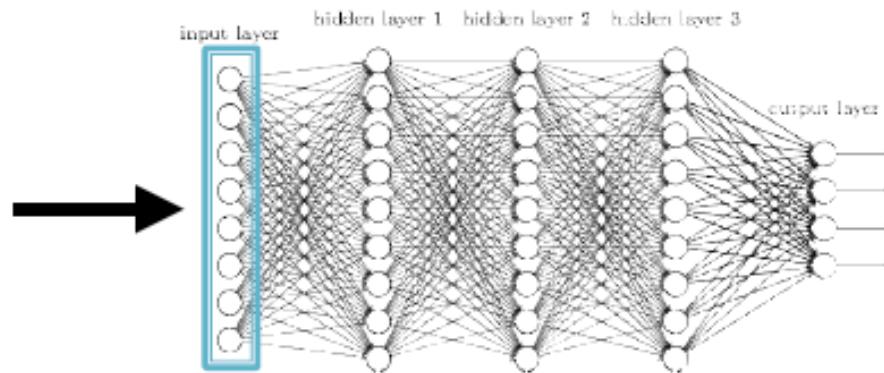
- Assume we have a graph  $G$ :
  - $V$  is the **vertex set**
  - $A$  is the **adjacency matrix** (assume binary)
  - $X \in \mathbb{R}^{m \times |V|}$  is a matrix of **node features**
  - $v$ : a node in  $V$ ;  $N(v)$ : the set of neighbors of  $v$ .
  - **Node features:**
    - Social networks: User profile, User image
    - Biological networks: Gene expression profiles, gene functional information
    - When there is no node feature in the graph dataset:
      - Indicator vectors (one-hot encoding of a node)
      - Vector of constant 1:  $[1, 1, \dots, 1]$

# Naïve Approach

- Join adjacency matrix and features
- Feed them into a deep neural net:

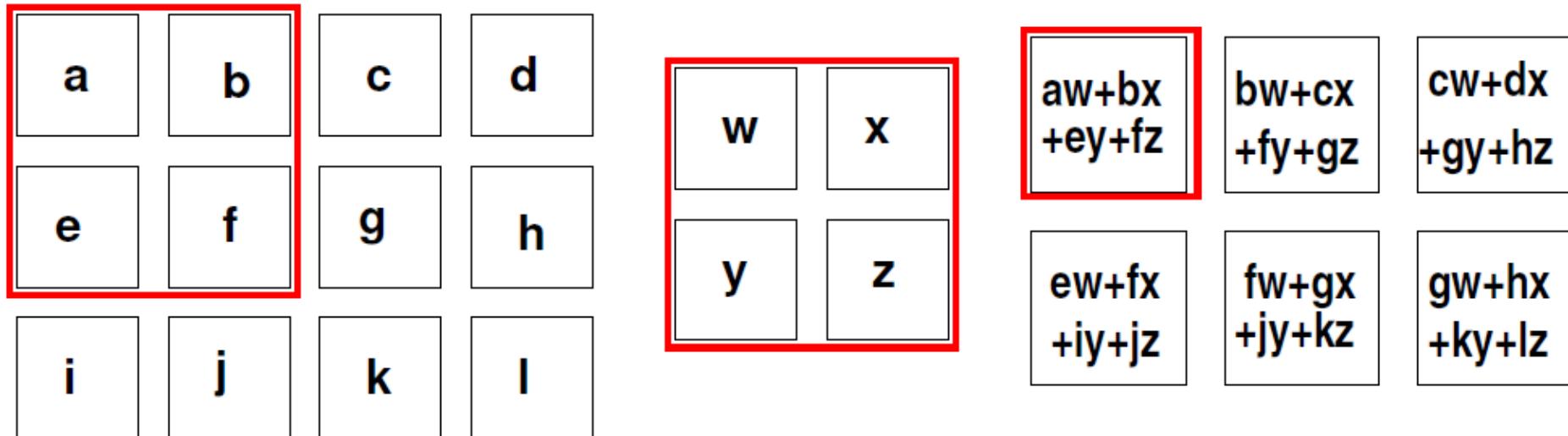


	A	B	C	D	E	Feat
A	0	1	1	1	0	1 0
B	1	0	0	1	1	0 0
C	1	0	0	1	0	0 1
D	1	1	1	0	1	1 1
E	0	1	0	1	0	1 0

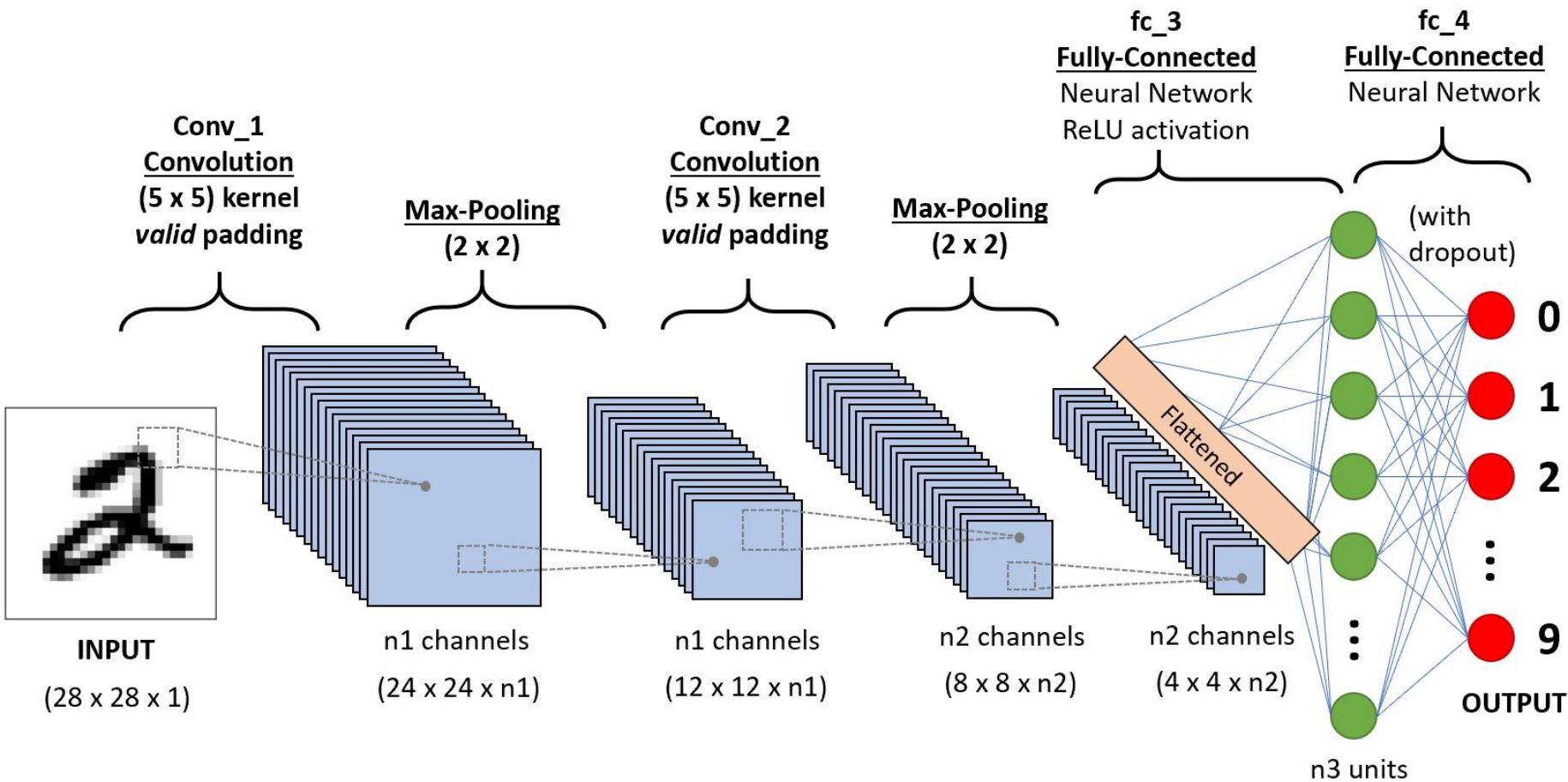


- Issues with this idea:
  - $O(|V|)$  parameters
  - Not applicable to graphs of different sizes
  - Sensitive to node ordering

# Idea – Use Convolution Networks



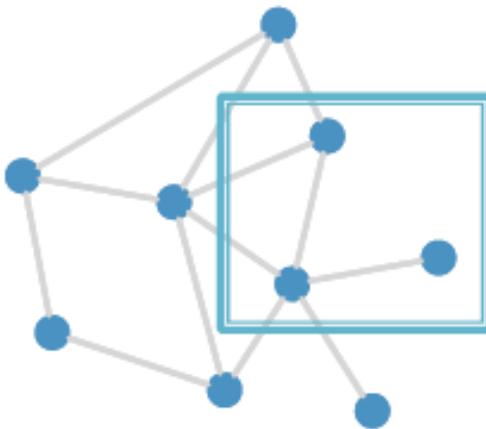
# Idea – Use Convolution Neural Networks



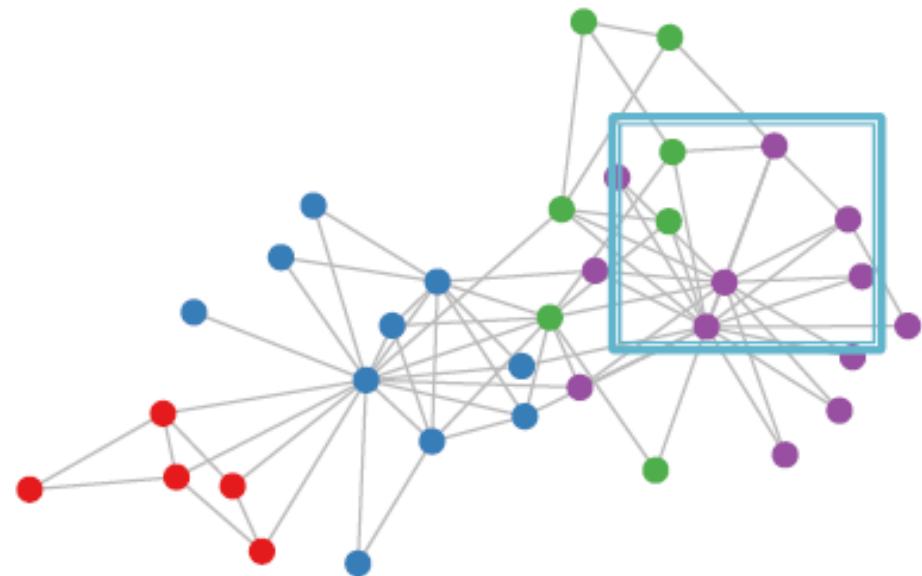
- Goal is to generalize convolutions beyond simple lattices
  - Leverage node feature attributes (texts and Images)

# Real World Graphs

**But our graphs look like this:**



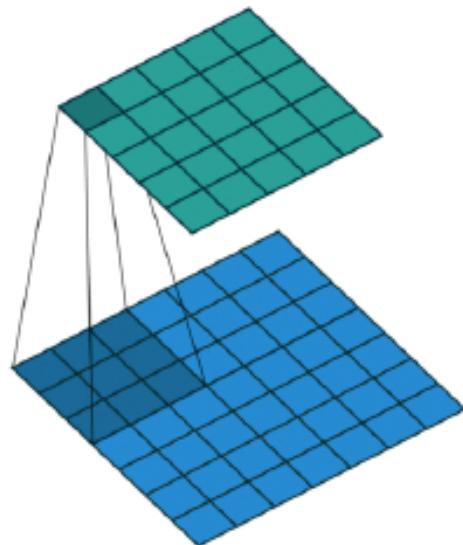
or this:



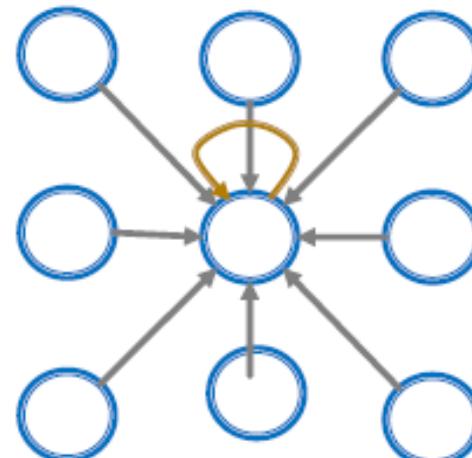
- There is no fixed notion of locality or sliding window on the graph
- Graph is permutation invariant

# Applying CNNs for graphs

Single Convolutional neural network (CNN) layer  
with 3x3 filter:



Image



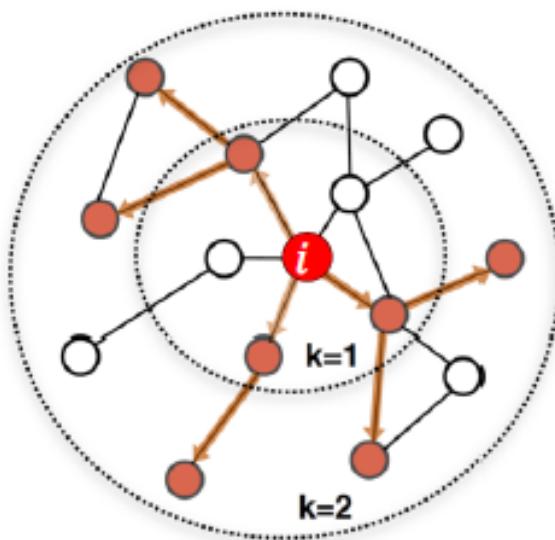
Graph

**Idea:** transform information at the neighbors and combine it:

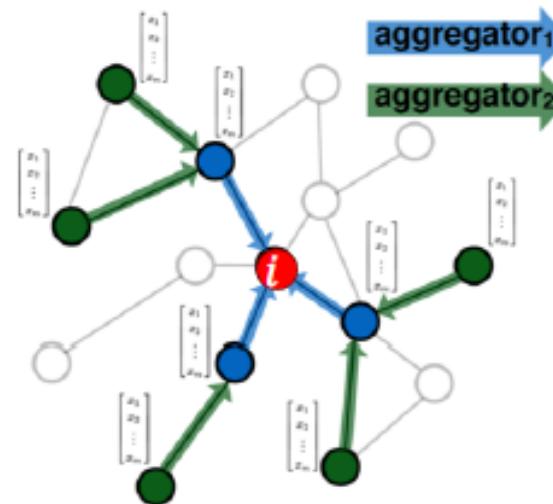
- Transform “messages”  $h_i$  from neighbors:  $W_i h_i$
- Add them up:  $\sum_i W_i h_i$

# Graph Convolution Networks

Idea: Node's neighborhood defines a computation graph



Determine node  
computation graph

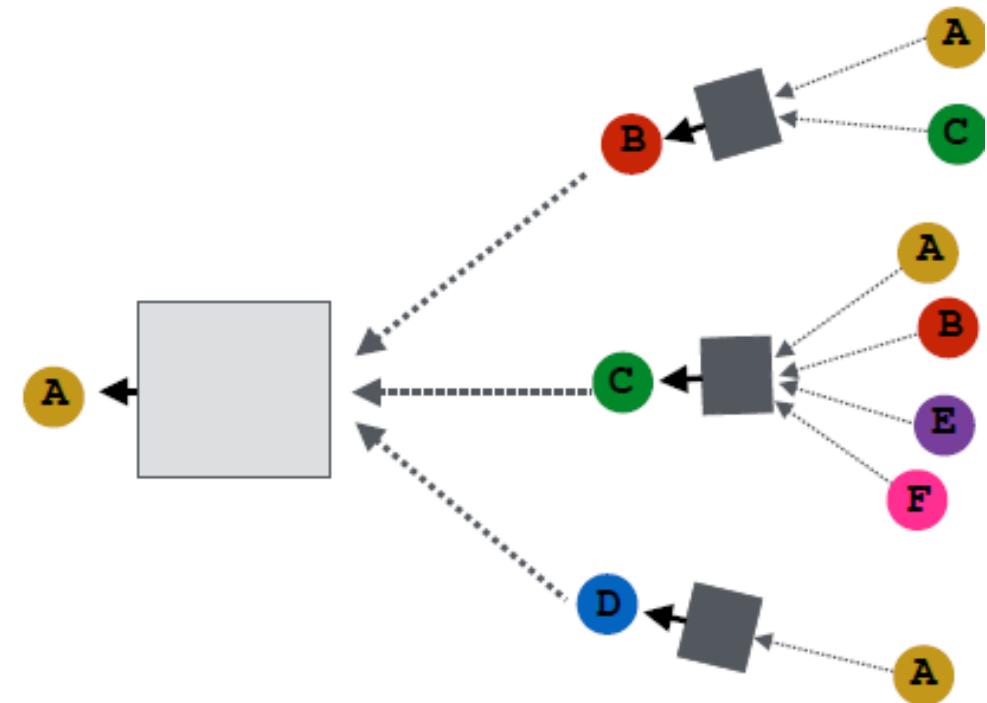
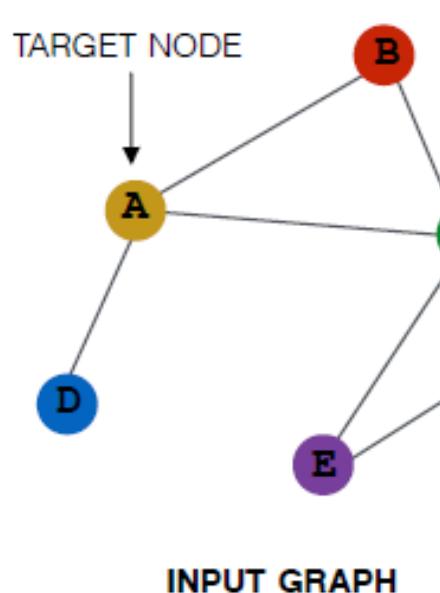


Propagate and  
transform information

Learn how to propagate information across the graph to compute node features

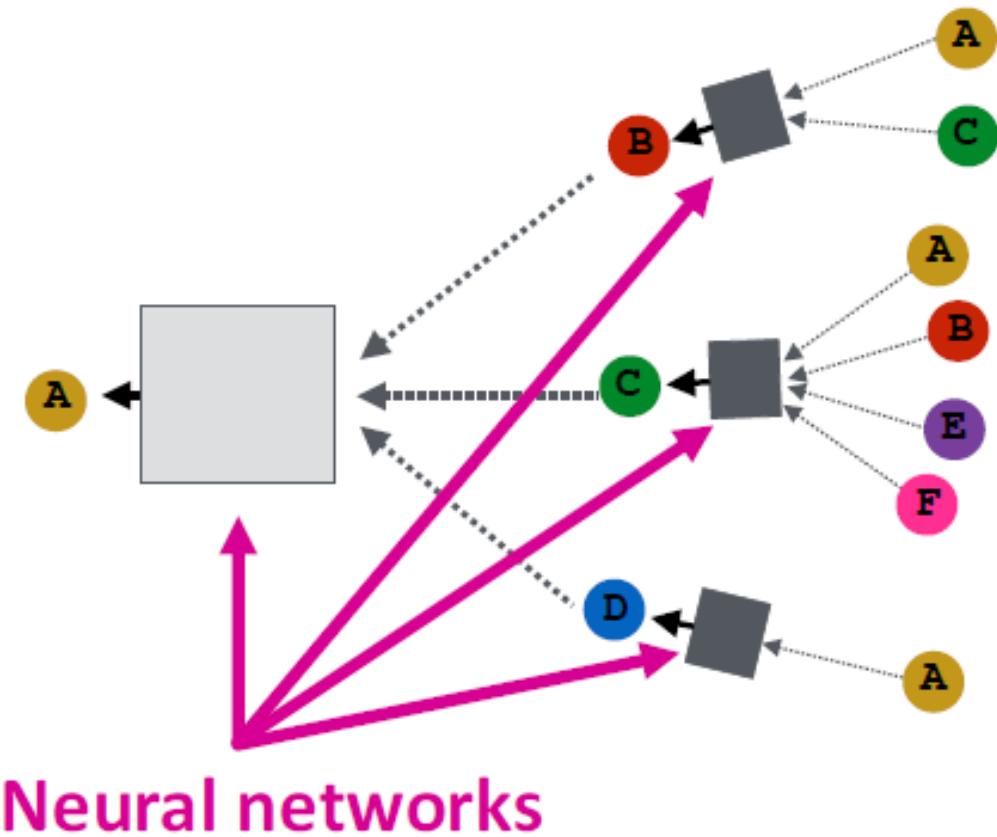
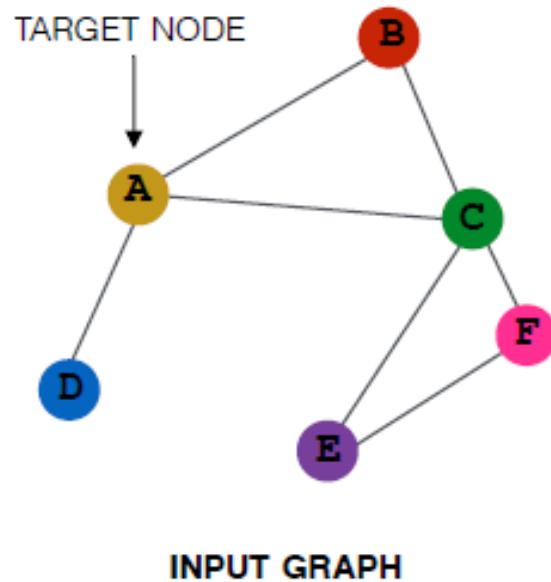
# Idea – Aggregate Neighbors

- **Key idea:** Generate node embeddings based on **local network neighborhoods**



# Idea – Aggregate Neighbors

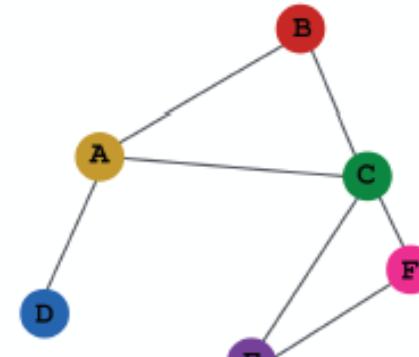
- **Intuition:** Nodes aggregate information from their neighbors using neural networks



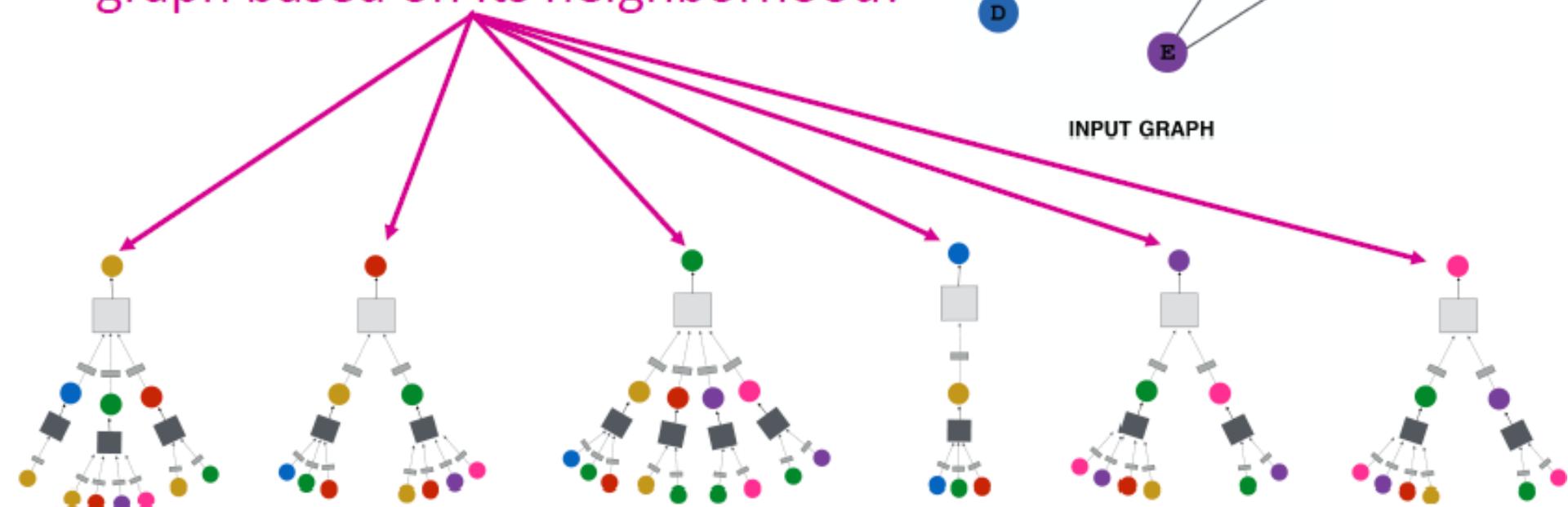
# Key Idea – Aggregate Neighbors

- **Intuition:** Network neighborhood defines a computation graph

Every node defines a computation graph based on its neighborhood!

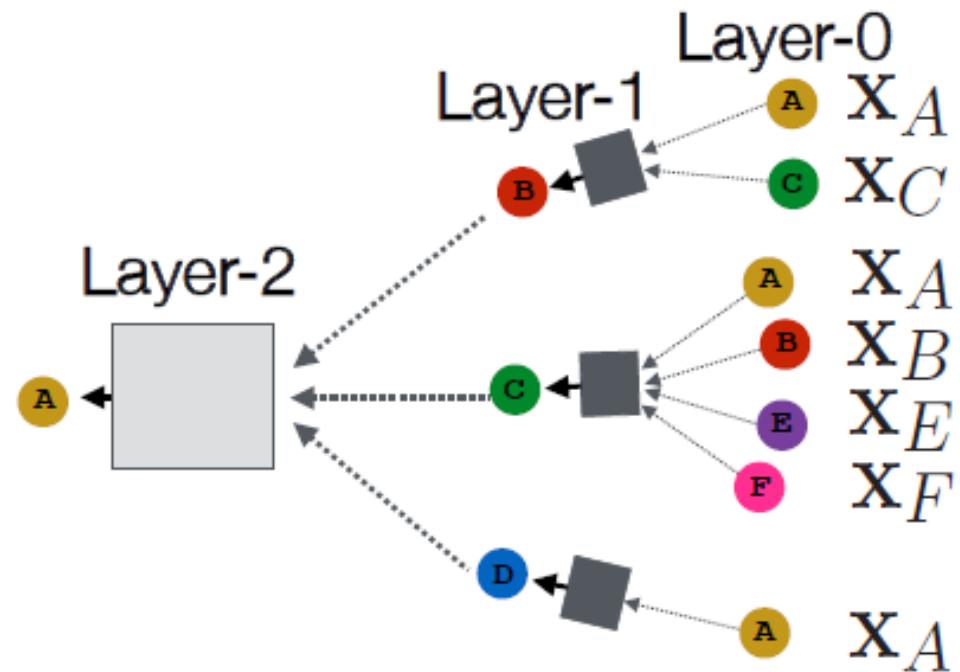
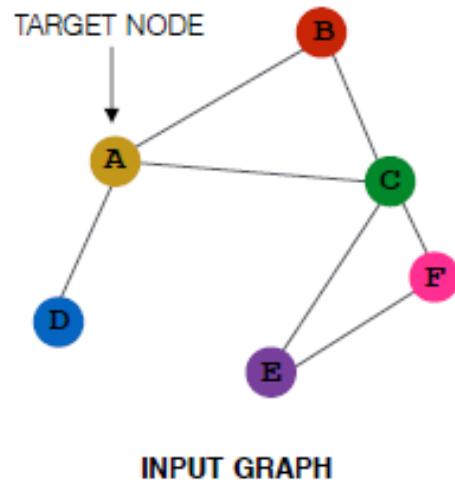


INPUT GRAPH



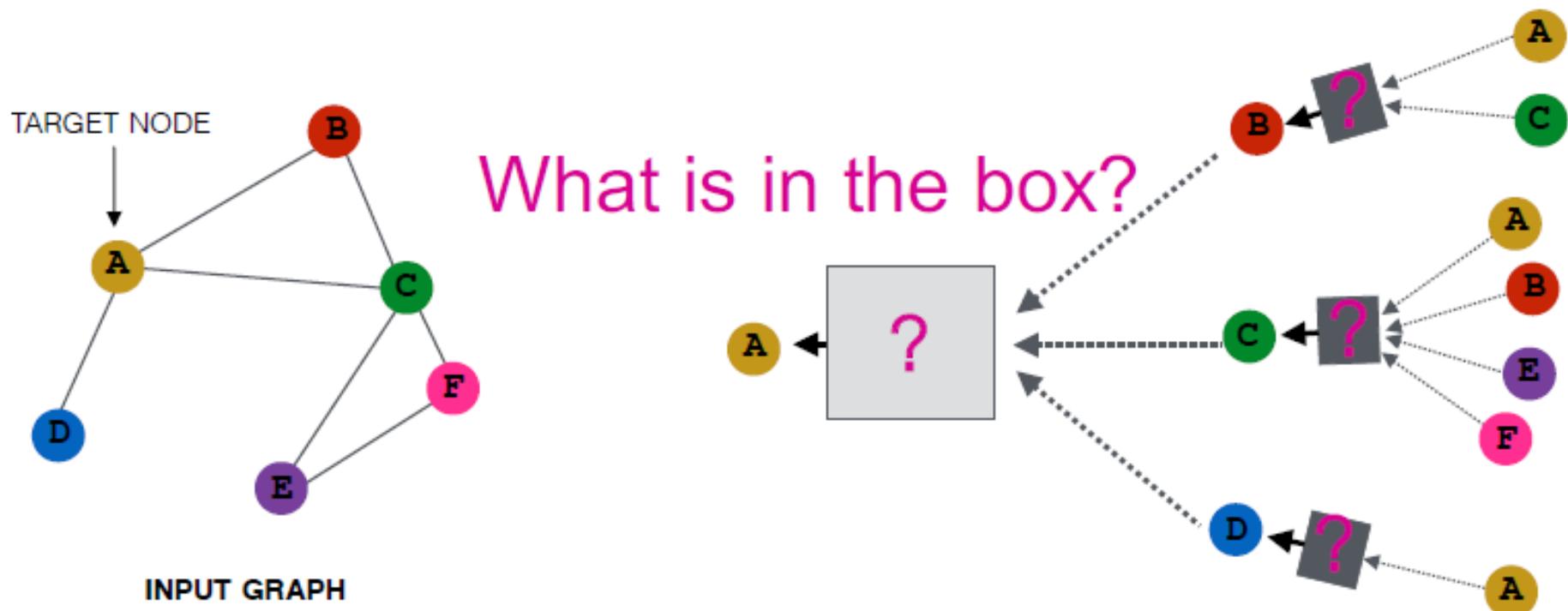
# Deep Model – Many Layers

- 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



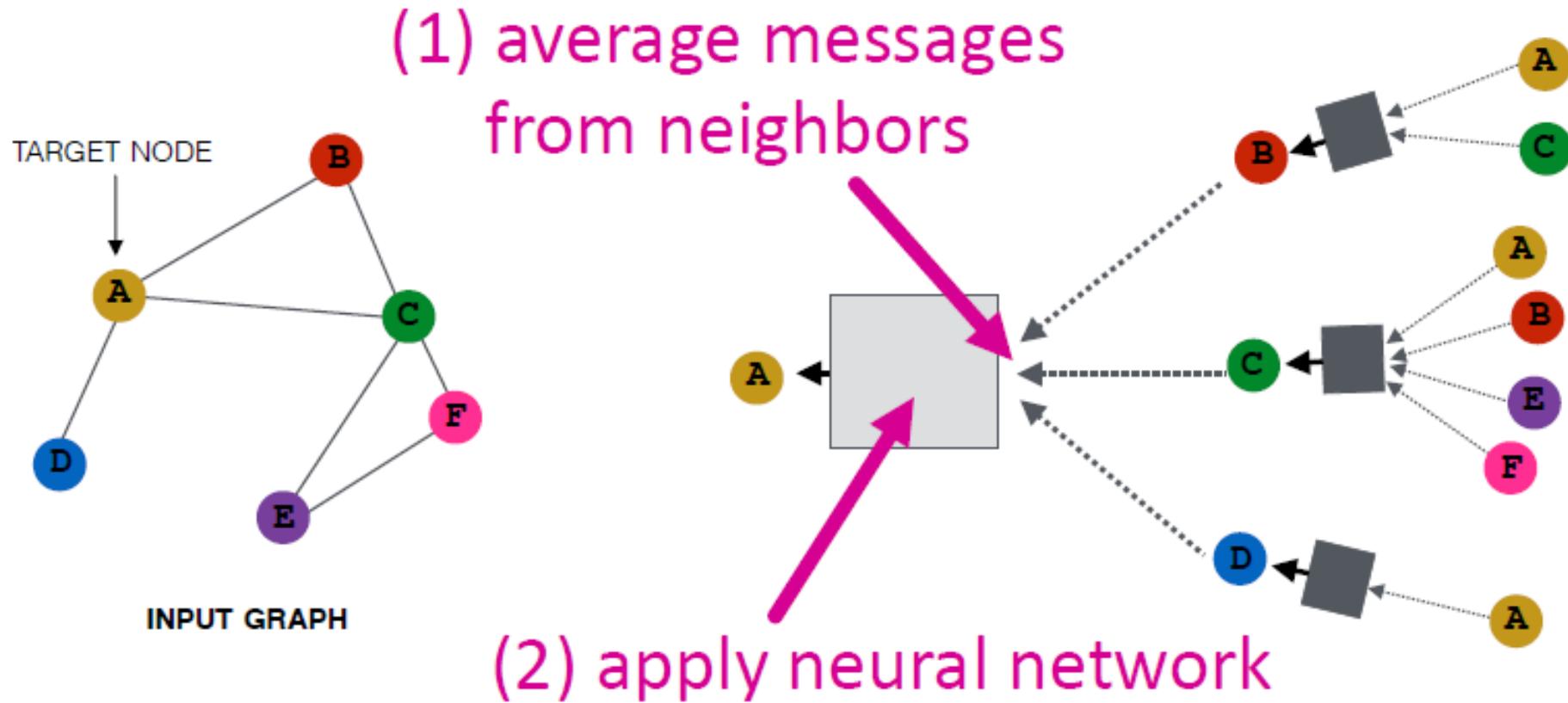
# Neighborhood Aggregation

- **Neighborhood aggregation:** Key distinctions are in how different approaches aggregate information across the layers



# Neighborhood Aggregation

- **Basic approach:** Average information from neighbors and apply a neural network



# The maths behind the deep encoder

- **Basic approach:** Average neighbor messages and apply a neural network

Initial 0-th layer embeddings are equal to node features

$$h_v^0 = x_v$$

embedding of  $v$  at layer  $l$

$$h_v^{(l+1)} = \sigma(W_l \sum_{u \in N(v)} \frac{h_u^{(l)}}{|N(v)|} + B_l h_v^{(l)}), \forall l \in \{0, \dots, L-1\}$$

$z_v = h_v^{(L)}$

Embedding after  $L$  layers of neighborhood aggregation

Non-linearity (e.g., ReLU)

Average of neighbor's previous layer embeddings

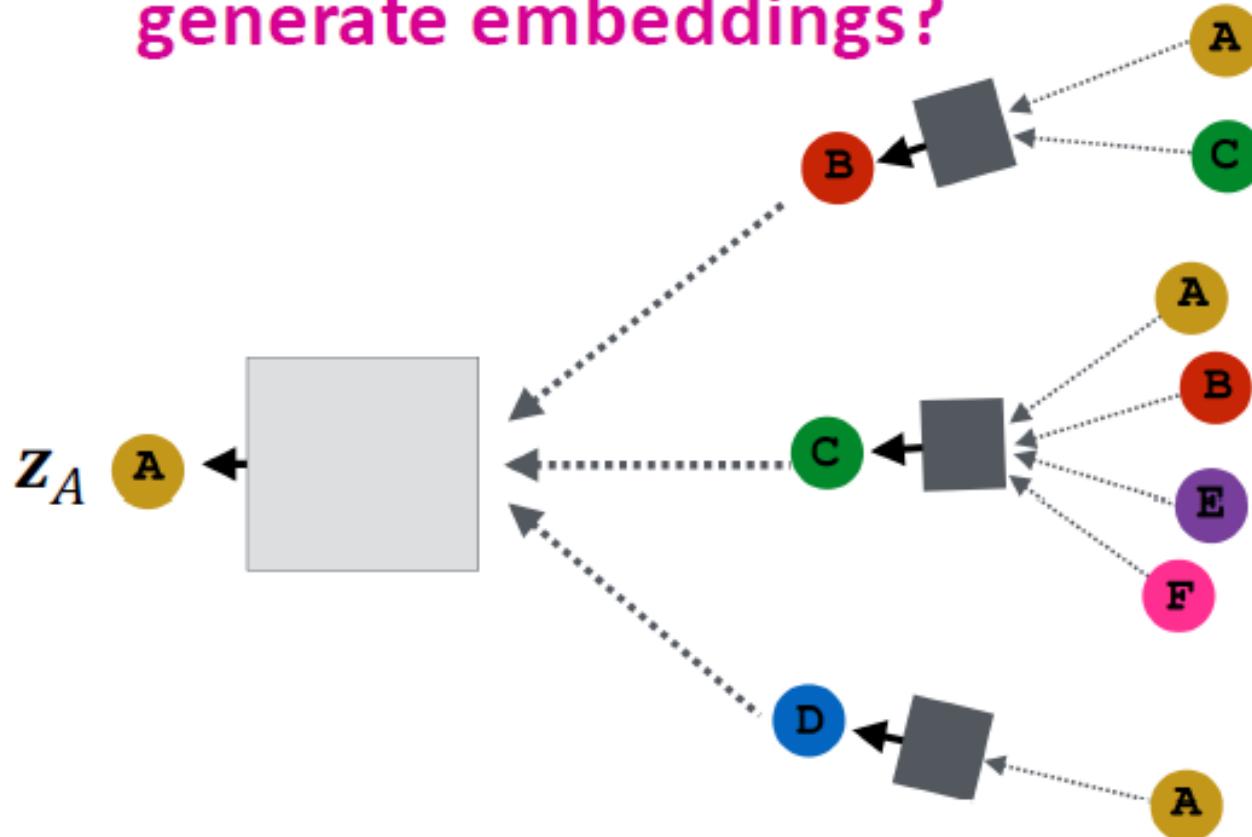
Total number of layers

Annotations:

- Initial 0-th layer embeddings are equal to node features
- embedding of  $v$  at layer  $l$
- embedding after  $L$  layers of neighborhood aggregation
- Non-linearity (e.g., ReLU)
- Average of neighbor's previous layer embeddings
- Total number of layers

# Training the model

How do we train the model to generate embeddings?



Need to define a loss function on the embeddings

# Model Parameters

Trainable weight matrices  
(i.e., what we learn)

$$\begin{aligned} h_v^{(0)} &= x_v \\ h_v^{(l+1)} &= \sigma(W_l \sum_{u \in N(v)} \frac{h_u^{(l)}}{|N(v)|} + B_l h_v^{(l)}), \forall l \in \{0, \dots, L-1\} \\ z_v &= h_v^{(L)} \end{aligned}$$

Final node embedding

We can feed these **embeddings into any loss function** and run SGD to **train the weight parameters**

$h_v^l$ : the hidden representation of node  $v$  at layer  $l$

- $W_k$ : weight matrix for neighborhood aggregation
- $B_k$ : weight matrix for transforming hidden vector of self

# Matrix Formulation

- Many aggregations can be performed efficiently by (sparse) matrix operations

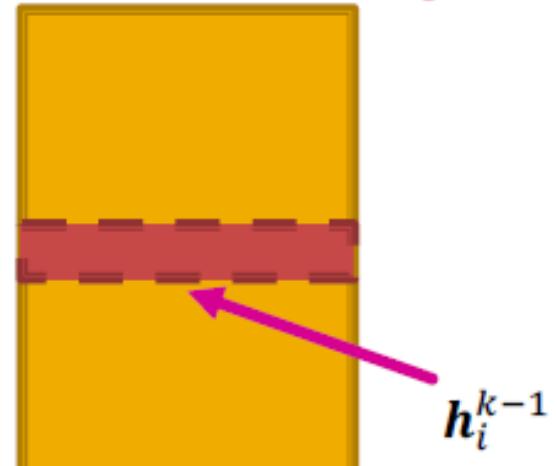
- Let  $H^{(l)} = [h_1^{(l)} \dots h_{|V|}^{(l)}]^T$
- Then:  $\sum_{u \in N(v)} h_u^{(l)} = A_{v,:} H^{(l)}$
- Let  $D$  be diagonal matrix where  $D_{v,v} = \text{Deg}(v) = |N(v)|$ 
  - The inverse of  $D$ :  $D^{-1}$  is also diagonal:  
$$D_{v,v}^{-1} = 1/|N(v)|$$
- Therefore,

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



$$H^{(l+1)} = D^{-1} A H^{(l)}$$

Matrix of hidden embeddings  $H^{k-1}$

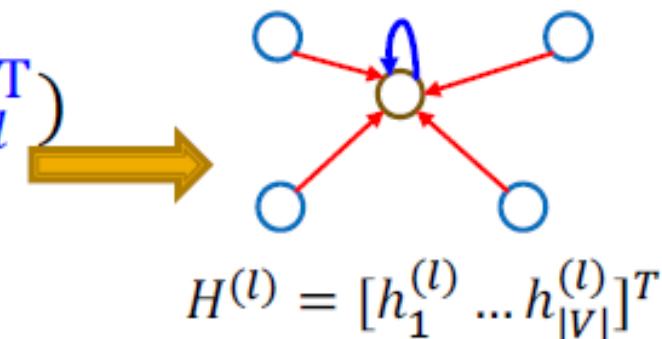


# Matrix Formulation

- Re-writing update function in matrix form:

$$H^{(l+1)} = \sigma(\tilde{A}H^{(l)}W_l^T + H^{(l)}B_l^T)$$

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



- Red: neighborhood aggregation
- Blue: self transformation
- In practice, this implies that efficient sparse matrix multiplication can be used ( $\tilde{A}$  is sparse)
- **Note:** not all GNNs can be expressed in matrix form, when aggregation function is complex

# Training the GNN

- Node embedding  $\mathbf{z}_v$  is a function of input graph
- **Supervised setting**: we want to minimize the loss  $\mathcal{L}$  (see also slide 15):

$$\min_{\Theta} \mathcal{L}(\mathbf{y}, f(\mathbf{z}_v))$$

- $\mathbf{y}$ : node label
- $\mathcal{L}$  could be L2 if  $\mathbf{y}$  is real number, or cross entropy if  $\mathbf{y}$  is categorical
- **Unsupervised setting**:
  - No node label available
  - Use the graph structure as the supervision!

# Unsupervised Training

- “Similar” nodes have similar embeddings

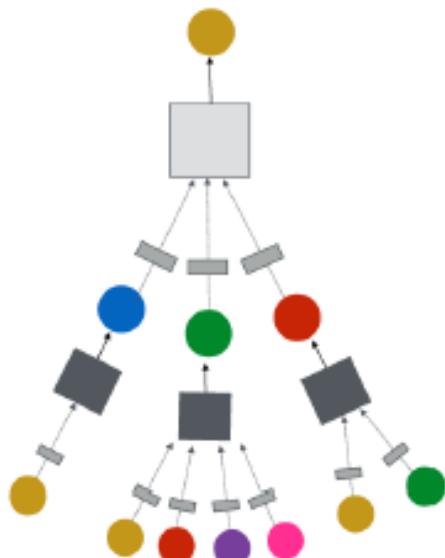
$$\mathcal{L} = \sum_{z_u, z_v} \text{CE}(y_{u,v}, \text{DEC}(z_u, z_v))$$

- Where  $y_{u,v} = 1$  when node  $u$  and  $v$  are **similar**
- **CE** is the cross entropy (slide 16)
- **DEC** is the decoder such as inner product (lecture 4)
- **Node similarity** can be anything from lecture 3, e.g., a loss based on:
  - **Random walks** (node2vec, DeepWalk, struc2vec)
  - **Matrix factorization**
  - **Node proximity in the graph**

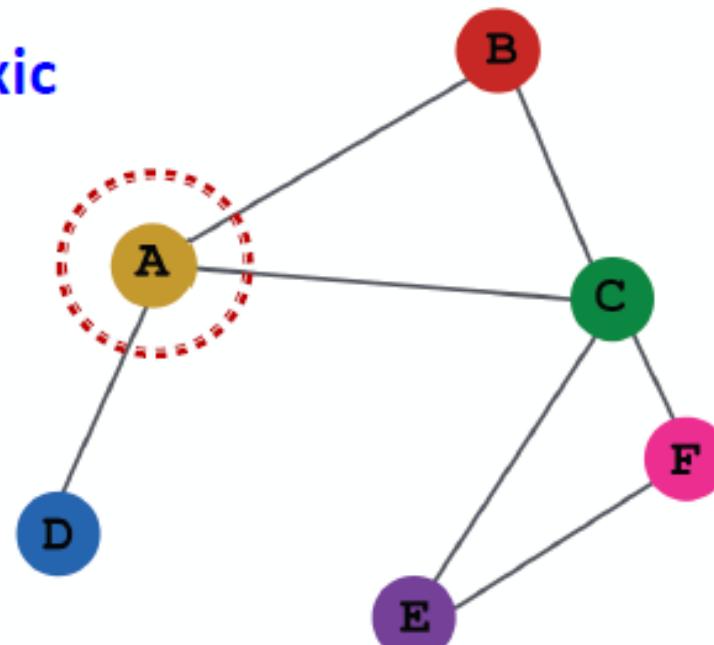
# Supervised Training

**Directly train** the model for a supervised task  
(e.g., node classification)

Safe or toxic  
drug?



Safe or toxic  
drug?

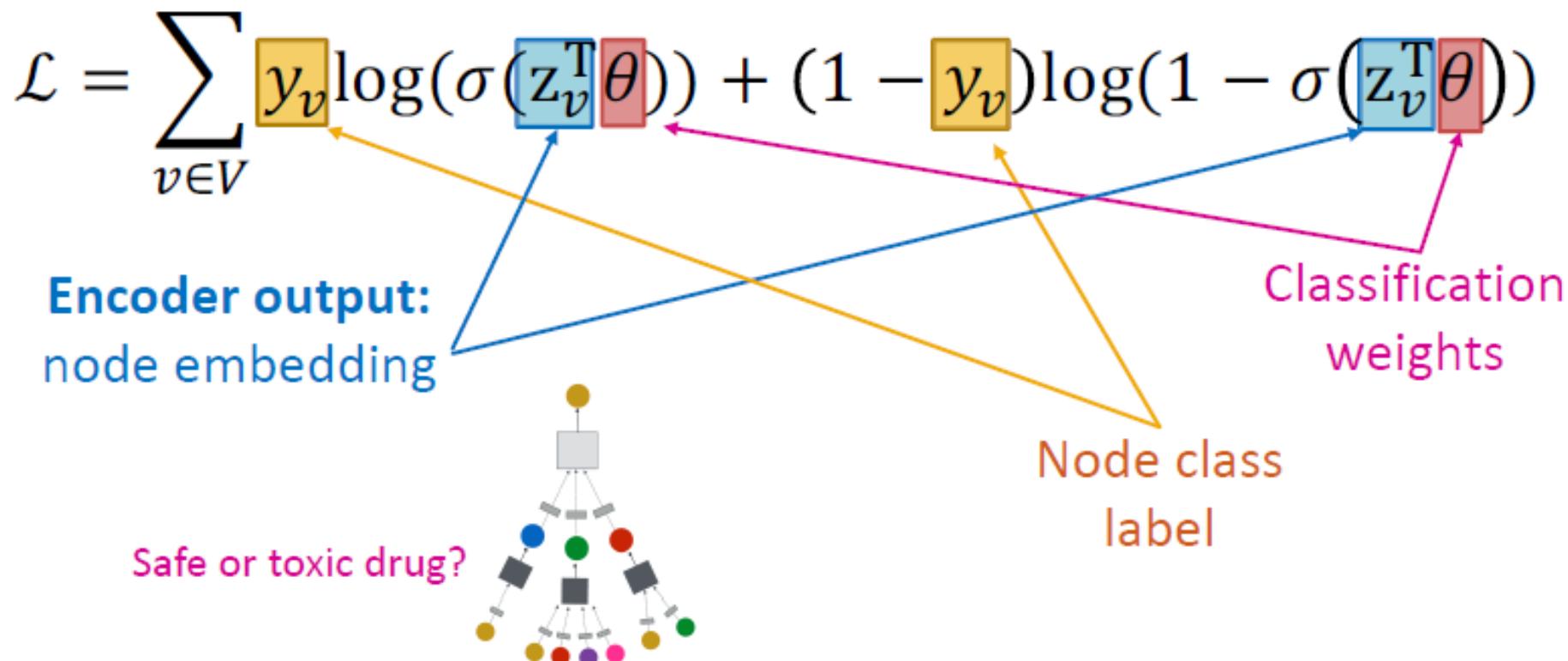


E.g., a drug-drug  
interaction network

# Supervised Training

**Directly train** the model for a supervised task  
(e.g., **node classification**)

- Use cross entropy loss (slide 16)



# Evaluation Metrics: Classification

- Evaluate classification tasks on graphs:
- (1) Multi-class classification
  - We simply report the accuracy

$$\frac{1[\operatorname{argmax}(\hat{\mathbf{y}}^{(i)}) = \mathbf{y}^{(i)}]}{N}$$

- (2) Binary classification
  - Metrics sensitive to classification threshold
    - Accuracy
    - Precision / Recall
    - If the range of prediction is [0,1], we will use 0.5 as threshold
  - Metric Agnostic to classification threshold
    - ROC AUC

# Metrics for Binary Classification

- **Accuracy:**

$$\frac{TP + TN}{TP + TN + FP + FN} = \frac{TP + TN}{|\text{Dataset}|}$$

- **Precision (P):**

$$\frac{TP}{TP + FP}$$

**Confusion matrix**

- **Recall (R):**

$$\frac{TP}{TP + FN}$$

- **F1-Score:**

$$\frac{2P * R}{P + R}$$

	Actually Positive (1)	Actually Negative (0)
Predicted Positive (1)	True Positives (TPs)	False Positives (FPs)
Predicted Negative (0)	False Negatives (FNs)	True Negatives (TNs)

[Sklearn Classification Report](#)

# Evaluation Metrics

- **ROC Curve:** Captures the tradeoff in TPR and FPR **as the classification threshold is varied for a binary classifier.**

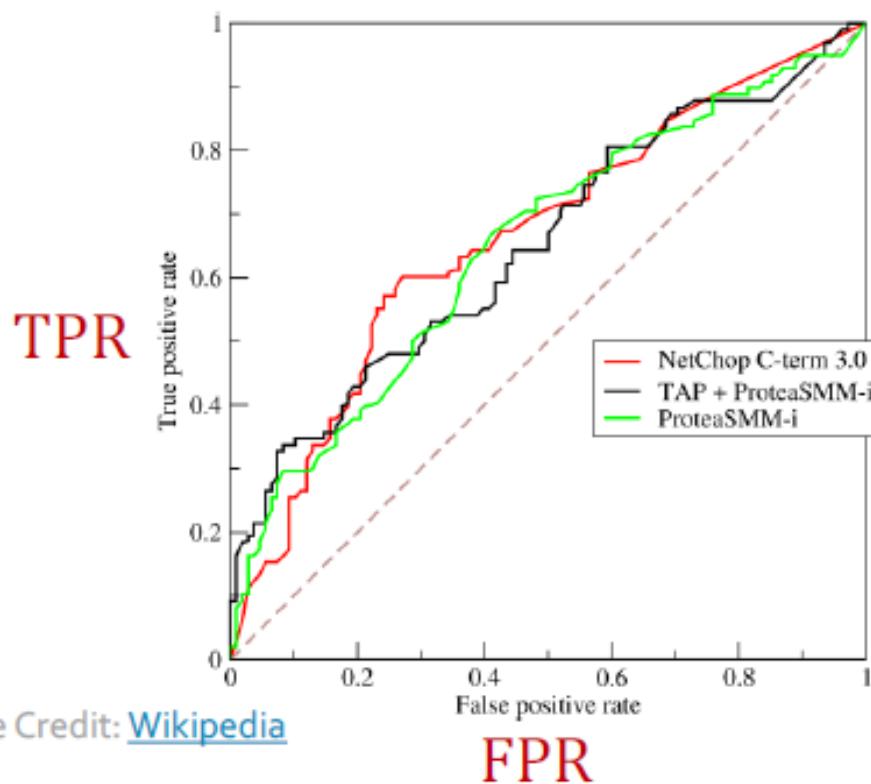


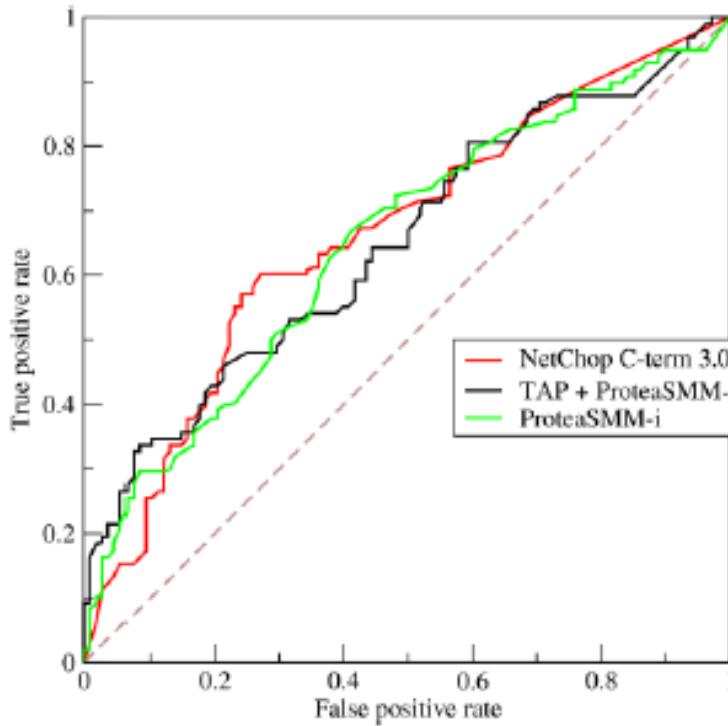
Image Credit: [Wikipedia](#)

$$\text{TPR} = \text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

**Note:** the dashed line represents performance of a random classifier

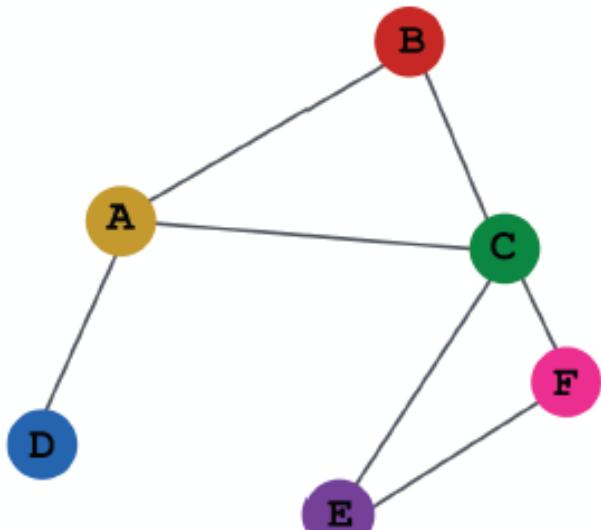
# Evaluation Metrics



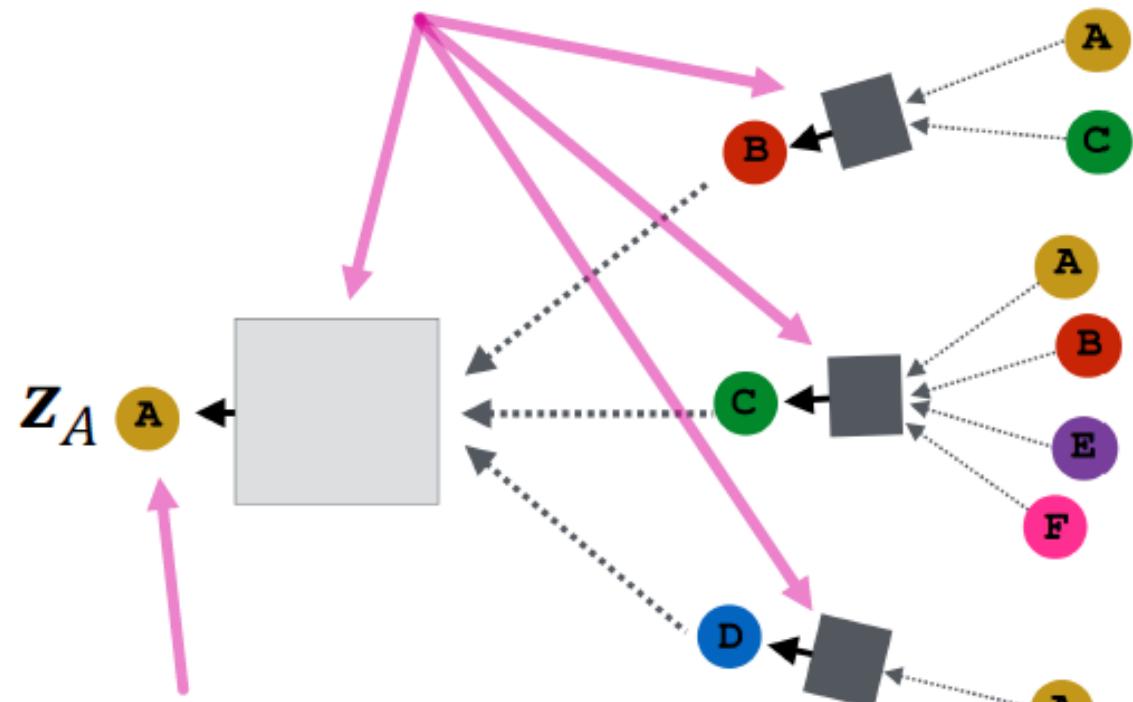
Content Credit: [Wikipedia](#)

- **ROC AUC: Area under the ROC Curve.**
- **Intuition:** The probability that a classifier will rank a randomly chosen positive instance higher than a randomly chosen negative one

# Model Design Overview

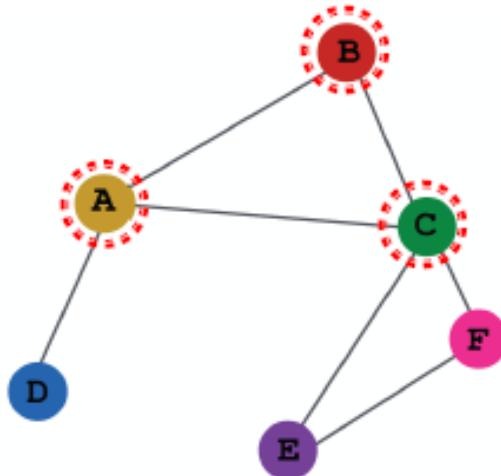


(1) Define a neighborhood aggregation function



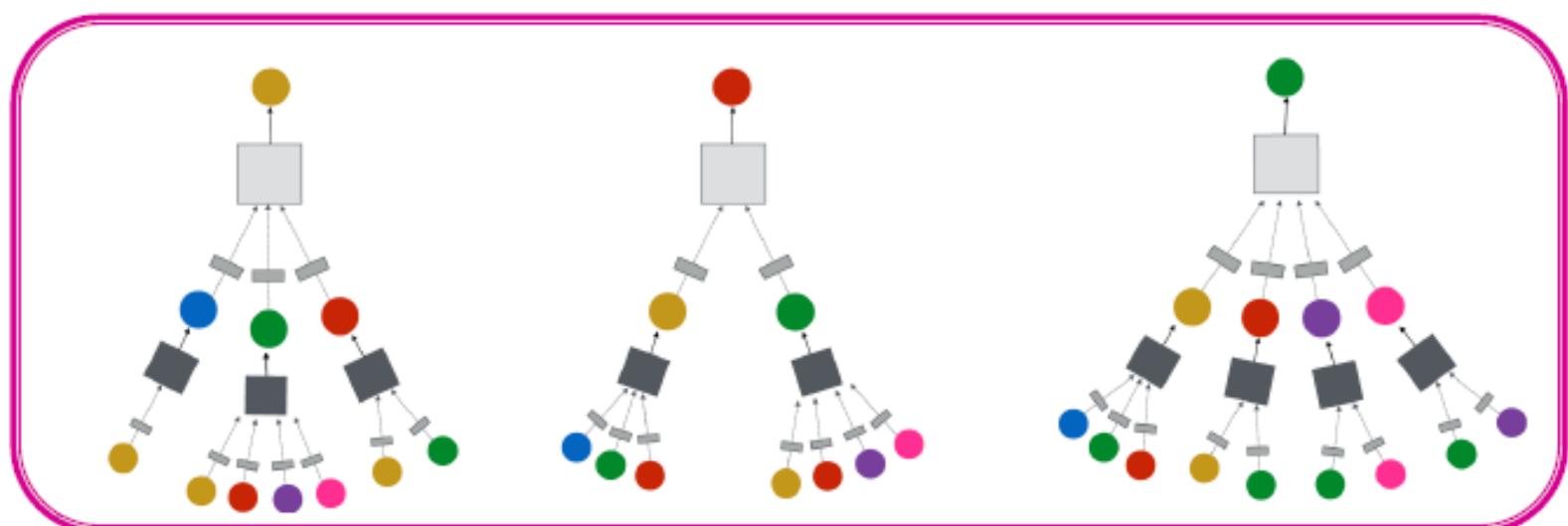
(2) Define a loss function on the embeddings

# Model Design Overview

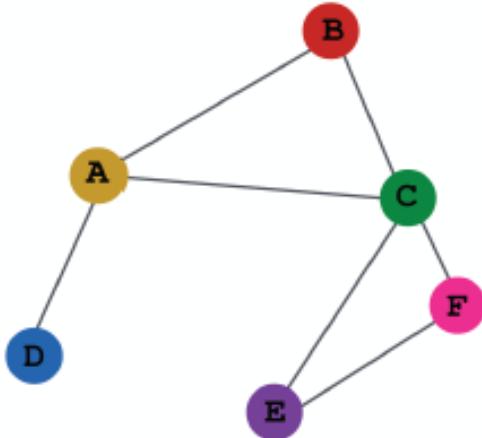


INPUT GRAPH

(3) Train on a set of nodes, i.e.,  
a batch of compute graphs



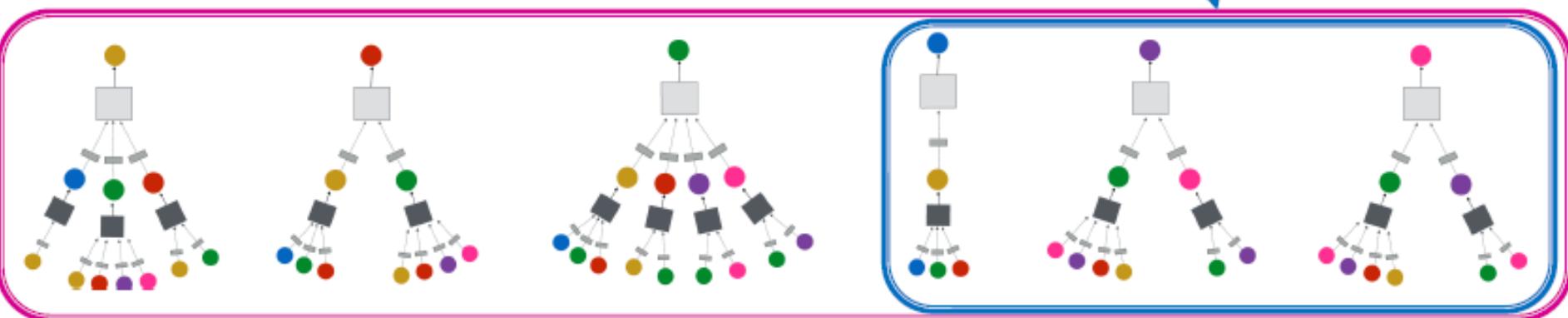
# Model Design Overview



INPUT GRAPH

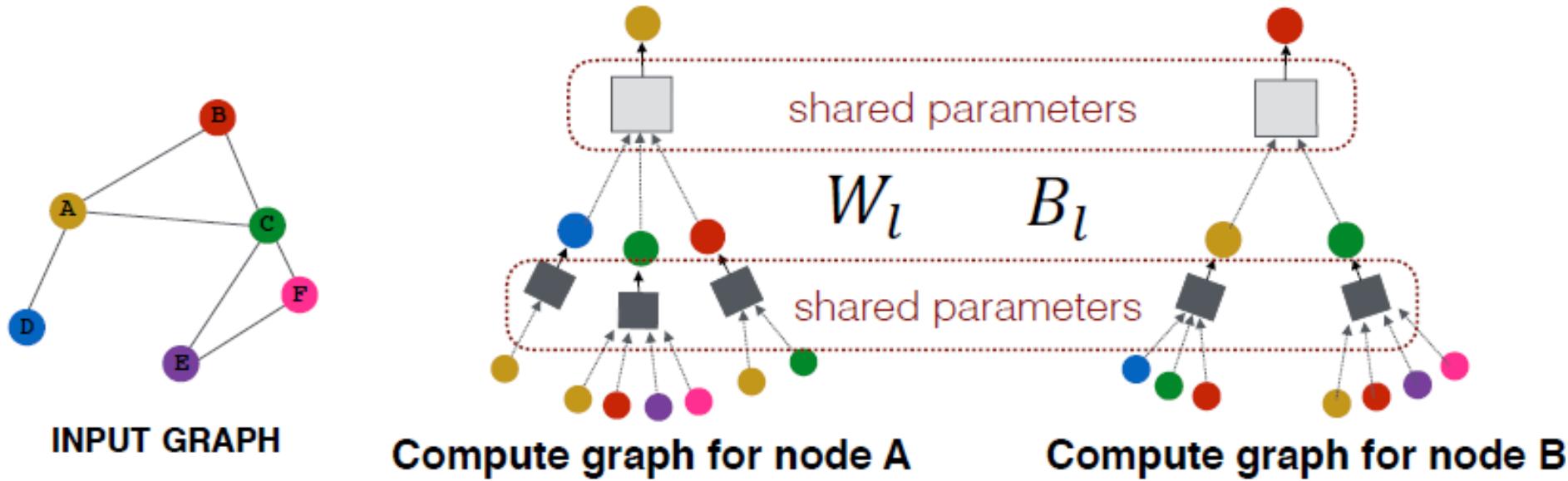
(4) Generate embeddings  
for nodes as needed

Even for nodes we never  
trained on!

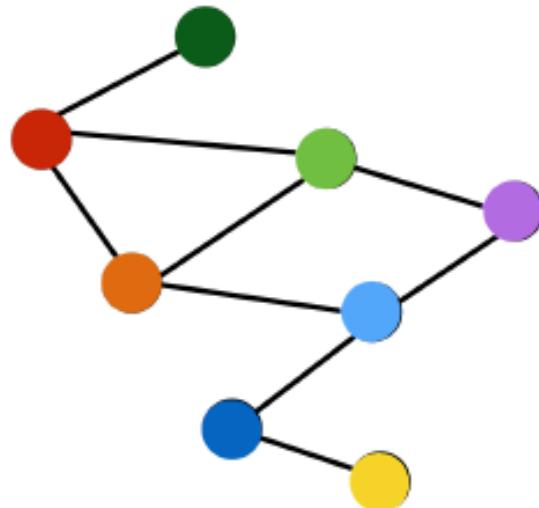


# Inductive Capability

- The same aggregation parameters are shared for all nodes:
  - The number of model parameters is sublinear in  $|V|$  and we can **generalize to unseen nodes!**



# Inductive Capability on New Graphs



Train on one graph

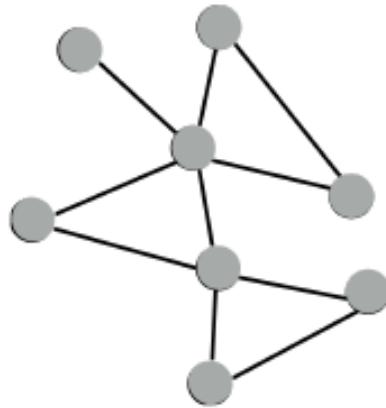


Generalize to new graph

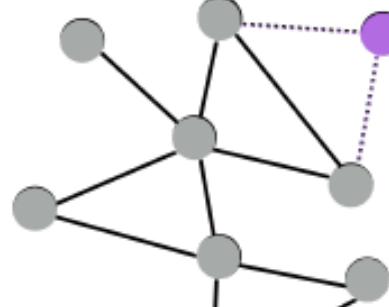
Inductive node embedding → Generalize to entirely unseen graphs

E.g., train on protein interaction graph from model organism A and generate embeddings on newly collected data about organism B

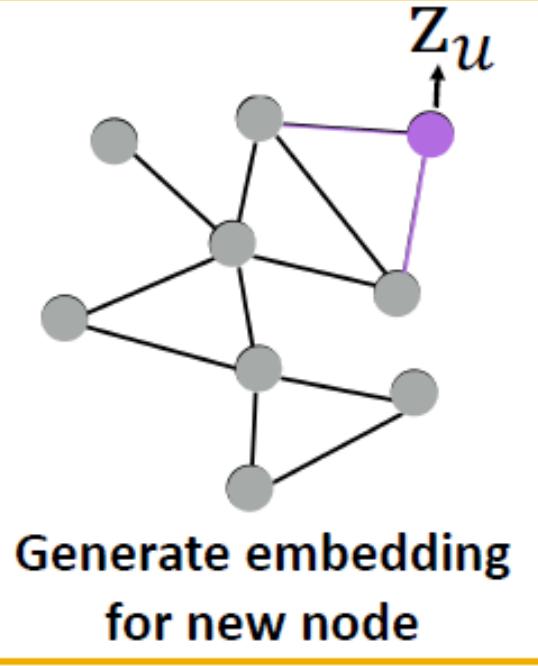
# Inductive Capability on New Nodes



Train with snapshot



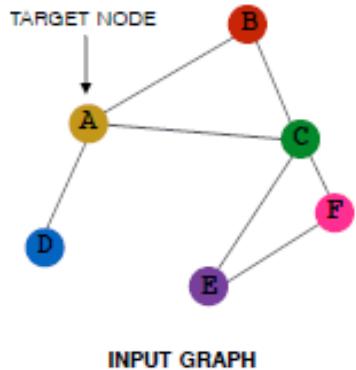
New node arrives



Generate embedding  
for new node

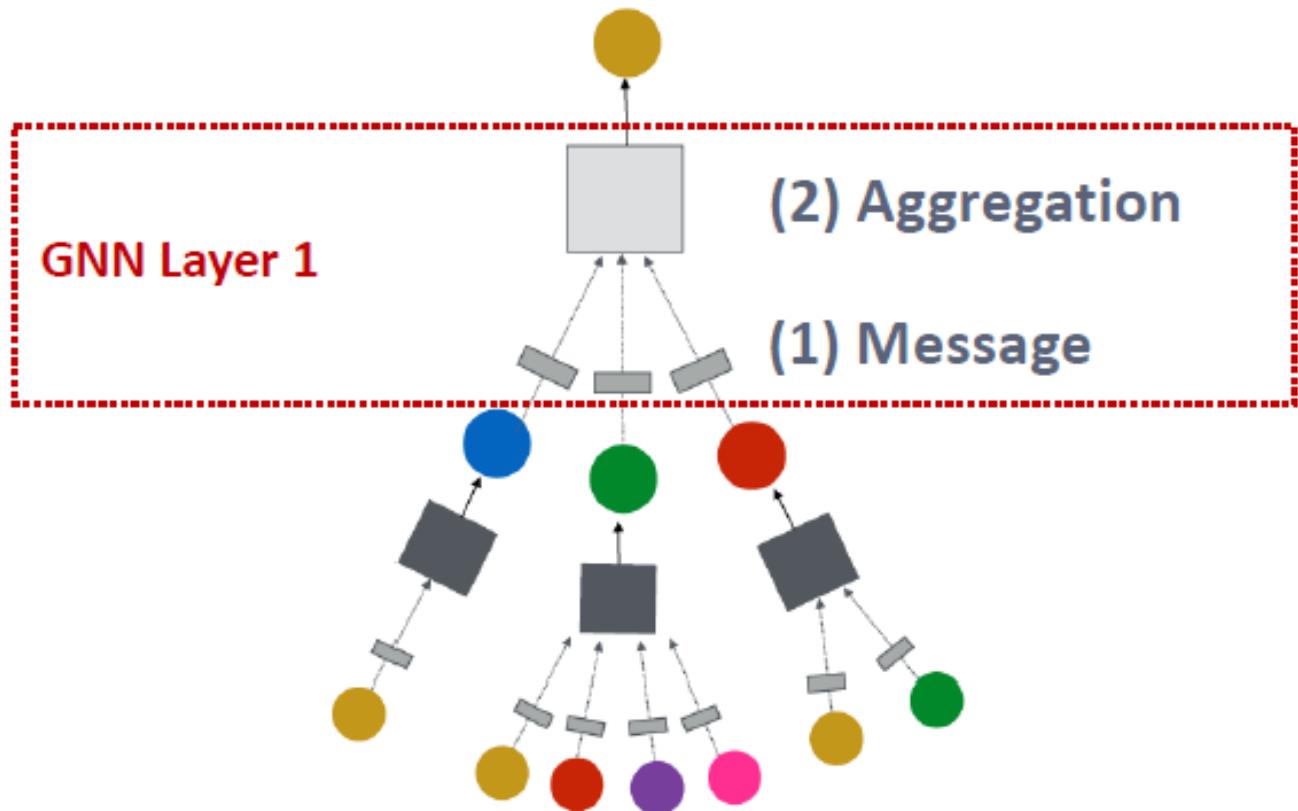
- Many application settings constantly encounter previously unseen nodes:
  - E.g., Reddit, YouTube, Google Scholar
- Need to generate new embeddings “on the fly”

# A Single GNN Layer



**GNN Layer = Message + Aggregation**

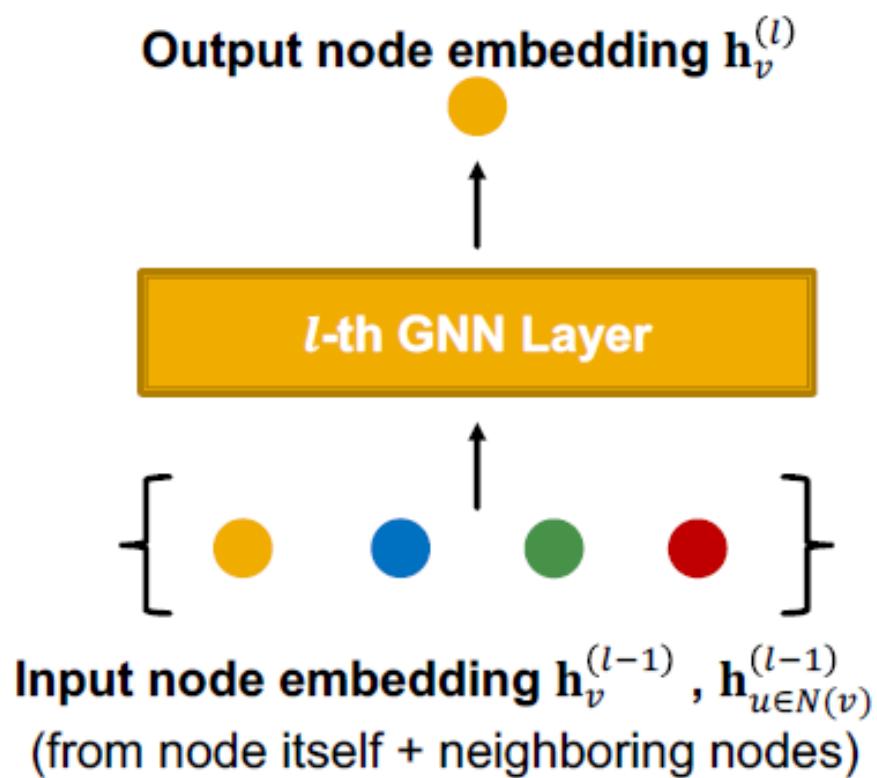
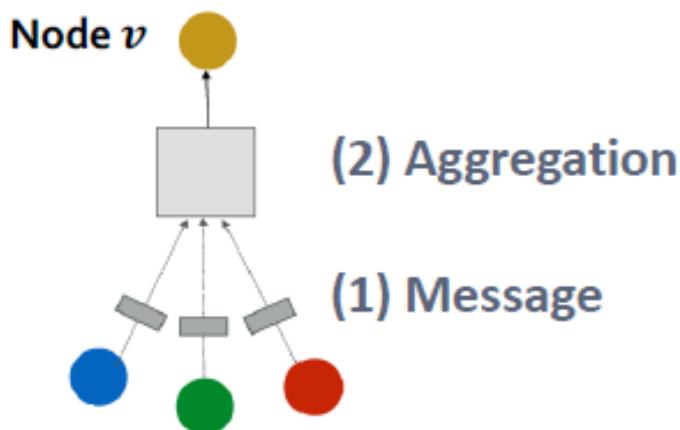
- Different instantiations under this perspective
- GCN, GraphSAGE, GAT, ...



# A Single GNN Layer

## ■ Idea of a GNN Layer:

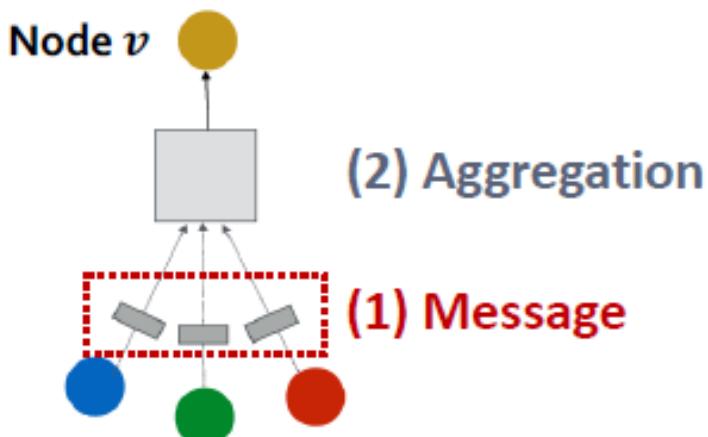
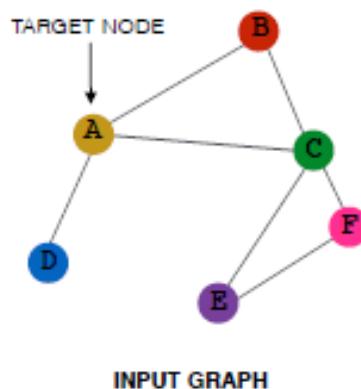
- Compress a set of vectors into a single vector
- Two step process:
  - (1) Message
  - (2) Aggregation



# Message Computation

## ■ (1) Message computation

- **Message function:**  $\mathbf{m}_u^{(l)} = \text{MSG}^{(l)}(\mathbf{h}_u^{(l-1)})$
- **Intuition:** Each node will create a message, which will be sent to other nodes later
- **Example:** A Linear layer  $\mathbf{m}_u^{(l)} = \mathbf{W}^{(l)}\mathbf{h}_u^{(l-1)}$ 
  - Multiply node features with weight matrix  $\mathbf{W}^{(l)}$



# Message Aggregation

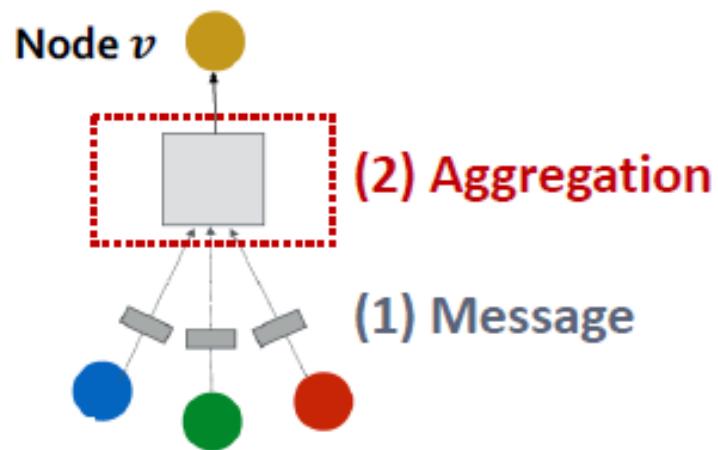
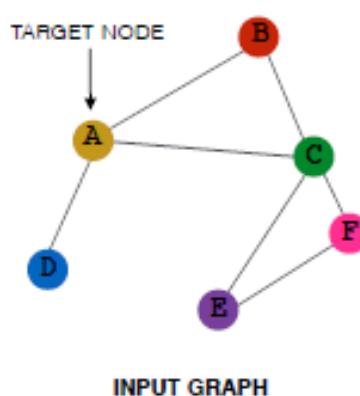
## ■ (2) Aggregation

- **Intuition:** Each node will aggregate the messages from node  $v$ 's neighbors

$$\mathbf{h}_v^{(l)} = \text{AGG}^{(l)} \left( \left\{ \mathbf{m}_u^{(l)}, u \in N(v) \right\} \right)$$

- **Example:** Sum( $\cdot$ ), Mean( $\cdot$ ) or Max( $\cdot$ ) aggregator

- $\mathbf{h}_v^{(l)} = \text{Sum}(\{\mathbf{m}_u^{(l)}, u \in N(v)\})$



# Message Aggregation Issue

- **Issue:** Information from node  $v$  itself **could get lost**

- Computation of  $\mathbf{h}_v^{(l)}$  does not directly depend on  $\mathbf{h}_v^{(l-1)}$

- **Solution:** Include  $\mathbf{h}_v^{(l-1)}$  when computing  $\mathbf{h}_v^{(l)}$

- **(1) Message:** compute message from node  $v$  itself

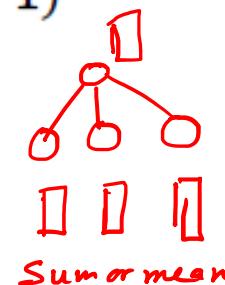
- Usually, a **different message computation** will be performed



$$\mathbf{m}_u^{(l)} = \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)}$$



$$\mathbf{m}_v^{(l)} = \mathbf{B}^{(l)} \mathbf{h}_v^{(l-1)}$$



- **(2) Aggregation:** After aggregating from neighbors, we can **aggregate the message from node  $v$  itself**

- Via **concatenation** or **summation**

Then aggregate from node itself

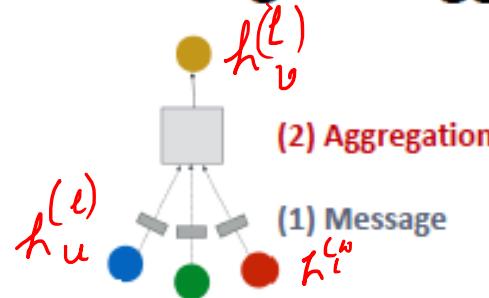
$$\mathbf{h}_v^{(l)} = \text{CONCAT} \left( \text{AGG} \left( \left\{ \mathbf{m}_u^{(l)}, u \in N(v) \right\} \right), \boxed{\mathbf{m}_v^{(l)}} \right)$$

First aggregate from neighbors

# A Single GNN Layer

## ■ Putting things together:

- **(1) Message**: each node computes a message  
$$\mathbf{m}_u^{(l)} = \text{MSG}^{(l)}\left(\mathbf{h}_u^{(l-1)}\right), u \in \{N(v) \cup v\}$$
- **(2) Aggregation**: aggregate messages from neighbors  
$$\mathbf{h}_v^{(l)} = \text{AGG}^{(l)}\left(\left\{\mathbf{m}_u^{(l)}, u \in N(v)\right\}, \mathbf{m}_v^{(l)}\right)$$
- **Nonlinearity (activation)**: Adds expressiveness
  - Often written as  $\sigma(\cdot)$ : ReLU( $\cdot$ ), Sigmoid( $\cdot$ ) , ...
  - Can be added to **message or aggregation**



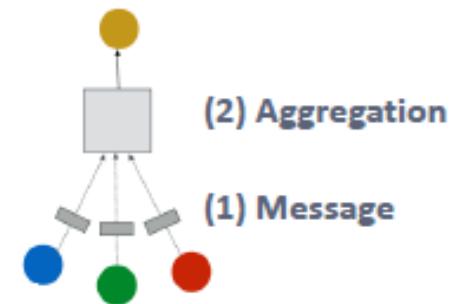
# Different GNN layers

## ■ (1) Graph Convolutional Networks (GCN)

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

## ■ How to write this as Message + Aggregation?

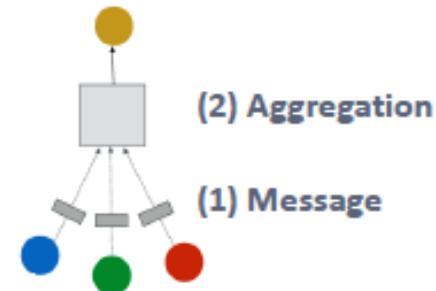
$$\mathbf{h}_v^{(l)} = \sigma \left( \sum_{u \in N(v)} \underbrace{W^{(l)} \frac{\mathbf{h}_u^{(l-1)}}{|N(v)|}}_{\text{Aggregation}} \underbrace{\mathbf{h}_u^{(l-1)}}_{\text{Message}} \right)$$



# Different GNN Layers

## ■ (1) Graph Convolutional Networks (GCN)

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



### ■ Message:

- Each Neighbor:  $\mathbf{m}_u^{(l)} = \frac{1}{|N(v)|} \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)}$

Normalized by node degree  
(In the GCN paper they use a slightly different normalization)

### ■ Aggregation:

- **Sum** over messages from neighbors, then apply activation
- $\mathbf{h}_v^{(l)} = \sigma \left( \text{Sum} \left( \{\mathbf{m}_u^{(l)}, u \in N(v)\} \right) \right)$

# Different GNN Layers: GraphSage

## ■ (2) GraphSAGE

$$\mathbf{h}_v^{(l)} = \sigma \left( \mathbf{w}^{(l)} \cdot \text{CONCAT} \left( \mathbf{h}_v^{(l-1)}, \text{AGG} \left( \left\{ \mathbf{h}_u^{(l-1)}, \forall u \in N(v) \right\} \right) \right) \right)$$

## ■ How to write this as Message + Aggregation?

- **Message** is computed within the **AGG( $\cdot$ )**

- **Two-stage aggregation**

- **Stage 1:** Aggregate from node neighbors

$$\mathbf{h}_{N(v)}^{(l)} \leftarrow \text{AGG} \left( \left\{ \mathbf{h}_u^{(l-1)}, \forall u \in N(v) \right\} \right)$$

- **Stage 2:** Further aggregate over the node itself

$$\mathbf{h}_v^{(l)} \leftarrow \sigma \left( \mathbf{w}^{(l)} \cdot \text{CONCAT}(\mathbf{h}_v^{(l-1)}, \mathbf{h}_{N(v)}^{(l)}) \right)$$

# GraphSAGE Neighbor Aggregation

- **Mean:** Take a weighted average of neighbors

$$\text{AGG} = \underbrace{\sum_{u \in N(v)} \frac{\mathbf{h}_u^{(l-1)}}{|N(v)|}}_{\text{Aggregation}} \quad \text{Message computation}$$

- **Pool:** Transform neighbor vectors and apply symmetric vector function  $\text{Mean}(\cdot)$  or  $\text{Max}(\cdot)$

$$\text{AGG} = \underbrace{\text{Mean}(\{\text{MLP}(\mathbf{h}_u^{(l-1)}), \forall u \in N(v)\})}_{\text{Aggregation}} \quad \text{Message computation}$$

- **LSTM:** Apply LSTM to reshuffled of neighbors

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

# GraphSAGE L2 normalization

- **$\ell_2$  Normalization:**

- **Optional:** Apply  $\ell_2$  normalization to  $\mathbf{h}_v^{(l)}$  at every layer
- $\mathbf{h}_v^{(l)} \leftarrow \frac{\mathbf{h}_v^{(l)}}{\|\mathbf{h}_v^{(l)}\|_2} \quad \forall v \in V \text{ where } \|u\|_2 = \sqrt{\sum_i u_i^2} \text{ ( $\ell_2$ -norm)}$
- Without  $\ell_2$  normalization, the embedding vectors have different scales ( $\ell_2$ -norm) for vectors
- In some cases (not always), normalization of embedding results in performance improvement
- After  $\ell_2$  normalization, all vectors will have the same  $\ell_2$ -norm

# Graph Attention (GAT)

## ■ (3) Graph Attention Networks

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

Attention weights

## ■ In GCN / GraphSAGE

- $\alpha_{vu} = \frac{1}{|N(v)|}$  is the **weighting factor (importance)** of node  $u$ 's message to node  $v$
- $\Rightarrow \alpha_{vu}$  is defined **explicitly** based on the **structural properties** of the graph (node degree)
- $\Rightarrow$  All neighbors  $u \in N(v)$  are **equally important** to node  $v$

# Graph Attention (GAT)

## ■ (3) Graph Attention Networks

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

Attention weights

**Not all node's neighbors are equally important**

- **Attention** is inspired by cognitive attention.
- The **attention**  $\alpha_{vu}$  focuses on the important parts of the input data and fades out the rest.
  - **Idea:** the NN should devote more computing power on that small but important part of the data.
  - Which part of the data is more important depends on the context and is learned through training.

# Graph Attention Networks

Can we do better than simple neighborhood aggregation?

Can we let weighting factors  $\alpha_{vu}$  to be learned?

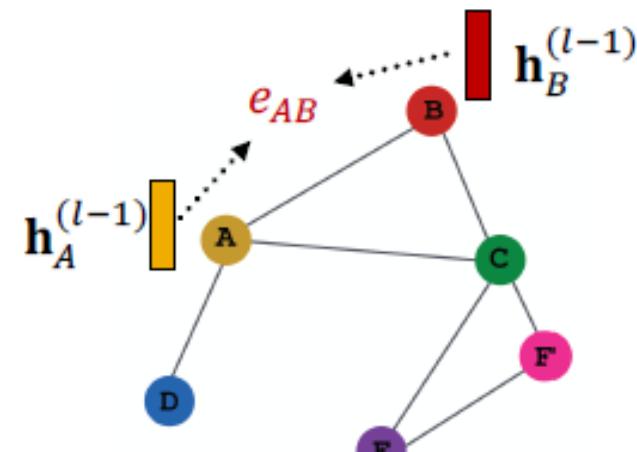
- **Goal:** Specify **arbitrary importance** to different neighbors of each node in the graph
- **Idea:** Compute embedding  $\mathbf{h}_v^{(l)}$  of each node in the graph following an **attention strategy**:
  - Nodes attend over their neighborhoods' message
  - Implicitly specifying different weights to different nodes in a neighborhood

# Attention Mechanism

- Let  $\alpha_{vu}$  be computed as a byproduct of an **attention mechanism**  $a$ :
  - (1) Let  $a$  compute **attention coefficients**  $e_{vu}$  across pairs of nodes  $u, v$  based on their messages:

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

- $e_{vu}$  indicates the importance of  $u$ 's message to node  $v$



$$e_{AB} = a(\mathbf{W}^{(l)} \mathbf{h}_A^{(l-1)}, \mathbf{W}^{(l)} \mathbf{h}_B^{(l-1)})$$

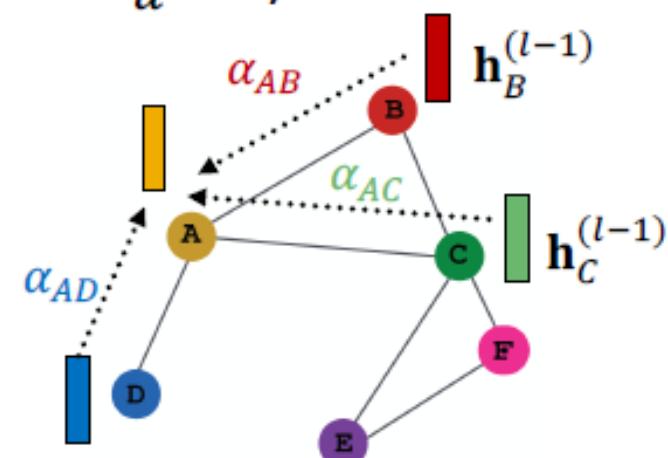
# Attention Mechanism

- **Normalize**  $e_{vu}$  into the **final attention weight**  $\alpha_{vu}$ 
  - Use the **softmax** function, so that  $\sum_{u \in N(v)} \alpha_{vu} = 1$ :
$$\alpha_{vu} = \frac{\exp(e_{vu})}{\sum_{k \in N(v)} \exp(e_{vk})}$$
- **Weighted sum** based on the **final attention weight**  $\alpha_{vu}$

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

**Weighted sum using**  $\alpha_{AB}$ ,  $\alpha_{AC}$ ,  $\alpha_{AD}$ :

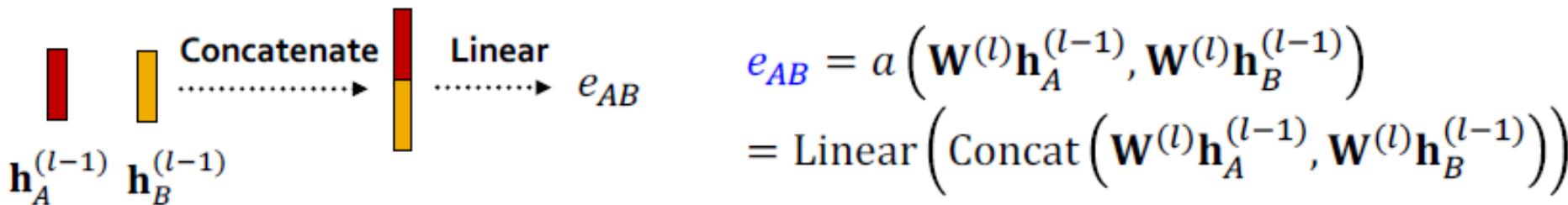
$$\mathbf{h}_A^{(l)} = \sigma(\alpha_{AB} \mathbf{W}^{(l)} \mathbf{h}_B^{(l-1)} + \alpha_{AC} \mathbf{W}^{(l)} \mathbf{h}_C^{(l-1)} + \alpha_{AD} \mathbf{W}^{(l)} \mathbf{h}_D^{(l-1)})$$



# Attention Mechanism

## ■ What is the form of attention mechanism $a$ ?

- The approach is agnostic to the choice of  $a$ 
  - E.g., use a simple single-layer neural network
    - $a$  have trainable parameters (weights in the Linear layer)



- Parameters of  $a$  are trained jointly:
  - Learn the parameters together with weight matrices (i.e., other parameter of the neural net  $\mathbf{W}^{(l)}$ ) in an end-to-end fashion

# Multihead Attention

- **Multi-head attention:** Stabilizes the learning process of attention mechanism

- Create **multiple attention scores** (each replica with a different set of parameters):

$$\mathbf{h}_v^{(l)}[1] = \sigma(\sum_{u \in N(v)} \alpha_{vu}^1 \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)})$$

$$\mathbf{h}_v^{(l)}[2] = \sigma(\sum_{u \in N(v)} \alpha_{vu}^2 \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)})$$

$$\mathbf{h}_v^{(l)}[3] = \sigma(\sum_{u \in N(v)} \alpha_{vu}^3 \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)})$$

- **Outputs are aggregated:**

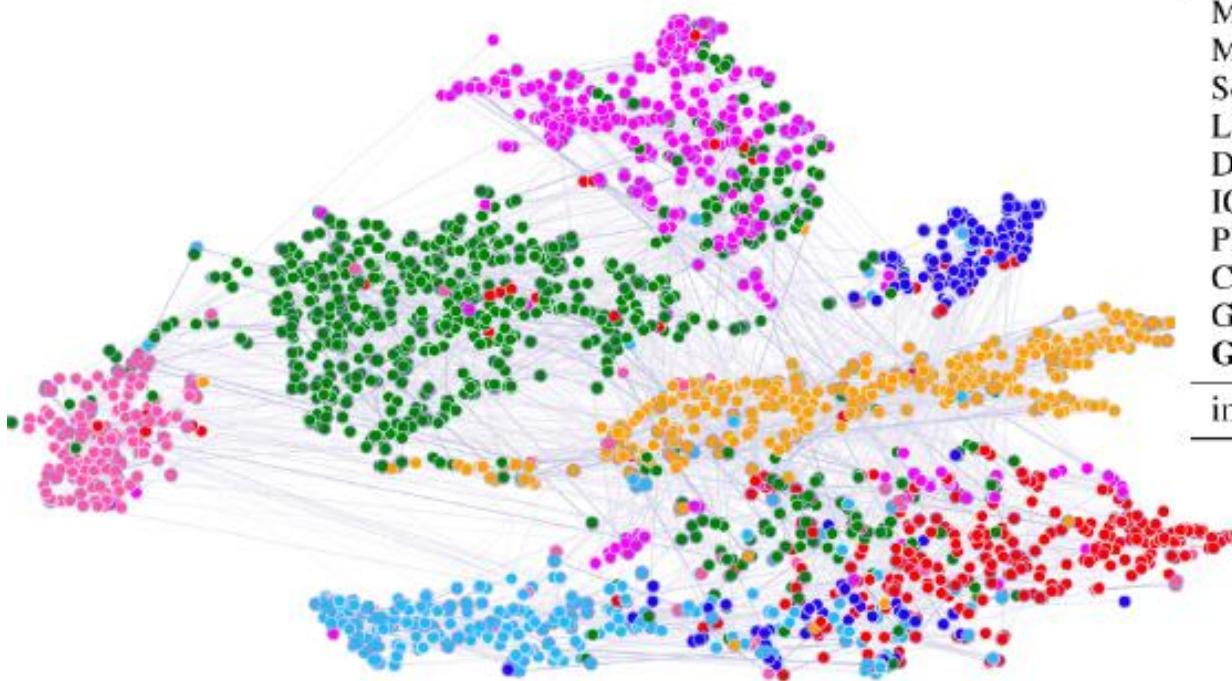
- By concatenation or summation

- $\mathbf{h}_v^{(l)} = \text{AGG}(\mathbf{h}_v^{(l)}[1], \mathbf{h}_v^{(l)}[2], \mathbf{h}_v^{(l)}[3])$

# Benefits of GAT

- **Key benefit:** Allows for (implicitly) specifying **different importance values ( $\alpha_{vu}$ ) to different neighbors**
- **Computationally efficient:**
  - Computation of attentional coefficients can be parallelized across all edges of the graph
  - Aggregation may be parallelized across all nodes
- **Storage efficient:**
  - Sparse matrix operations do not require more than  $O(V + E)$  entries to be stored
  - **Fixed** number of parameters, irrespective of graph size
- **Localized:**
  - Only **attends over local network neighborhoods**
- **Inductive capability:**
  - It is a shared *edge-wise* mechanism
  - It does not depend on the global graph structure

# GAT Example Cora Citation Net



Method	Cora
MLP	55.1%
ManiReg (Belkin et al., 2006)	59.5%
SemiEmb (Weston et al., 2012)	59.0%
LP (Zhu et al., 2003)	68.0%
DeepWalk (Perozzi et al., 2014)	67.2%
ICA (Lu & Getoor, 2003)	75.1%
Planetoid (Yang et al., 2016)	75.7%
Chebyshev (Defferrard et al., 2016)	81.2%
GCN (Kipf & Welling, 2017)	81.5%
<b>GAT</b>	<b>83.3%</b>
improvement w.r.t GCN	1.8%

Attention mechanism can be used with many different graph neural network models

In many cases, attention leads to performance gains

- **t-SNE plot of GAT-based node embeddings:**
  - Node color: 7 publication classes
  - Edge thickness: Normalized attention coefficients between nodes  $i$  and  $j$ , across eight attention heads,  $\sum_k (\alpha_{ij}^k + \alpha_{ji}^k)$

# GNN Layers in practice

- Many modern deep learning modules can be incorporated into a GNN layer

- **Batch Normalization:**

- Stabilize neural network training

- **Dropout:**

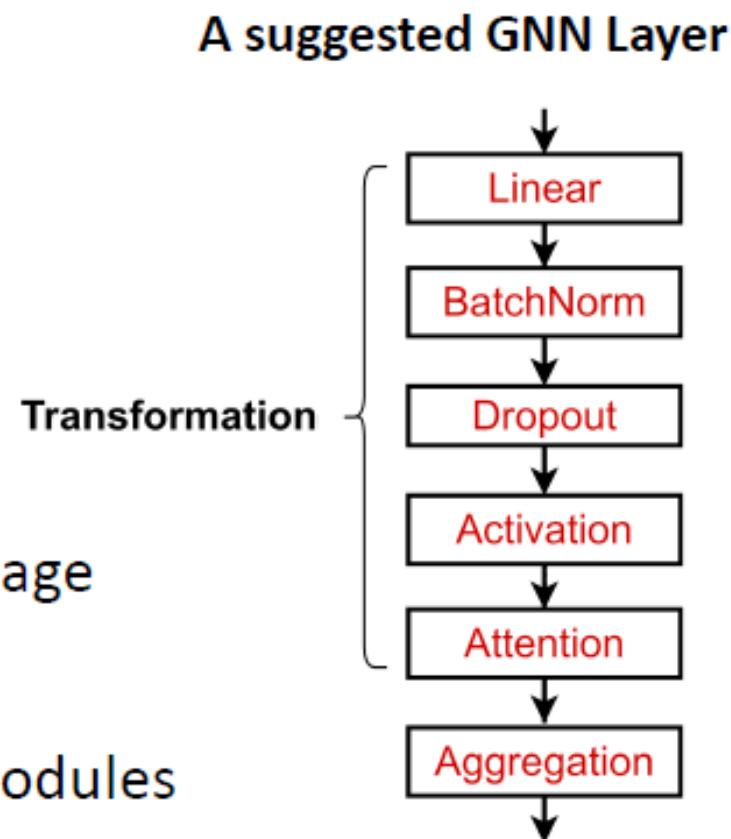
- Prevent overfitting

- **Attention/Gating:**

- Control the importance of a message

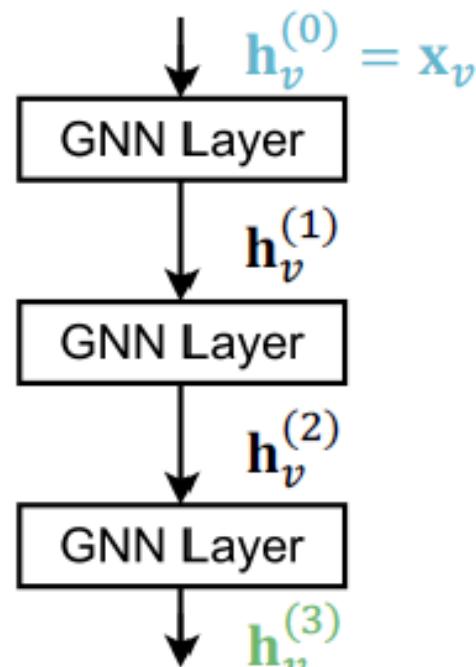
- **More:**

- Any other useful deep learning modules



# Stacking GNN layers

- **How to construct a Graph Neural Network?**
  - **The standard way:** Stack GNN layers sequentially
  - **Input:** Initial raw node feature  $\mathbf{x}_v$
  - **Output:** Node embeddings  $\mathbf{h}_v^{(L)}$  after  $L$  GNN layers



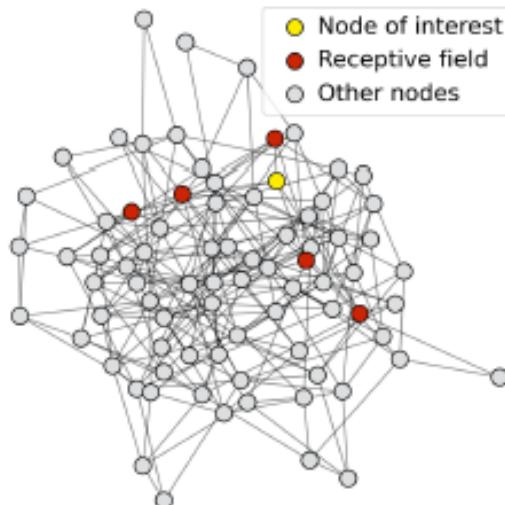
# The over-smoothing problem

- **The Issue of stacking many GNN layers**
  - GNN suffers from **the over-smoothing problem**
- **The over-smoothing problem: all the node embeddings converge to the same value**
  - This is bad because we **want to use node embeddings to differentiate nodes**
- **Why does the over-smoothing problem happen?**

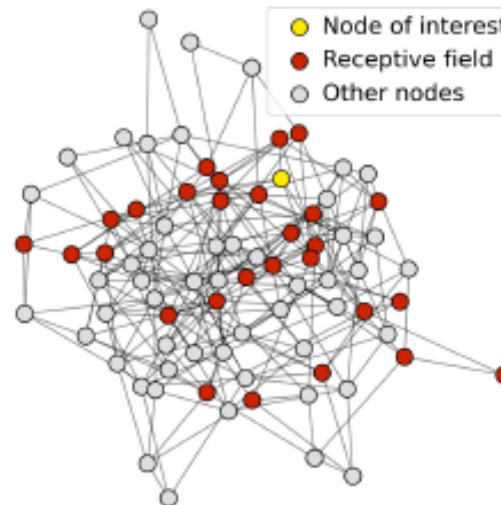
# Receptive fields of GNNs

- **Receptive field:** the set of nodes that determine the embedding of a node of interest
  - In a  $K$ -layer GNN, each node has a receptive field of  $K$ -hop neighborhood

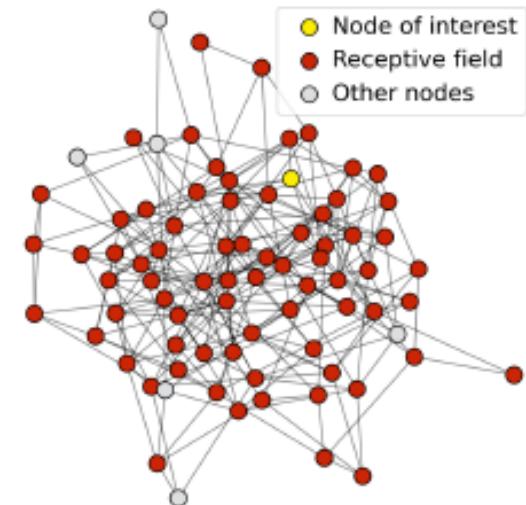
Receptive field for  
1-layer GNN



Receptive field for  
2-layer GNN



Receptive field for  
3-layer GNN

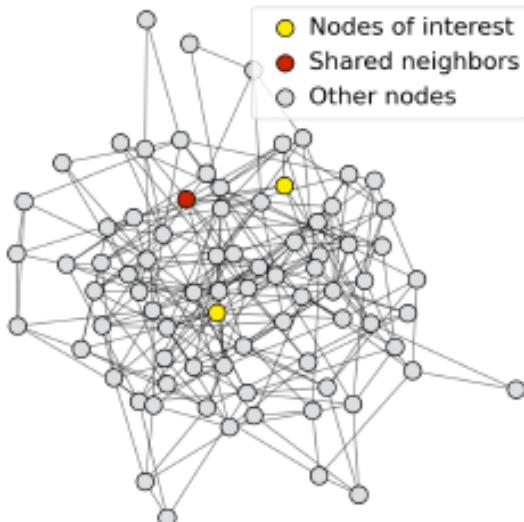


# Receptive fields of GNNs

- **Receptive field overlap for two nodes**
  - **The shared neighbors quickly grows** when we increase the number of hops (num of GNN layers)

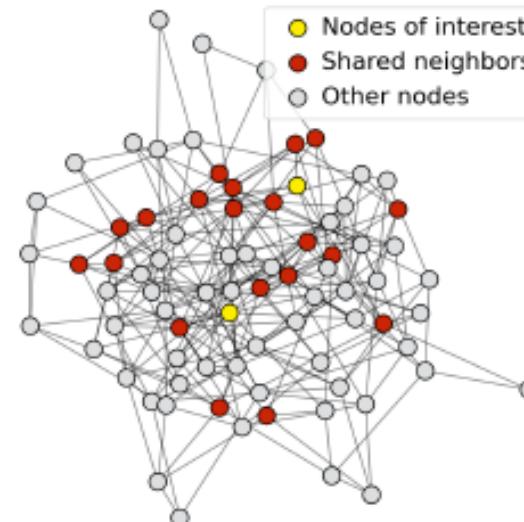
**1-hop neighbor overlap**

Only 1 node



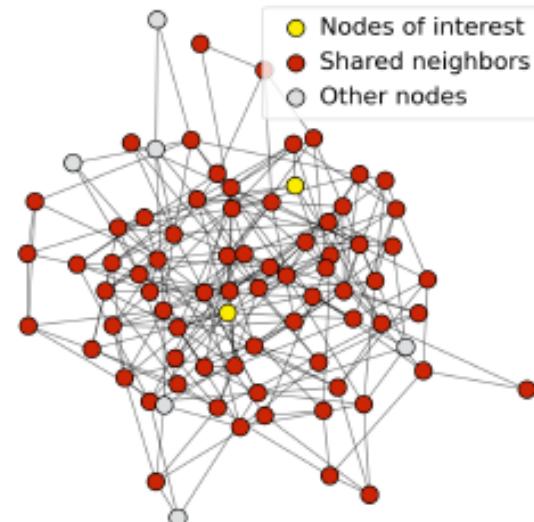
**2-hop neighbor overlap**

About 20 nodes



**3-hop neighbor overlap**

Almost all the nodes!



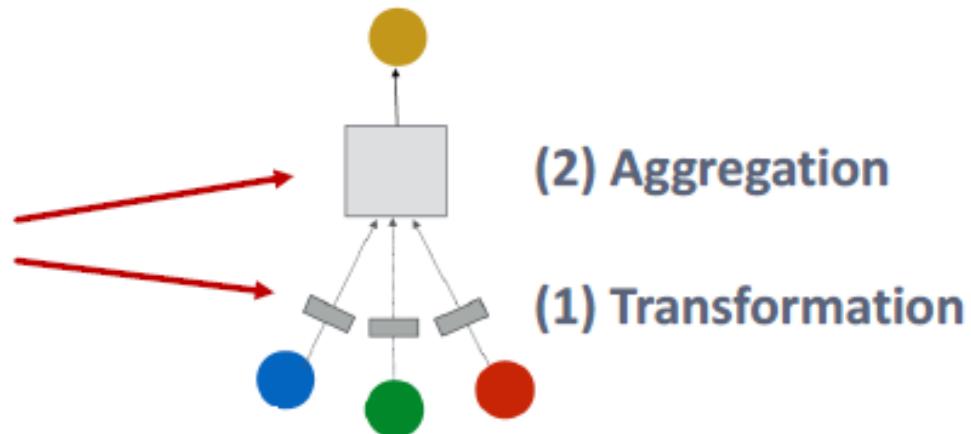
# Mitigating Over-Smoothing Problem

- **What do we learn from the over-smoothing problem?**
- **Lesson 1: Be cautious when adding GNN layers**
  - Unlike neural networks in other domains (CNN for image classification), **adding more GNN layers do not always help**
  - **Step 1:** Analyze the necessary receptive field to solve your problem. E.g., by computing the diameter of the graph
  - **Step 2:** Set number of GNN layers  $L$  to be a bit more than the receptive field we like. **Do not set  $L$  to be unnecessarily large!**
- **Question:** How to enhance the expressive power of a GNN, if the number of GNN layers is small?

# Solution - 1

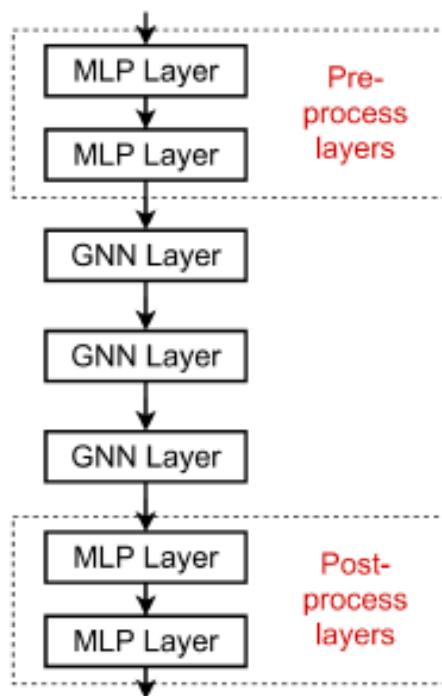
- How to make a shallow GNN more expressive?
- Solution 1: Increase the expressive power within each GNN layer
  - In our previous examples, each transformation or aggregation function only include one linear layer
  - We can make aggregation / transformation become a deep neural network!

If needed, each box could include a 3-layer MLP



# Solution - 2

- **How to make a shallow GNN more expressive?**
- **Solution 2:** Add layers that do not pass messages
  - A GNN does not necessarily only contain GNN layers
    - E.g., we can add **MLP layers** (applied to each node) before and after GNN layers, as **pre-process layers** and **post-process layers**



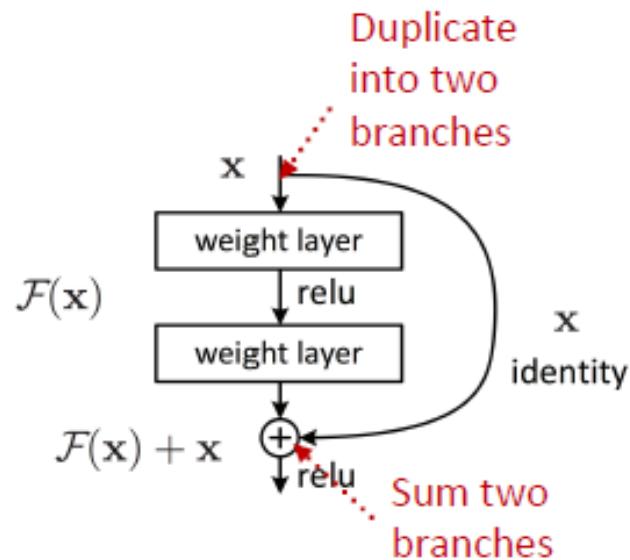
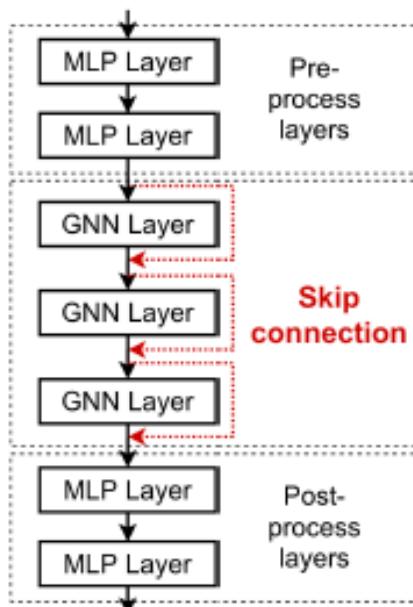
**Pre-processing layers:** Important when encoding node features is necessary.  
E.g., when nodes represent images/text

**Post-processing layers:** Important when reasoning / transformation over node embeddings are needed  
E.g., graph classification, knowledge graphs

In practice, adding these layers works great!

# Solution - 3

- What if my problem still requires many GNN layers?
- Lesson 2: Add skip connections in GNNs
  - Observation from over-smoothing: Node embeddings in earlier GNN layers can sometimes better differentiate nodes
  - Solution: We can increase the impact of earlier layers on the final node embeddings, by adding shortcuts in GNN



**Idea of skip connections:**  
Before adding shortcuts:  
 $F(x)$   
After adding shortcuts:  
 $F(x) + x$

# Scaling up GCNs

- How to implement GCNs for large graphs

# Problem with scale

- **Large-scale:**
  - #nodes ranges from 10M to 10B.
  - #edges ranges from 100M to 100B.
- **Tasks**
  - **Node-level:** User/item/paper classification.
  - **Link-level:** Recommendation, completion.
- **Todays' lecture**
  - Scale up GNNs to large graphs!

# Why it is hard

- **Recall:** How we usually train an ML model on large data ( $N=\#\text{data}$  is large)
- **Objective:** Minimize the averaged loss

$$\ell(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=0}^{N-1} \ell_i(\boldsymbol{\theta})$$

- $\boldsymbol{\theta}$ : model parameters,  $\ell_i(\boldsymbol{\theta})$ : loss for  $i$ -th data point.
- We perform **Stochastic Gradient Descent (SGD)**.
  - Randomly sample  $M$  ( $\ll N$ ) data (**mini-batches**).
  - Compute the  $\ell_{\text{sub}}(\boldsymbol{\theta})$  over the  $M$  data points.
  - Perform SGD:  $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \nabla \ell_{\text{sub}}(\boldsymbol{\theta})$

# Methods

We introduce **three methods for scaling up GNNs**:

- Two methods perform message-passing over **small subgraphs in each mini-batch**; only the subgraphs need to be loaded on a GPU at a time.
  - **Neighbor Sampling** [Hamilton et al. NeurIPS 2017]
  - **Cluster-GCN** [Chiang et al. KDD 2019]
- One method **simplifies a GNN into feature-preprocessing operation** (can be efficiently performed even on a CPU)
  - **Simplified GCN** [Wu et al. ICML 2019]

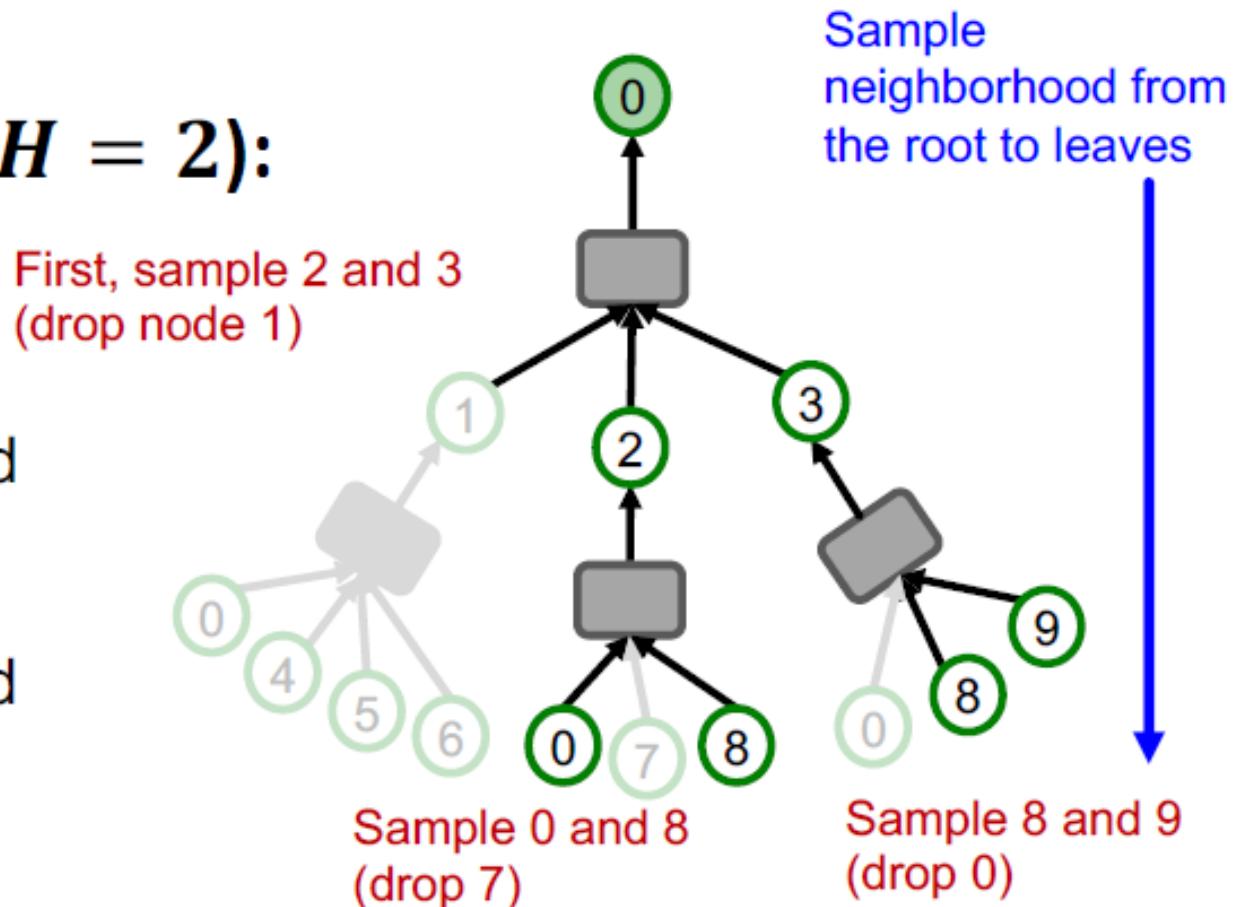
# Neighbor Sampling

**Key idea:** Construct the computational graph by (randomly) sampling at most  $H$  neighbors at each hop.

■ Example ( $H = 2$ ):

1<sup>st</sup>-hop neighborhood

2<sup>nd</sup>-hop neighborhood



# Issues

- **Remark 1: Trade-off in sampling number  $H$** 
  - Smaller  $H$  leads to more efficient neighbor aggregation, but results in more unstable training **due to the larger variance** in neighbor aggregation.
- **Remark 2: Computational time**
  - Even with neighbor sampling, **the size of the computational graph is still exponential with respect to number of GNN layers  $K$ .**
  - Increasing one GNN layer would make computation  $H$  times more expensive.

# How to sample nodes?

## ■ Remark 3: How to sample the nodes

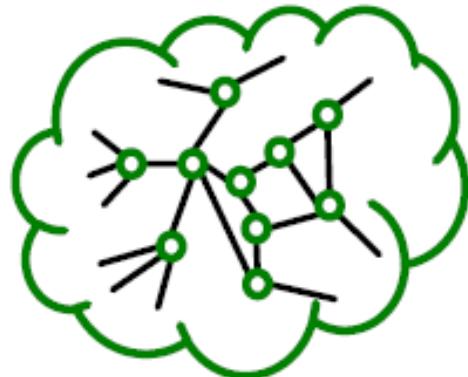
- **Random sampling:** fast but many times not optimal (may sample many “unimportant” nodes)
- **Random Walk with Restarts:**
  - Natural graphs are “scale free”, sampling random neighbors, samples many low degree “leaf” nodes.
  - Strategy to sample important nodes:
    - Compute Random Walk with Restarts score  $R_i$  starting at the **green** node
    - At each level sample  $H$  neighbors  $i$  with the highest  $R_i$
  - This strategy works much better in practice.



# Subgraph Sampling

- **Key idea:** We can **sample a small subgraph of the large graph** and then perform the efficient **layer-wise** node embeddings update over the subgraph.

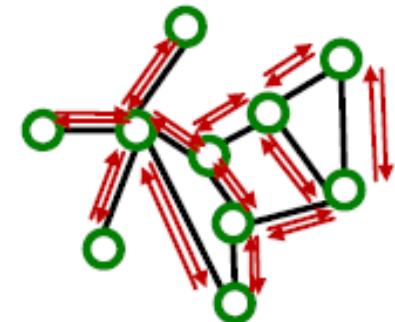
Large graph



Sampled subgraph  
(small enough to  
be put on a GPU)



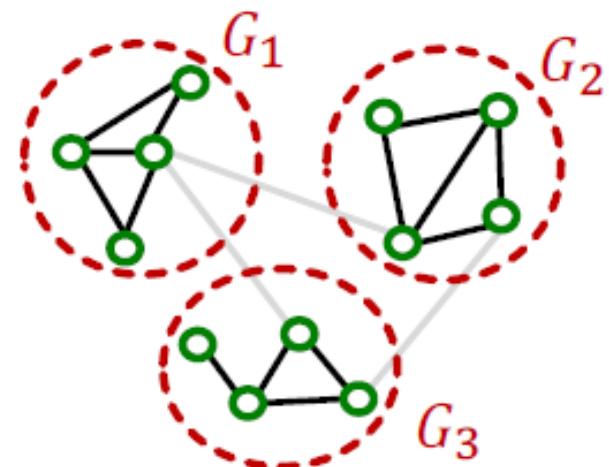
Layer-wise  
node embeddings  
update on the GPU



# Cluster GCN

- Given a large graph  $G = (V, E)$ , **partition its nodes  $V$  into  $C$  groups:  $V_1, \dots, V_C$ .**
  - We can use any scalable community detection methods, e.g., Louvain, METIS [Karypis et al. SIAM 1998].
- $V_1, \dots, V_C$  induces  $C$  subgraphs,  $G_1, \dots, G_C$ ,
  - Recall:  $G_c \equiv (V_c, E_c)$ ,
  - where  $E_c = \{(u, v) | u, v \in V_c\}$

**Notice: Between-group edges are *not* included in  $G_1, \dots, G_C$ .**



# Simplifying GCN

- We start from **Graph Convolutional Network (GCN)** [Kipf & Welling ICLR 2017].
- We simplify GCN by **removing the non-linear activation** from the GCN [Wu et al. ICML 2019].
  - *Wu et al.* demonstrated that the performance on benchmark is not much lower by the simplification.
- Simplified GCN turns out to be extremely scalable by the model design.

# GCN (Mean Pool)

- **Given:** Graph  $G = (V, E)$  with input node features  $X_v$  for  $v \in V$ , where  **$E$  includes the self-loop:**

- $(v, v) \in E$  for all  $v \in V$ .



- Set input node embeddings:  $h_v^{(0)} = X_v$  for  $v \in V$ .
- For  $k \in \{0, \dots, K - 1\}$ :
  - For all  $v \in V$ , aggregate neighboring information as

$$h_v^{(k+1)} = \text{ReLU} \left( \mathbf{W}_k \left[ \frac{1}{|N(v)|} \sum_{u \in N(v)} h_u^{(k)} \right] \right)$$

Trainable weight matrices  
(i.e., what we learn)

Mean-pooling

- **Final node embedding:**  $z_v = h_v^{(K)}$

# Simplifying GCN

- Simplify GCN by removing ReLU non-linearity:

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

- The final node embedding matrix is given as

$$H^{(K)} = \tilde{A} \underbrace{H^{(K-1)}}_{\text{pink bracket}} W_{K-1}^T$$

$$= \tilde{A} (\underbrace{\tilde{A} H^{(K-2)} W_{K-2}^T}_{\text{green bracket}}) W_{K-1}^T$$

$$\dots = \underbrace{\tilde{A} (\tilde{A} (\dots (\underbrace{\tilde{A} H^{(0)} W_0^T}_{\text{blue bracket}} \dots) W_{K-2}^T) W_{K-1}^T}_{\text{orange bracket}}$$

$$= \tilde{A}^K X \underbrace{(W_0^T \dots W_{K-1}^T)}_{\text{pink bracket}}$$

Composition of linear transformation is still linear!

$$= \tilde{A}^K X W^T \quad \text{where } W \equiv W_{K-1} \dots W_0$$