

# GRAPH NEURAL NETWORK

Alireza Akhavanpour  
<http://Class.vision>

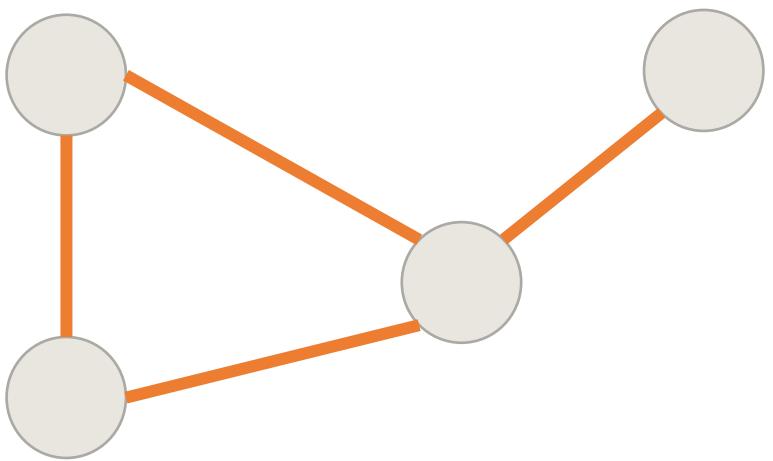
# GRAPH TERMINOLOGY

What is Node, Edge, and ...

How we can store graphs?

...

# GRAPH DEFINITION

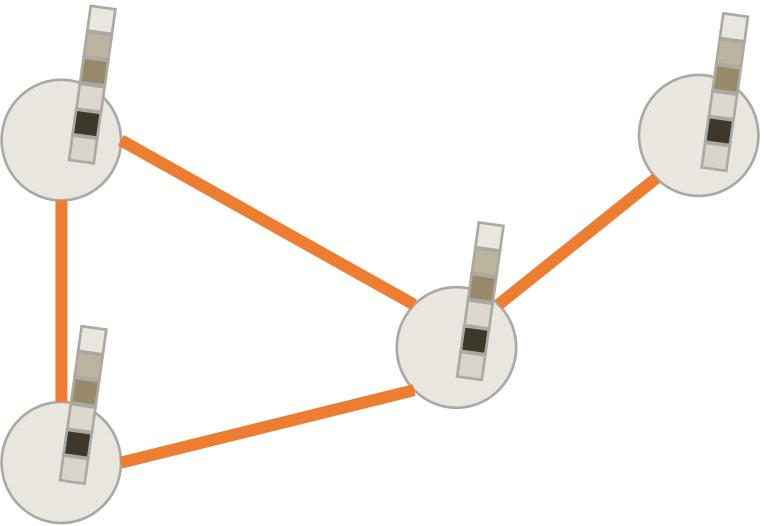


Nodes



Edges

# GRAPH DEFINITION



Nodes



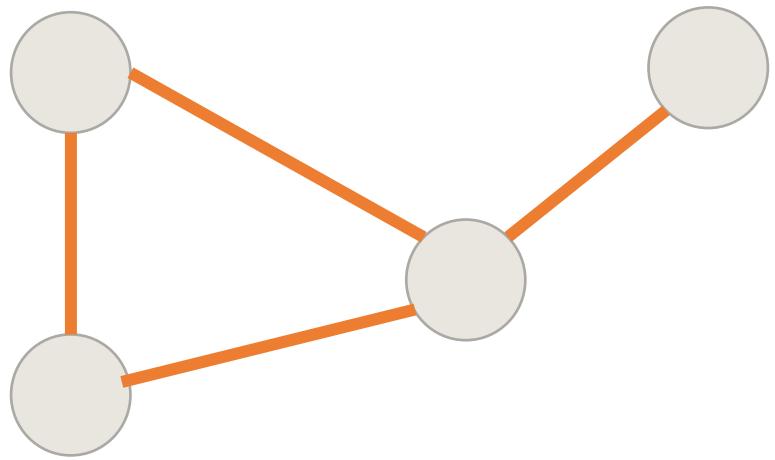
Edges



Node Features

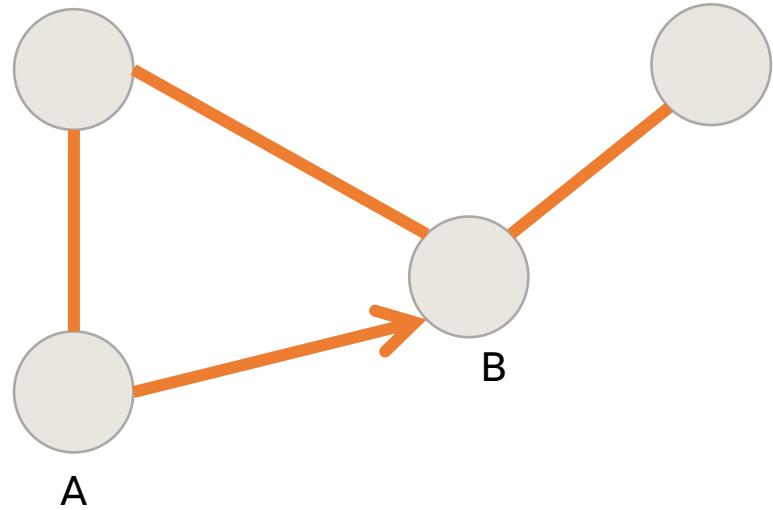
# TYPES OF GRAPH

- Undirected graph
- Directed graph



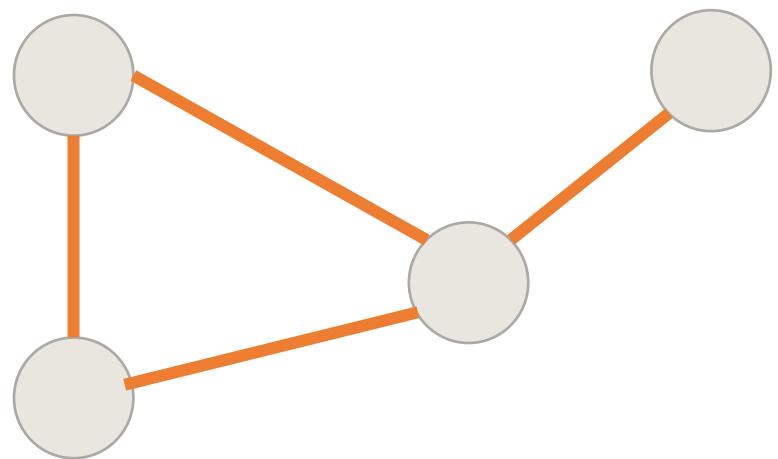
# TYPES OF GRAPH

- Undirected graph  or 
- Directed graph 



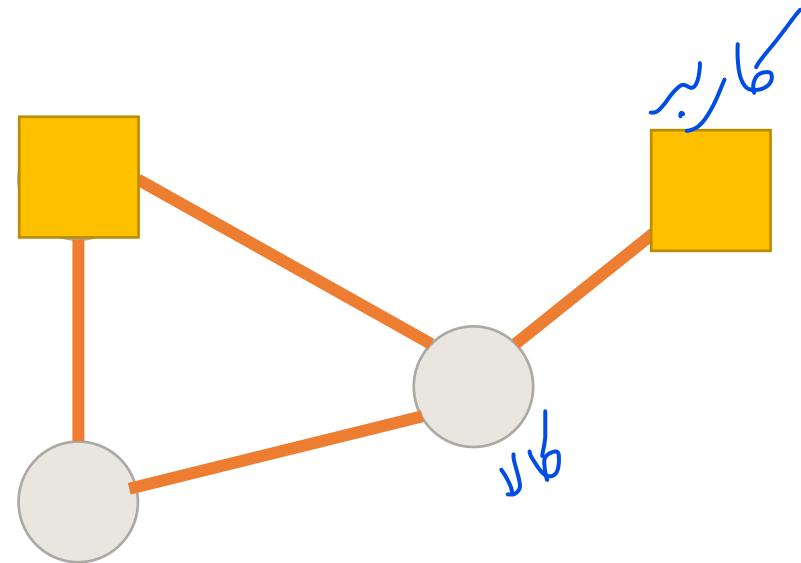
# TYPES OF GRAPH

- Homogeneous graph
- Heterogeneous graph



Homogeneous

نحو از مکانیک



Heterogeneous

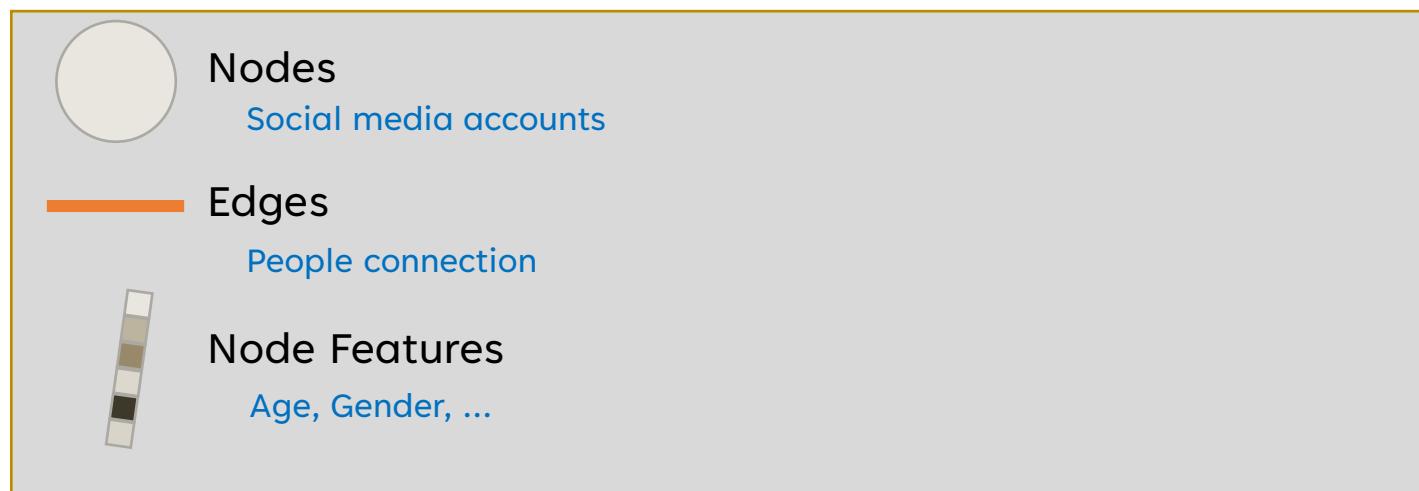
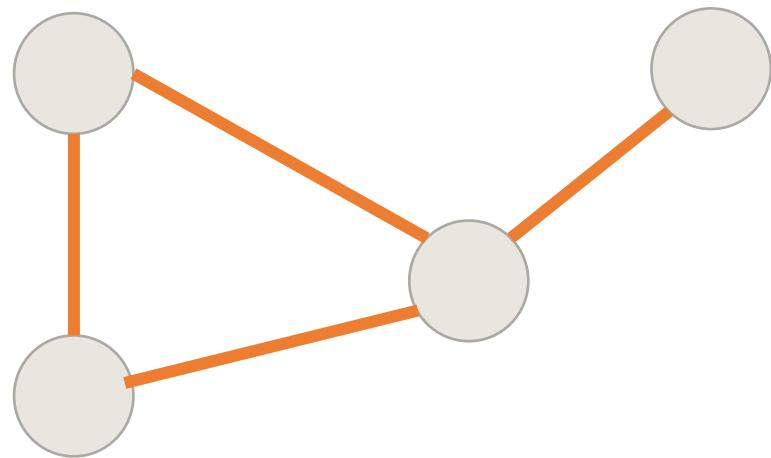
# GRAPH EXAMPLE

$G = (V, E, u)$

FEATUE VECTORS

EDGES (Adjacency, Weight) =  $(A, W)$

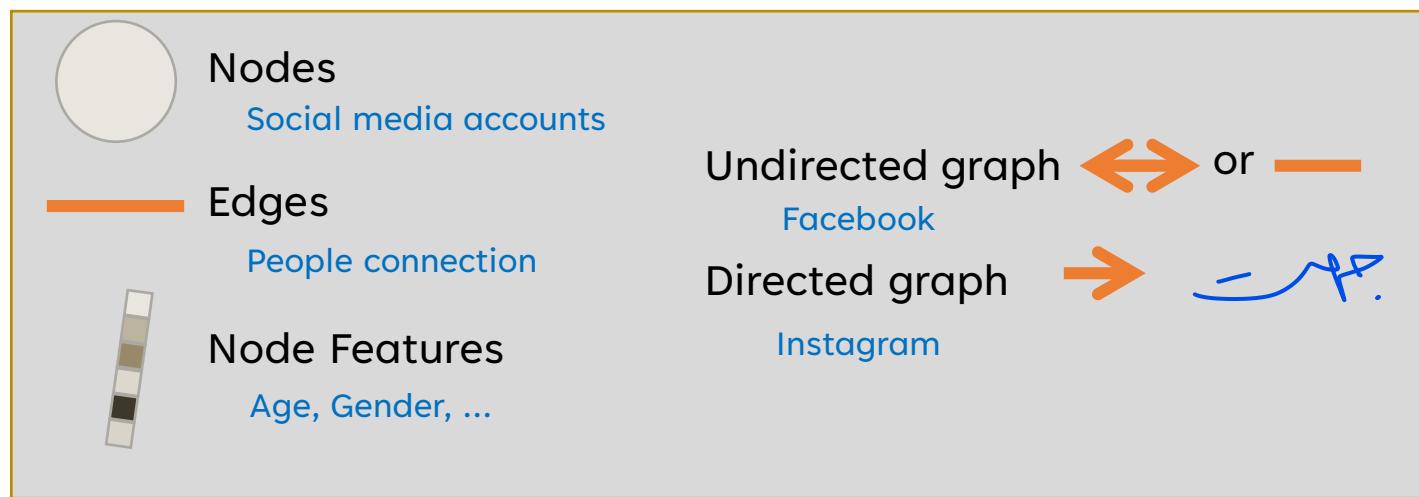
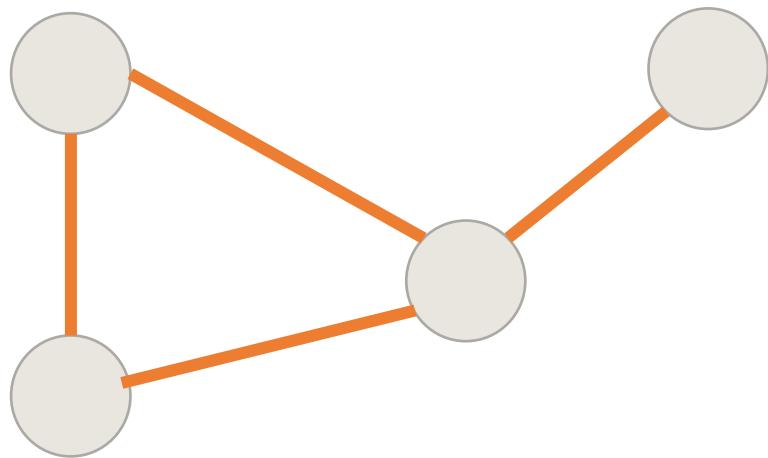
VERTECIES (NODES)

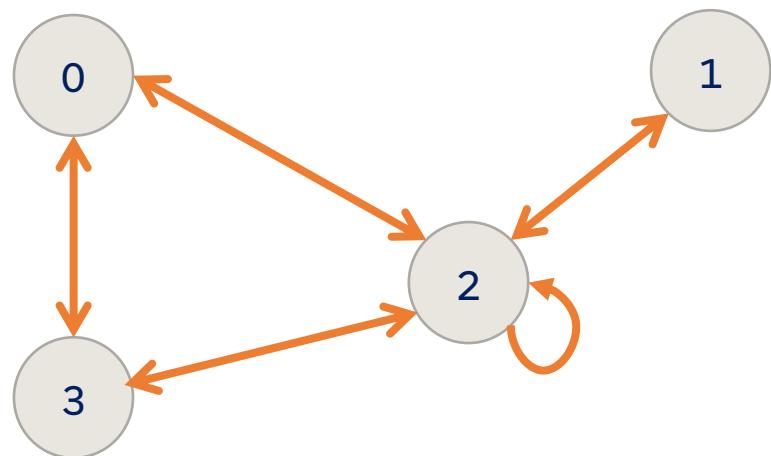


# GRAPH EXAMPLE

$G = (V, E, u)$

- FEATURE VECTORS
- EDGES (Adjacency, Weight) = (A,W)
- VERTECIES (NODES)





**Homogeneous**

## STORING GRAPH

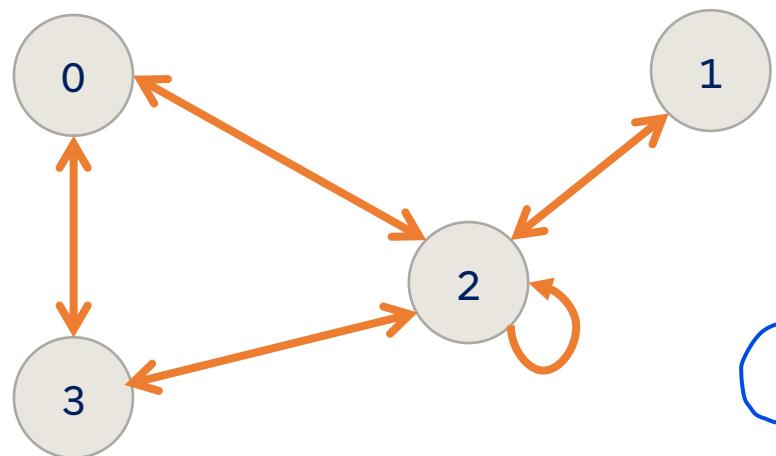
Source Node, Target Node

Source Node, Target Node
(0,1)
(0,2)
(0,3)
(1,0)
(2,0)
(2,2)
(2,3)
(3,0)
(3,2)



Edge list:

جملہ  
Eij



## STORING GRAPH

②

Adjacency Matrix:

ماتریس ارجاعی را در ماتریس می‌دانیم.

مقدار فهرست

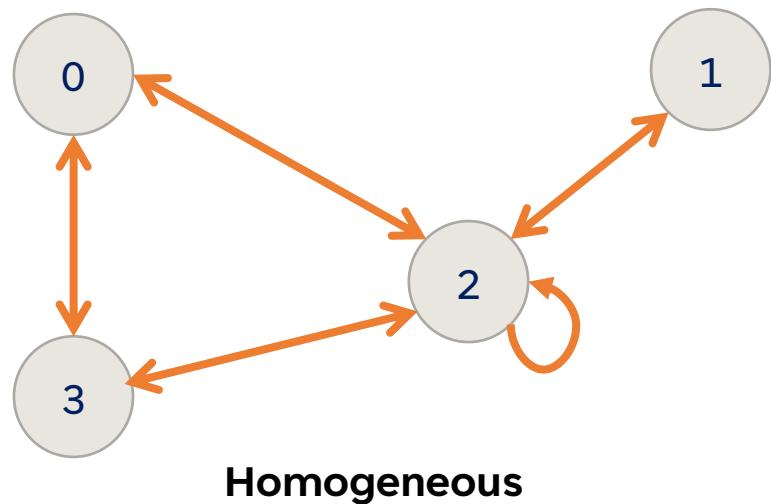
$(j_i - l_i)$

$\sum_{i=1}^n l_i$

	0	1	2	3
0	0	1	1	1
1	1	0	0	0
2	1	0	1	1
3	1	0	0	1

$V \times V$

# STORING GRAPH



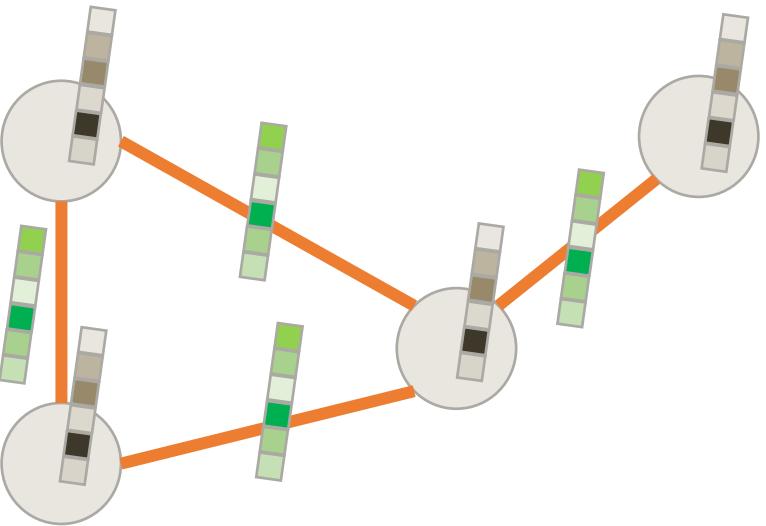
We can use **weight** instead of Boolean!  
To show how strong the connection is!

## Adjacency Matrix:

	0	1	2	3
0	0	2	1.5	4
1	5	0	0	0
2	1.5	0	1	1
3	12	0	0	1

لهم اذْعُنْ لِجَنَاحِي

# EDGE FEATURES



Nodes



Edges



Node Features



Edge Features

# YOU CAN MODEL COMPLEX SYSTEMS, DEPENDING ON HOW YOU CHOOSE TO DEFINE THE GRAPH

- Edge type:**

- weighted vs binary

- Edge directionality:**

- undirected vs directed

- Features:**

- None, node-based, edge-based

- Temporal Aspects:**

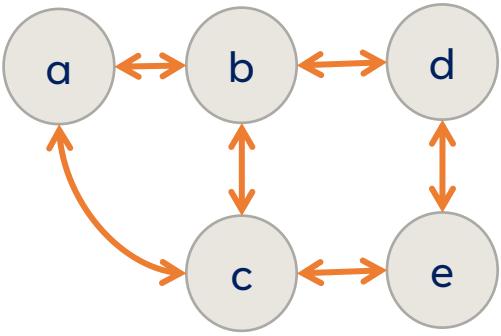
- Features, topology

- Others:**

- Multi-graphs, hypergraphs, complex networks

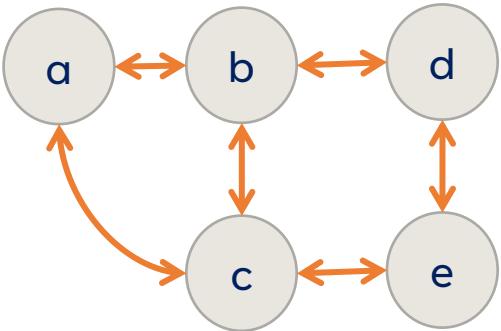
زنگنه

# GRAPH DEGREE



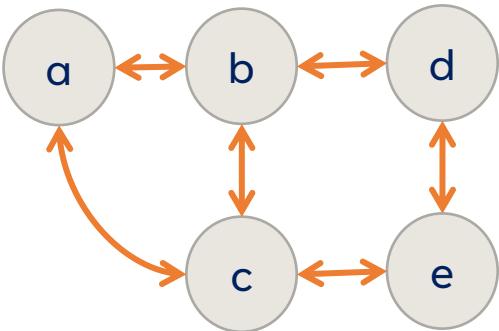
$$x = \begin{bmatrix} a \\ b \\ c \\ d \\ e \end{bmatrix}$$

## GRAPH DEGREE



$$x = \begin{bmatrix} a \\ b \\ c \\ d \\ e \end{bmatrix} \quad A = \begin{bmatrix} 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 \\ 1 & 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 1 & 0 \end{bmatrix}$$

# GRAPH DEGREE



$$x = \begin{bmatrix} a \\ b \\ c \\ d \\ e \end{bmatrix} \quad A = \begin{bmatrix} 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 \\ 1 & 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 1 & 0 \end{bmatrix}$$

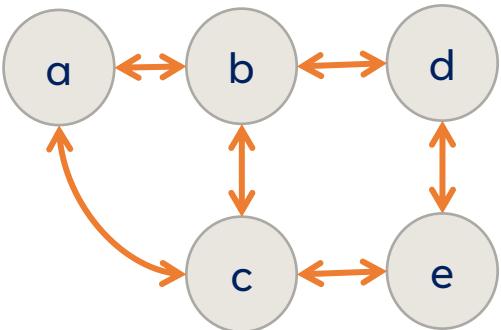
Degree matrix (D) is a diagonal matrix defining number of connection per node

*bw tf*  $\rightarrow$

$$D = \begin{bmatrix} 2 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 & 2 \end{bmatrix}$$

Degree matrix shows influence of each node on the whole graph

# LAPLACIAN OF GRAPH

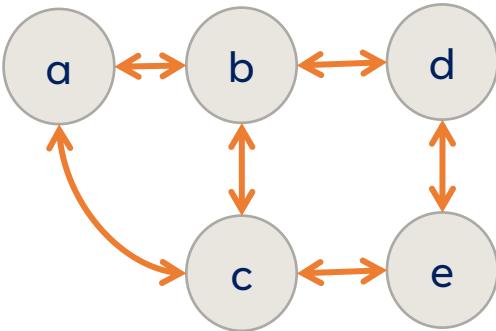


$$x = \begin{bmatrix} a \\ b \\ c \\ d \\ e \end{bmatrix}$$
$$A = \begin{bmatrix} 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 \\ 1 & 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 1 & 0 \end{bmatrix}$$
$$D = \begin{bmatrix} 2 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 & 2 \end{bmatrix}$$

Laplacian matrix ( $L$ ) is a  $L = D - A$  OR  $L = D - W$  in weighted matrix

$$L = \begin{bmatrix} 2 & -1 & -1 & 0 & 0 \\ -1 & 3 & -1 & -1 & 0 \\ -1 & -1 & 3 & 0 & -1 \\ 0 & -1 & 0 & 2 & -1 \\ 0 & 0 & -1 & -1 & 2 \end{bmatrix}$$

## NORMALIZED GRAPH

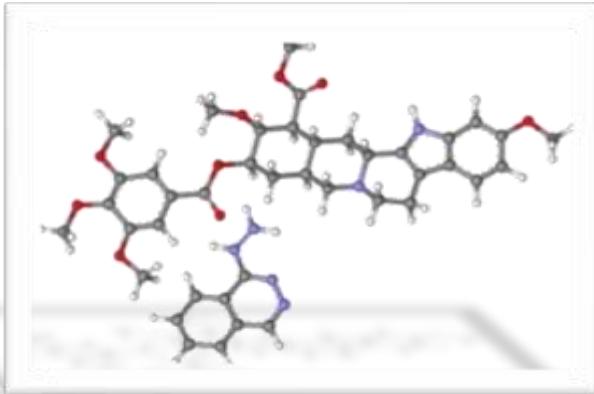
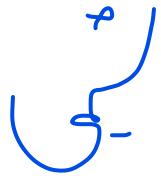


We can decide to show the relation between the nodes, with any of the following matrices:

$$A, L, \bar{A}, \bar{L}$$

# GRAPH USAGE AND APPLICATIONS

# GRAPH DATA IS EVERYWHERE



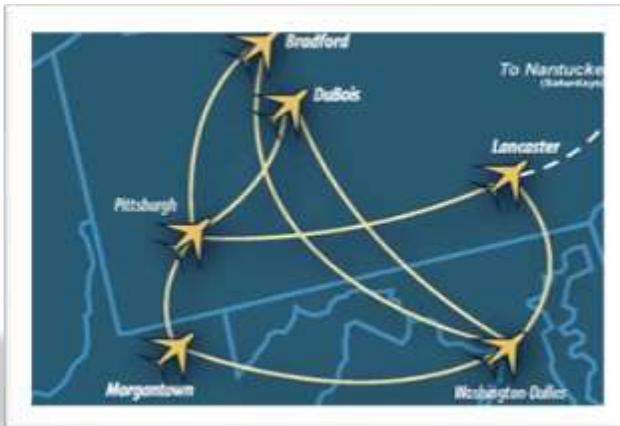
Medicine/ pharmacy



Recommender system



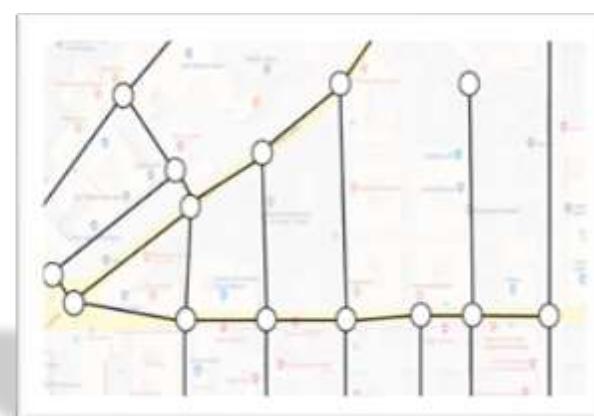
Social Networks



Airports connection



Brain cortex

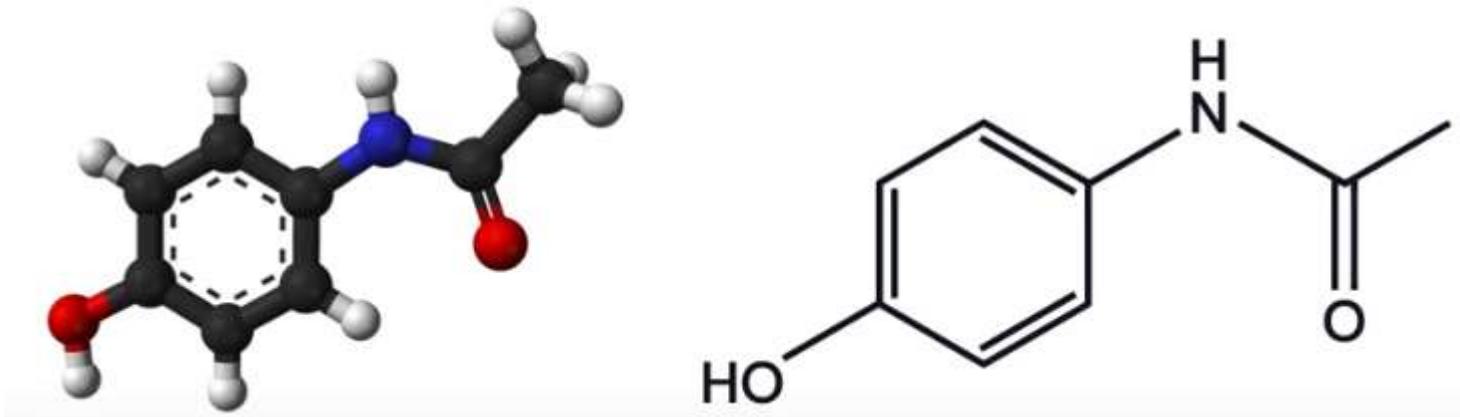


Traffic map

# MOLECULES ARE GRAPHS!

A very natural way to represent molecules is as a **graph**

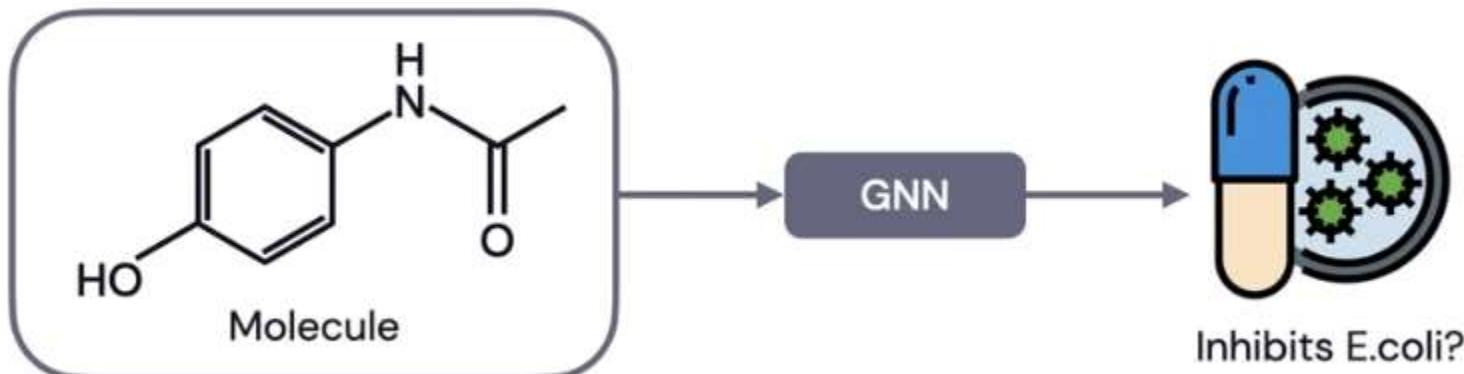
- **Atoms** as nodes, **bonds** as edges
- Features such as atom type, charge, bond type..



# GNNS FOR MOLECULE CLASSIFICATION

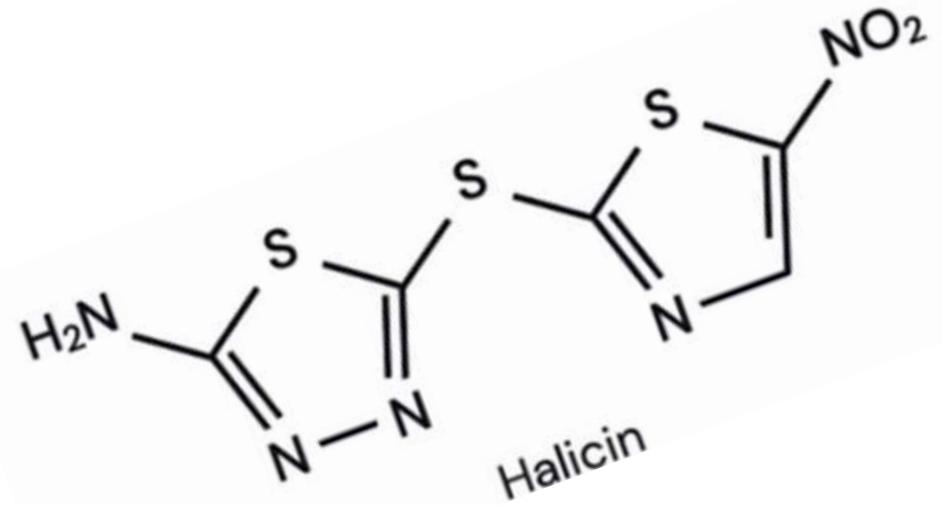
Interesting task to predict is, for example, whether the molecule is a potent drug

- Can do binary classification on whether the drug will inhibit certain bacteria. (E.coli)
- Train on a curated dataset for compounds where response is known.



# FOLLOW-UP STUDY

- Once trained, the model can be applied to any molecule.
  - Execute on a large dataset of known candidate molecules.
  - Select the —top-100 candidates from your GNN model.
  - Have chemists thoroughly investigate those (after some additional filtering).
- Discover a previously overlooked compound that is a highly potent antibiotic!



# SUCCESS STORY!

Cell

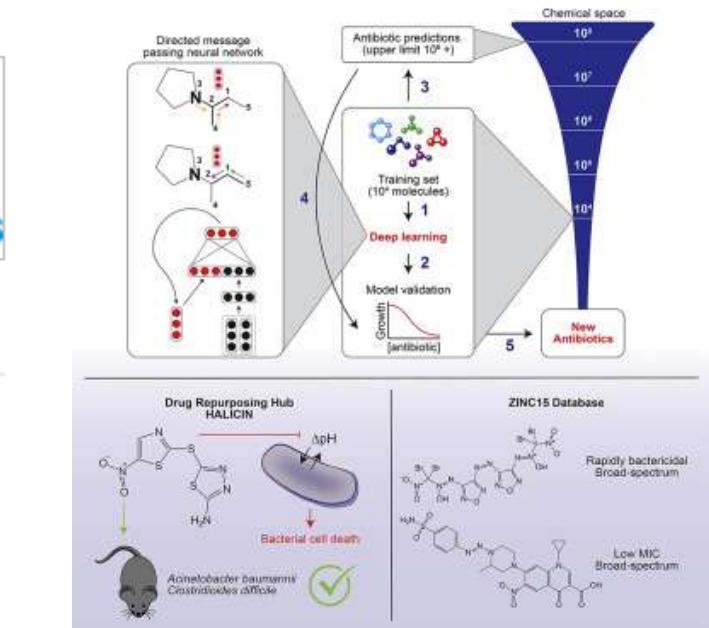
Volume 180, Issue 4, 20 February 2020, Pages 688-702.e13



Article

## A Deep Learning Approach to Antibiotic Discovery

Jonathan M. Stokes<sup>1 2 3</sup>, Kevin Yang<sup>3 4 10</sup>, Kyle Swanson<sup>3 4 10</sup>, Wengong Jin<sup>3 4</sup>,  
Andres Cubillos-Ruiz<sup>1 2 5</sup>, Nina M. Donghia<sup>1 5</sup>, Craig R. MacNair<sup>6</sup>, Shawn French<sup>6</sup>,  
Lindsey A. Carfrae<sup>6</sup>, Zohar Bloom-Ackermann<sup>2 7</sup>, Victoria M. Tran<sup>2</sup>, Anush Chiappino-Pepe<sup>5 7</sup>,  
Ahmed H. Badran<sup>2</sup>, Ian W. Andrews<sup>1 2 5</sup>, Emma J. Chory<sup>1 2</sup>, George M. Church<sup>5 7 8</sup>,  
Eric D. Brown<sup>6</sup>, Tommi S. Jaakkola<sup>3 4</sup>, Regina Barzilay<sup>3 4 9</sup> , James J. Collins<sup>1 2 5 8 9 11</sup>



SUCCESS STORY!

**nature**

Subscribe

NEWS · 20 FEBRUARY 2020

## Powerful antibiotics discovered using AI

Machine learning spots molecules that work even against ‘untreatable’ strains of bacteria.

# SUCCESS STORY!

**FINANCIAL TIMES**

S COMPANIES TECH MARKETS GRAPHICS OPINION WORK & CAREERS LIFE & ARTS HOW TO SPEND IT

**CORONAVIRUS BUSINESS UPDATE**

Get 30 days' complimentary access to our Coronavirus Business Update newsletter



**POW** Artificial intelligence

Robotics  'Death of the office' homeworking claims exaggerated

Anti-social robots harr increase social distanc

**AI**

Artificial intelligence + Add to myFT

## AI discovers antibiotics to treat drug-resistant diseases

Machine learning uncovers potent new drug able to kill 35 powerful bacteria

# SUCCESS STORY!

The screenshot shows the BBC News homepage with a sidebar on the left containing various news snippets. The main content area features a prominent blue banner for 'BBC WORKLIFE' with the text 'Our new guide for getting ahead'. Below this, a large, bold headline reads 'Scientists discover powerful antibiotic using AI'. A timestamp indicates it was published on 21 February 2020. A 'Share' button is visible in the bottom right corner of the article area.

BBC Sign in News Sport Reel Worklife Travel Future

na NEWS

Home Video World UK Business Tech Science Stories Entertainment & Arts

NEWS

Powerful intelligence

Robotics

Machines

bacteria

Artificial

AI

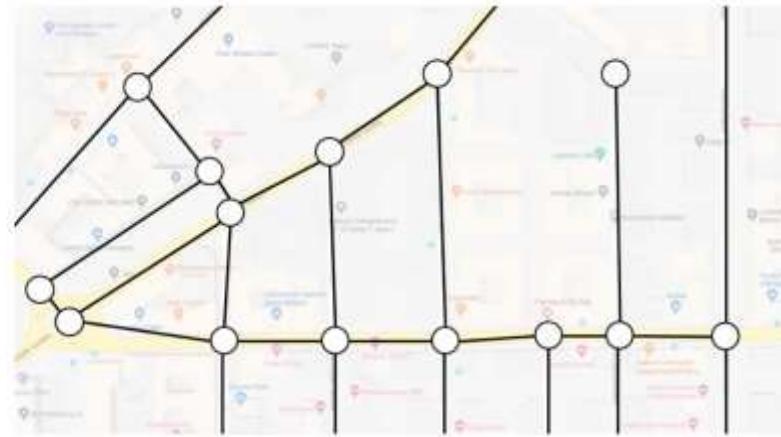
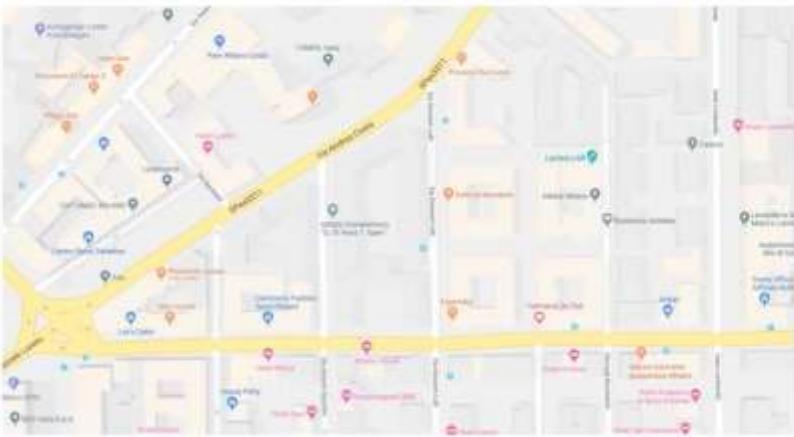
dis

Machine learning uncovers potent new drug able to kill off powerful bacteria

21 February 2020 Share

# TRAFFIC MAPS ARE GRAPHS!

Transportation maps (e.g. the ones found on Google Maps)  
naturally modeled as graphs.



**Nodes** could be **intersections**, and **edges** could be **roads**.

(Relevant **node features**: road length, current speeds, historical speeds)

# DEEPMIND'S ETA PROBLEM!

Partition candidate route into super-segments, sampled proportionally to (est.) traffic density.

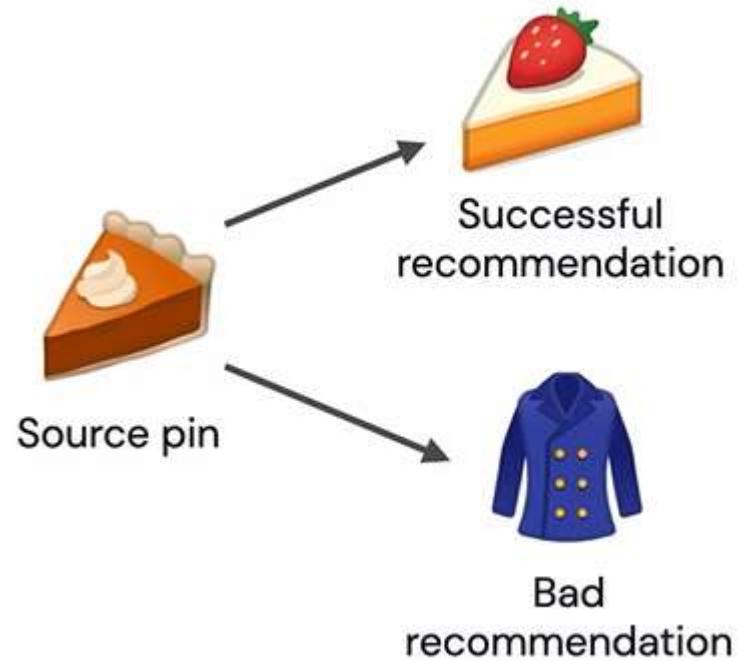
Run GNN on super-segment graph to estimate estimated time of arrival (ETA) (graph regression).

[گوگل-مپ-ترافیک-شبکه-عصبی-گرافی/](https://class.vision/blog/)

# RECOMMENDER SYSTEMS

A common task on **social networks** is **recommendation**.

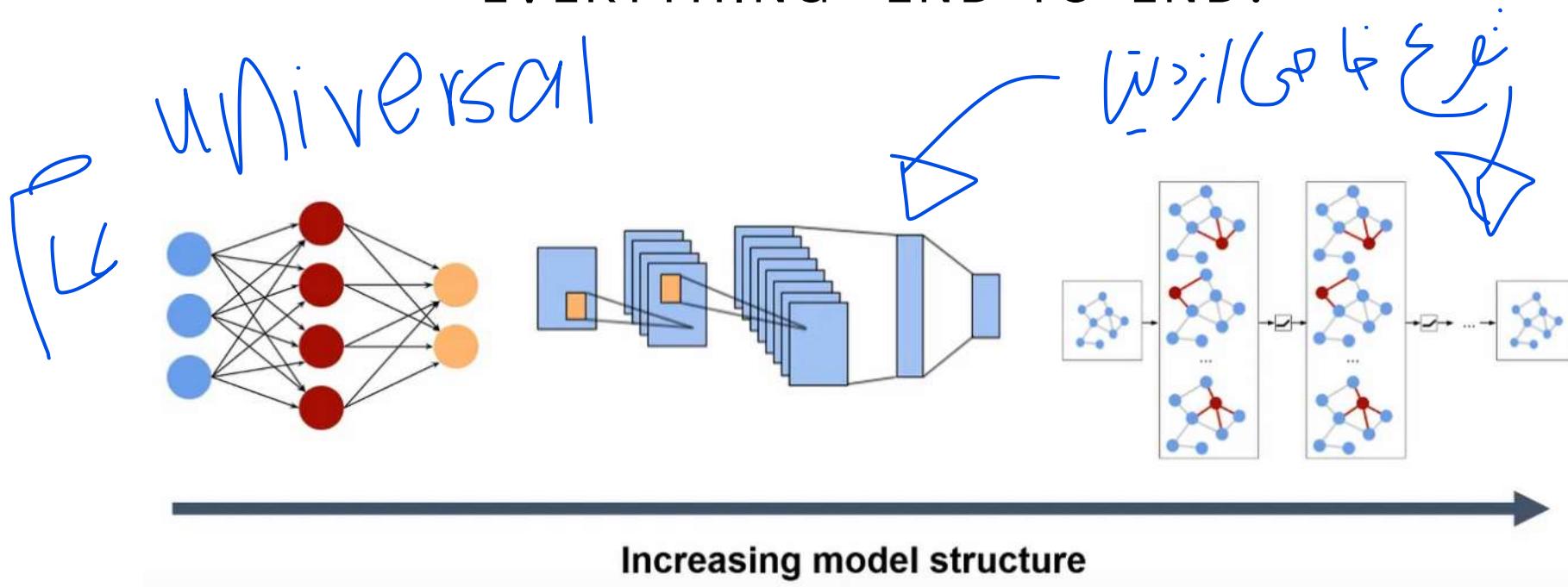
- Based on a user's preferences, recommend new content
- Can leverage existing links as adjacency input to a (link-prediction) GNN!
- Major issue: our methods (so far) assume the graph is processed all-at- once! (one solution is GraphSAGE)



# GRAPH CHALLENGE AND PROBLEMS



# WHY USE GRAPHS? WHY NOT JUST USE MLP OR ATTENTION AND LEARN “EVERYTHING” END-TO-END?



# PROBLEM: GRAPH DATA IS DIFFERENT

Challenge 1: Data size and shape

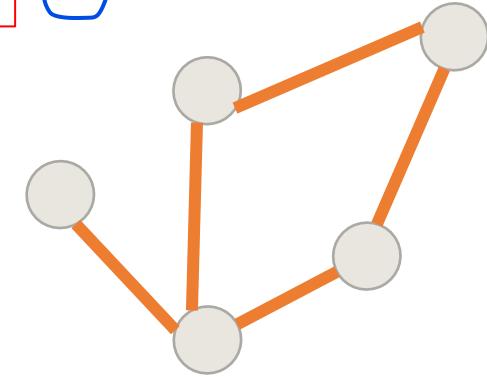
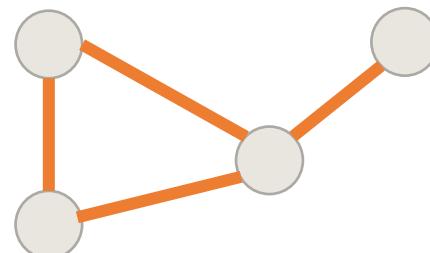
Feature



It should be Size independent

~~DATA~~

GNN ✓



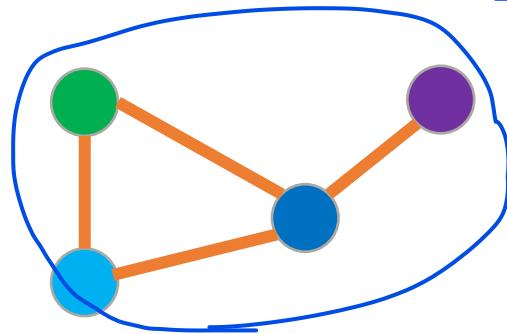
# PROBLEM: GRAPH DATA IS DIFFERENT

Flip

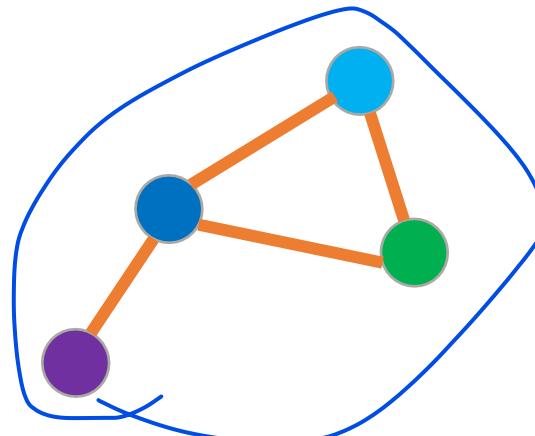
Challenge 2: Isomorphism



It should be **Permutation invariance**



UniqR

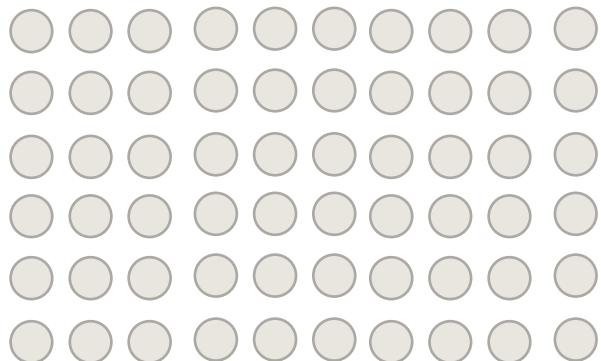
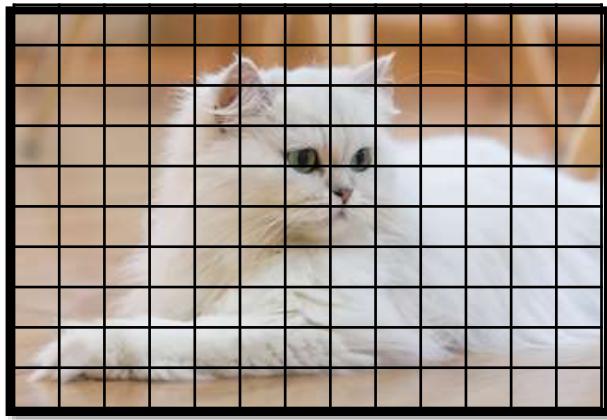


We cannot feed adjacency matrix to MLP

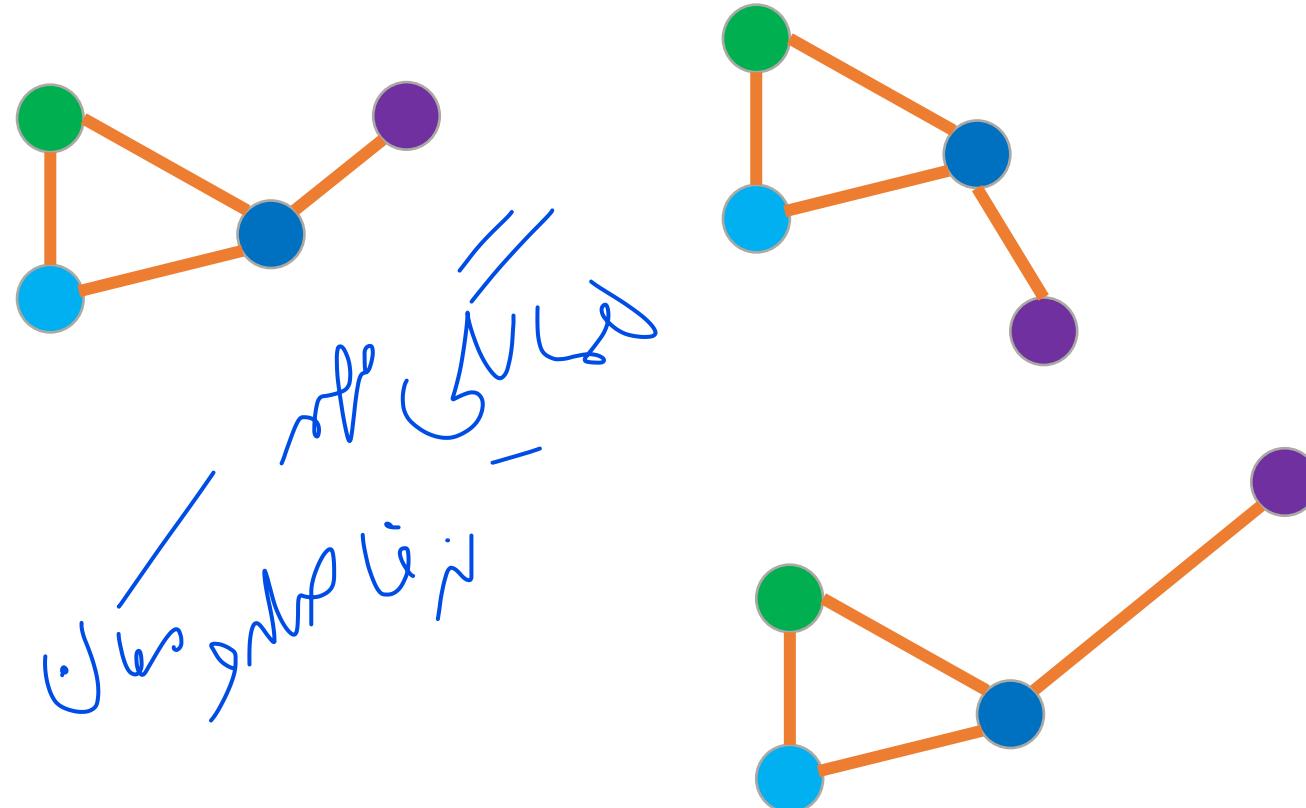
# PROBLEM: GRAPH DATA IS DIFFERENT

فهذا  
المعنى  
أي

Challenge 3: Grid structure



It should be in **Non-Euclidean space**



# OTHER CHALLENGES WITH GRAPH CONVOLUTIONS

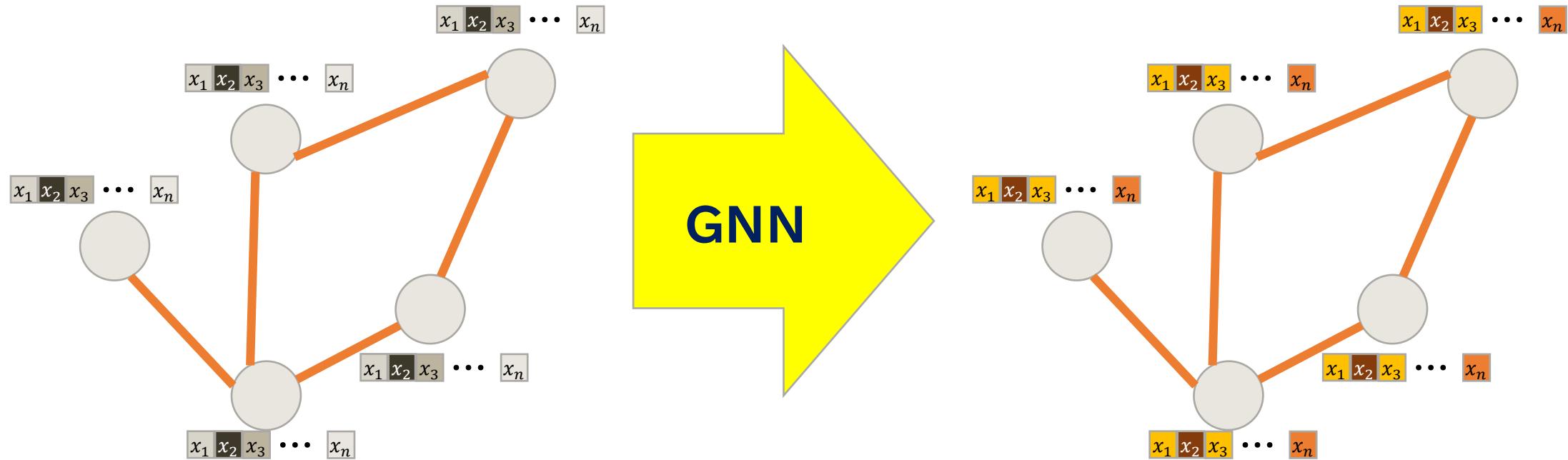
Desirable properties for a graph convolutional layer:

- Computational and storage efficiency ( $\sim O(V + E)$ )
- Fixed number of parameters (independent of input size)
- Localisation (acts on a local neighbourhood of a node)
- Specifying different importances to different neighbours
- Applicability to inductive problems.

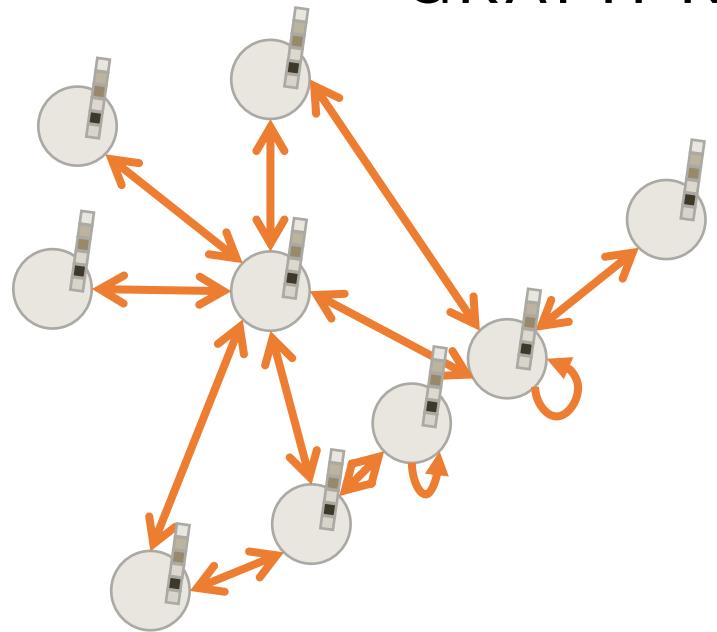
الآن نحن في درجات الحرارة في التعلم المترافق  
inductive GNN

# LEARNING IN GRAPH

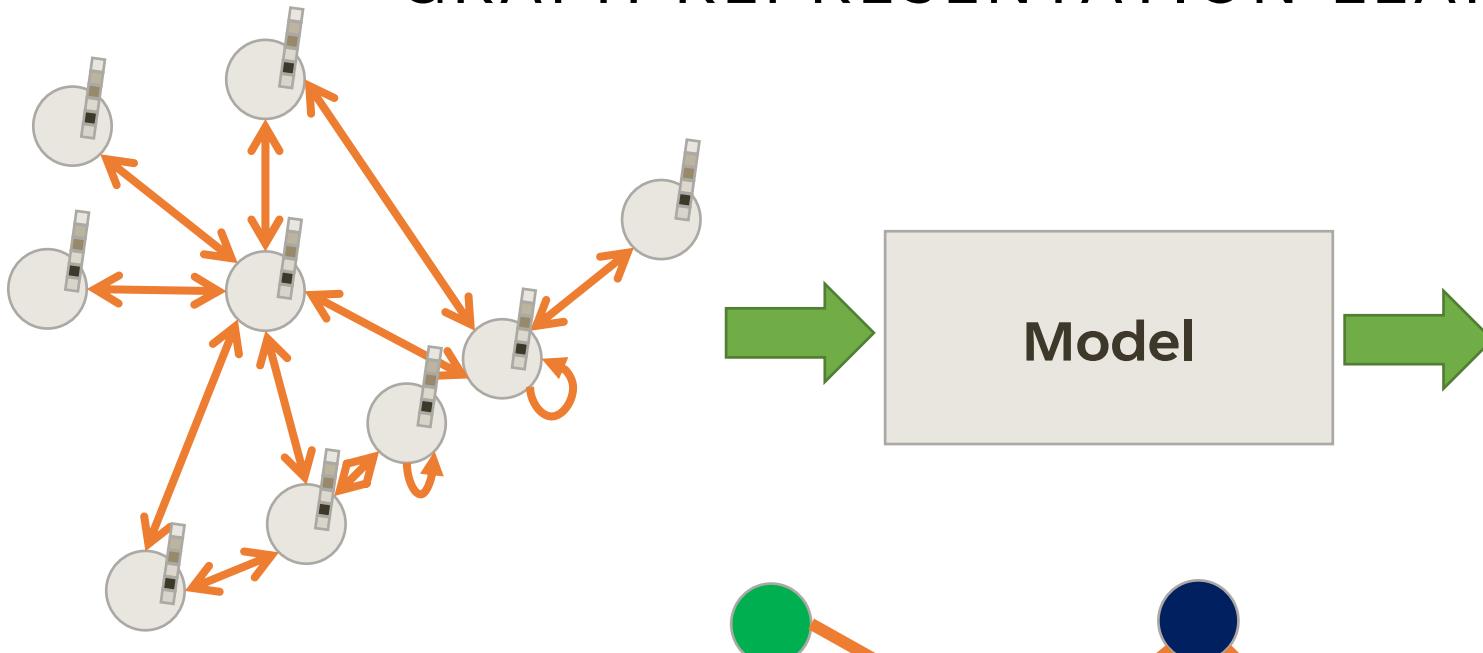
# REPRESENTATION LEARNING



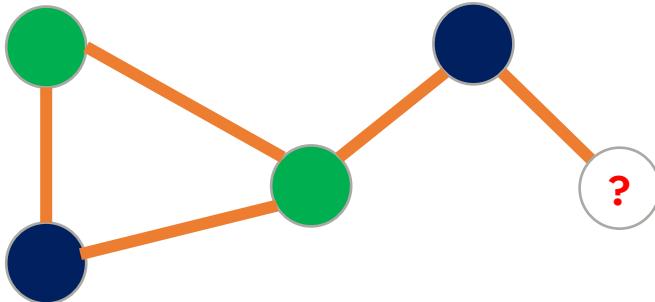
# LEARNING IN GRAPH REPRESENTATION LEARNING



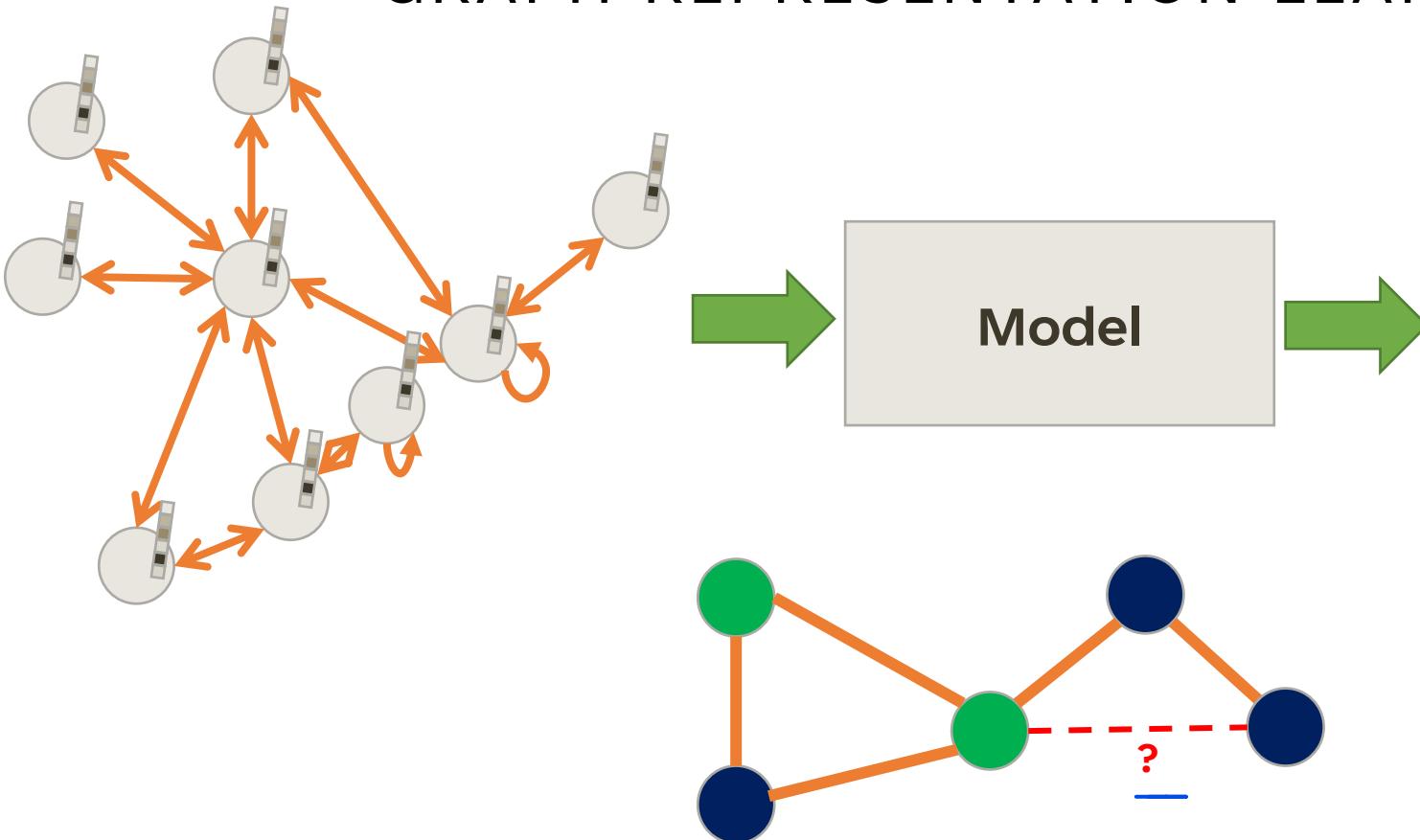
# LEARNING IN GRAPH REPRESENTATION LEARNING



- ✓ Node Prediction  
(Node-level Prediction)

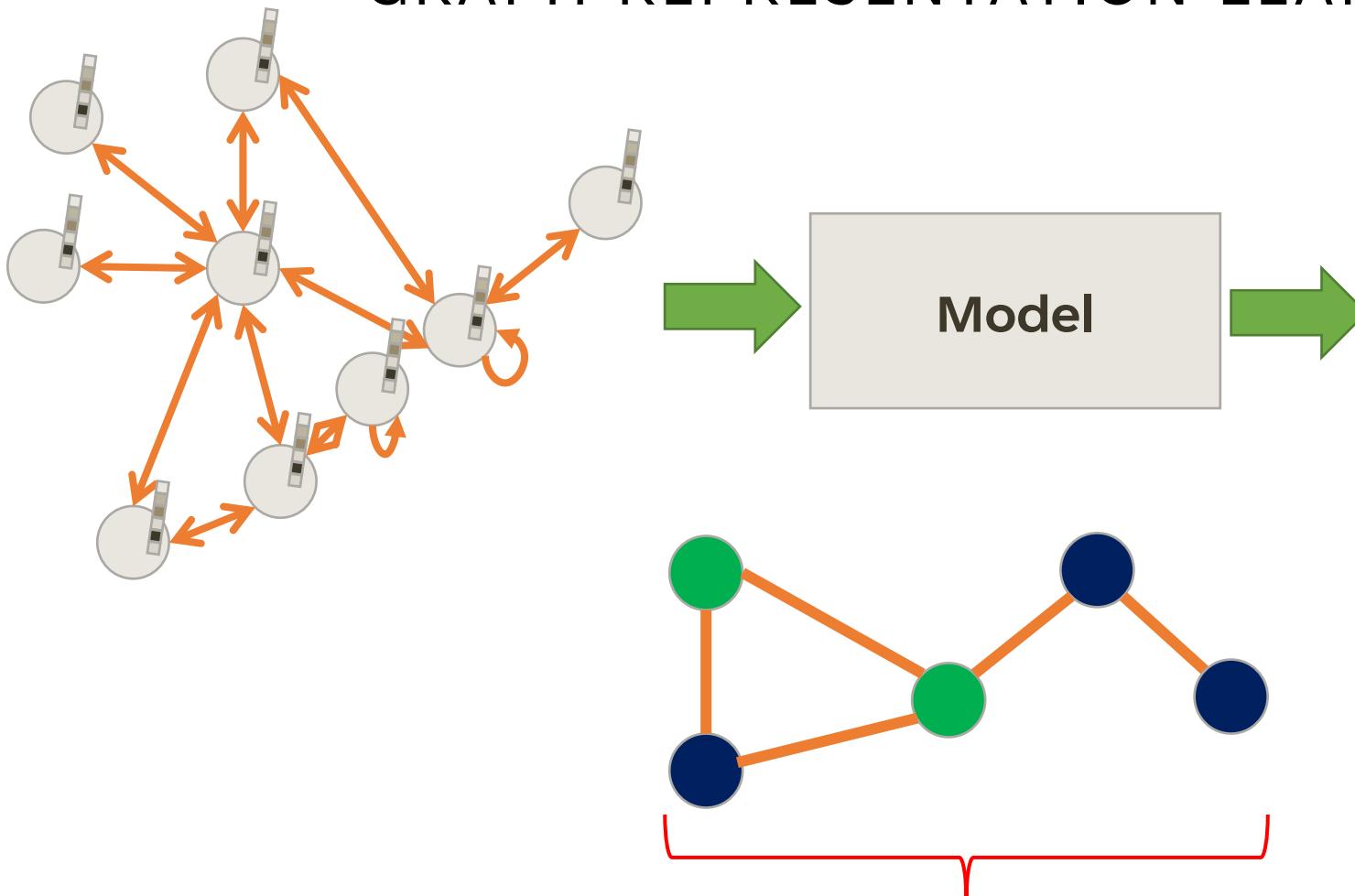


# LEARNING IN GRAPH REPRESENTATION LEARNING



- ✓ Node Prediction  
(Node-level Prediction)
- ✓ Link Prediction  
(Edge-level prediction)

# LEARNING IN GRAPH REPRESENTATION LEARNING



- ✓ Node Prediction  
(Node-level Prediction)
- ✓ Link Prediction  
(Edge-level prediction)
- ✓ Graph representation  
(Graph-level prediction)

کل ترا فکر

# WHAT TYPES OF PROBLEMS CAN GNNS SOLVE?

## Unsupervised

- Node, Edge, or Graph clustering
  - Use embeddings to find “similar” nodes, edges, or graphs
- Link Prediction
- Graph Generation

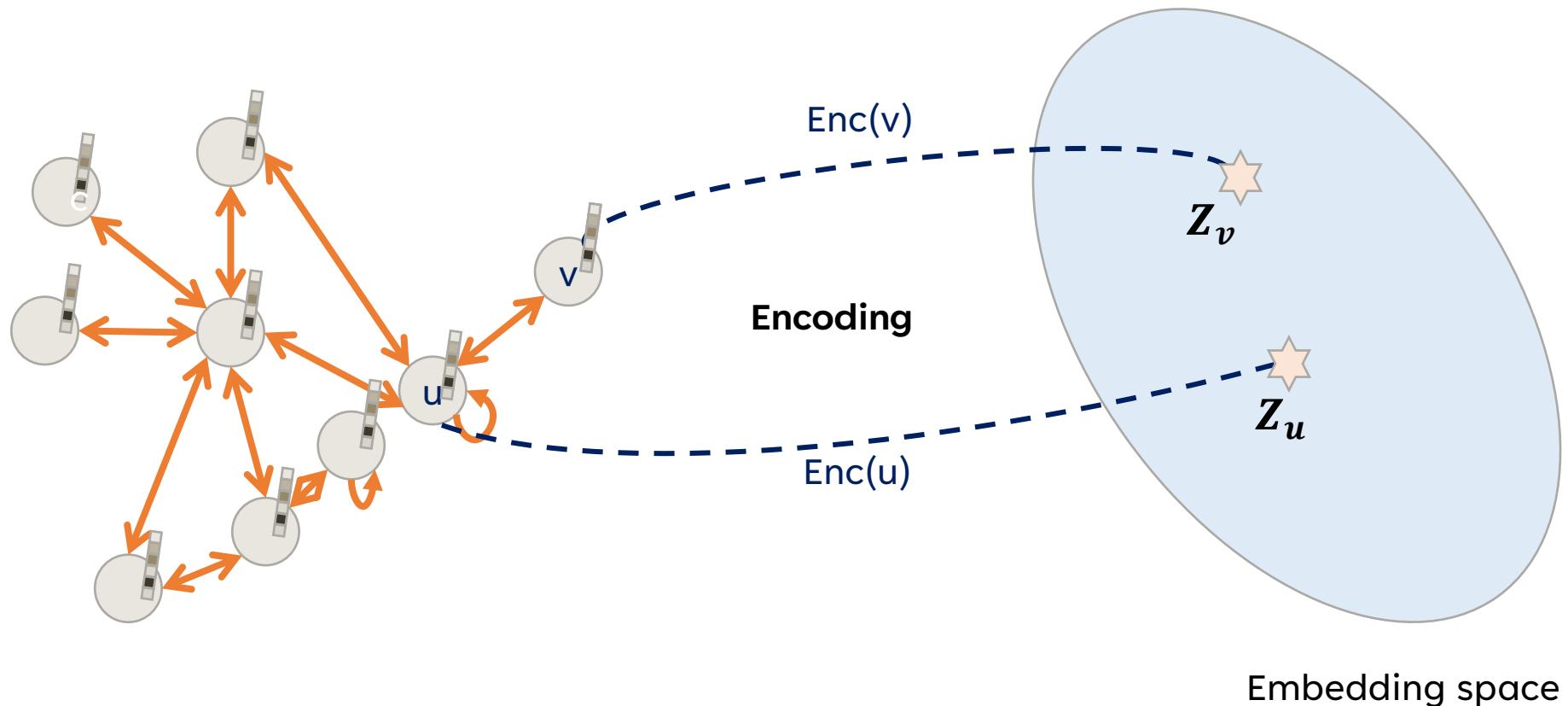
## Supervised

- Node, Edge, or Graph classification / regression
  - Use embeddings to predict based on known data

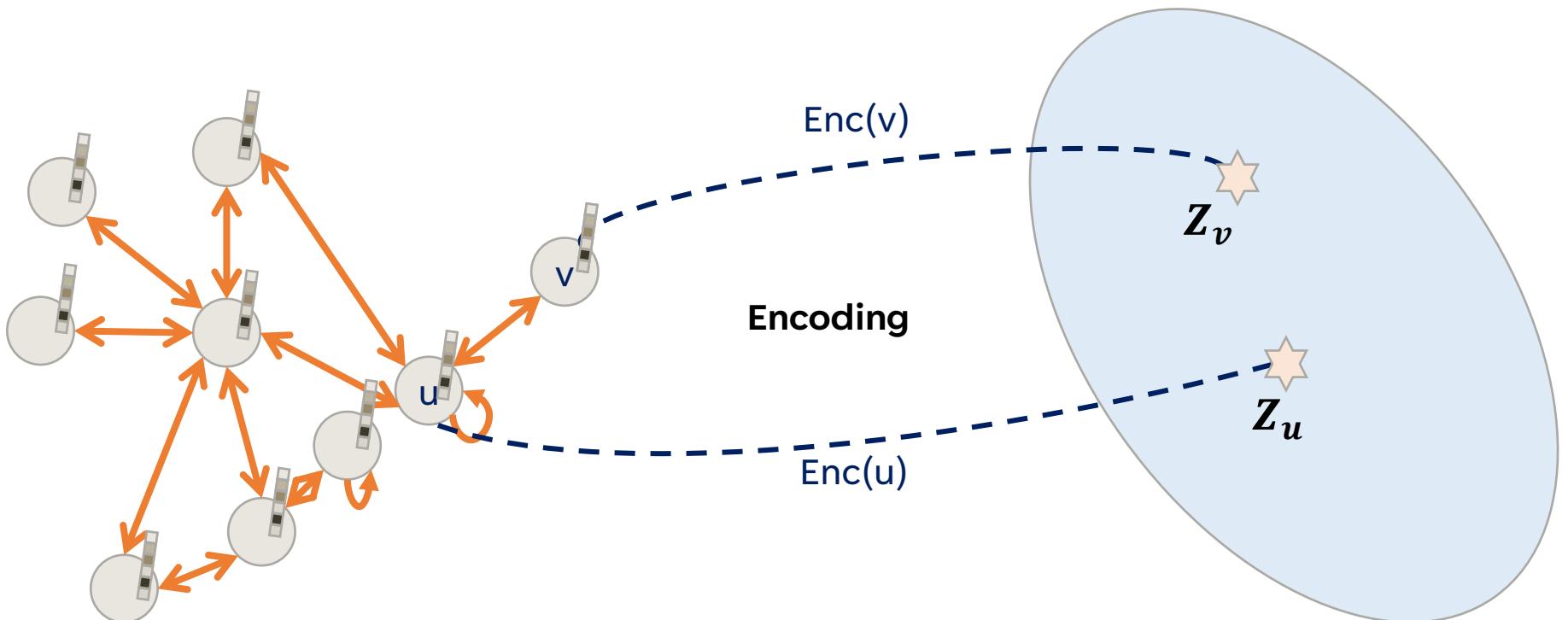
“A Fair Comparison of Graph Neural Networks for Graph Classification”, ICLR 2020

“Revisiting Graph Neural Networks for Link Prediction” (2020)

# LEARNING IN GRAPH REPRESENTATION LEARNING



# LEARNING IN GRAPH REPRESENTATION LEARNING

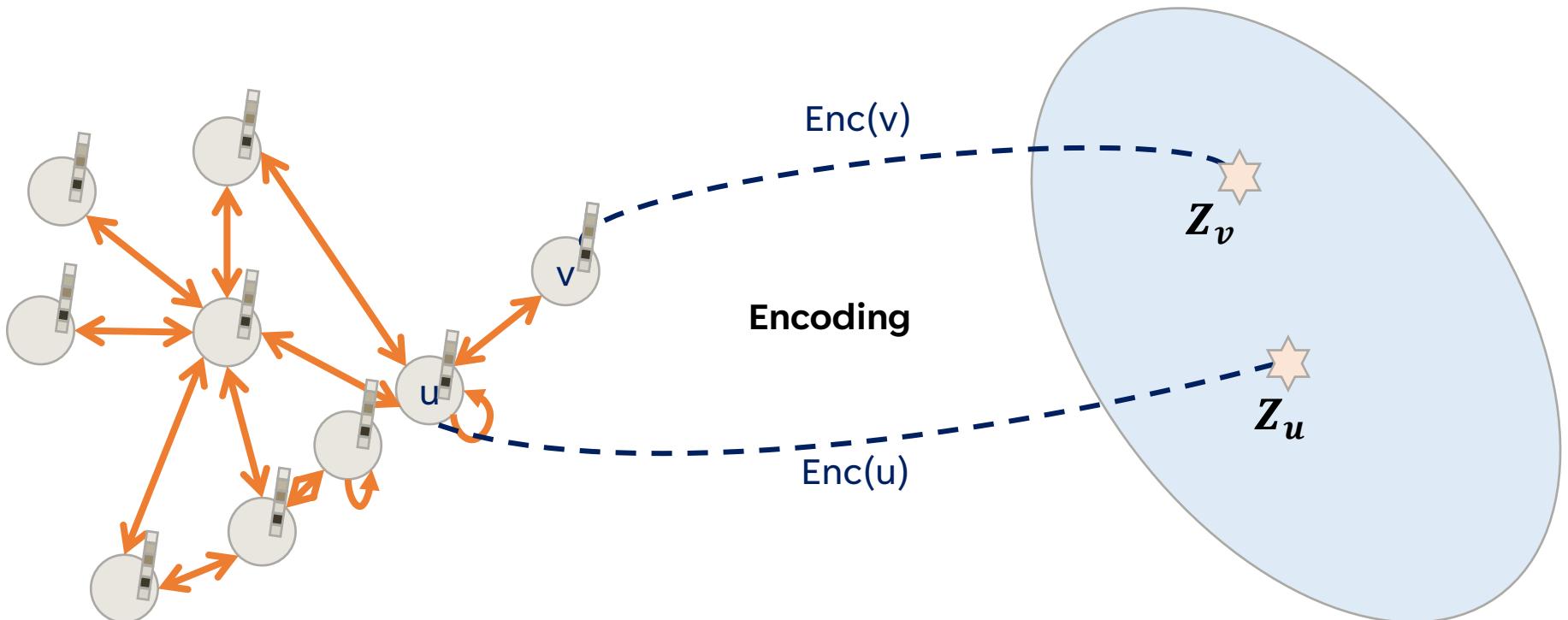


**Goal:**  $\text{Similarity}(u, v) \approx \text{Similarity}(Z_u, Z_v)$

$$S_G(u, v) \approx S_V(Z_u, Z_v)$$

Embedding space

# LEARNING IN GRAPH REPRESENTATION LEARNING



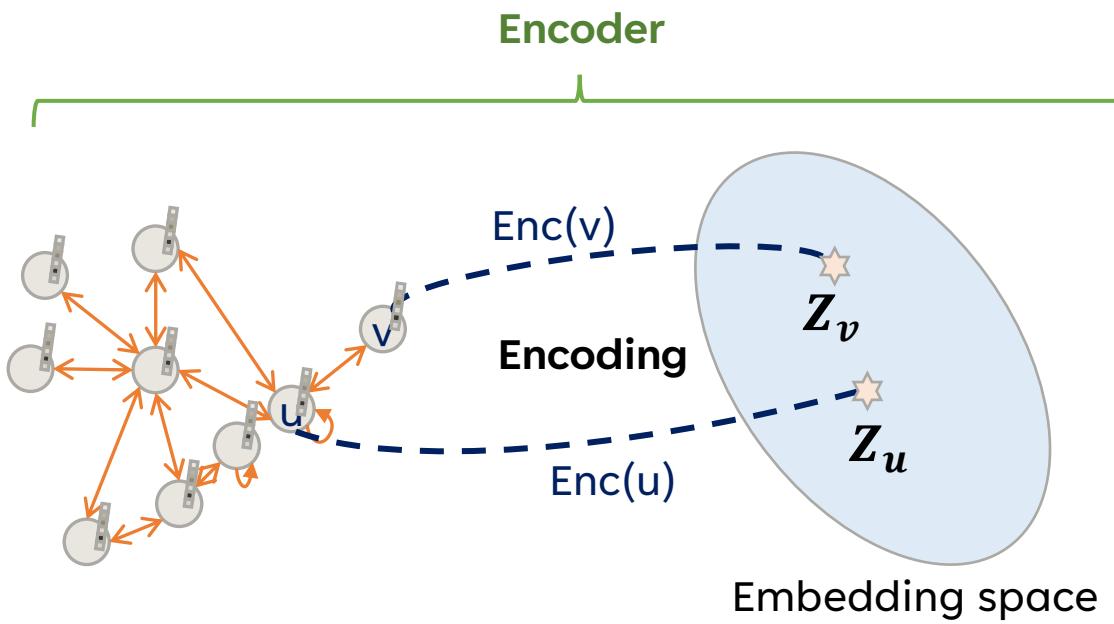
**Goal:**  $S_G(u, v) \approx S_V(Z_u, Z_v)$



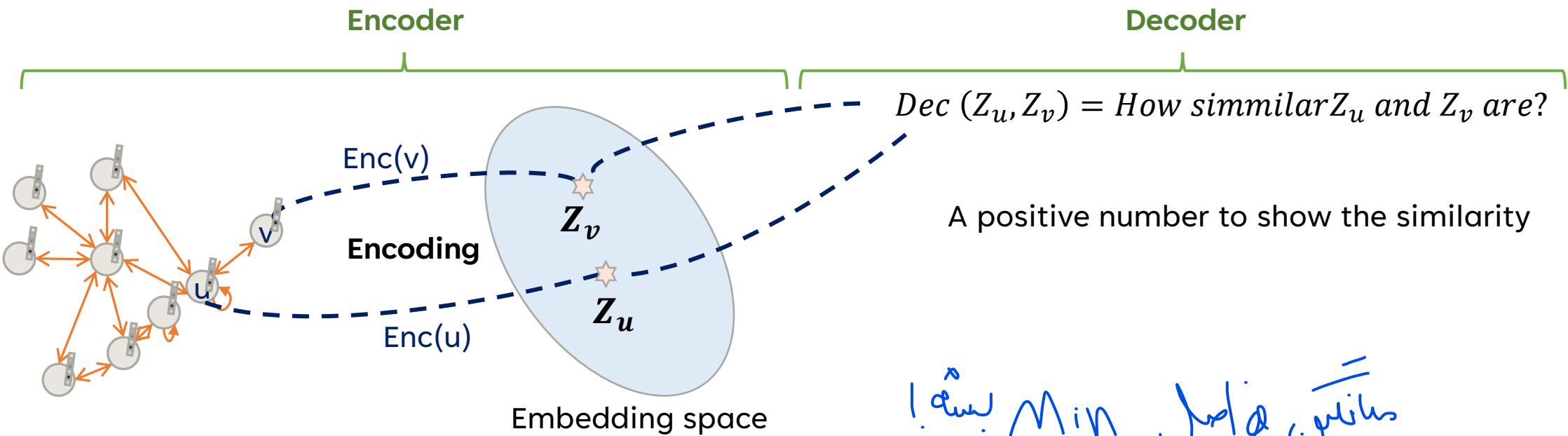
- How to perform **Encoding**?
- What is the **meaning of similarity** ?

Embedding space

# HOW TO ENCODE AND DECODE?



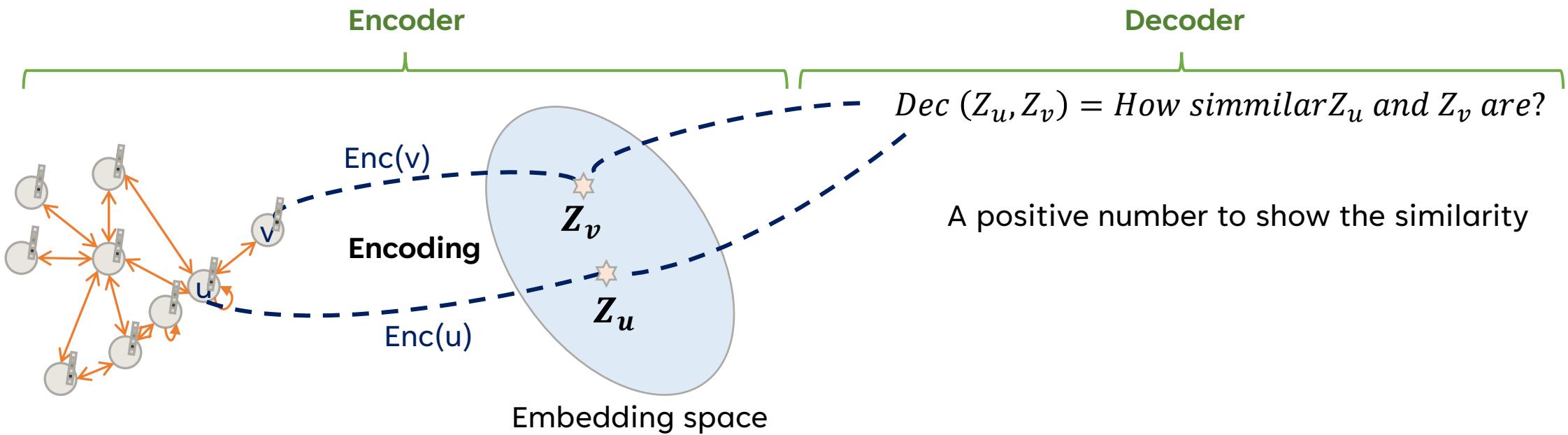
# HOW TO ENCODE AND DECODE?



$$S_G(u, v) \approx S_V(Z_u, Z_v)$$

$$\ell(z_v, z_u) = \sum_{(v,u) \in V} \left\| S_E(z_v, z_u) - S_G(v, u) \right\|_2^2$$

# HOW TO ENCODE AND DECODE?



Matrix factorization

Look-up table

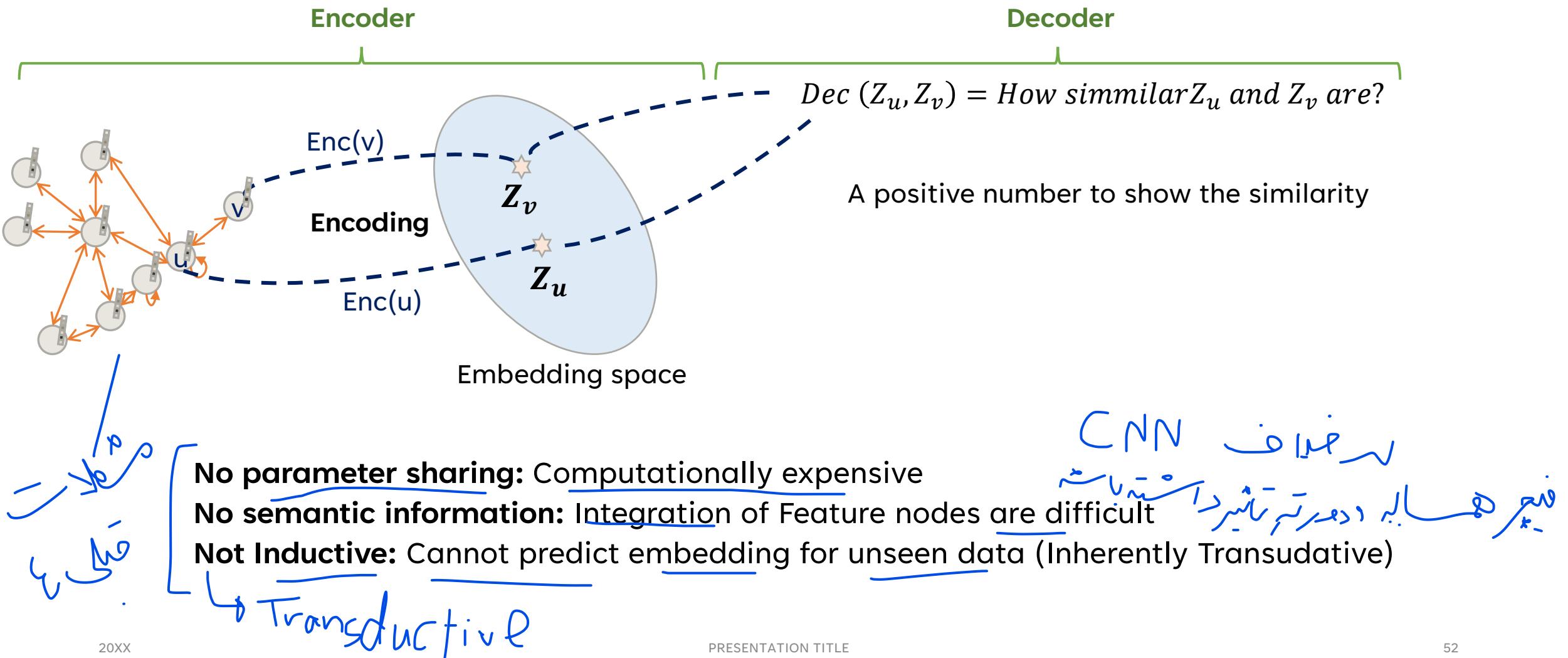
Random Walk

Inner product  $Z_u^T Z_v$

Inner product  $Z_u^T Z_v$

Decode statistic of random walk

# DRAWBACKS



# DEEP VS SHALLOW

Older methods (“shallow”, non-neural network models)  
Deepwalk, node2vec

Generally fallen out of favor with researchers because:

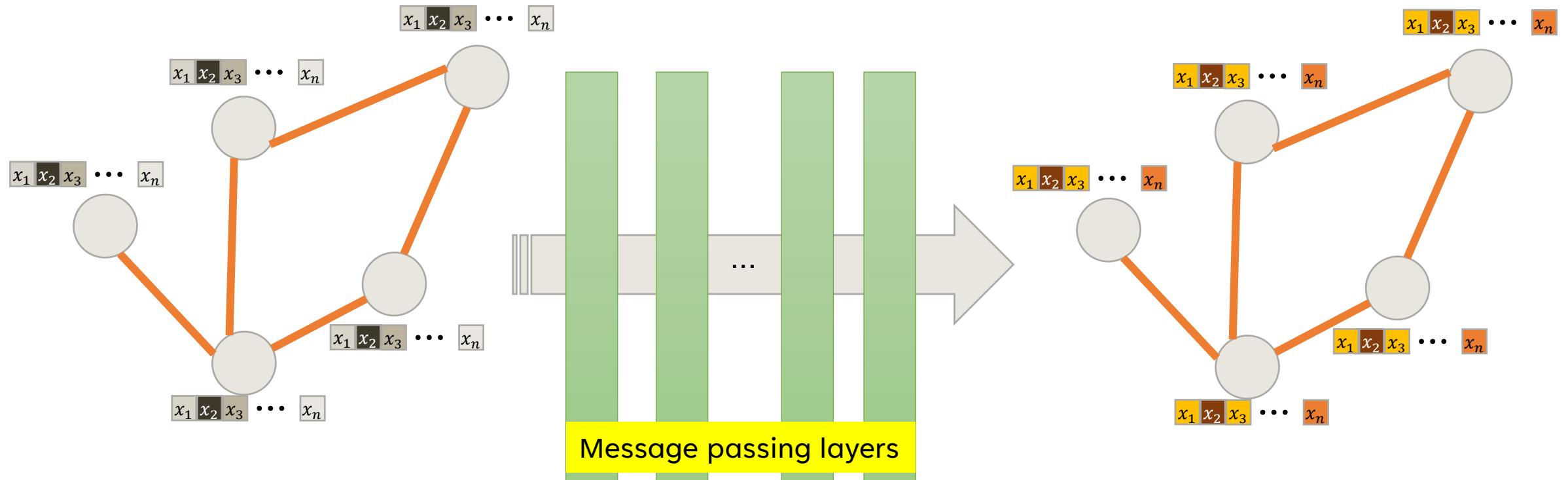
- **No parameter sharing** (bad scaling, overfitting)
- Transductive (only work with nodes present during training)

GNNs solve these problems, they can

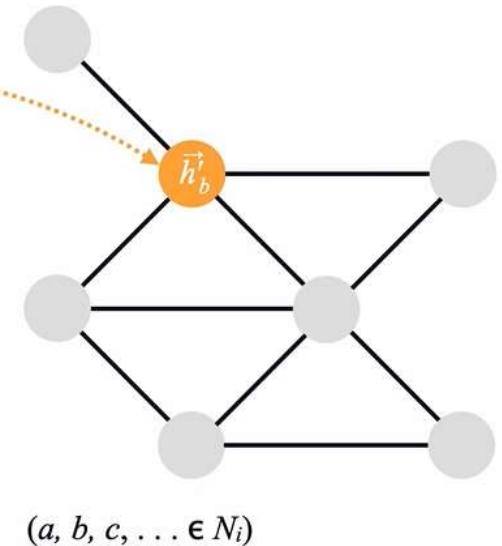
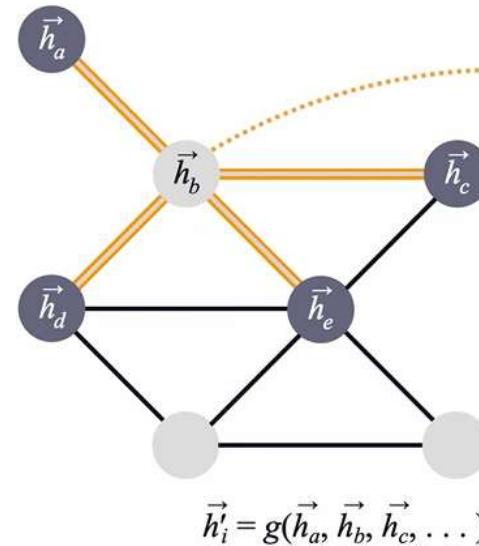
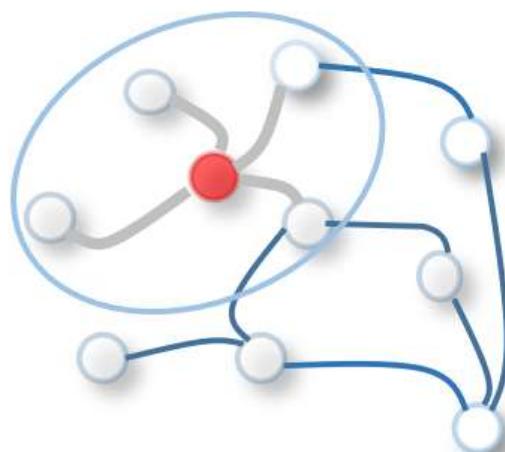
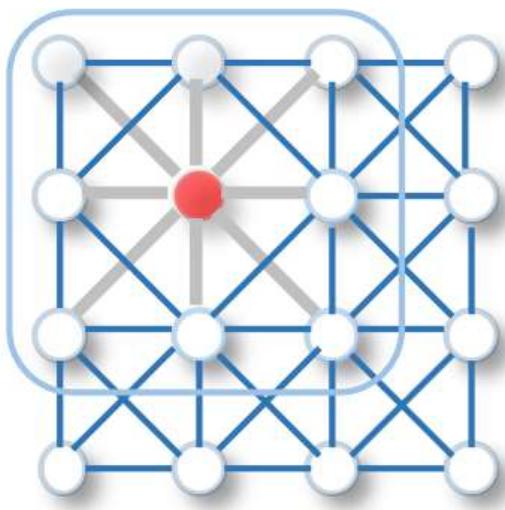
- ✓ Share parameters
- ✓ Can generalize to inductive tasks

inductive and transductive!

# REPRESENTATION LEARNING



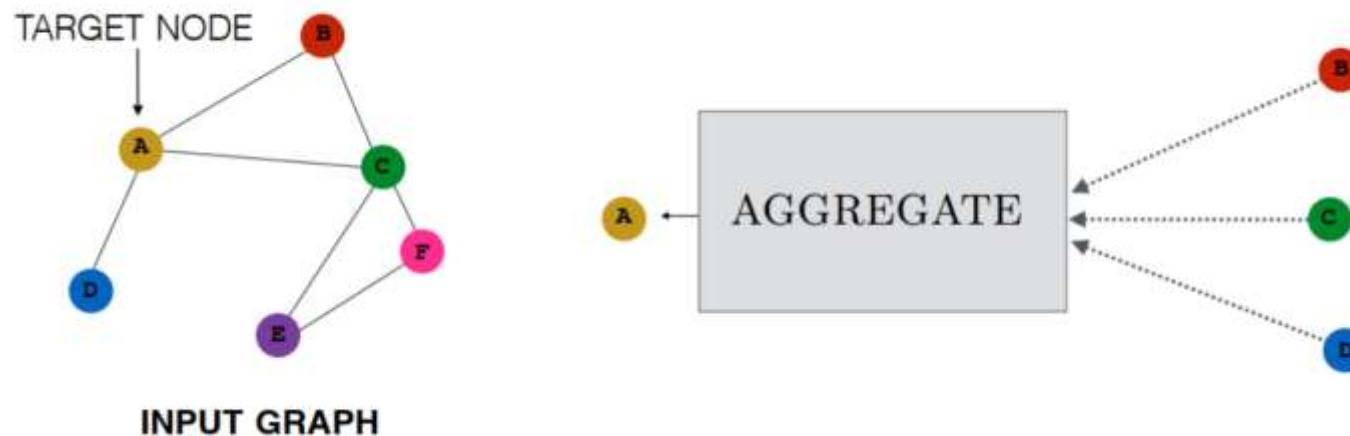
# GRAPH CONVOLUTIONAL NETWORK



[A Comprehensive Survey on Graph Neural Networks](#)

# UNDERSTANDING GRAPH NEURAL NETWORKS

GNNs were originally based on 2-step message passing



## 1. Aggregate :

Pass information (the “message”) from a target node’s neighbors to the target node

## 2. Update:

Update each node’s features based on “message” to form an embedded representation

# MESSAGE PASSING

$$h_u = \text{UPDATE}(h_u, \text{AGGREGATE}(\{h_v, \forall v \in N(u)\}))$$

$h$  = node features / embeddings

Aggregate function operates over sets, must be permutation invariant or permutation equivariant

# MESSAGE PASSING

$$h_u = \text{UPDATE}(h_u, \text{AGGREGATE}(\{h_v, \forall v \in N(u)\}))$$

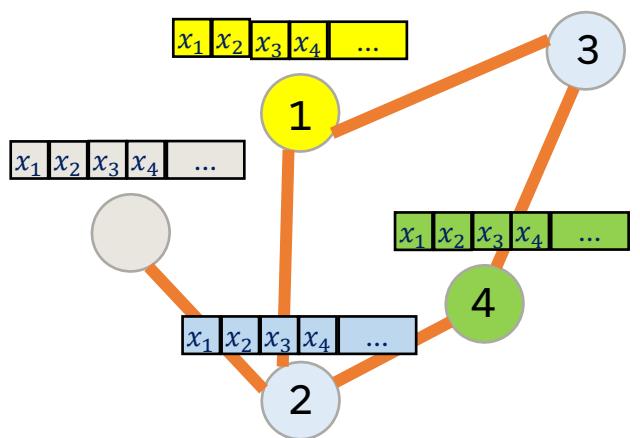
$h$  = node features / embeddings

Aggregate function operates over sets, must be permutation invariant or permutation equivariant

$$h_u = \sigma \left( W_{self} h_u + W_{neigh} \sum_{v \in N(u)} h_v \right)$$

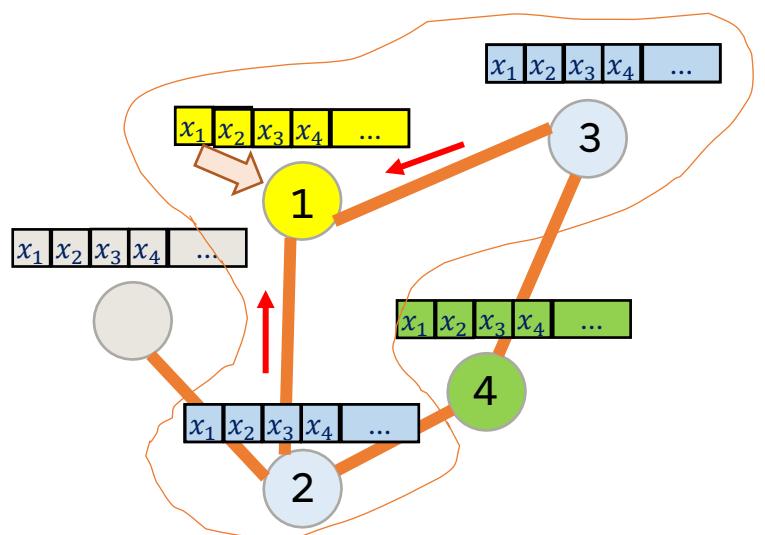
# MESSAGE PASSING

$$h_u = \sigma \left( W_{self} h_u + W_{neigh} \sum_{v \in N(u)} h_v \right)$$



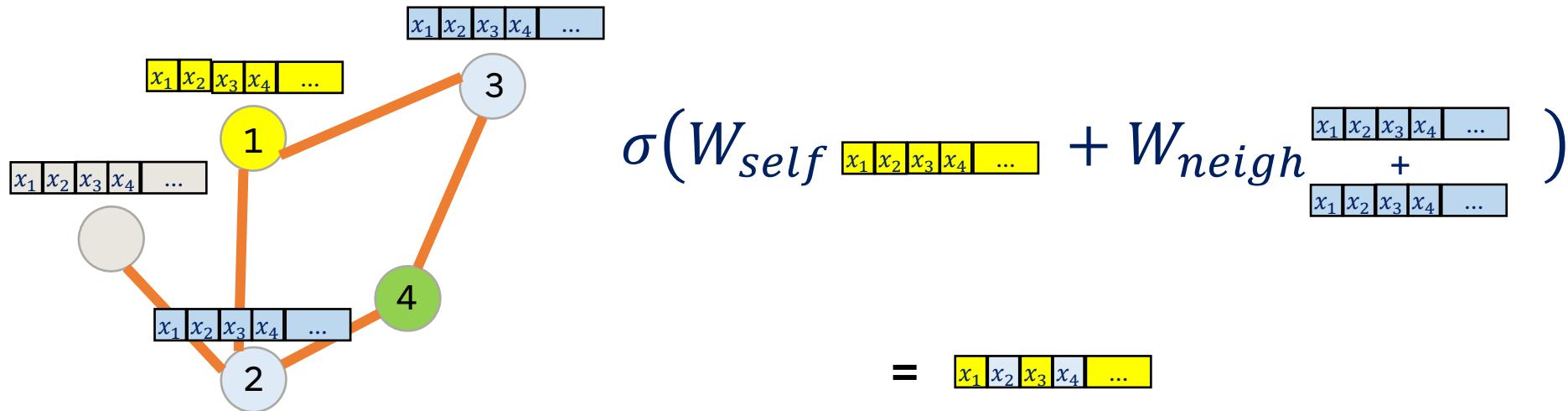
# MESSAGE PASSING

$$h_u = \sigma \left( W_{self} h_u + W_{neigh} \sum_{v \in N(u)} h_v \right)$$



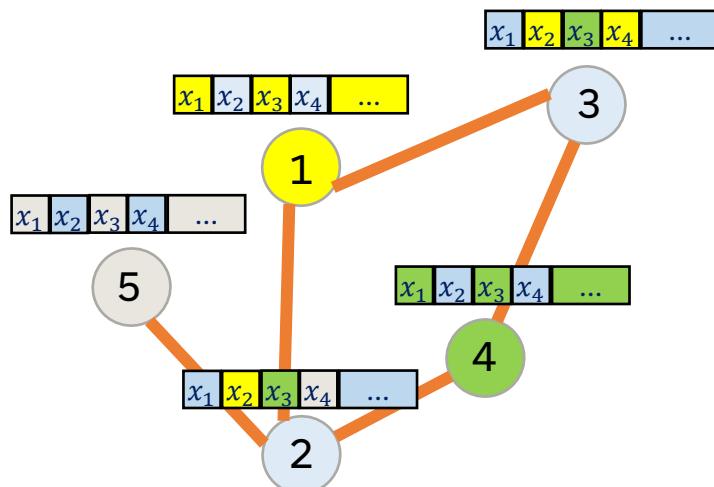
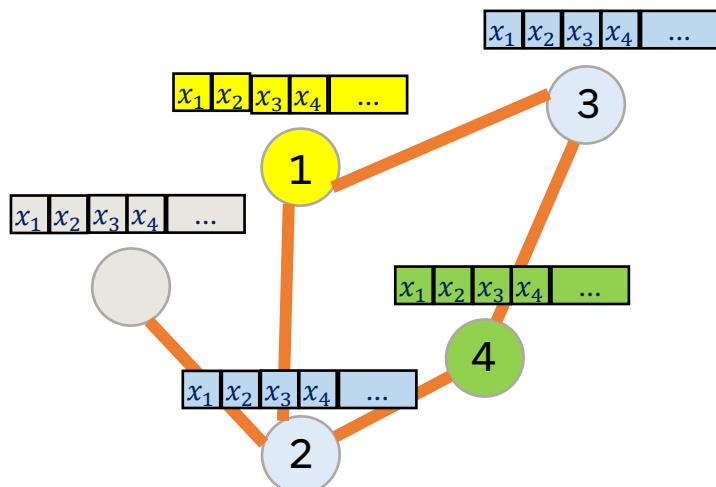
# MESSAGE PASSING

$$h_u = \sigma \left( W_{self} h_u + W_{neigh} \sum_{v \in N(u)} h_v \right)$$



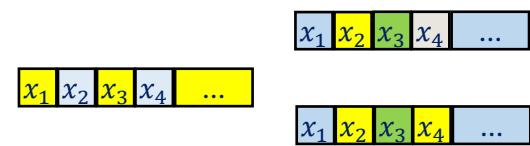
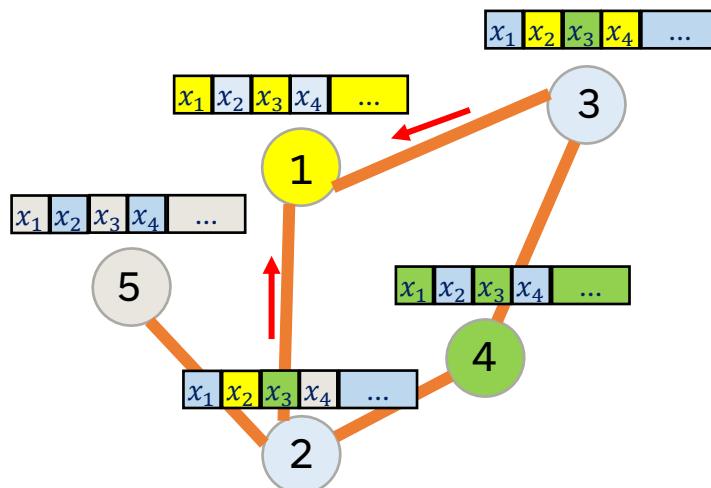
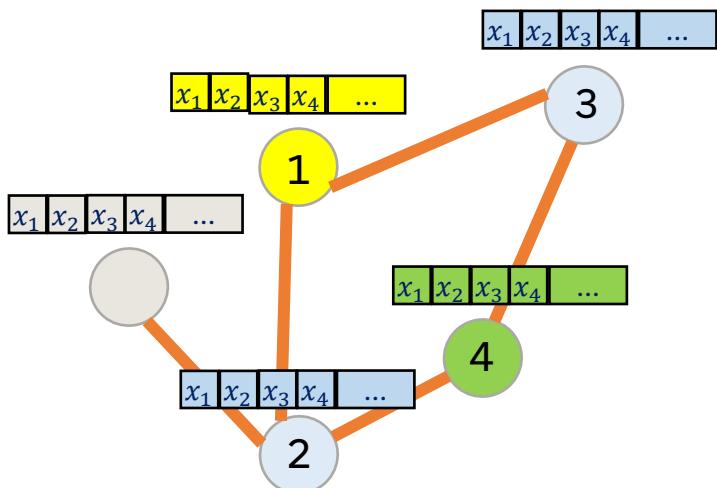
# MESSAGE PASSING

$$h_u = \sigma \left( W_{self} h_u + W_{neigh} \sum_{v \in N(u)} h_v \right)$$



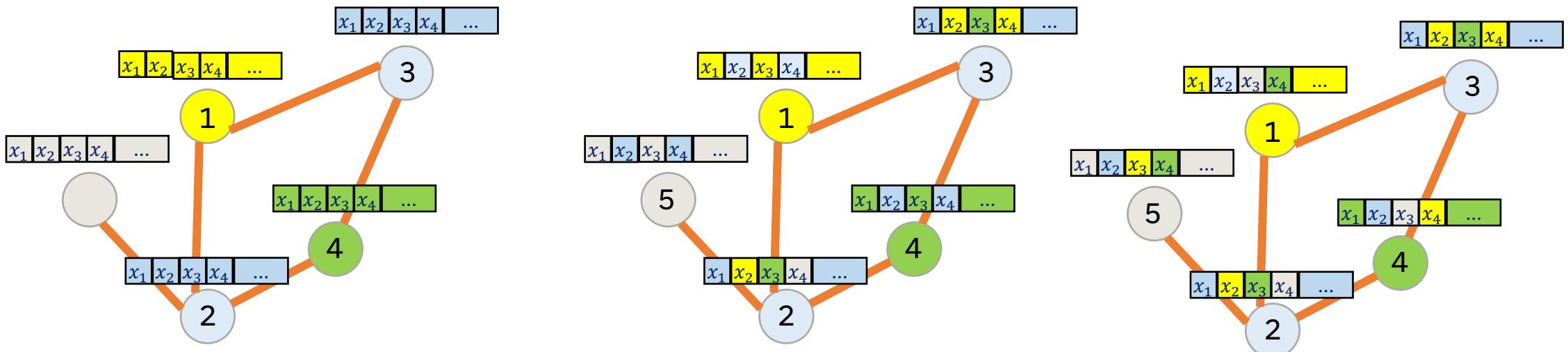
# MESSAGE PASSING

$$h_u = \sigma \left( W_{self} h_u + W_{neigh} \sum_{v \in N(u)} h_v \right)$$



# MESSAGE PASSING

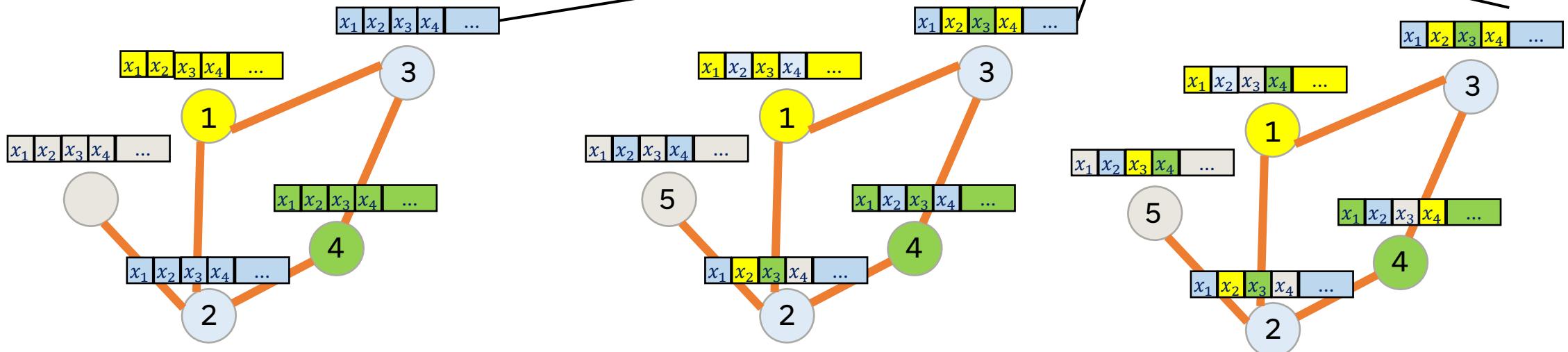
$$h_u = \sigma \left( W_{self} h_u + W_{neigh} \sum_{v \in N(u)} h_v \right)$$



# MESSAGE PASSING

$$h_u^{(k+1)} = \sigma \left( W_{self}^{(k+1)} h_u^k + W_{neigh}^{(k+1)} \sum_{v \in N(u)} h_v^k \right)$$

The dimensions can be different  
 $\text{Len}(h_u^k) \neq \text{len}(h_u^{k+1})$



- ✓ The **local feature aggregation** can be compared to **learnable CNN kernels**:  
<https://distill.pub/2021/gnn-intro/>

# MESSAGE PASSING

$$\mathbf{h}_u^{(k+1)} = \text{UPDATE}^{(k)} \left( \mathbf{h}_u^{(k)}, \text{AGGREGATE}^{(k)}(\{\mathbf{h}_v^{(k)}, \forall v \in \mathcal{N}(u)\}) \right)$$

$$h_u^{(k+1)} = \sigma \left( W_{\text{self}}^{(k+1)} h_u^k + W_{\text{neigh}}^{(k+1)} \sum_{v \in \mathcal{N}(u)} h_v^{(k)} \right)$$

- $\mathbf{h}$  = node features / embeddings
- $k$  = number of hops

Each node's updated value becomes a weighting of its previous value + a weighting of its neighbor's values

The choice to sum over neighboring nodes isn't the only valid choice, other choices include mean, max, concatenation, etc.

# MESSAGE PASSING

$$\mathbf{h}_u^{(k+1)} = \text{UPDATE}^{(k)} \left( \mathbf{h}_u^{(k)}, \text{AGGREGATE}^{(k)}(\{\mathbf{h}_v^{(k)}, \forall v \in \mathcal{N}(u)\}) \right)$$

$$h_u^{(k+1)} = \sigma \left( W_{\text{self}}^{(k+1)} h_u^k + W_{\text{neigh}}^{(k+1)} \sum_{v \in \mathcal{N}(u)} h_v^{(k)} \right)$$

- ❑ Collapse  $\mathbf{W}_{\text{self}}$  and  $\mathbf{W}_{\text{neigh}}$  into  $\mathbf{W}$  by adding self-loops to the adjacency matrix  $\mathbf{A}$

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

This method reduces message passing to relatively simple matrix multiplication

# THE MEAN-POOLING UPDATE RULE

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

❑ **Problem:** Multiplication by  $A+I$  may increase the scale of the output features.

✓ **Solution:** We need to normalize appropriately:

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

We arrive at the mean-pooling update rule:

$$h^{(k+1)} = \sigma \sum_{j \in N_i} \frac{1}{|N_i|} Wh_j^k$$

which is simple but versatile (common for inductive problems!).

# GCN GRAPH CONVOLUTIONAL NETWORK

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

“Original” GNN

(Merkwirth, 2005 + Scarselli et al., 2009)

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

GCN

(Kipf + Welling, 2016)

$$\tilde{\mathbf{A}} = (\mathbf{D} + \mathbf{I})^{-\frac{1}{2}}(\mathbf{I} + \mathbf{A})(\mathbf{D} + \mathbf{I})^{-\frac{1}{2}}$$

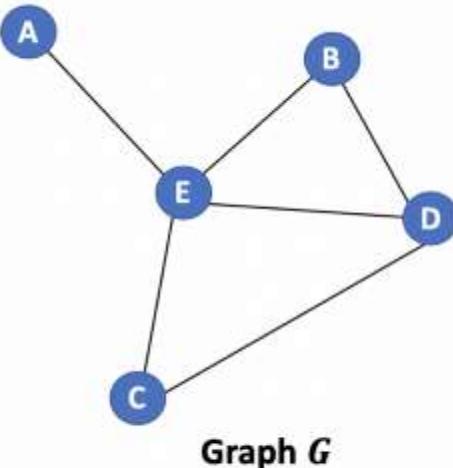
Normalizes by # of nodes in neighborhood

Node-wise, this can be written as follows:

$$h^{(k+1)} = \sigma\left(\sum_{j \in N_i} \frac{1}{\sqrt{|N_i||N_j|}} Wh_j^k\right)$$

Most commonly  
cited GNN paper

# INTUITION AND THE MATH'S BEHIND



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

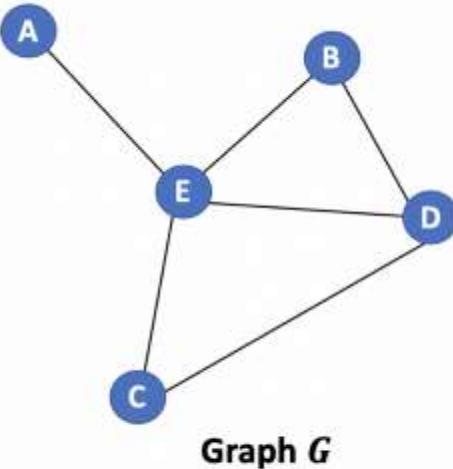
Adjacency matrix  $A$

	A	B	C
A	-1.1	3.2	4.2
B	0.4	5.1	-1.2
C	1.2	1.3	2.1
D	1.4	-1.2	2.5
E	1.4	2.5	4.5

Feature vector  $X$

<https://www.topbots.com/graph-convolutional-networks/>

# INTUITION AND THE MATH'S BEHIND



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

Adjacency matrix  $A$

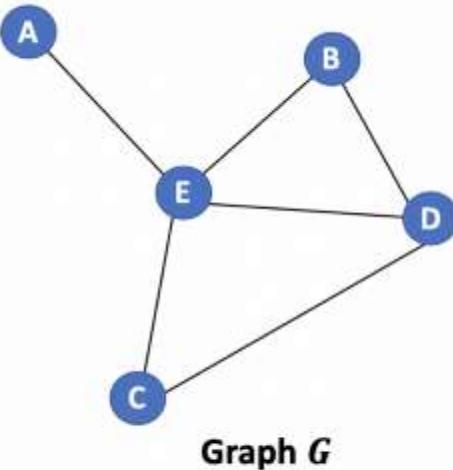
	A	B	C
A	-1.1	3.2	4.2
B	0.4	5.1	-1.2
C	1.2	1.3	2.1
D	1.4	-1.2	2.5
E	1.4	2.5	4.5

Feature vector  $X$

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

$$h^{(k+1)} = \sigma \sum_{j \in N_i} Wh_j^k$$

# INTUITION AND THE MATH'S BEHIND



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

Adjacency matrix  $A$

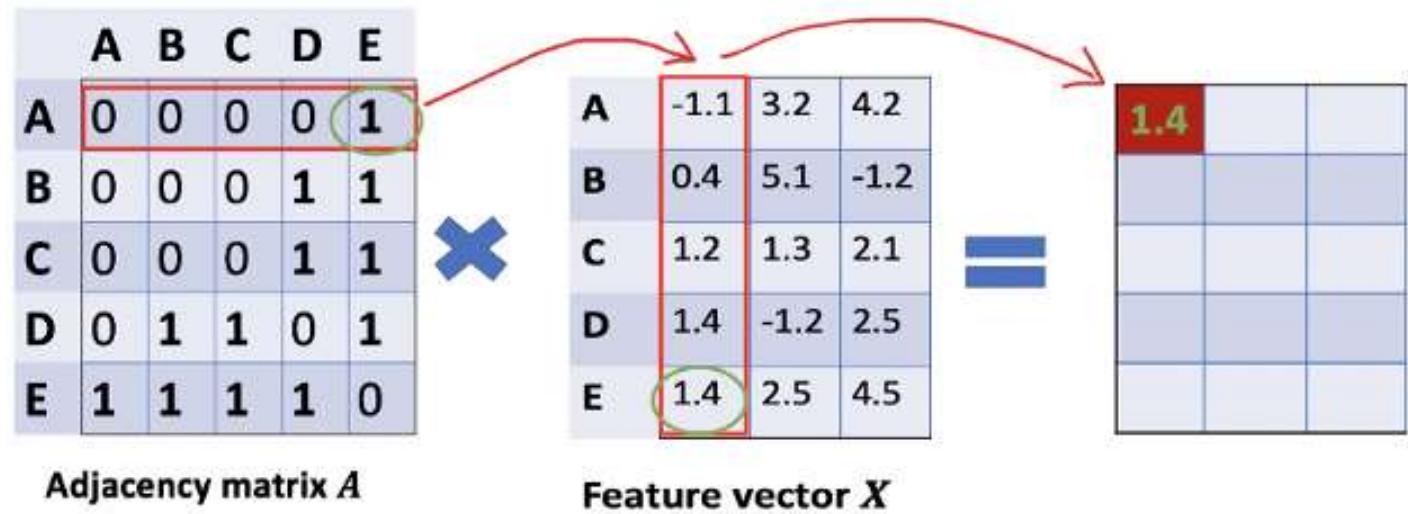
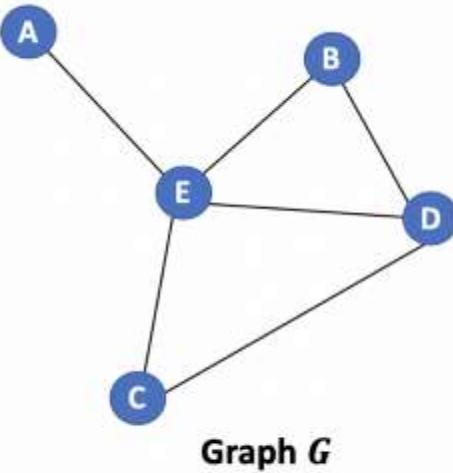
	A	B	C
A	-1.1	3.2	4.2
B	0.4	5.1	-1.2
C	1.2	1.3	2.1
D	1.4	-1.2	2.5
E	1.4	2.5	4.5

Feature vector  $X$

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

$$h^{(k+1)} = \sigma \sum_{j \in N_i} Wh_j^k$$

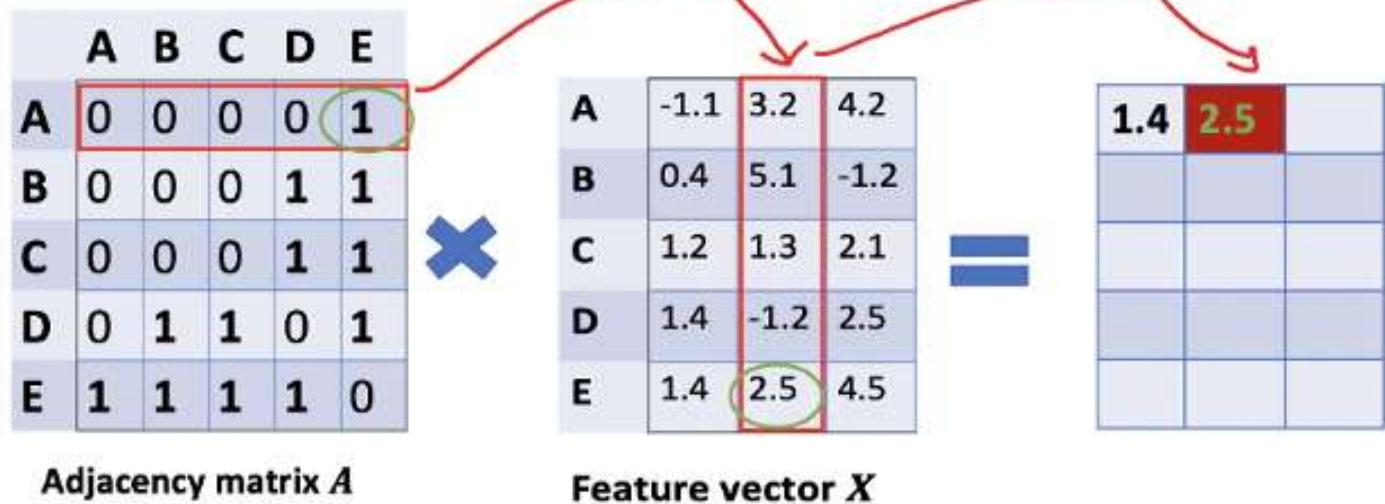
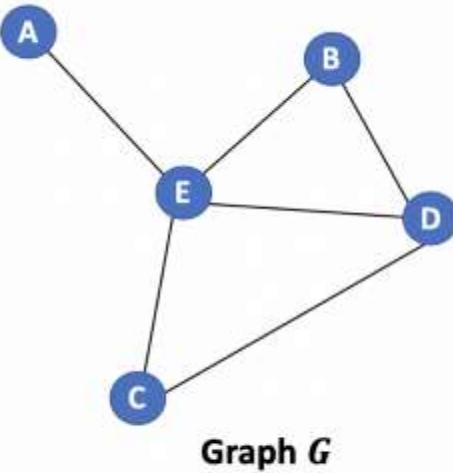
# INTUITION AND THE MATH'S BEHIND



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

$$h^{(k+1)} = \sigma \sum_{j \in N_i} Wh_j^k$$

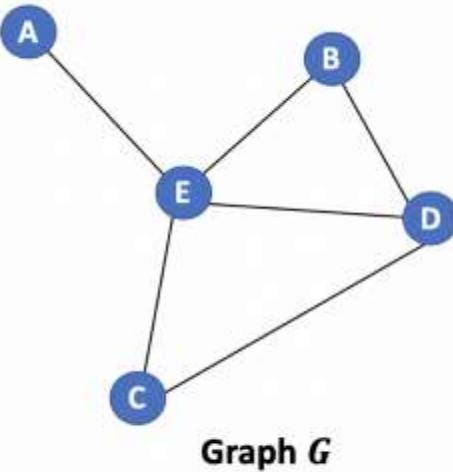
# INTUITION AND THE MATH'S BEHIND



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

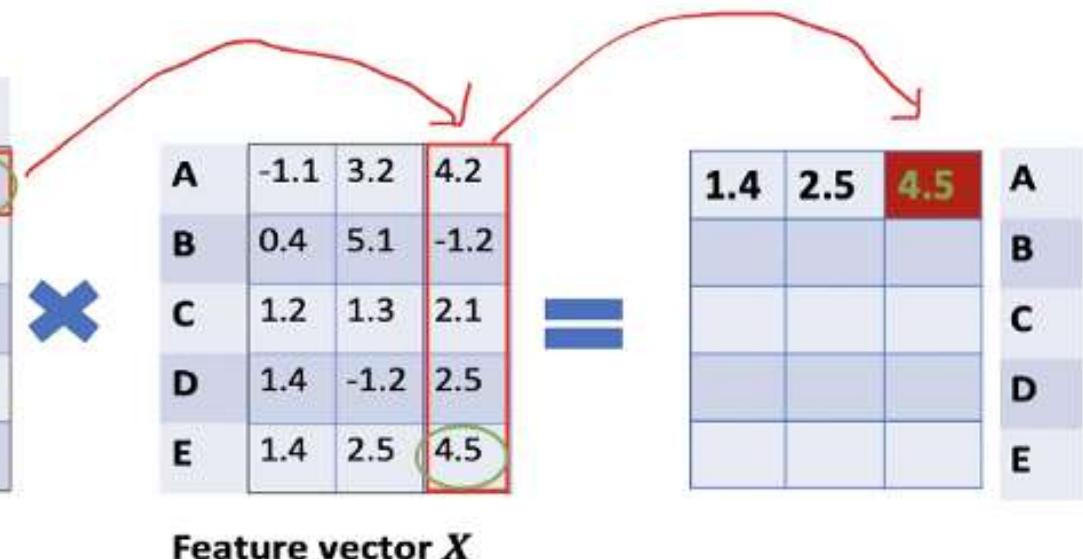
$$h^{(k+1)} = \sigma \sum_{j \in N_i} Wh_j^k$$

# INTUITION AND THE MATH'S BEHIND



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

Adjacency matrix  $A$



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

$$h^{(k+1)} = \sigma \sum_{j \in N_i} Wh_j^k$$

# PROBLEMS!

1. We miss the **feature of the node itself**. For example, the first row of the result matrix should contain features of node A too.

# PROBLEMS!

1. We miss the **feature of the node itself**. For example, the first row of the result matrix should contain features of node A too.

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

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



Adjacency matrix  $A$

1	0	0	0	0
0	1	0	0	0
0	0	1	0	0
0	0	0	1	0
0	0	0	0	1

Identity matrix  $I$

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

New Adjacency matrix  $\tilde{A}$

# PROBLEMS!

1. We miss the **feature of the node itself**. For example, the first row of the result matrix should contain features of node A too.
2. Instead of sum() function, we need to take the average, or even better, the weighted average of neighbors' feature vectors. **Why don't we use the sum() function?** The reason is that when using the sum() function, high-degree nodes are likely to have huge v vectors, while low-degree nodes tend to get small aggregate vectors, which may later cause **exploding or vanishing gradients** (e.g., when using sigmoid). Besides, Neural networks seem to be **sensitive to the scale of input data**. Thus, we need to normalize these vectors to get rid of the potential issues.

# PROBLEMS!

2. Instead of sum() function, we need to take the average, or even better, the weighted average of neighbors' feature vectors. **Why don't we use the sum() function?** The reason is that when using the sum() function, high-degree nodes are likely to have huge v vectors, while low-degree nodes tend to get small aggregate vectors, which may later cause **exploding or vanishing gradients** (e.g., when using sigmoid). Besides, Neural networks seem to be **sensitive to the scale of input data**. Thus, we need to normalize these vectors to get rid of the potential issues.

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

*New adjacency matrix  $\tilde{A}$*



2	0	0	0	0
0	3	0	0	0
0	0	3	0	0
0	0	0	4	0
0	0	0	0	5

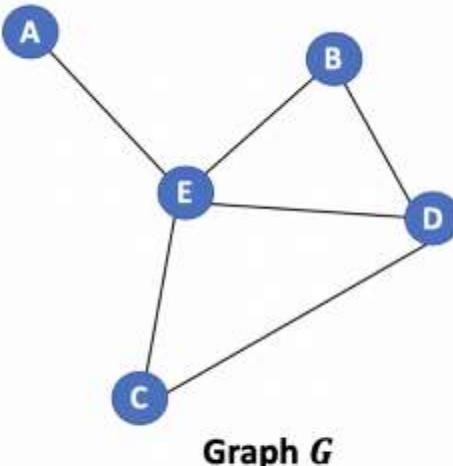
*New degree matrix  $\tilde{D}$*



1/2	0	0	0	0
0	1/3	0	0	0
0	0	1/3	0	0
0	0	0	1/4	0
0	0	0	0	1/5

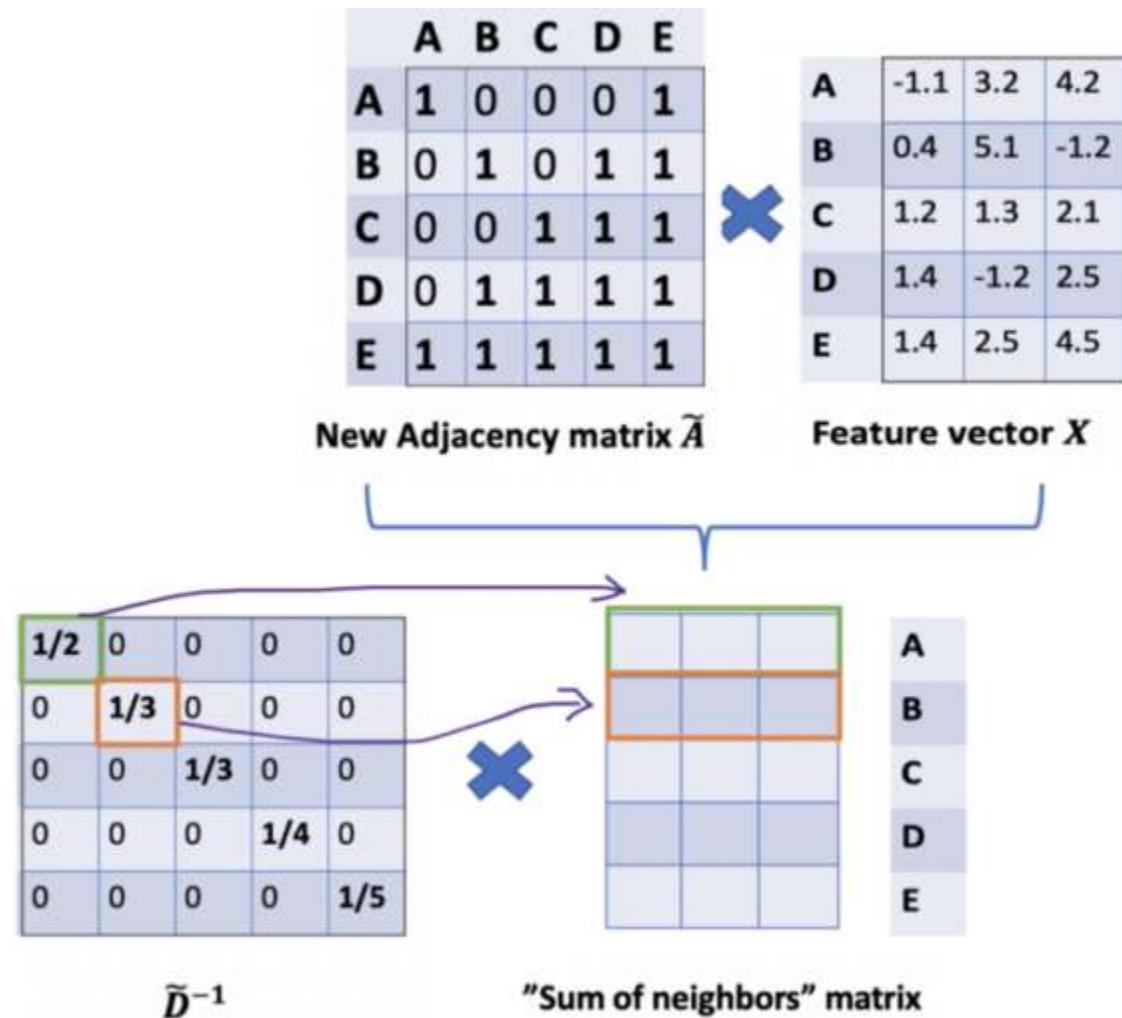
*$\tilde{D}^{-1}$*

# INTUITION AND THE MATH'S BEHIND



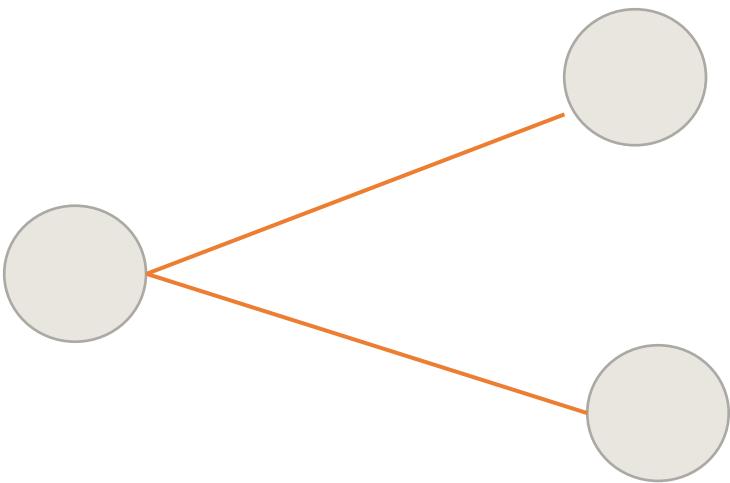
$$H^{(k+1)} = \sigma(D^{-1}(A + I)H^{(k)}W^{(k+1)})$$

$$h^{(k+1)} = \sigma \sum_{j \in N_i} \frac{1}{|N_i|} Wh_j^k$$



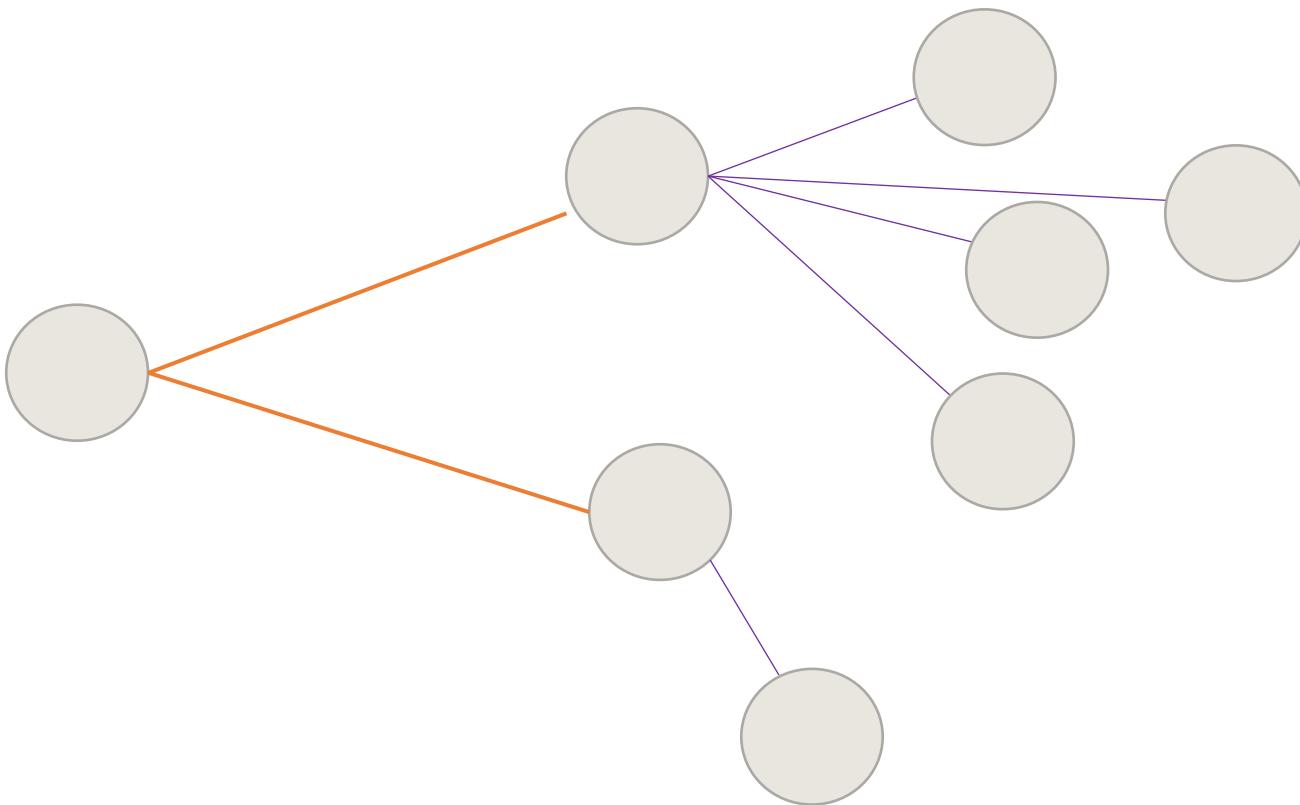
# INTUITION AND THE MATH'S BEHIND

- So far, so good!

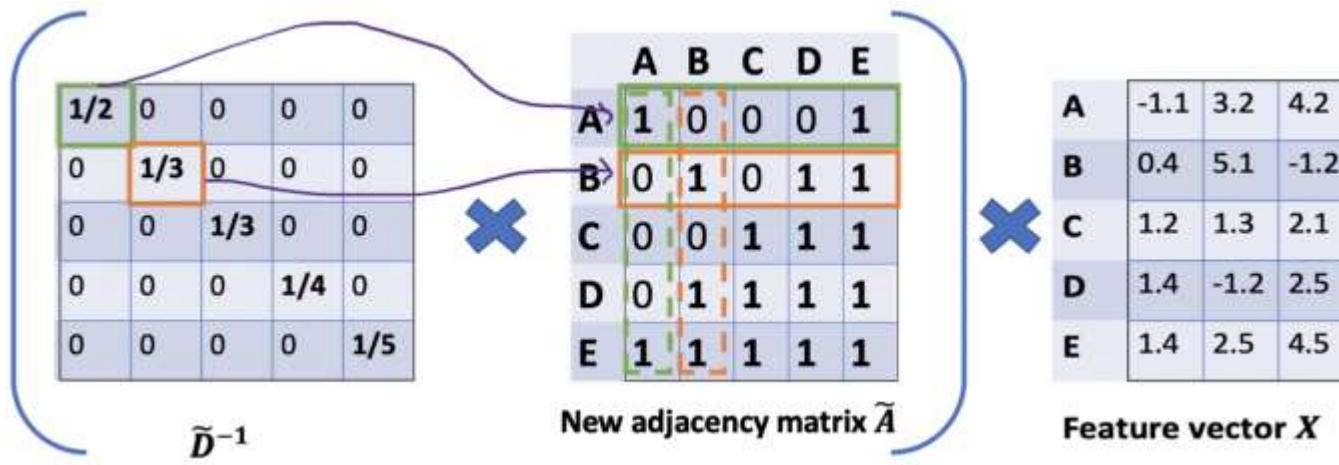


# INTUITION AND THE MATH'S BEHIND

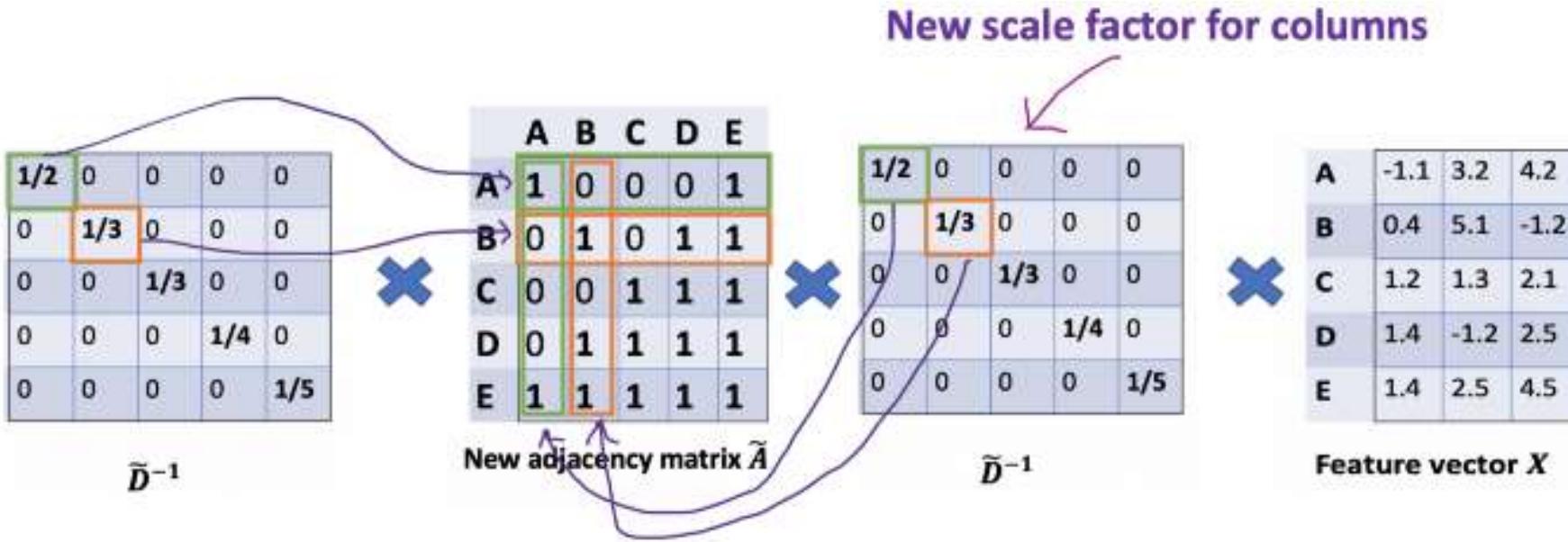
- So far, so good!
- Intuitively, it should be better if we treat high and low degree nodes differently.



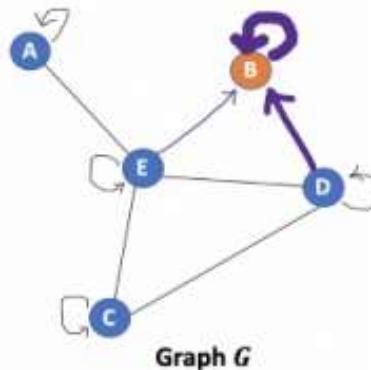
# INTUITION AND THE MATH'S BEHIND



# INTUITION AND THE MATH'S BEHIND



The new scalar gives us the “weighted” average.  
What we are doing here is to put more weights on  
the nodes that have low-degree and reduce the  
impact of high-degree nodes.



# INTUITION AND THE MATH'S BEHIND

**One more minor note:** When using two scalers ( $\tilde{D}_{ii}$  and  $\tilde{D}_{jj}$ ), we actually normalize twice, one time for the row as before, and another time for the column. It would make sense if we rebalance by modifying  $\tilde{D}_{ii}\tilde{D}_{jj}$  to  $\sqrt{\tilde{D}_{ii}\tilde{D}_{jj}}$ . In other words, instead of using  $\tilde{D}^{-1}$ , we use  $\tilde{D}^{-1/2}$ . So, we further alter the formula to  $\tilde{D}^{-1/2}\tilde{A}\tilde{D}^{-1/2}X$ , which is exactly used in the paper.

2	0	0	0	0
0	3	0	0	0
0	0	3	0	0
0	0	0	4	0
0	0	0	0	5

$\tilde{D}$



1/2	0	0	0	0
0	1/3	0	0	0
0	0	1/3	0	0
0	0	0	1/4	0
0	0	0	0	1/5

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

1/ $\sqrt{2}$	0	0	0	0
0	1/ $\sqrt{3}$	0	0	0
0	0	1/ $\sqrt{3}$	0	0
0	0	0	1/2	0
0	0	0	0	1/ $\sqrt{5}$

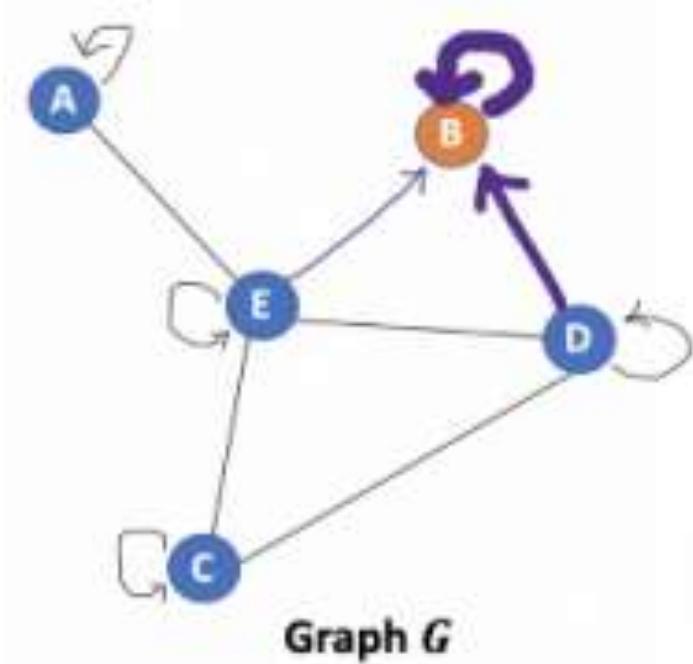
$\tilde{D}^{-1/2}$

# INTUITION AND THE MATH'S BEHIND

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

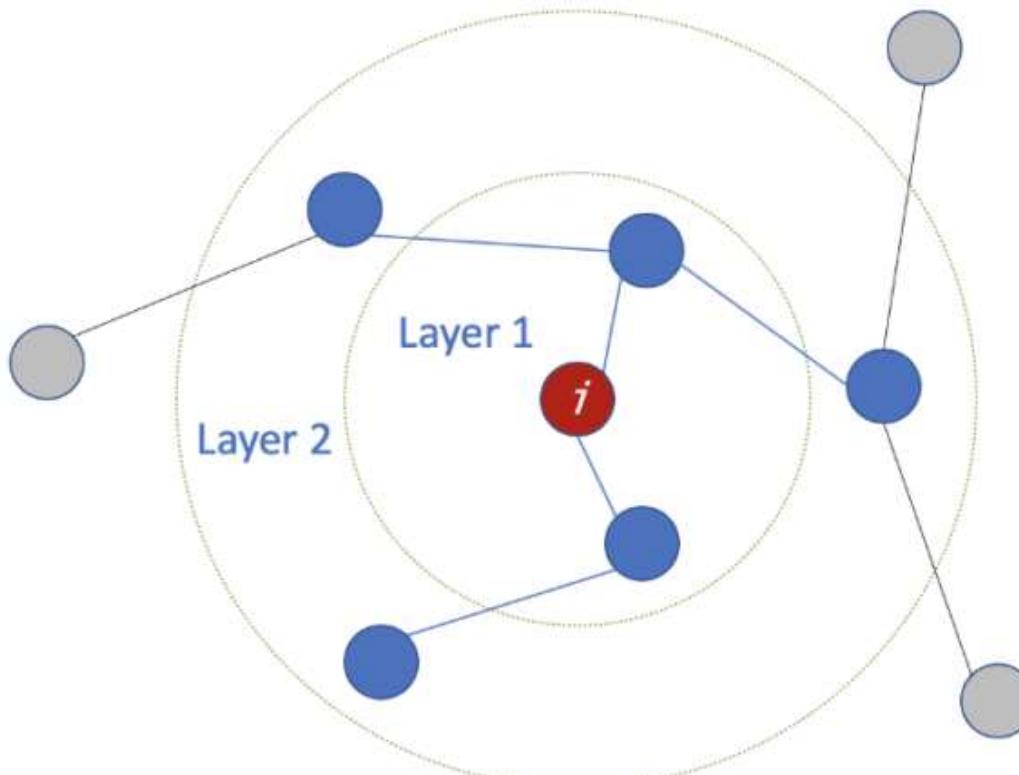
$$\tilde{\mathbf{A}} = (\mathbf{D} + \mathbf{I})^{-\frac{1}{2}}(\mathbf{I} + \mathbf{A})(\mathbf{D} + \mathbf{I})^{-\frac{1}{2}}$$

$$h^{(k+1)} = \sigma\left(\sum_{j \in N_i} \frac{1}{\sqrt{|N_i||N_j|}} Wh_j^k\right)$$

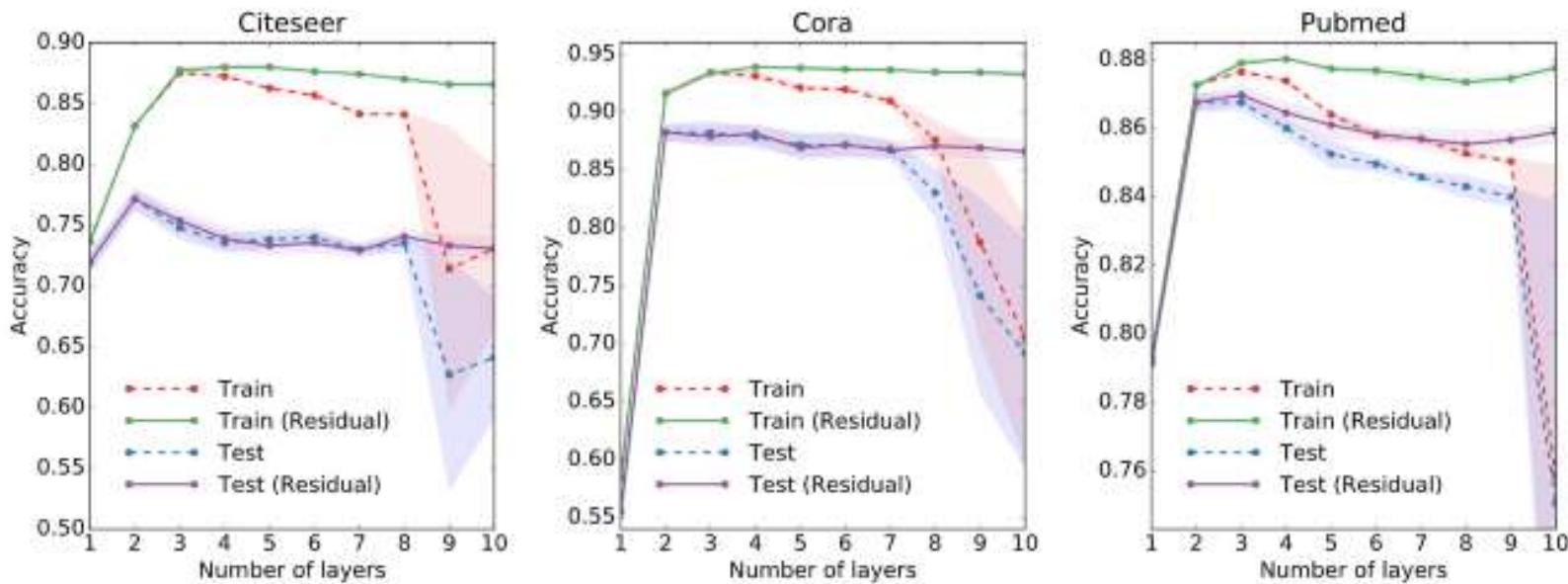


# THE NUMBER OF LAYERS

- ❑ The number of layers is the farthest distance that node features can travel.
- ❑ Normally we don't want to go too far. With 6–7 hops, we almost get the entire graph which makes the aggregation less meaningful.



# HOW MANY LAYERS SHOULD WE STACK THE GCN?



# GNN VARIANTS

$$h_u = \text{UPDATE}\left(h_u, \text{AGREGATE}(\{h_v, \forall v \in N(u)\})\right)$$

Graph Convolutional Networks,  
Kipf and Welling [2016]

$$\mathbf{h}_v^{(k)} = \sigma \left( \mathbf{W}^{(k)} \sum_{v \in \mathcal{N}(u) \cup \{u\}} \frac{\mathbf{h}_v}{\sqrt{|\mathcal{N}(u)||\mathcal{N}(v)|}} \right)$$

Sum of normalized  
neighbor embeddings

Multi-Layer-Perceptron as  
Aggregator, Zaheer et al. [2017]

Aggregated message

$$\mathbf{m}_{\mathcal{N}(u)} = \underset{\text{trainable!}}{\boxed{\text{MLP}_\theta}} \left( \sum_{v \in \mathcal{N}(u)} \text{MLP}_\phi(\mathbf{h}_v) \right)$$

Send states through a MLP

Graph Attention Networks,  
Veličković et al. [2017]

$$\mathbf{m}_{\mathcal{N}(u)} = \sum_{v \in \mathcal{N}(u)} \alpha_{u,v} \mathbf{h}_v$$

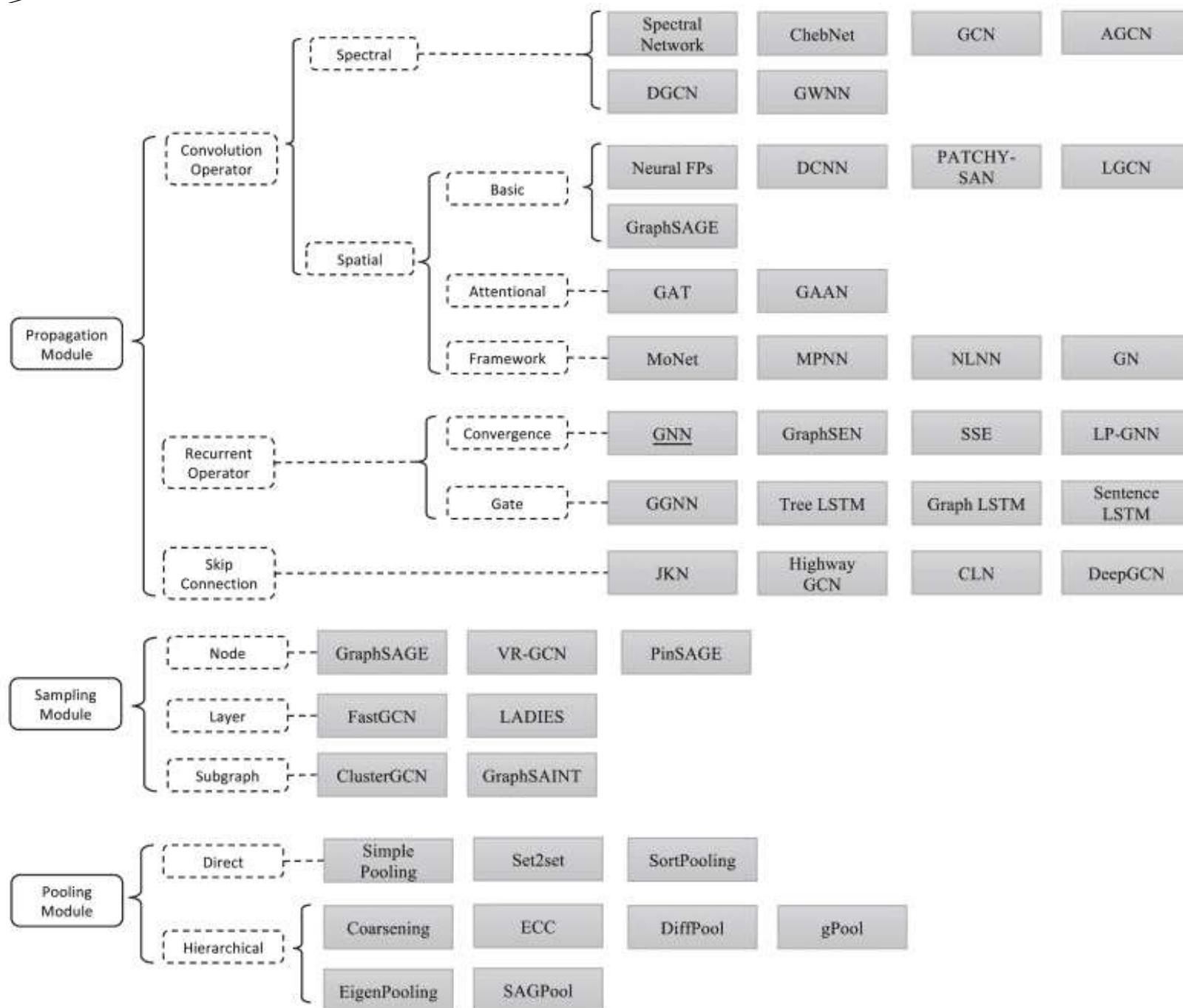
Attention weights

$$\alpha_{u,v} = \frac{\exp(\mathbf{a}^\top [\mathbf{W}\mathbf{h}_u \oplus \mathbf{W}\mathbf{h}_v])}{\sum_{v' \in \mathcal{N}(u)} \exp(\mathbf{a}^\top [\mathbf{W}\mathbf{h}_u \oplus \mathbf{W}\mathbf{h}_{v'}])}$$

Gated Graph Neural Networks,  
Li et al. [2015]

$$\mathbf{h}_u^{(k)} = \text{GRU}(\mathbf{h}_u^{(k-1)}, \mathbf{m}_{\mathcal{N}(u)}^{(k)})$$

Recurrent update of the state



---

# GRAPH REPRESENTATION LEARNING

---

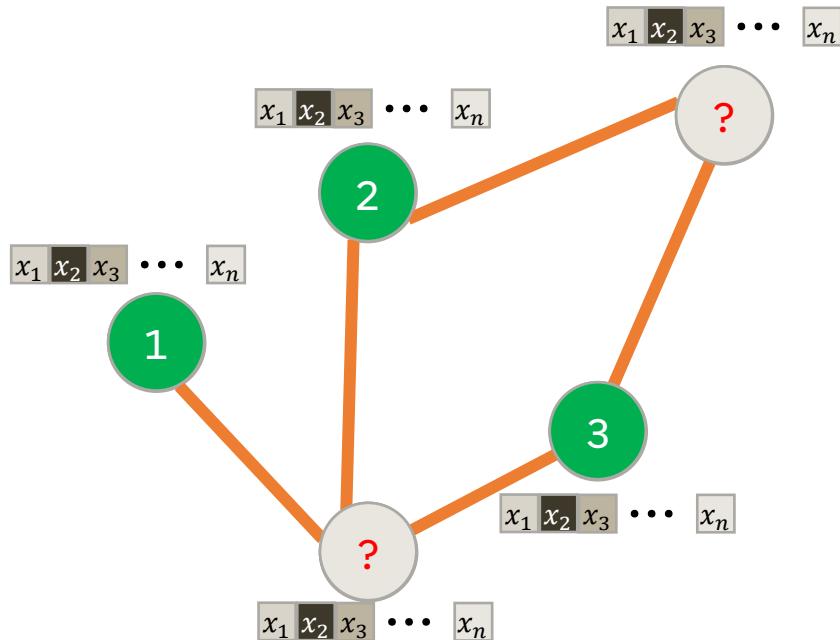
WILLIAM L. HAMILTON

*McGill University*

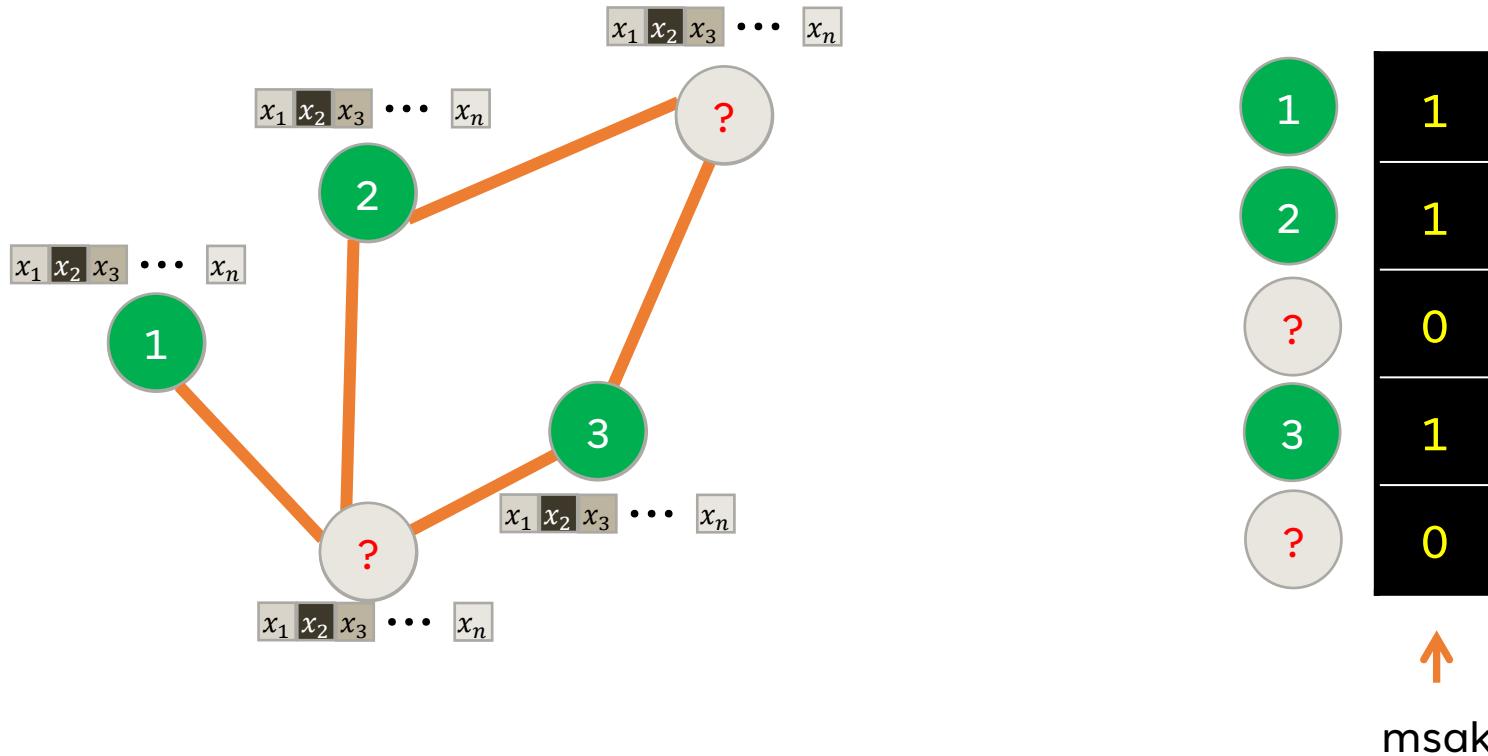
2020

[https://www.cs.mcgill.ca/~wlh/grl\\_book/files/GRL\\_Book.pdf](https://www.cs.mcgill.ca/~wlh/grl_book/files/GRL_Book.pdf)

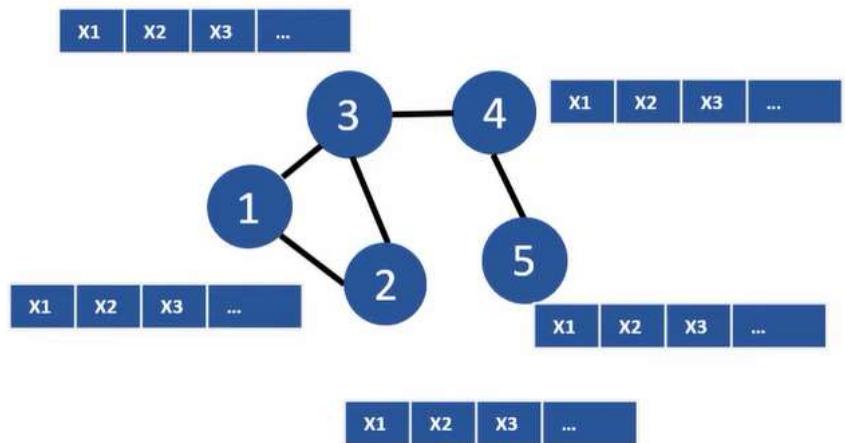
# BINARY MASKS FOR NODE-LEVEL PREDICTION



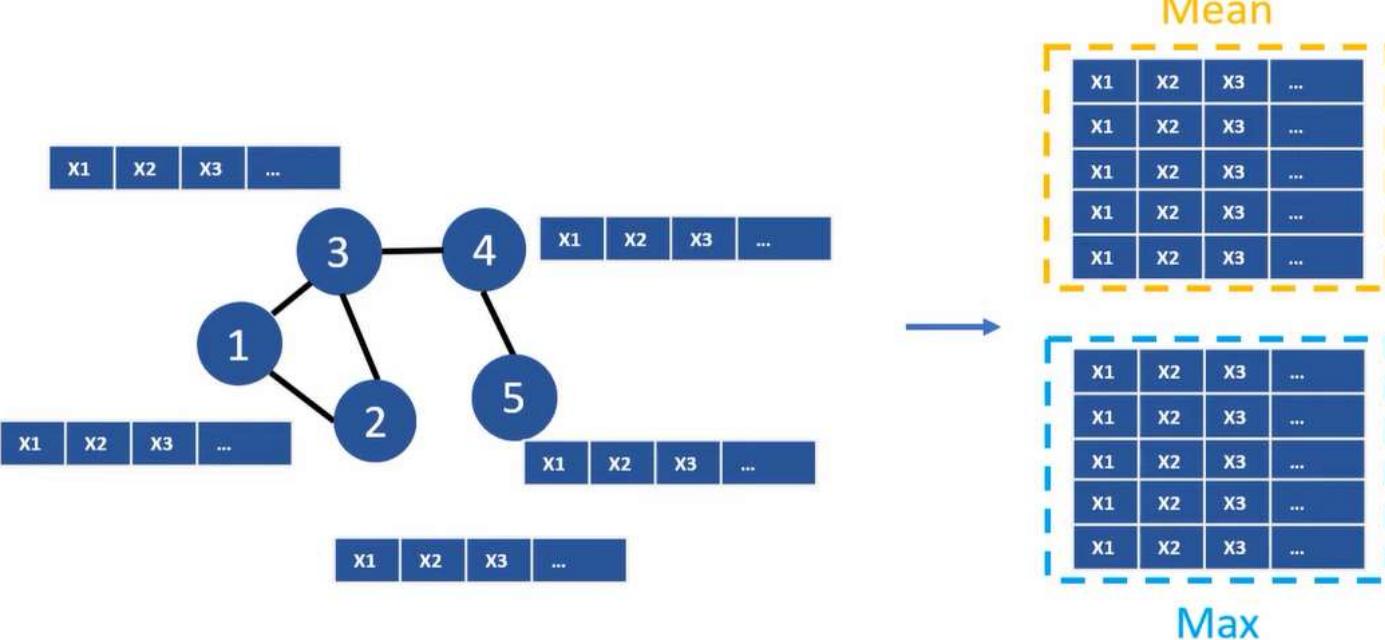
# BINARY MASKS FOR NODE-LEVEL PREDICTION



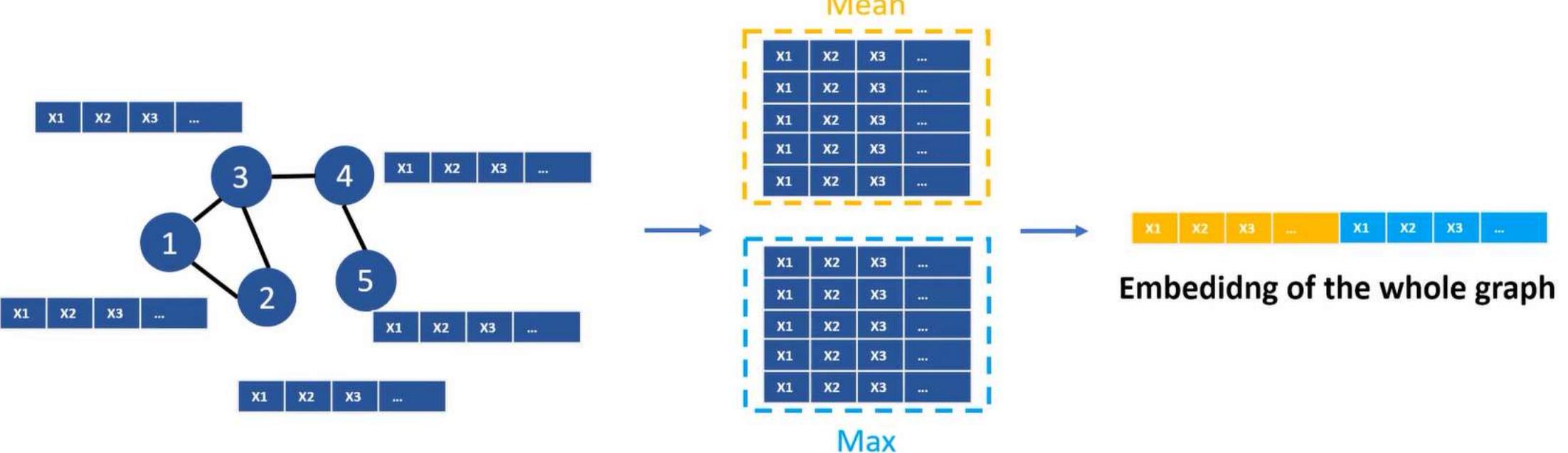
# GLOBAL GRAPH POOLING



# GLOBAL GRAPH POOLING



# GLOBAL GRAPH POOLING



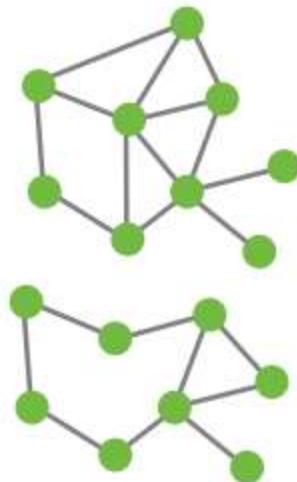
# BATCHING WITH GRAPHS

In the image or language domain:  
**rescaling or padding**

WHAT ABOUT Graphs?

## BATCHING WITH GRAPHS

$$\mathcal{G}_1 = (\mathbf{X}_1, \mathbf{A}_1)$$

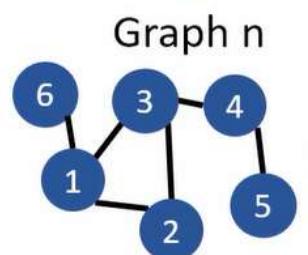
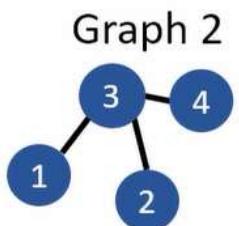
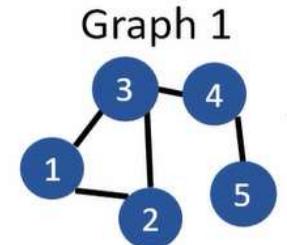


$$\mathcal{G}_2 = (\mathbf{X}_2, \mathbf{A}_2)$$

GNN

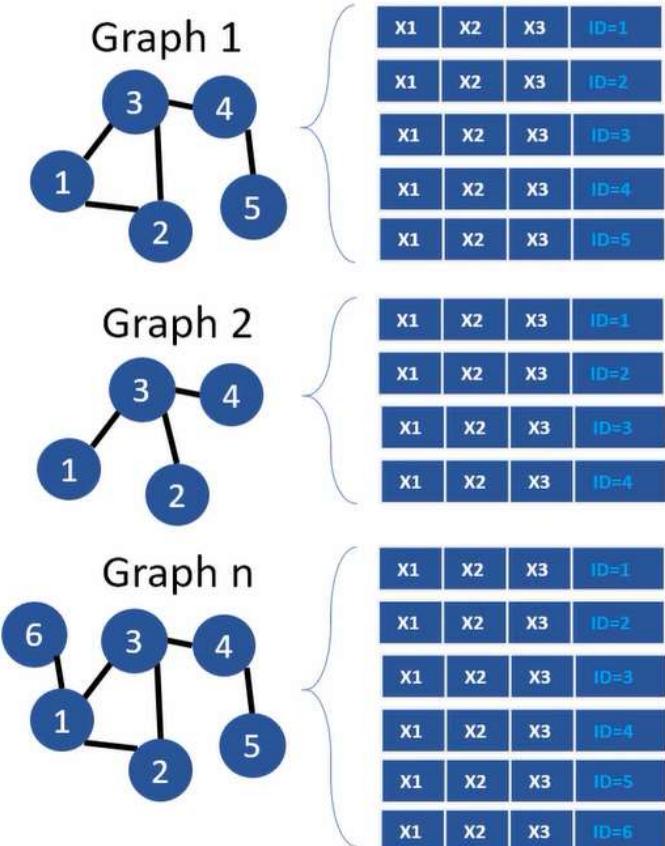
$$\left( \begin{array}{|c|c|} \hline \mathbf{A}_1 & \mathbf{A}_2 \\ \hline \mathbf{X}_1 & \mathbf{X}_2 \\ \hline \end{array} \right) = \left( \begin{array}{|c|} \hline \mathbf{X}'_1 \\ \hline \mathbf{X}'_2 \\ \hline \end{array} \right)$$

# BATCHING WITH GRAPHS



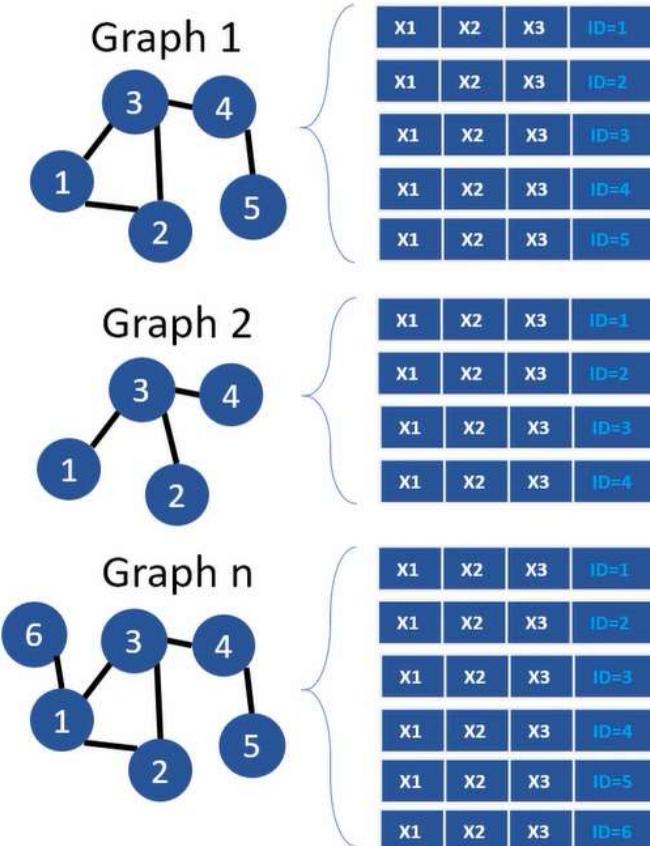
$n$  = Batch Size

# BATCHING WITH GRAPHS



n = Batch Size

# BATCHING WITH GRAPHS

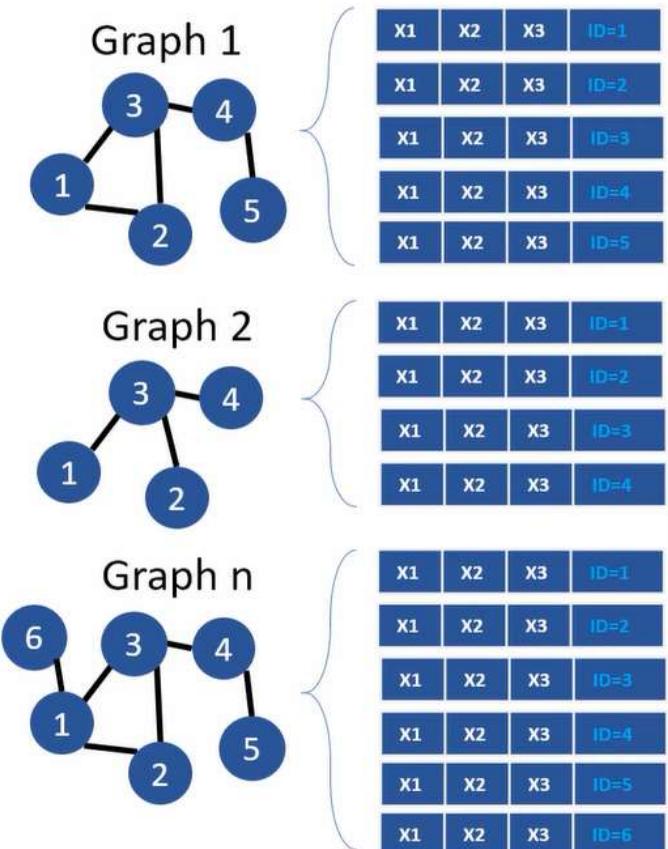


n = Batch Size

Large Adjacency Matrix

0	1	0	1	0	0	0	...
1							
0							
1							
...							

# BATCHING WITH GRAPHS



n = Batch Size

Large Adjacency Matrix

The diagram shows the final output of the GNN, which are the embeddings for each graph node. There are two sets of tables labeled "Embeddings" to the right of the GNN block.

X1	X2	X3	ID
X1	X2	X3	ID=1
X1	X2	X3	ID=2
X1	X2	X3	ID=3
X1	X2	X3	ID=4
X1	X2	X3	ID=5

X1	X2	X3	ID
X1	X2	X3	ID=1
X1	X2	X3	ID=2
X1	X2	X3	ID=3
X1	X2	X3	ID=4
X1	X2	X3	ID=5
X1	X2	X3	ID=6

Embeddings

# **SCALING UP GRAPH NEURAL NETWORKS TO LARGE GRAPHS**

# GRAPHS IN MODERN APPLICATIONS

## Recommender systems:

- Amazone
- YouTube
- Pinterest
- Instagram



**amazon**

## Tasks:

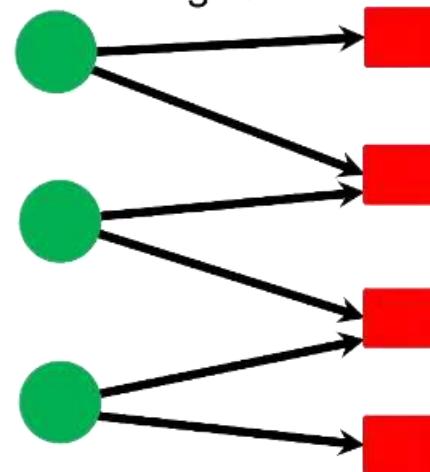
- Recommend Items (Link Prediction)
- Classify users/Items (Node Classification)

Users:       Products / Videos:

100M ~ 1B

10M~ 1B

Bought/saw



# GRAPHS IN MODERN APPLICATIONS

## Social Networks

- Facebook
- Twitter
- Instagram

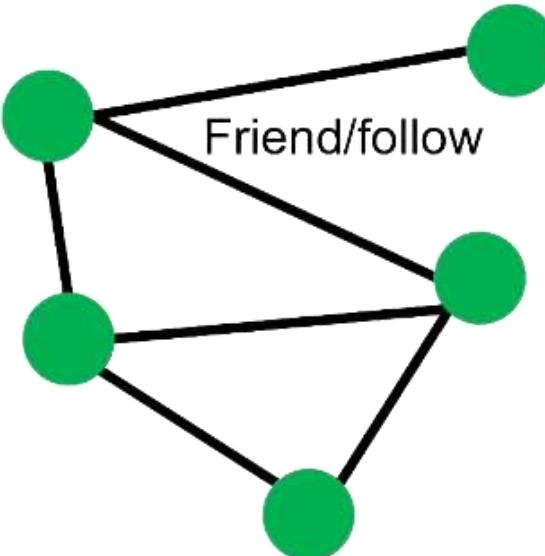


## Tasks:

- Friend Recommend(Link Prediction)
- User property recommendation (Node-Level)

## ❑ Users:

300M ~ 3B



# GRAPHS IN MODERN APPLICATIONS

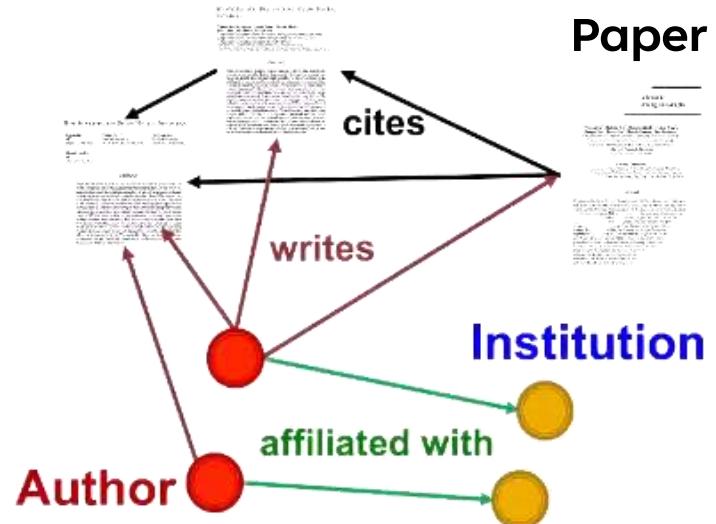
## Academic Graph

- Microsoft Academic Graph/

### Tasks:

- Paper categorization  
(node classification)
- Author collaboration recommendation
- Paper citation recommendation  
(Link prediction)

Papers      Authors  
120M      120 M



# GRAPHS IN MODERN APPLICATIONS

## Knowledge Graphs (KGs)

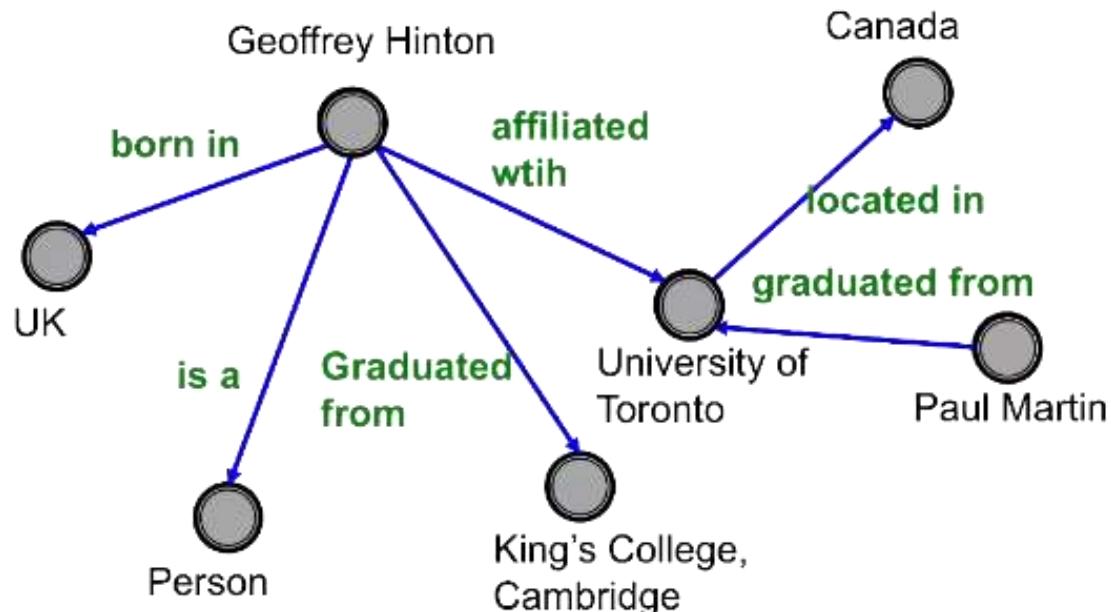
- Wikipedia
- Freebase

## Tasks:

- KG completion
- Reasoning

### ❑ Entities:

80M ~ 90M



# WHAT IS IN COMMON?!

## ❑ 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

# PROBLEM!

**Full-batch** implementation is **not feasible** for a large graphs

## Time inefficiency

- In CPU takes too much time!

## Memory Limitations

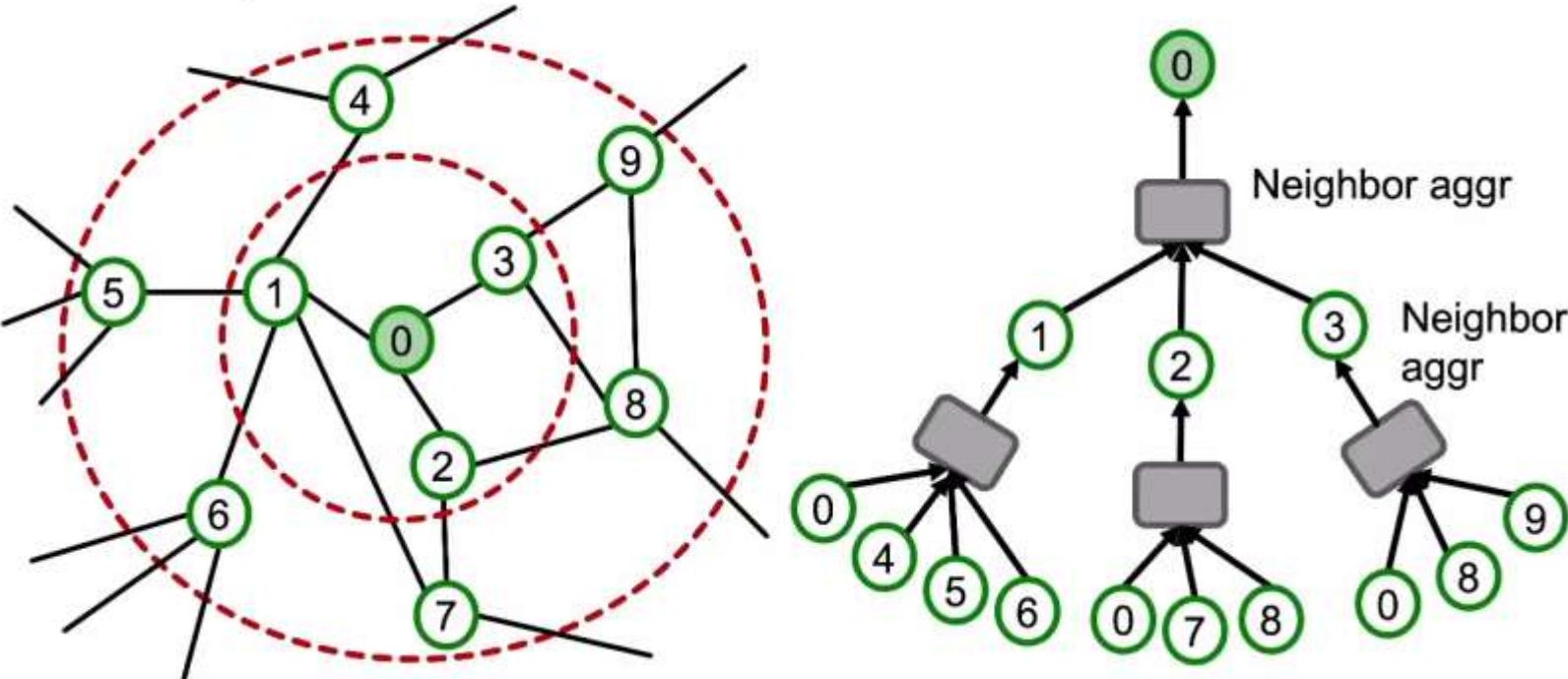
- GPU memory is extremely limited
- We cannot load entire dataset into memory

# SOLUTIONS! SOME METHODS FOR SCALING UP GNNs

- Perform message-passing over **small subgraphs in each mini-batch**
  - ❖ Only the subgraphs need to be loaded on a GPU at a time.
    - Neighbour Sampling [Hamilton NeurIPS 2017]
    - Cluster-GCN [Chiang et al. KDD 2019]
- Simplifies a GNN into feature-preprocessing operation
  - ❖ Can be efficiently performed even on a CPU
    - Simplified GCN [Wu et al. ICML2019]

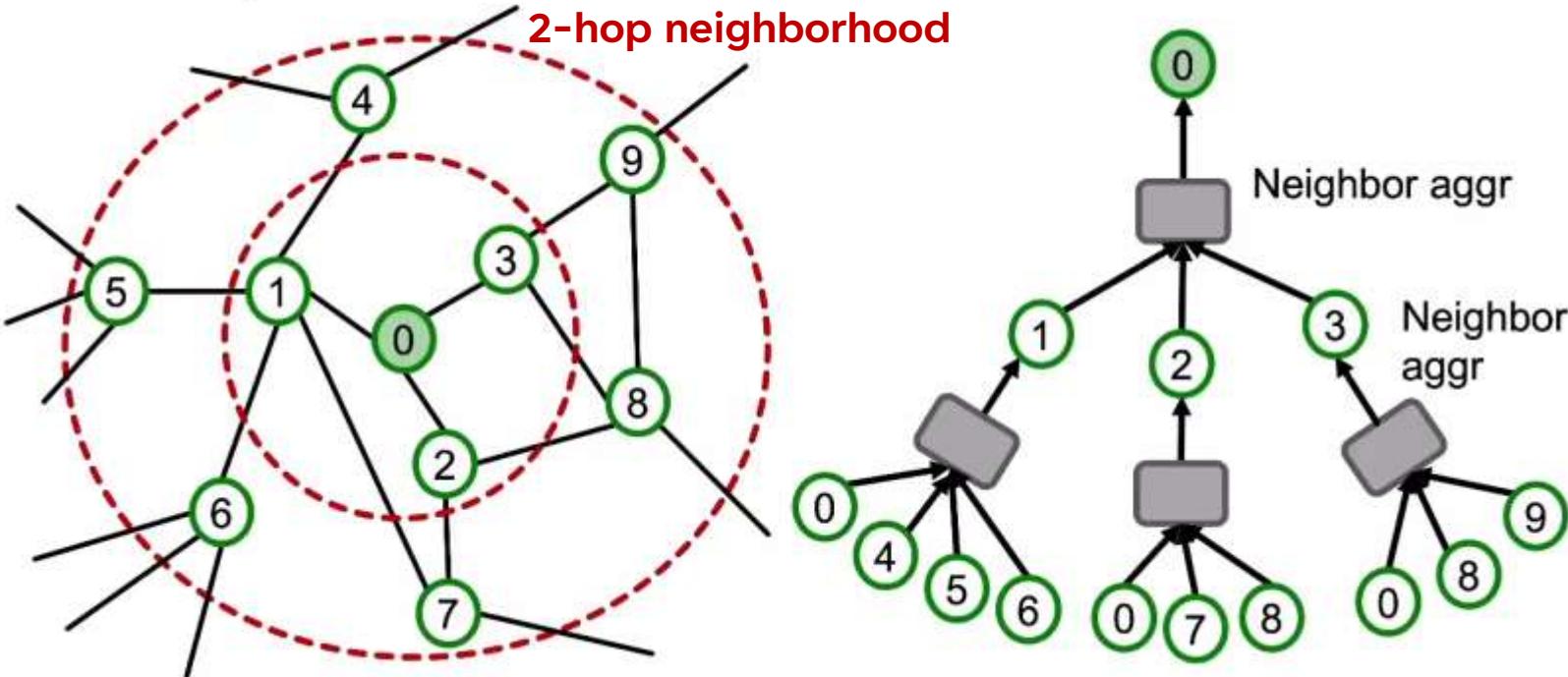
# GRAPHSAGE NEIGHBOR SAPLING

GNNs generate node embeddings via neighbour aggregation.



# GRAPHSAGE NEIGHBOR SAPLING

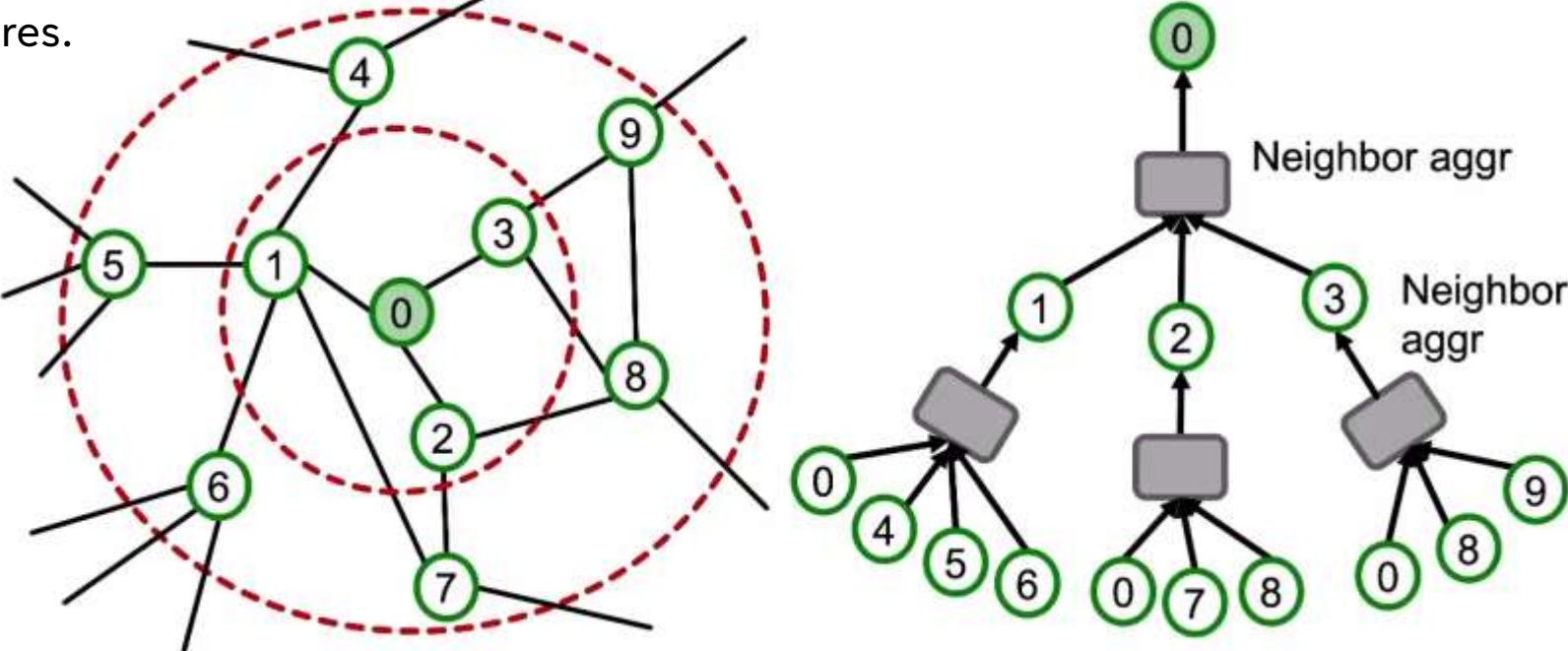
**Observation:** A 2-layer GNN generates embedding of node "0" using 2-hop neighborhood structure and features.



# GRAPHSAGE NEIGHBOR SAPLING

**Observation:** A 2-layer GNN generates embedding of node "0" using 2-hop neighborhood structure and features.

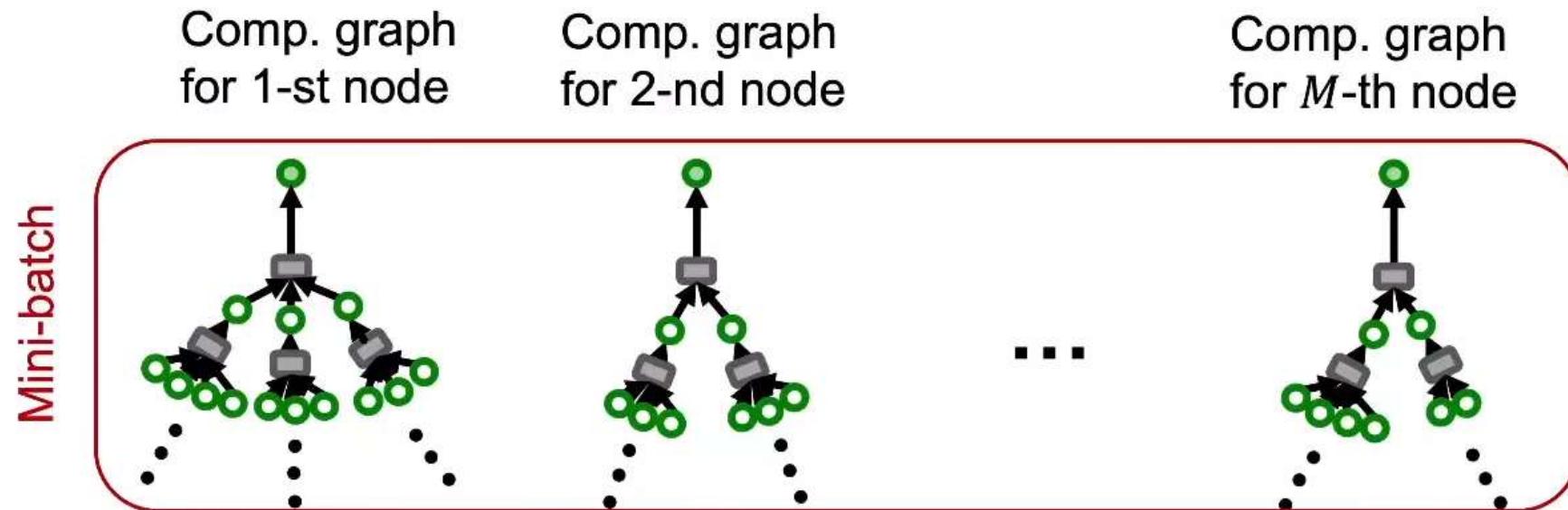
More generally, K-layer GNNs generate embedding of a node using K-hop neighborhood structure and features.



# GRAPHSAGE NEIGHBOR SAPLING

**Key insight:** To compute embedding of a single node, all we need is the **K-hop neighborhood** (which defines the computation graph).

- Given a set of **M different nodes in a mini-batch**, we can generate their embeddings using **M computational graphs. Can be computed on GPU!**



# STOCHASTIC TRAINING OF GNNS

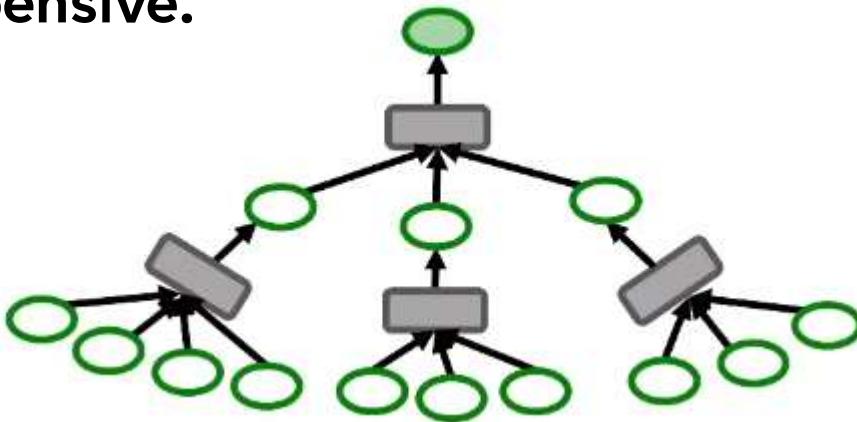
We can now consider the following SGD strategy for training K-layer GNNs:

- Randomly sample M ( $<< N$ ) nodes.
- For each sampled node  $v$ :
  - Get **k-hop neighbourhood**, and construct the **computation graph**.
  - Use the above to generate  $v$ 's embedding.
- Compute the loss  $l_{sub}(\theta)$  averaged over the M nodes.
- Perform SGD:  $\theta \leftarrow \theta - \nabla l_{sub}(\theta)$



# ISSUE STOCHASTIC TRAINING

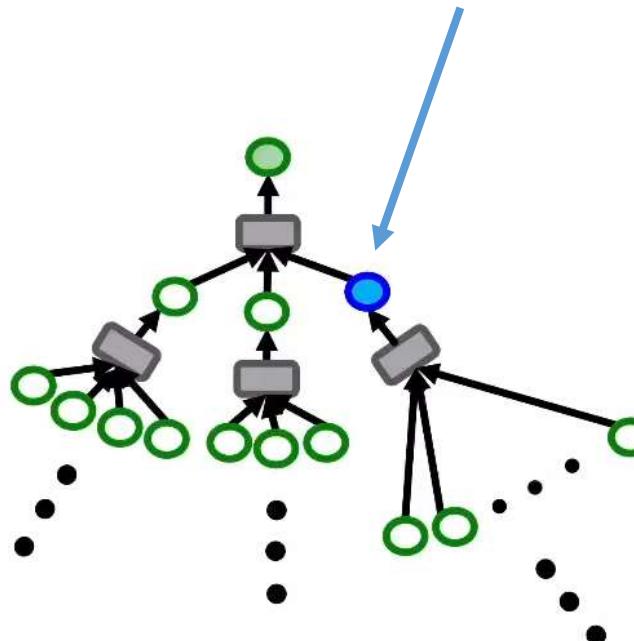
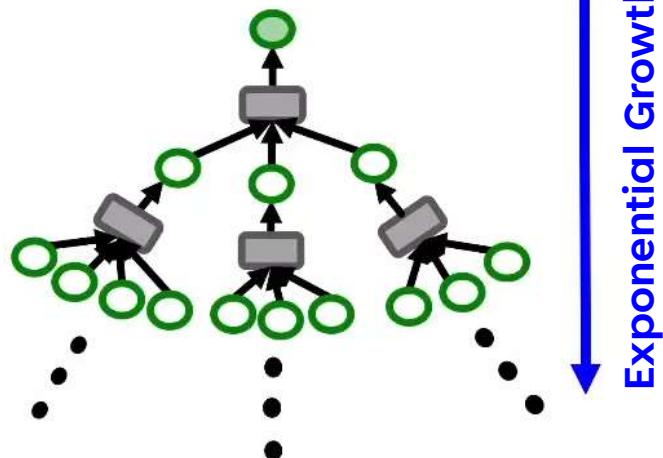
- For each node, we need to get **the entire K-hop neighborhood** and pass it through the computation graph.
- We need to aggregate lot of information just to compute one node embedding.
- **Computationally expensive.**



# ISSUE STOCHASTIC TRAINING

**More details:**

- Computation graph becomes **exponentially large** with respect to the layer size K.
- Computation graph explodes when it hits a **hub node** (high-degree node).

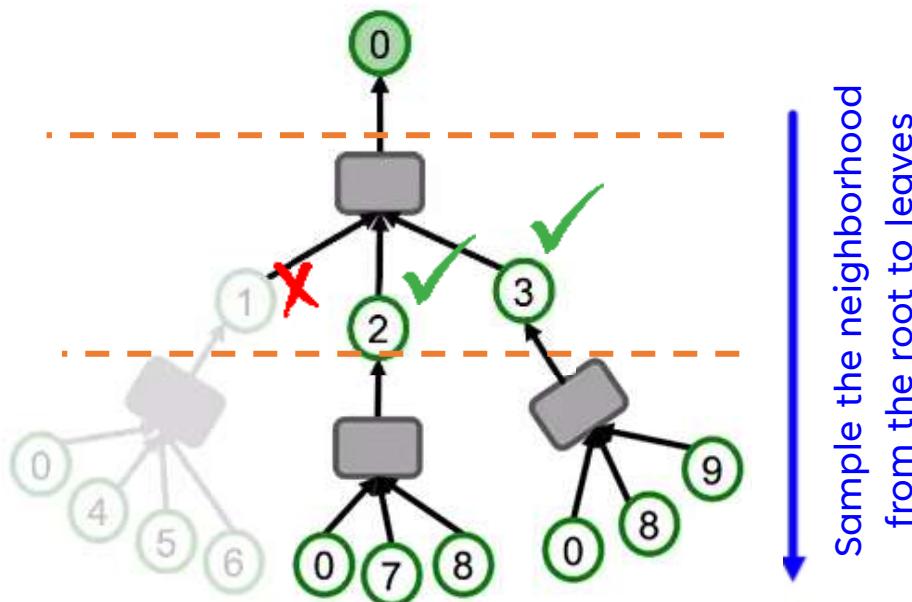


# NEIGHBOR SAMPLING

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

❑ Example:

1<sup>st</sup> hub neighborhood  
Sample 2, 3 | Drop 1



# NEIGHBOR SAMPLING

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

❑ Example:

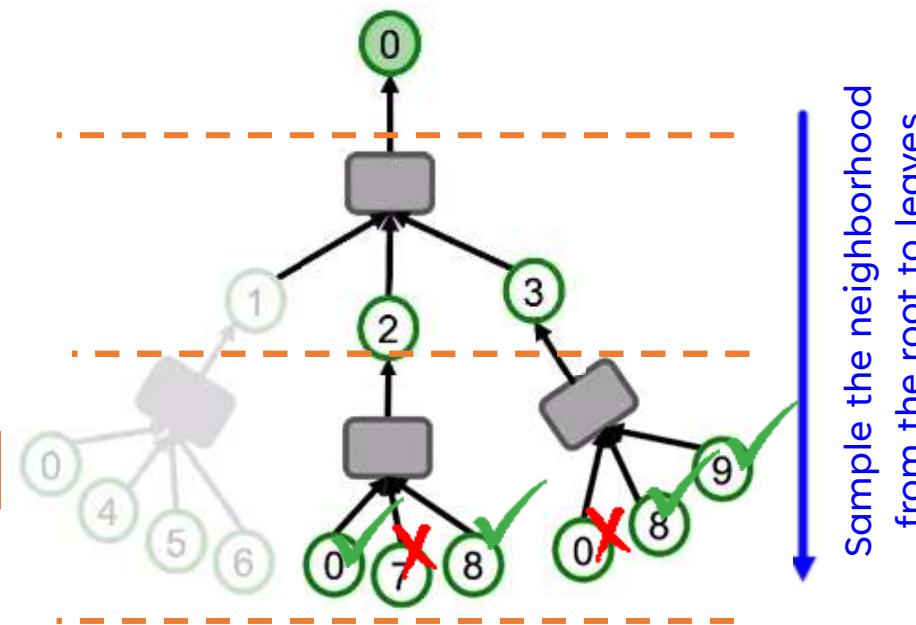
1<sup>st</sup> hub neighborhood

Sample 2, 3 | Drop 1

2<sup>nd</sup> hub neighborhood

Sample 0, 8 | Drop 7

Sample 8, 9 | Drop 0



# NEIGHBOR SAMPLING

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

❑ Example:

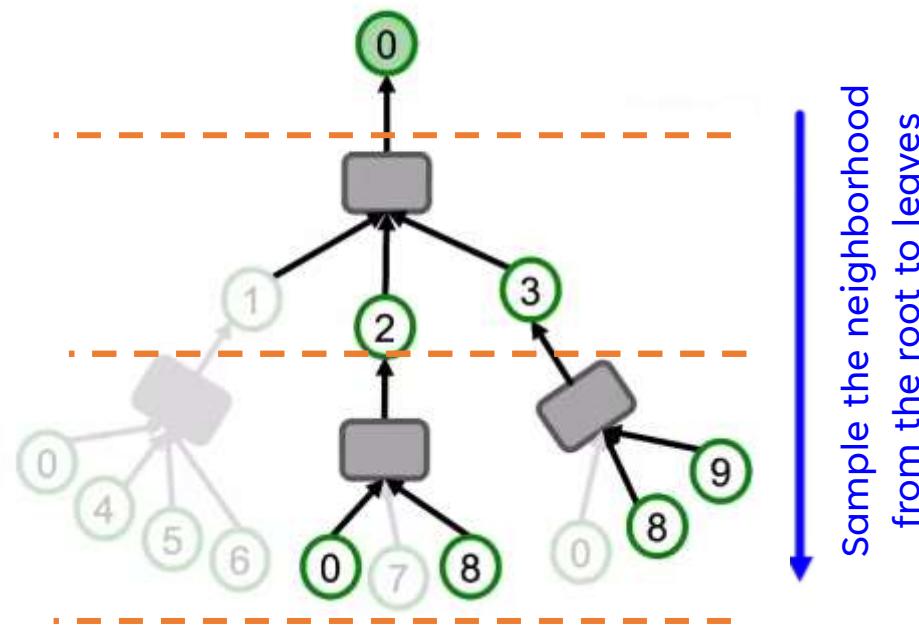
1<sup>st</sup> hub neighborhood

Sample 2, 3 | Drop 1

2<sup>nd</sup> hub neighborhood

Sample 0, 8 | Drop 7

Sample 8, 9 | Drop 0



❖ K-layer GNN will at most involve  $\prod_{k=1}^K H_k$  leaf nodes in computation graph.

# REMARKS ON NEIGHBOR SAMPLING

## ❑ Remark 1: Trade-off in sampling number H

- ❖ Smaller  $H$  leads to more efficient neighbour aggregation, but results in more unstable training due to the larger variance in neighbour aggregation.

## ❑ Remark 2: Computational time

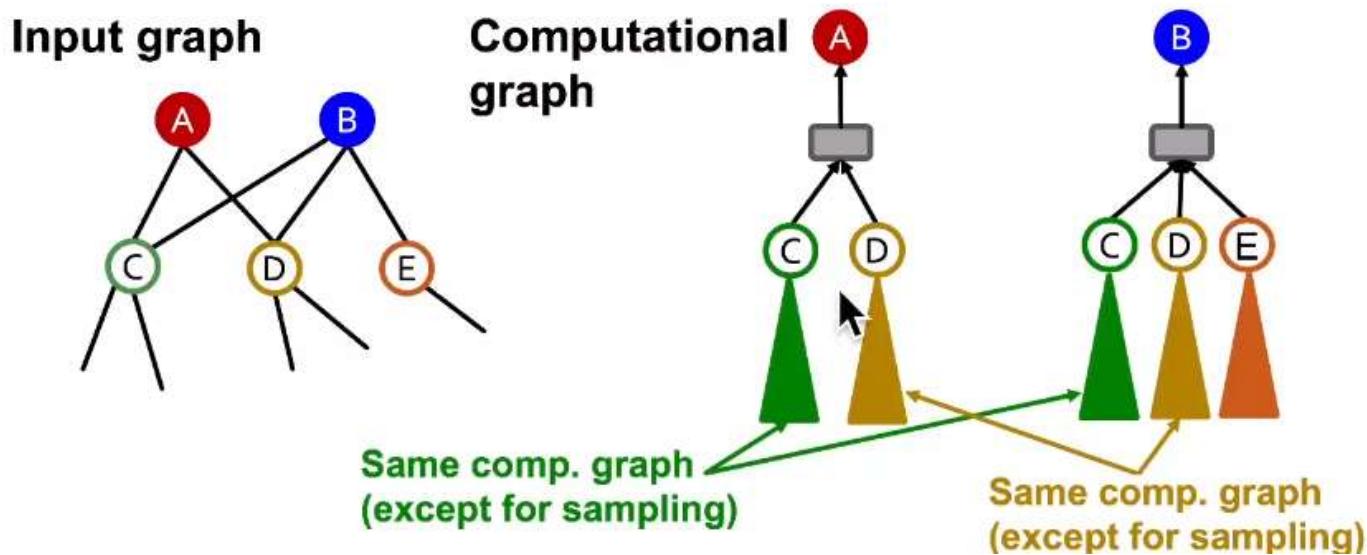
- ❖ Even with neighbour sampling, the size of the computational graph is still exponential with respect to number of GNN layers  $K$ .
- ❖ Increasing one GNN layer would make computation  $H$  times more expensive.

## ❑ Remark 3: How to sample the nodes

- ❖ Random sampling: fast but many times not optimal!
- ❖ Random walk with restart

# ISSUE WITH NEIGHBOUR SAMPLING

- Issue with neighbour sampling:
  - The size of computational graph becomes exponentially large w.r.t. the #GNN layers.
  - **Computation is redundant**, especially when nodes in a mini-batch share many neighbours.



# conv.SAGEConv qq

```
class SAGEConv ( in_channels: Union[int, Tuple[int, int]], out_channels: int, aggr: Optional[Union[str, List[str], Aggregation]] = 'mean', normalize: bool = False, root_weight: bool = True, project: bool = False, bias: bool = True, **kwargs )
```

[\[source\]](#)

Bases: [MessagePassing](#)

The GraphSAGE operator from the “[Inductive Representation Learning on Large Graphs](#)” paper

$$\mathbf{x}'_i = \mathbf{W}_1 \mathbf{x}_i + \mathbf{W}_2 \cdot \text{mean}_{j \in \mathcal{N}(i)} \mathbf{x}_j$$

# Redundancy-Free Computation for Graph Neural Networks

Zhihao Jia  
Stanford University  
zhihao@cs.stanford.edu

Jiaxuan You  
Stanford University  
jiaxuan@stanford.edu

Sina Lin  
Microsoft  
silin@microsoft.com

Jure Leskovec  
Stanford University  
jure@cs.stanford.edu

Rex Ying  
Stanford University  
rexying@stanford.edu

Alex Aiken  
Stanford University  
aiken@cs.stanford.edu

## ABSTRACT

Graph Neural Networks (GNNs) are based on repeated aggregations of information from nodes' neighbors in a graph. However, because nodes share many neighbors, a naive implementation leads to repeated and inefficient aggregations and represents significant computational overhead. Here we propose *Hierarchically Aggregated computation Graphs* (HAGs), a new GNN representation technique that explicitly avoids redundancy by managing intermediate aggregation results hierarchically and eliminates repeated computations and unnecessary data transfers in GNN training and inference. HAGs perform the same computations and give the same models/accuracy as traditional GNNs, but in a much shorter time due

## 1 INTRODUCTION

Graph Neural Network models (GNNs) generalize deep representation learning to graph data [3, 9, 23] and have achieved state-of-the-art performance across a number of graph-based tasks, such as node classification, link prediction, and graph classification and recommender systems [8, 14, 24, 27].

GNNs are based on a recursive neighborhood aggregation scheme, where within a single layer of a GNN each node aggregates its neighbors' activations and uses the aggregated value to update its own activation [23]. Such updated activations are then recursively propagated multiple times (multiple layers). In the end, every node in a GNN collects information from other nodes that are in its k-

One approach to solve the redundancy problem!

<https://dl.acm.org/doi/pdf/10.1145/3394486.3403142>

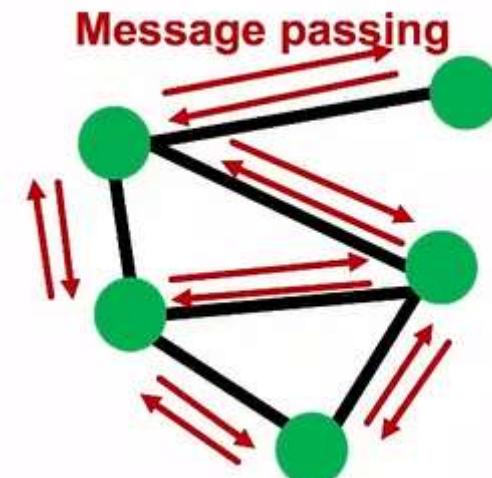
# CLUSTER-GCN: REVIEW FULL-BATCH GNN

- In full-batch GNN implementation, all the node embeddings are updated together using embeddings of the previous layer

**Update for all  $v \in V$**

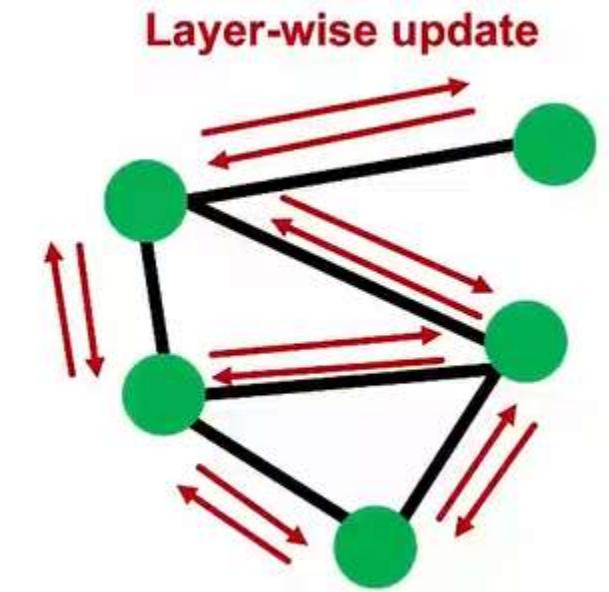
$$h_v^{(\ell)} = \text{COMBINE} \left( h_v^{(\ell-1)}, \text{AGGR} \left( \left\{ \textcolor{red}{h_u^{(\ell-1)}} \right\}_{u \in N(v)} \right) \right)$$

- In each layer, only **2\*(edges)** messages need to be computed.
- For K-layer GNN, only  $2K * \#(\text{edges})$  messages need to be computed.
- GNN's entire computation is only **linear** in  $\#(\text{edges})$  and  $\#(\text{GNN layers})$ . **Fast!**



# CLUSTER-GCN: INSIGHT FROM FULL-BATCH GNN

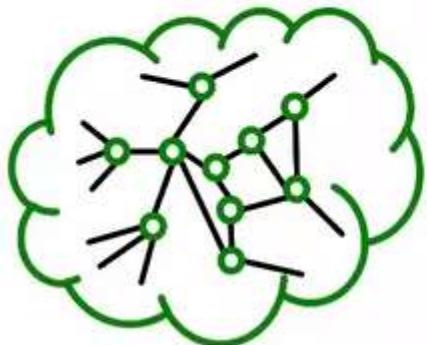
- The **layer-wise** node embedding update allows the re-use of embeddings from the previous layer.
- This significantly **reduces the computational redundancy of neighbour sampling**.
  - ❖ Of course, the **layer-wise update is not feasible** for a large graph due to **limited GPU memory**.



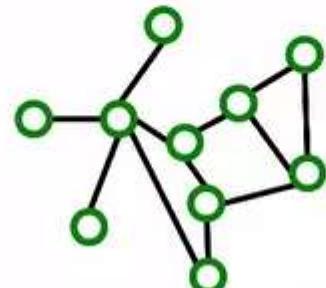
# CLUSTER-GCN: SUB-GRAPH SAMPLING

- ✓ **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.

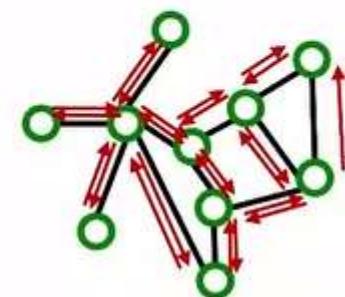
**Large graph**



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



**Layer-wise  
node embeddings  
update on the GPU**



## CLUSTER-GCN: SUB-GRAPH SAMPLING

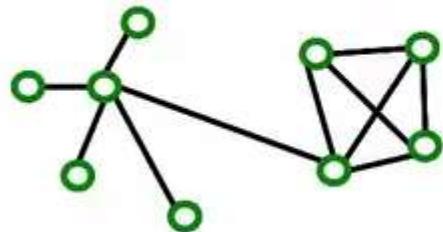
**Key question:** What subgraphs are good for training GNNs?

- Recall: GNN performs node embedding by passing messages [via the edges](#).
  - Subgraphs **should retain edge connectivity structure of the original graph as much as possible.**
  - This way, the GNN over the subgraph generates embeddings closer to the GNN over the original graph.

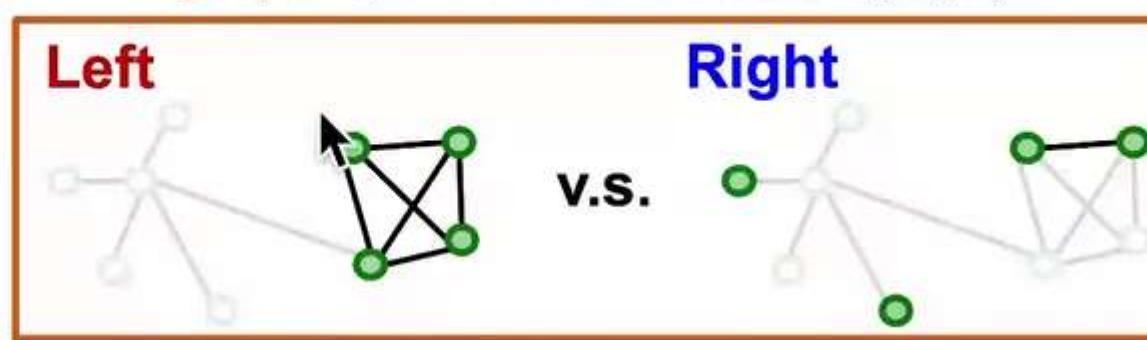
# CLUSTER-GCN: SUB-GRAPH SAMPLING

Which subgraph is good for training GNN?

Original graph



Subgraphs (both 4-node induced subgraph)



- **Left subgraph:**  
retains the essential community structure among the 4 nodes → **Good** ✓
- **Right subgraph:**  
drops many connectivity patterns, even leading to isolated nodes → **Bad** ✗

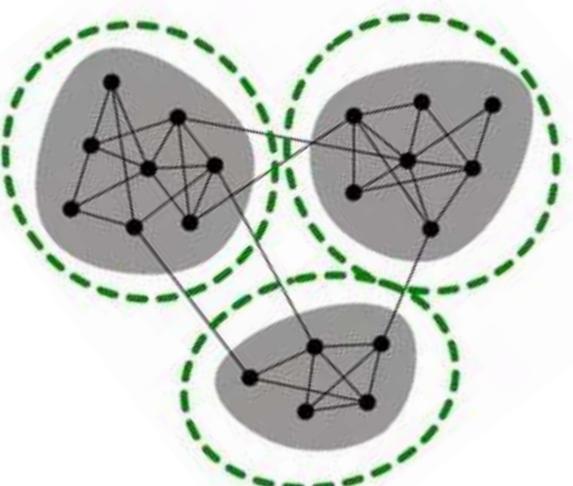
# CLUSTER-GCN: EXPLOITING COMMUNITY STRUCTURE

**Real-world graph exhibits community structure**

- A large graph can be decomposed into many small communities.

**Key insight** [Chiang et al. KDD 2019]:

- Sample a community as a subgraph.
- Each subgraph retains essential local connectivity pattern of the original graph.



# CLUSTER-GCN: OVERVIEW

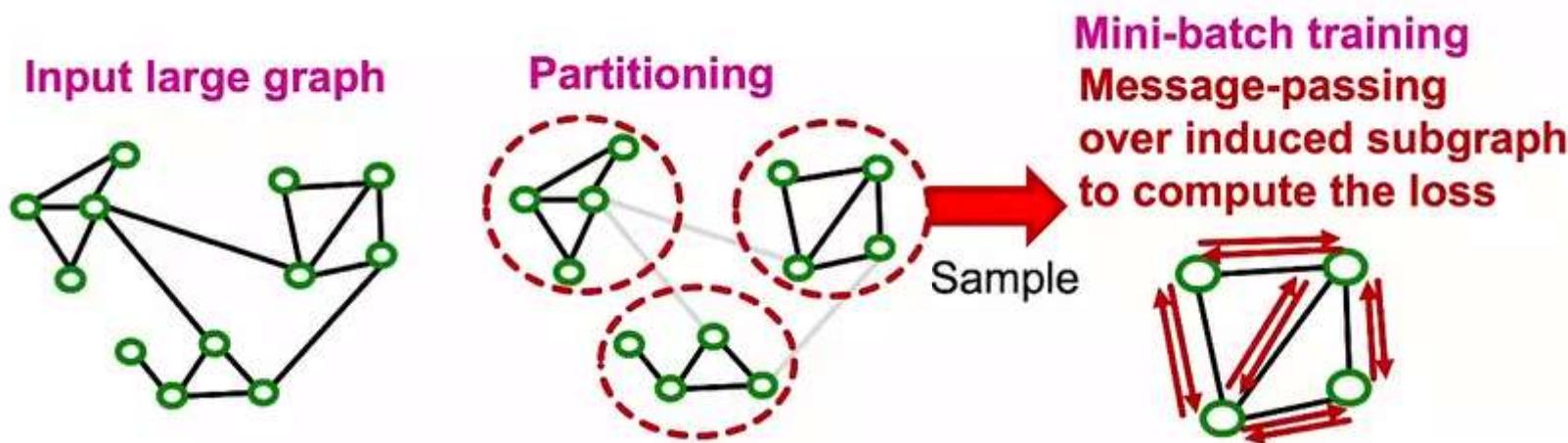
**Cluster-GCN consists of two steps:** It is a Vanilla cluster-GCN

## 1. Pre-processing:

Given a large graph, partition it into groups of nodes (i.e., subgraphs).

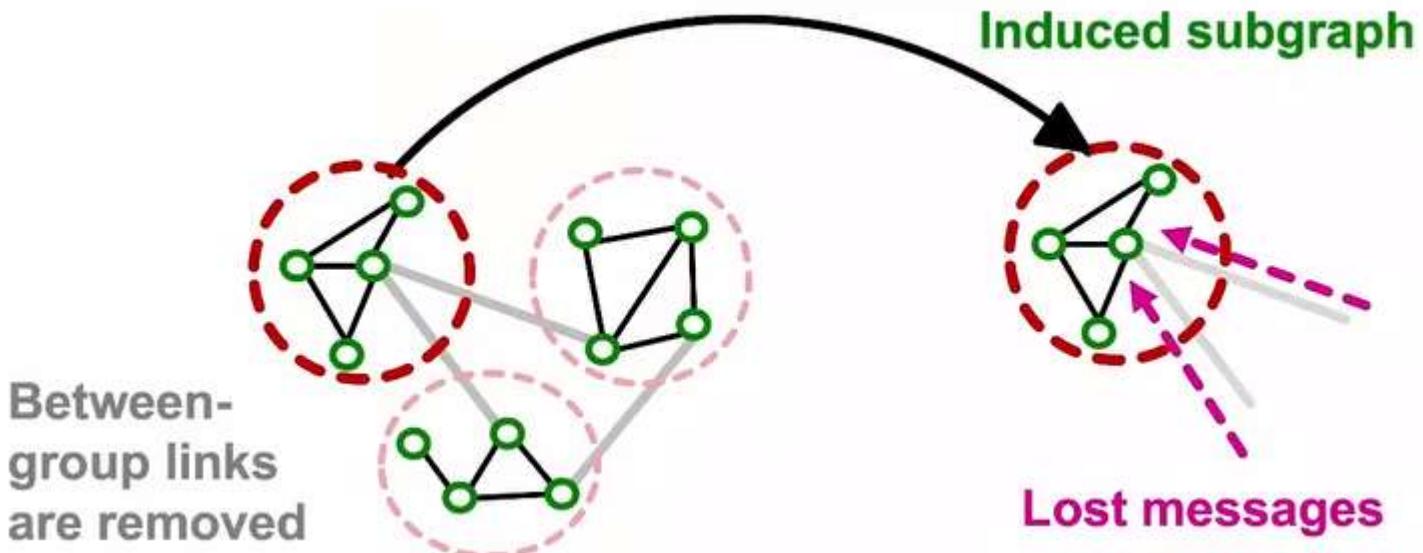
## 2. Mini-batch training:

Sample one node group at a time. Apply GNN's message passing over the induced subgraph.



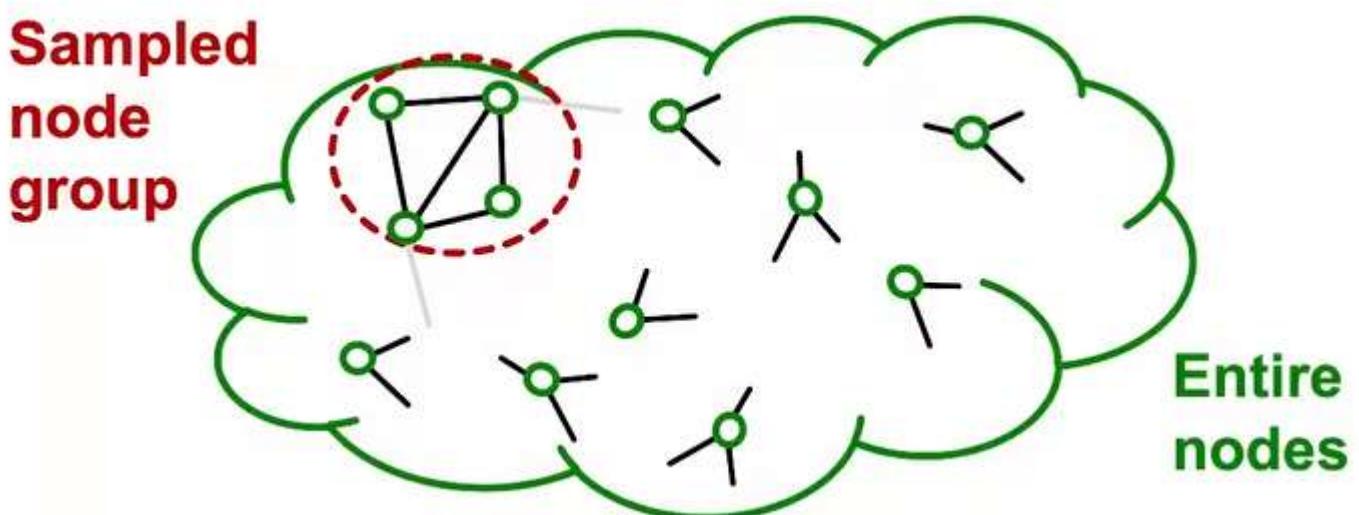
## CLUSTER-GCN: ISSUES(1)

- ❑ The induced subgraph **removes** between-group links.
- ❑ As a result, messages from other groups will be lost during message passing, which could hurt the GNN's performance.



## CLUSTER-GCN: ISSUES(2)

- ❑ Graph community detection algorithm puts similar nodes together in the same group.
- ❑ Sampled node group tends to only cover the small-concentrated portion of the entire data.



## ADVANCED CLUSTER-GCN: ISSUES(3)

**Sampled nodes are not diverse enough to be represent the graph structure:**

- ❑ As a result, the gradient averaged over the sampled nodes,  $\frac{1}{|V_c|} \sum_{v \in V_c} \nabla l_v(\theta)$ , becomes unreliable.
  - Fluctuates a lot from a node group to another.
  - In other words, the gradient has high variance.
- ❑ Leads to slow convergence of SGD

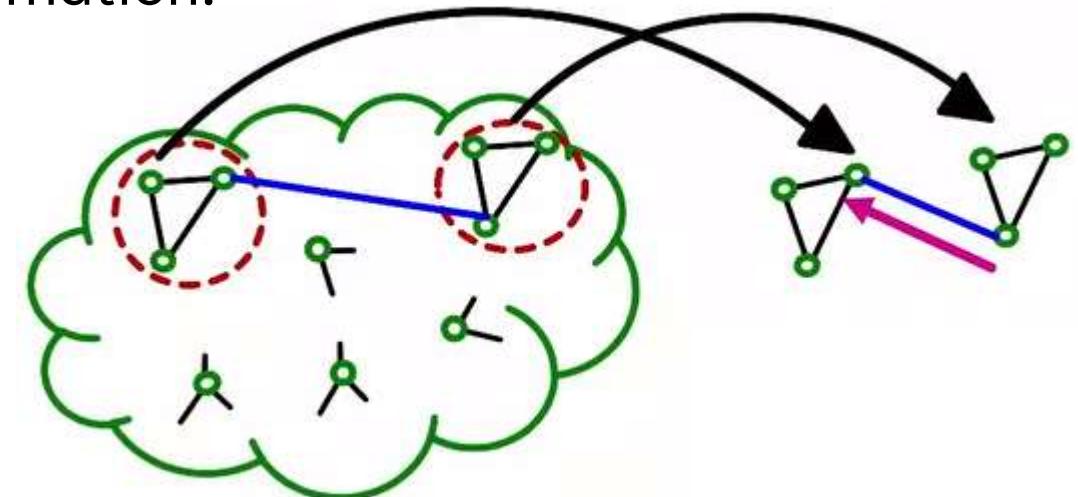
## ADVANCED CLUSTER-GCN

- ✓ **Solution:** Aggregate multiple node groups per mini-batch.
- ❑ Partition the graph into **relatively-small groups of nodes**.
- ❑ For each mini-batch:
  1. Sample and aggregate multiple node groups.
  2. Construct the induced subgraph of the aggregated node group.
  3. The rest is the same as vanilla Cluster-GCN (compute node embeddings and the loss, update parameters)

# ADVANCED CLUSTER-GCN

Why does the solution work?

- ❑ Sampling multiple node groups
  - Makes the sampled nodes more representative of the entire nodes.
  - Leads to less variance in gradient estimation.



- ❑ The induced subgraph over aggregated node groups
  - Includes between-group edges
  - Message can flow across groups.

## GRAPHSAGE VS CLUSTER-GCN

- Cluster-GCN is more computationally efficient than neighbour sampling, especially when #(GNN layers) is large.
- But Cluster-GCN leads to systematically biased gradient estimates (due to missing cross-community edges)

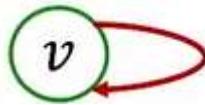
# SIMPLIFYING GNNS

- ❑ 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 all. 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.

## SIMPLIFYING GNNS: RECALL MEAN-POOL IN GCN

- ❑ 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 neighbouring 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

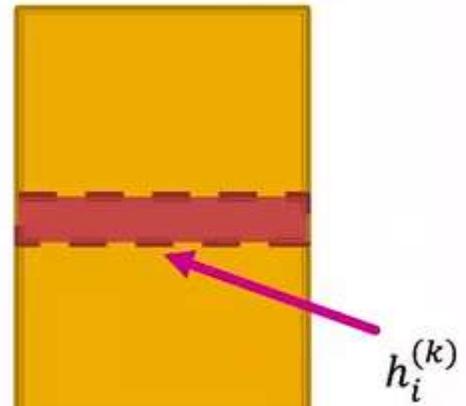
- ❑ Final node embedding:  $Z_v = h_v^{(k)}$

# SIMPLIFYING GNNS: RECALL MATRIX FORMULATION OF GCN

GCN aggregations can be formulated as matrix vector product:

Matrix of hidden embeddings  $\mathbf{H}^{(k)}$

- Let  $\mathbf{H}^{(k)} = [h_1^{(k)} \dots h_{|v|}^{(k)}]^T$
- Let  $\mathbf{A}$  be the adjacency matrix (w/ self-loop)
- Then:  $\sum_{u \in N(v)} h_u^{(k)} = \mathbf{A}_{v,:}\mathbf{H}^{(k)}$
- Let  $\mathbf{D}$  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,



$$\frac{1}{|N(v)|} \sum_{u \in N(v)} h_u^{(k)} \longrightarrow \mathbf{H}^{(l+1)} = \mathbf{D}^{-1} \mathbf{A} \mathbf{H}^{(l)}$$

## SIMPLIFYING GNNS: RECALL MATRIX FORMULATION OF GCN

GCN's neighbour aggregation:

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

In matrix form:

$$\mathbf{H}^{(k+1)} = \text{ReLU}(\tilde{\mathbf{A}} \mathbf{H}^{(k)} \mathbf{W}_k^T)$$

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

**Note:** The original GCN uses re-normalized version:  $\tilde{\mathbf{A}} = D^{-1/2} \mathbf{A} D^{-1/2}$

- Empirically, this version of  $\tilde{\mathbf{A}}$  often gives better performance than  $D^{-1} \mathbf{A}$

# SIMPLIFYING GNNS

Simplify GCN by removing ReLU non-linearity:

$$\mathbf{H}^{(k+1)} = \tilde{\mathbf{A}} \mathbf{H}^{(k)} \mathbf{W}_k^T$$

The final node embedding matrix is given as

$$\begin{aligned}\mathbf{H}^{(K)} &= \tilde{\mathbf{A}} \underbrace{\mathbf{H}^{(K-1)} \mathbf{W}_{K-1}^T}_{\text{Red bracket}} \\ &= \tilde{\mathbf{A}} (\tilde{\mathbf{A}} \underbrace{\mathbf{H}^{(K-2)} \mathbf{W}_{K-2}^T}_{\text{Green bracket}}) \mathbf{W}_{K-1}^T \\ &\dots = \tilde{\mathbf{A}} (\tilde{\mathbf{A}} (\dots (\tilde{\mathbf{A}} \underbrace{\mathbf{H}^{(0)} \mathbf{W}_0^T}_{\text{Blue bracket}} \dots) \mathbf{W}_{K-2}^T) \mathbf{W}_{K-1}^T \\ &= \tilde{\mathbf{A}}^K \mathbf{X} \underbrace{(\mathbf{W}_0^T \dots \mathbf{W}_{K-1}^T)}_{\text{Orange bracket}} \quad \text{Composition of linear transformation is still linear!} \\ &= \tilde{\mathbf{A}}^K \mathbf{X} \mathbf{W}^T \quad \text{where } \mathbf{W} \equiv \mathbf{W}_{K-1} \dots \mathbf{W}_0\end{aligned}$$

## SIMPLIFYING GNNS

- ❑ Removing ReLU significantly simplifies GCN!

$$H^{(K)} = \tilde{A}^K X W^T$$

- ❑ Notice  $\tilde{A}^K X$  does not contain any learnable parameters; hence,  
**it can be pre-computed.**
  - Efficiently computable as a sequence of sparse-matrix vector products:
  - Do  $X \leftarrow \tilde{A}X$  for K times.

## SIMPLIFYING GNNS

- Let  $\tilde{X} = \tilde{A}^K X$  be pre-computed matrix.

Simplified GCN's final embedding is

$$H^{(K)} = \tilde{X} W^T$$

- It's just a **linear transformation of pre-computed matrix!**
- Back to the node embedding form:

$$h_v^{(K)} = W \boxed{\tilde{X}_v}$$

Pre-computed feature vector for node  $v$

- Embedding of node  $v$  only depends on its own (pre-processed) feature!

## SIMPLIFYING GNNS

- Once  $\tilde{X}$  is pre-computed, embeddings of  $M$  nodes can be generated in time linear in  $M$ :
  - Given  $M$  nodes  $\{v_1, v_2, \dots, v_M\}$ , their embeddings are
    - $h_{v_1}^{(K)} = \mathbf{W}\tilde{\mathbf{X}}_{v_1},$
    - $h_{v_2}^{(K)} = \mathbf{W}\tilde{\mathbf{X}}_{v_2},$
    - $\dots$
    - $h_{v_M}^{(K)} = \mathbf{W}\tilde{\mathbf{X}}_{v_M}.$

## SIMPLIFYING GNNS

In summary, simplified GCN consists of **two steps**:

- **Pre-processing step:**

- Pre-compute  $\tilde{\mathbf{X}} = \tilde{\mathbf{A}}^K \mathbf{X}$ . Can be done on CPU.

- **Mini-batch training step:**

- For each mini-batch, randomly-sample  $M$  nodes  $\{v_1, v_2, \dots, v_M\}$ .
  - Compute their embeddings by
    - $h_{v_1}^{(K)} = \mathbf{W}\tilde{\mathbf{X}}_{v_1}, h_{v_2}^{(K)} = \mathbf{W}\tilde{\mathbf{X}}_{v_2}, \dots, h_{v_M}^{(K)} = \mathbf{W}\tilde{\mathbf{X}}_{v_M}$
  - Use the embeddings to make prediction and compute the loss averaged over the  $M$  data points.
  - Perform SGD parameter update.

# COMPARISON WITH OTHER MODELS

## ❑ Compared to neighbour sampling:

- Simplified GCN generates node embeddings much more efficiently (no need to construct the giant computational graph for each node).

## ❑ Compared to Cluster-GCN:

- Mini-batch nodes of simplified GCN can be sampled completely randomly from the entire nodes (no need to sample from multiple groups as Cluster-GCN does)
- Leads to lower SGD variance during training.

## ❑ But the model is much less expressive.

## COMPARISON WITH OTHER MODELS

**Compared to the original GN models, simplified GCN's expressive power is limited due to the lack of non-linearity in generating node embeddings.**

## COMPARISON WITH OTHER MODELS

**Compared to the original GN models, simplified GCN's expressive power is limited due to the lack of non-linearity in generating node embeddings.**

**Why the performance is good?**

<https://youtu.be/iTRW9Gh7yKI?list=PLoROMvodv4rPLKxIpqhjhPgdQy7imNkDn&t=880>



latest

Search docs

## INSTALL PYG

Installation

## GET STARTED

GINConv

Powerful are Graph Neural Networks?" paper

GINEConv

The modified `GINConv` operator from the "Strategies for Pre-training Graph Neural Networks" paper

ARMAConv

The ARMA graph convolutional operator from the "Graph Neural Networks with Convolutional ARMA Filters" paper

SGConv

The simple graph convolutional operator from the "Simplifying Graph Convolutional Networks" paper

SSGConv

The simple spectral graph convolutional operator from the "Simple Spectral Graph Convolution" paper

APPNP

The approximate personalized propagation of neural predictions layer from the "Predict then Propagate: Graph Neural Networks meet Personalized PageRank" paper

The graph neural network operator from the

<https://pytorch-geometric.readthedocs.io/en/latest/modules/nn.html>

# conv.SGConv

```
class SGConv ( in_channels: int, out_channels: int, K: int = 1, cached: bool =  
False, add_self_loops: bool = True, bias: bool = True, **kwargs )      [source]
```

Bases: `MessagePassing`

The simple graph convolutional operator from the “[Simplifying Graph Convolutional Networks](#)” paper

$$\mathbf{X}' = \left( \hat{\mathbf{D}}^{-1/2} \hat{\mathbf{A}} \hat{\mathbf{D}}^{-1/2} \right)^K \mathbf{X} \Theta,$$

where  $\hat{\mathbf{A}} = \mathbf{A} + \mathbf{I}$  denotes the adjacency matrix with inserted self-loops and  $\hat{D}_{ii} = \sum_{j=0} \hat{A}_{ij}$  its diagonal degree matrix. The adjacency matrix can include other values than `1` representing edge weights via the optional `edge_weight` tensor.

# SCALING UP GNNS VIA REMOTE BACKENDS

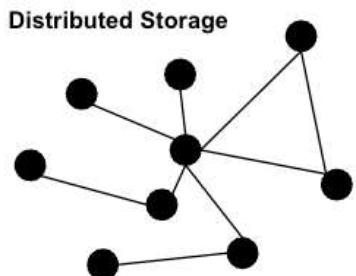
## ❑ Using key-value and graph database:

➤ Documentation:

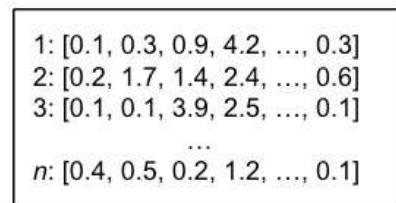
<https://pytorch-geometric.readthedocs.io/en/latest/advanced/remote.html>

➤ Example:

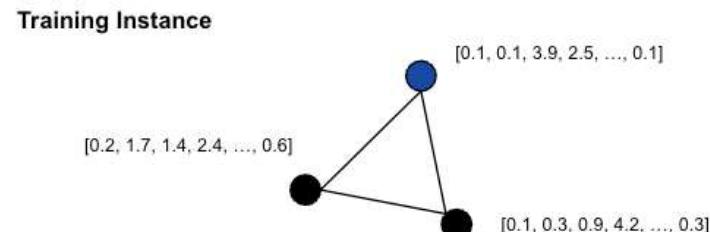
[https://github.com/pyg-team/pytorch\\_geometric/tree/master/examples/kuzu/papers\\_100M](https://github.com/pyg-team/pytorch_geometric/tree/master/examples/kuzu/papers_100M)



Graph Store: nodes and edges.



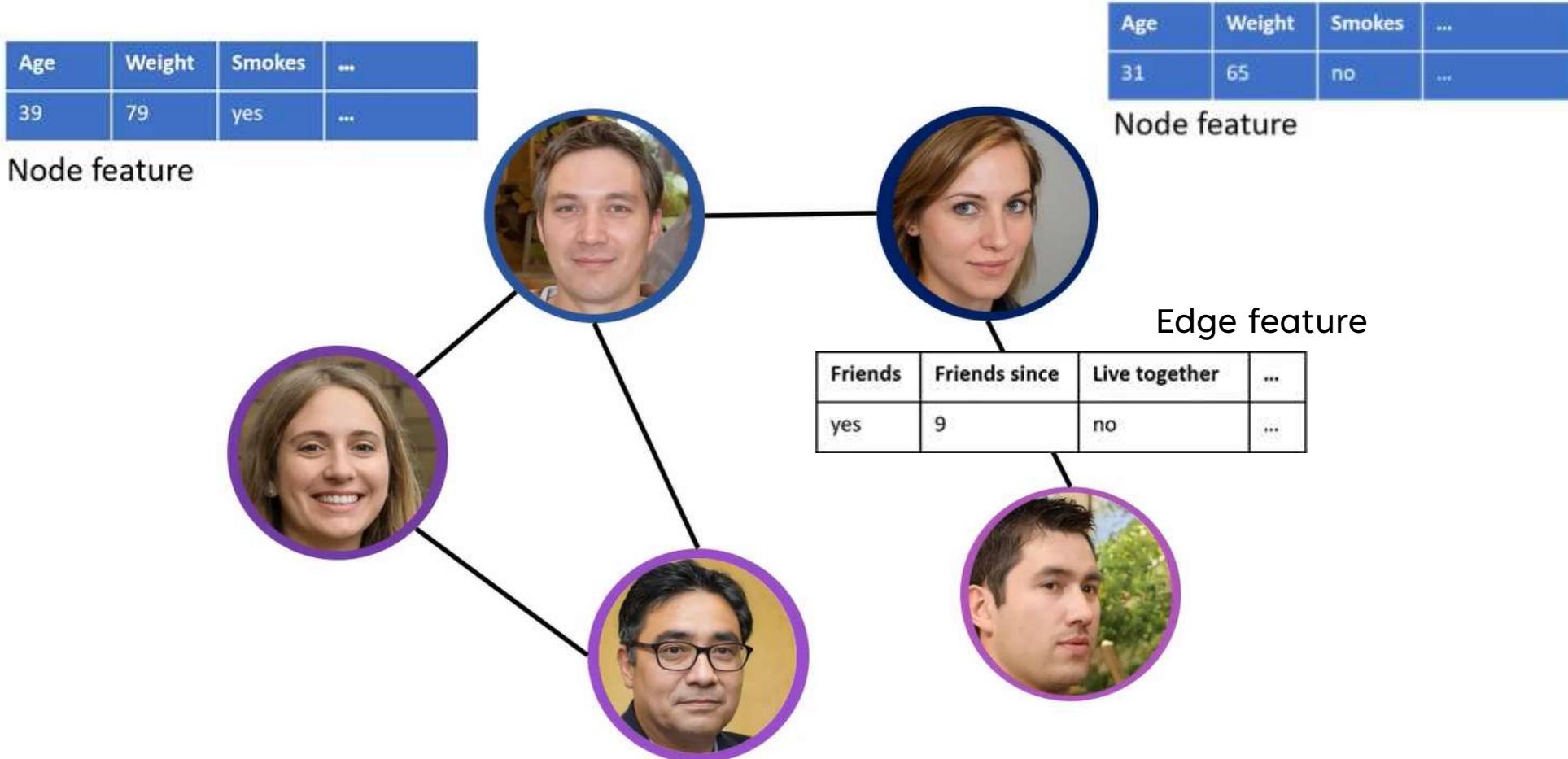
Feature store: node and edge tensors



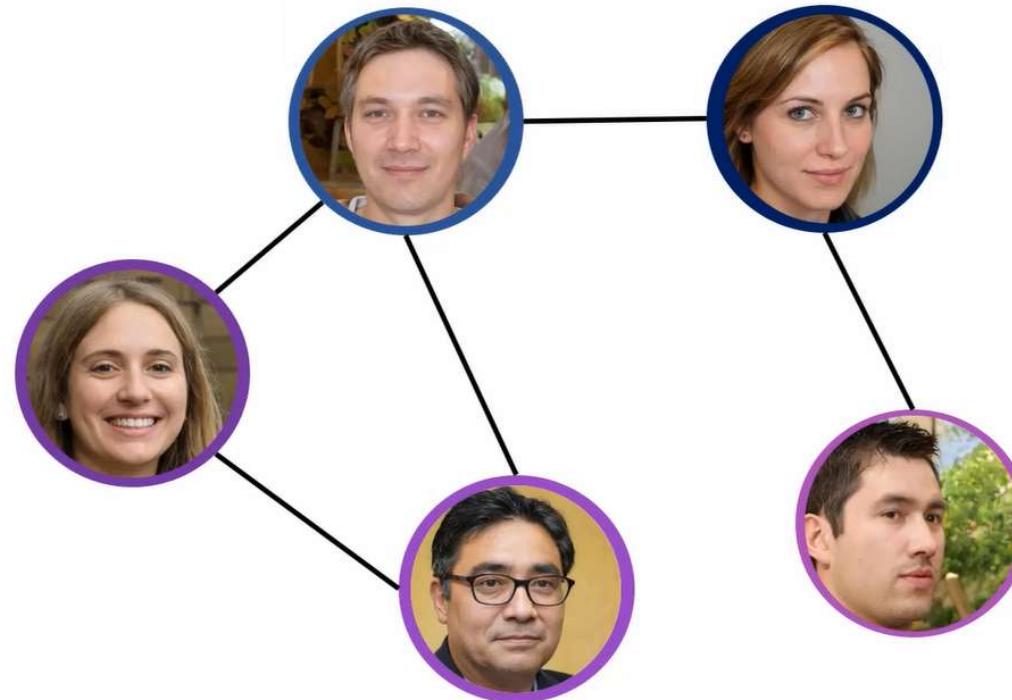
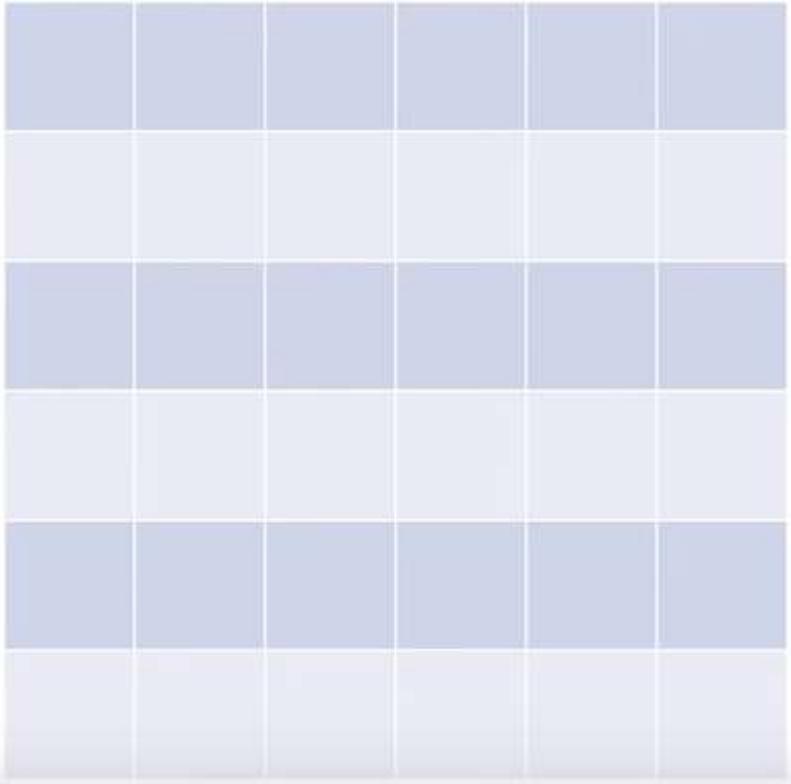
Sampled subgraph, joined with  
features; all that is necessary for  
forward/backward.

# EDGE FEATURES

# WHY ARE EDGE FEATURES ARE IMPORTANT?

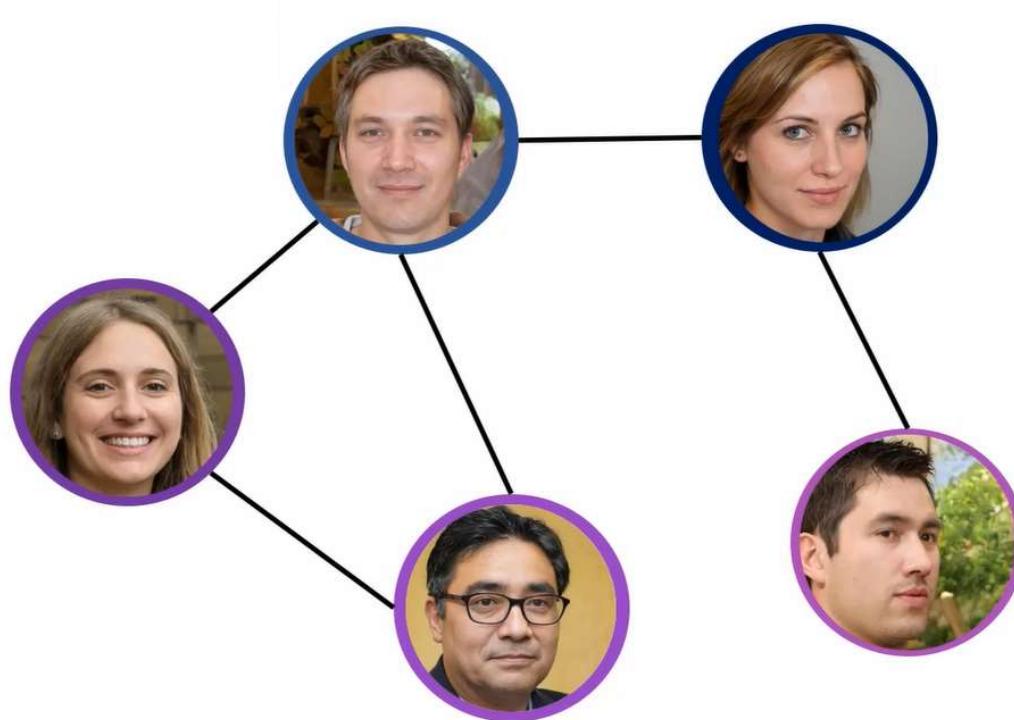


# WHY ARE EDGE FEATURES ARE IMPORTANT?



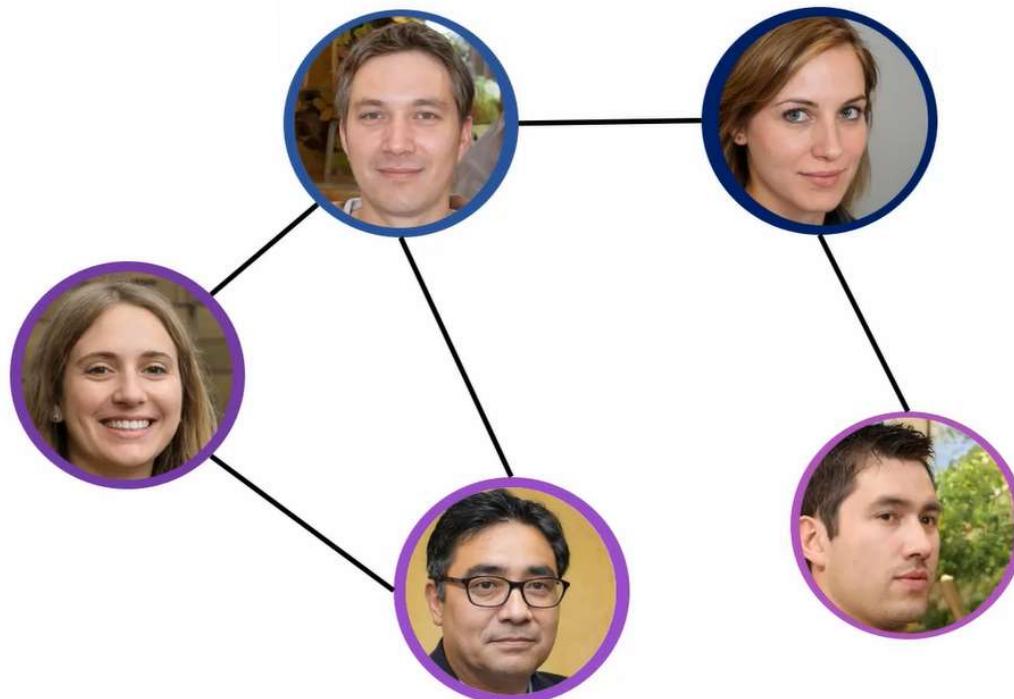
# WHY ARE EDGE FEATURES ARE IMPORTANT?

	1	2	3	4	5	6
1						
2	X	✓	X	X	✓	X
3	✓	X	✓	✓		X
4	X	✓	X	X		✓
5	✓	✓	X	X		X
6	X	X	✓	X		X

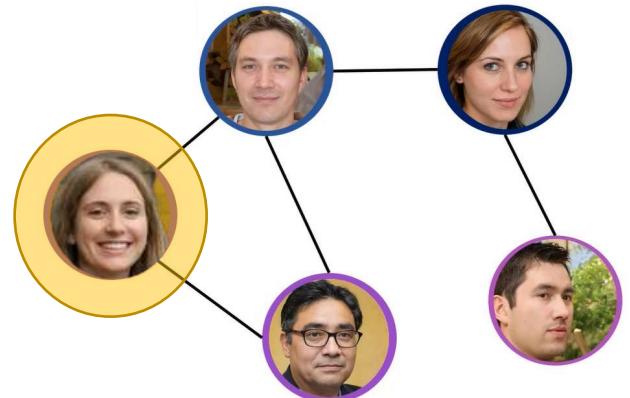
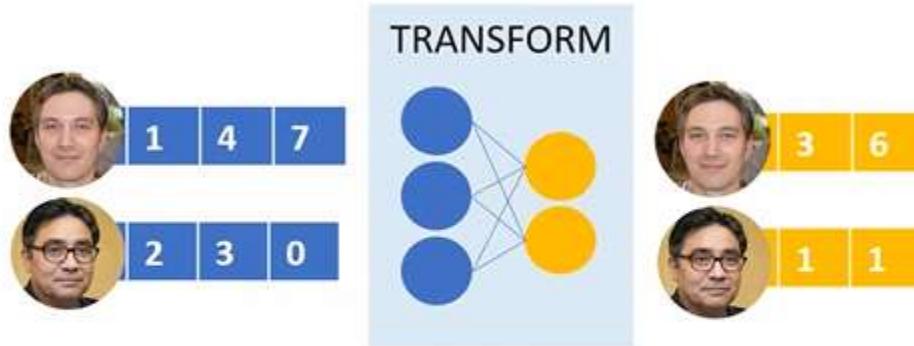


# WHY ARE EDGE FEATURES ARE IMPORTANT?

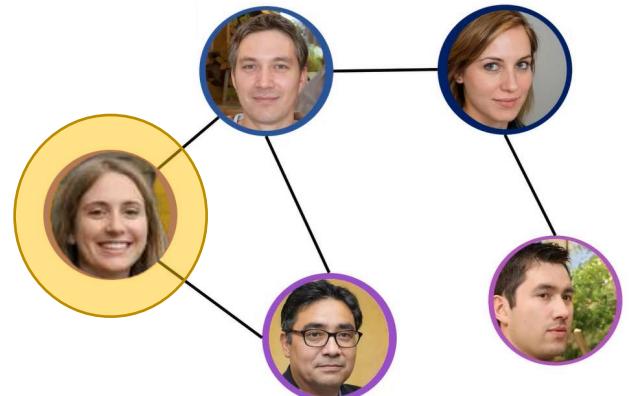
	0	1	0	1	0
	1	0	1	1	0
	0	1	0	0	1
	1	1	0	0	0
	0	0	1	0	0



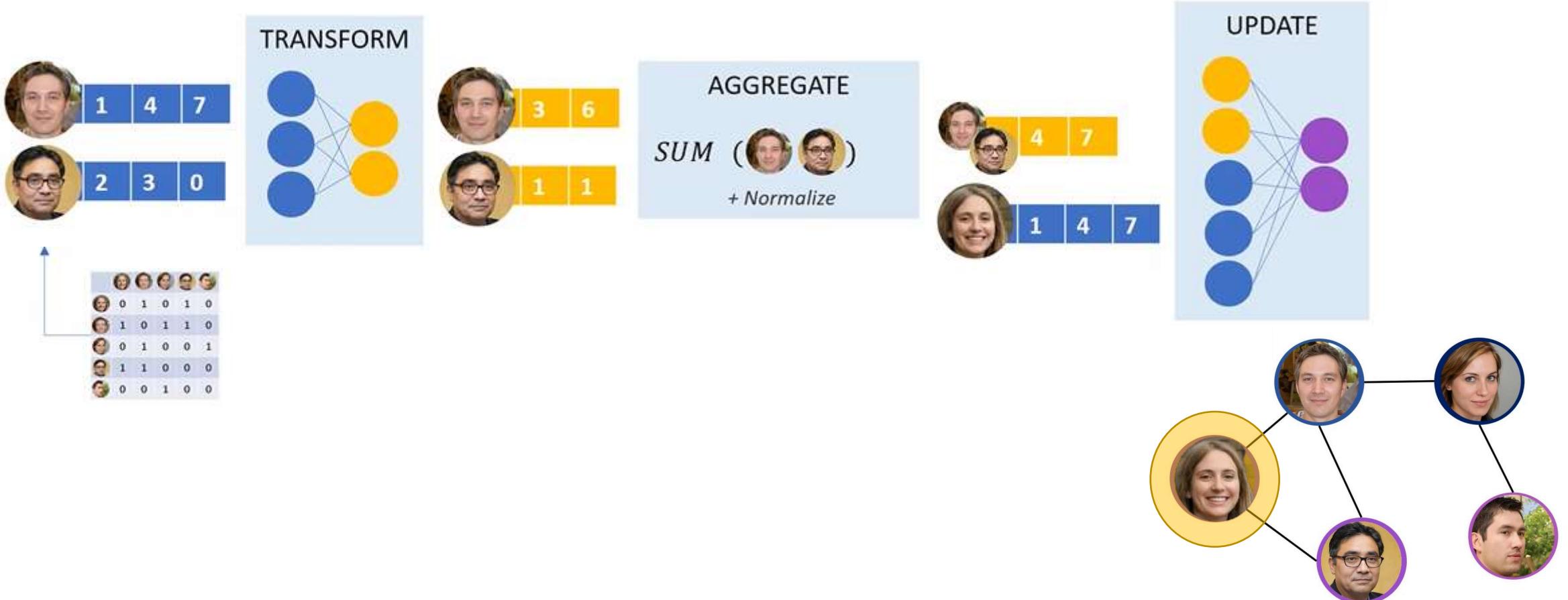
# THE GENERAL PROCESS IN GNNS



# THE GENERAL PROCESS IN GNNS



# THE GENERAL PROCESS IN GNNS



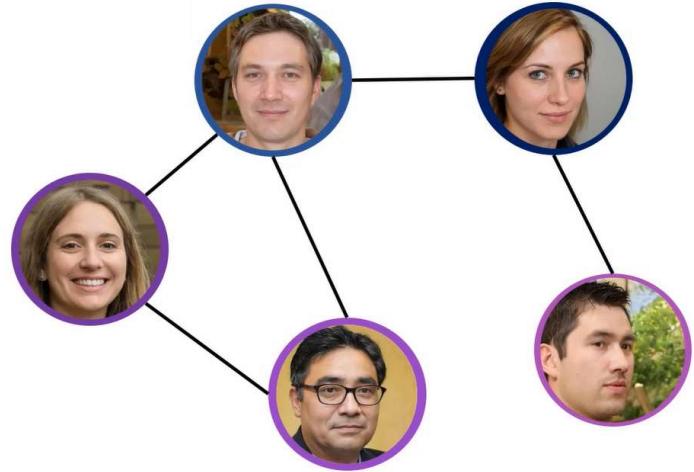
# USING EDGE WEIGHT

	1	2	3	4	5
1	0	1	0	1	0
2	1	0	1	1	0
3	0	1	0	0	1
4	1	1	0	0	0
5	0	0	1	0	0

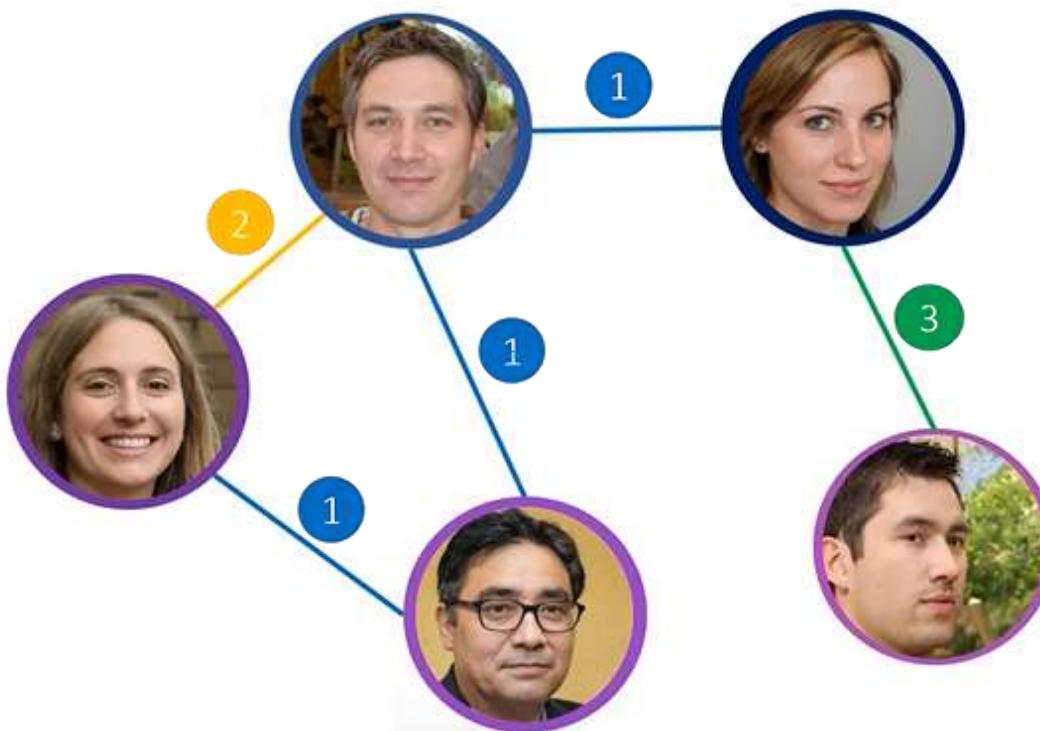
	1	2	3	4	5
1	0	0.4	0	0.5	0
2	0.9	0	1	1	0
3	0	1	0	0	0.5
4	1	0.7	0	0	0
5	0	0	0.1	0	0

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

Node Features/embeddings

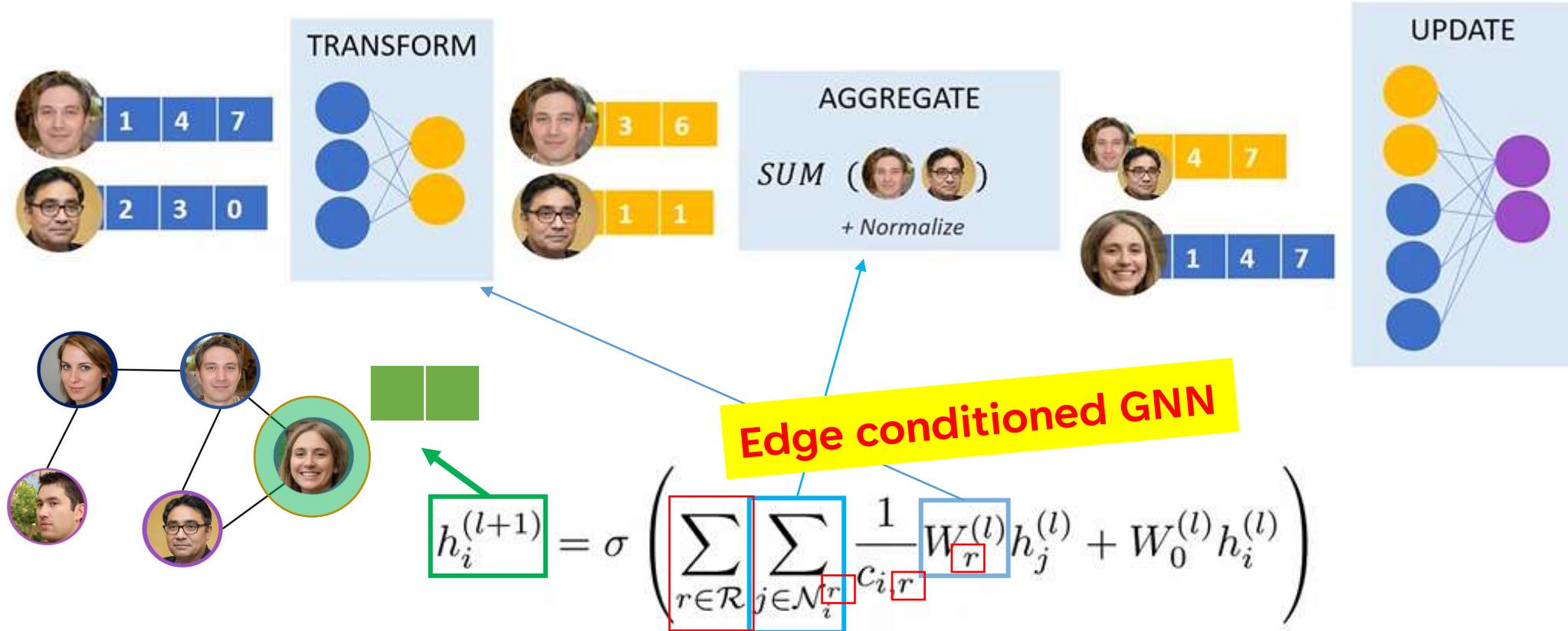


# DIFFERENT EDGE TYPES



1 = Friends  
2 = Couple  
3 = Colleagues

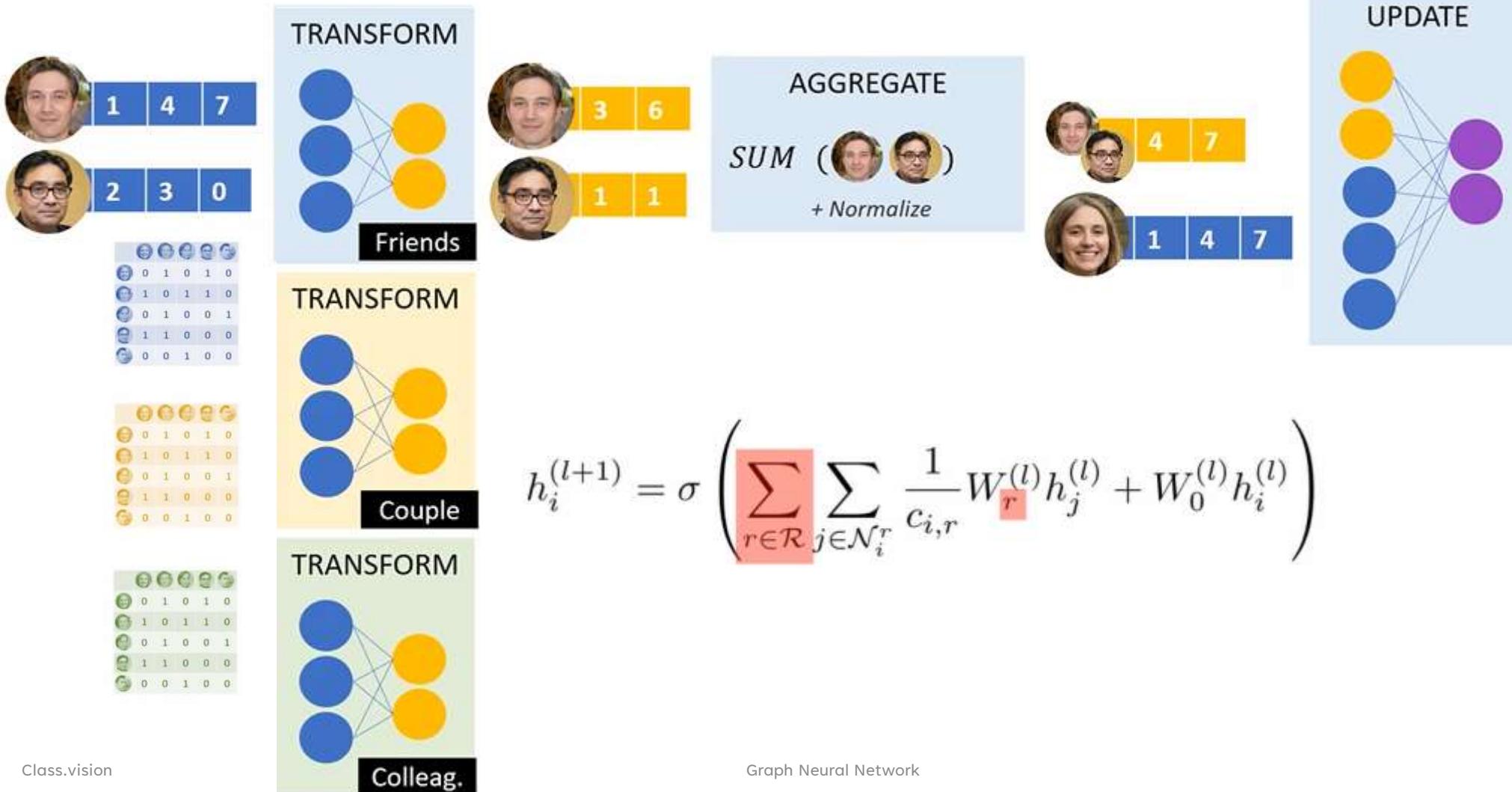
# DIFFERENT EDGE TYPES – RELATIONAL GCN



## Relational GCN

Modelling Relational Data with Graph Convolutional Networks, Schlichtkrull et al.

# DIFFERENT EDGE TYPES – RELATIONAL GCN



APPNP

The approximate personalized propagation of neural predictions layer from the "["Predict then Propagate: Graph Neural Networks meet Personalized PageRank"](#) paper

MFConv

The graph neural network operator from the "["Convolutional Networks on Graphs for Learning Molecular Fingerprints"](#) paper

RGCNConv

The relational graph convolutional operator from the "["Modeling Relational Data with Graph Convolutional Networks"](#) paper

FastRGCNConv

See [RGCNConv](#).

CuGraphRGCNConv

The relational graph convolutional operator from the "["Modeling Relational Data with Graph Convolutional Networks"](#) paper.

<https://pytorch-geometric.readthedocs.io/en/latest/modules/nn.html>

## conv.RGCNConv

```
class RGCNConv ( in_channels: Union[int, Tuple[int, int]], out_channels: int,
    num_relations: int, num_bases: Optional[int] = None, num_blocks: Optional[int] = None, aggr: str = 'mean', root_weight: bool = True, is_sorted: bool = False,
    bias: bool = True, **kwargs ) [source]
```

Bases: `MessagePassing`

The relational graph convolutional operator from the "Modeling Relational Data with Graph Convolutional Networks" paper

$$\mathbf{x}'_i = \Theta_{\text{root}} \cdot \mathbf{x}_i + \sum_{r \in \mathcal{R}} \sum_{j \in \mathcal{N}_r(i)} \frac{1}{|\mathcal{N}_r(i)|} \Theta_r \cdot \mathbf{x}_j,$$

where  $\mathcal{R}$  denotes the set of relations, i.e. edge types. Edge type needs to be a one-dimensional `torch.long` tensor which stores a relation identifier  $\in \{0, \dots, |\mathcal{R}| - 1\}$  for each edge.

[https://pytorch-geometric.readthedocs.io/en/latest/generated/torch\\_geometric.nn.conv.RGCNConv.html#torch\\_geometric.nn.conv.RGCNConv](https://pytorch-geometric.readthedocs.io/en/latest/generated/torch_geometric.nn.conv.RGCNConv.html#torch_geometric.nn.conv.RGCNConv)

## conv.FastRGCNConv

```
class FastRGCNConv ( in_channels: Union[int, Tuple[int, int]], out_channels: int,
    num_relations: int, num_bases: Optional[int] = None, num_blocks: Optional[int] = None, aggr: str = 'mean', root_weight: bool = True, is_sorted: bool = False,
    bias: bool = True, **kwargs ) [source]
```

Bases: `RGCNConv`

See `RGCNConv`.

```
forward ( x: Union[Tensor, None, Tuple[Optional[Tensor], Tensor]], edge_index:
    Union[Tensor, SparseTensor], edge_type: Optional[Tensor] = None ) [source]
```

Runs the forward pass of the module.

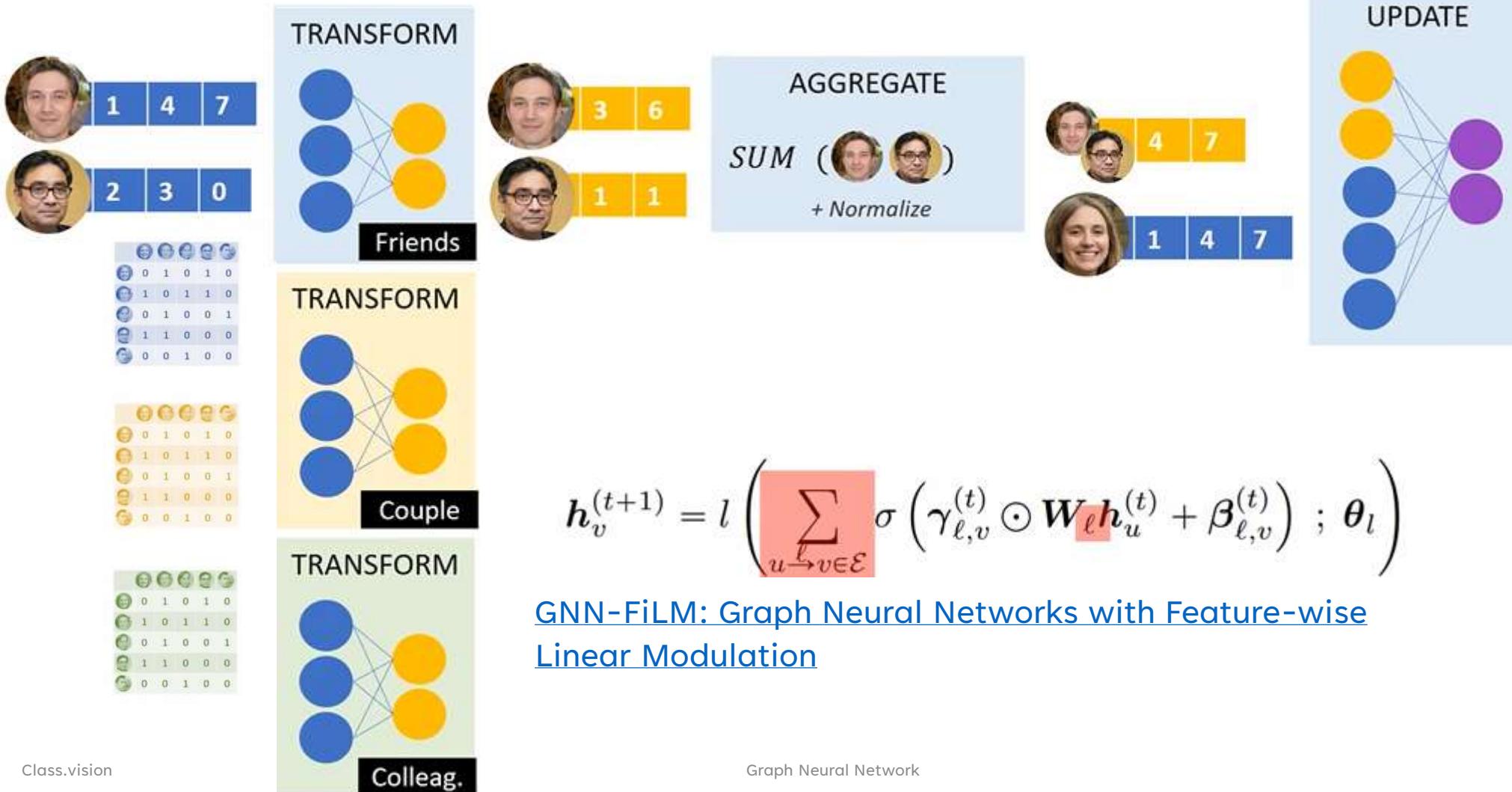
### PARAMETERS

- `x (torch.Tensor or tuple, optional)` – The input node features. Can be either a `[num_nodes, in_channels]` node feature matrix, or an optional one-dimensional node index tensor (in which case input features

[https://pytorch-geometric.readthedocs.io/en/latest/generated/torch\\_geometric.nn.conv.FastRGCNConv.html#torch\\_geometric.nn.conv.FastRGCNConv](https://pytorch-geometric.readthedocs.io/en/latest/generated/torch_geometric.nn.conv.FastRGCNConv.html#torch_geometric.nn.conv.FastRGCNConv)

Example: [https://github.com/pyg-team/pytorch\\_geometric/blob/master/examples/rgcn.py](https://github.com/pyg-team/pytorch_geometric/blob/master/examples/rgcn.py)

# DIFFERENT EDGE TYPES – GNN FiLM





## GNN-FiLM: Graph Neural Networks with Feature-wise Linear Modulation

Marc Brockschmidt

## GNN-FiLM

Graph Neural Networks with Feature-wise Linear Modulation

Marc Brockschmidt <mabrocks@microsoft.com>

@mmjb86



## GNN-FiLM: Graph Neural Networks with Feature-wise Linear Modulation

Jul 12, 2020

Soll diese Präsentation für 1000 Jahre gespeichert werden?  
Wie speichern wir Präsentationen?



<https://slideslive.com/38927627/gnnfilm-graph-neural-networks-with-featurewise-linear-modulation?ref=recommended>

NLConv

The Weisfeiler Lehman operator from the "["A Reduction of a Graph to a Canonical Form and an Algebra Arising During this Reduction"](#) paper, which iteratively refines node colorings:

NLConvContinuous

The Weisfeiler Lehman operator from the "["Wasserstein Weisfeiler-Lehman Graph Kernels"](#) paper.

FiLMConv

The FiLM graph convolutional operator from the "["GNN-FiLM: Graph Neural Networks with Feature-wise Linear Modulation"](#) paper

SuperGATConv

The self-supervised graph attentional operator from the "["How to Find Your Friendly Neighborhood: Graph Attention Design with Self-Supervision"](#) paper

FACConv

The Frequency Adaptive Graph Convolution operator from the "["Beyond Low-Frequency Information in Graph Convolutional Networks"](#) paper

# conv.FiLMConv

```
class FiLMConv ( in_channels: Union[int, Tuple[int, int]], out_channels: int,
    num_relations: int = 1, nn: Optional[Callable] = None, act: Optional[Callable] =
    ReLU(), aggr: str = 'mean', **kwargs )      [source]
```

Bases: `MessagePassing`

The FiLM graph convolutional operator from the “[GNN-FiLM: Graph Neural Networks with Feature-wise Linear Modulation](#)” paper

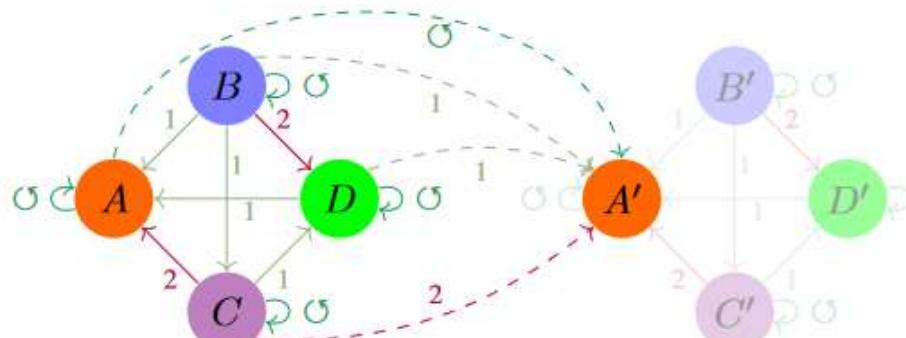
$$\mathbf{x}'_i = \sum_{r \in \mathcal{R}} \sum_{j \in \mathcal{N}(i)} \sigma(\boldsymbol{\gamma}_{r,i} \odot \mathbf{W}_r \mathbf{x}_j + \boldsymbol{\beta}_{r,i})$$

where  $\boldsymbol{\beta}_{r,i}, \boldsymbol{\gamma}_{r,i} = g(\mathbf{x}_i)$  with  $g$  being a single linear layer by default. Self-loops are automatically added to the input graph and represented as its own relation type.

[https://pytorch-geometric.readthedocs.io/en/latest/generated/torch\\_geometric.nn.conv.FiLMConv](https://pytorch-geometric.readthedocs.io/en/latest/generated/torch_geometric.nn.conv.FiLMConv.html#torch_geometric.nn.conv.FiLMConv)

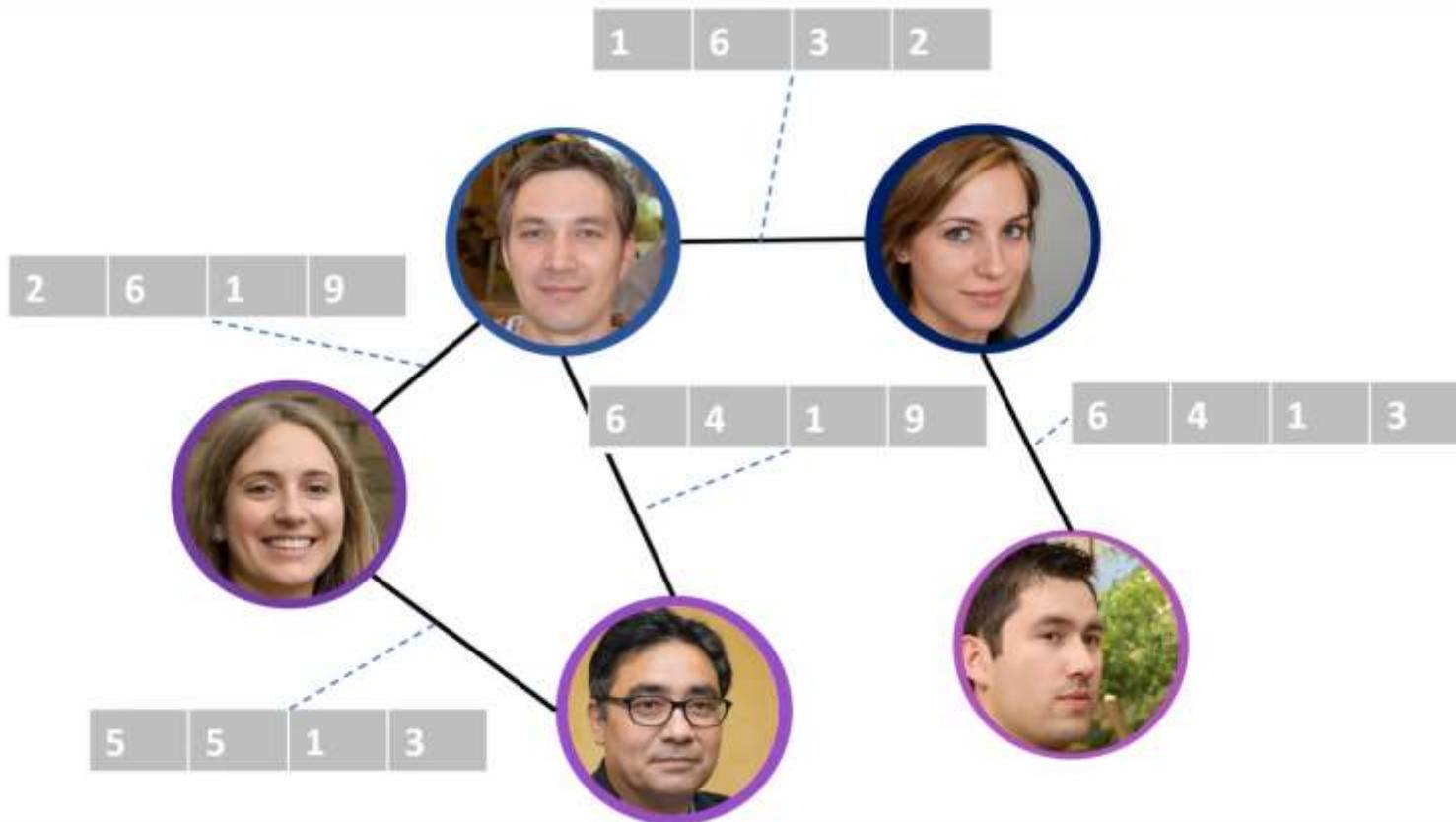
# DIFFERENT EDGE TYPES- OTHER VARIANTS

GGNN: $A' = GRU(A)$	$A$	$W_1 \cdot B + W_2 \cdot C + W_1 \cdot D$
R-GCN: $A' = \sigma(W_G \cdot A + W_B \cdot B + W_C \cdot C + W_D \cdot D)$	$W_G \cdot A + W_B \cdot B + W_C \cdot C + W_D \cdot D$	$W_2 \cdot C + W_1 \cdot D$
R-GAT: $A' = \sigma((a_{A'})_{A \xrightarrow{1} A} \cdot W_G \cdot A + (a_{A'})_{B \xrightarrow{1} A} \cdot W_1 \cdot B + (a_{A'})_{C \xrightarrow{2} A} \cdot W_2 \cdot C + (a_{A'})_{D \xrightarrow{1} A} \cdot W_1 \cdot D)$	$(a_{A'})_{A \xrightarrow{1} A} \cdot W_G \cdot A + (a_{A'})_{B \xrightarrow{1} A} \cdot W_1 \cdot B + (a_{A'})_{C \xrightarrow{2} A} \cdot W_2 \cdot C + (a_{A'})_{D \xrightarrow{1} A} \cdot W_1 \cdot D$	$W_2 \cdot C + W_1 \cdot D$
R-GIN: $A' = \sigma(MLP_G(A) + MLP_1(B) + MLP_2(C) + MLP_1(D))$	$MLP_G(A) + MLP_1(B) + MLP_2(C) + MLP_1(D)$	$MLP_1(D \parallel A) + MLP_1(D \parallel B) + MLP_1(D \parallel C) + MLP_1(D \parallel D)$
GNN-MLP: $A' = \sigma(MLP_G(A \parallel A) + MLP_1(B \parallel A) + MLP_2(C \parallel A) + MLP_1(D \parallel A))$	$MLP_G(A \parallel A) + MLP_1(B \parallel A) + MLP_2(C \parallel A) + MLP_1(D \parallel A)$	$MLP_1(D \parallel A) + MLP_1(D \parallel B) + MLP_1(D \parallel C) + MLP_1(D \parallel D)$
RGDCN: $A' = \sigma(W_{G,A} \cdot A + W_{1,A} \cdot B + W_{2,A} \cdot C + W_{1,D} \cdot D)$	$W_{G,A} \cdot A + W_{1,A} \cdot B + W_{2,A} \cdot C + W_{1,D} \cdot D$	$W_1 \cdot B + W_2 \cdot C + W_1 \cdot D$
GNN-FiLM: $A' = \sigma(\beta_{G,A} + \gamma_{G,A} \odot W_G \cdot A + \beta_{1,A} + \gamma_{1,A} \odot W_1 \cdot B + \beta_{2,A} + \gamma_{2,A} \odot W_2 \cdot C + \beta_{1,D} + \gamma_{1,D} \odot W_1 \cdot D)$	$\beta_{G,A} + \gamma_{G,A} \odot W_G \cdot A + \beta_{1,A} + \gamma_{1,A} \odot W_1 \cdot B + \beta_{2,A} + \gamma_{2,A} \odot W_2 \cdot C + \beta_{1,D} + \gamma_{1,D} \odot W_1 \cdot D$	$\beta_{G,D} + \gamma_{G,D} \odot W_G \cdot D + \beta_{1,D} + \gamma_{1,D} \odot W_1 \cdot D$



<https://arxiv.org/pdf/1906.12192.pdf>

# MULTIDIMENSIONAL EDGE FEATURES



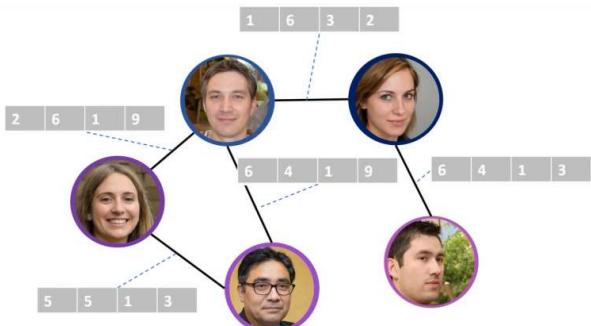
# MULTIDIMENSIONAL EDGE FEATURES

$$h_v = \gamma \left( x_v, \bigoplus_{w \in N(v)} \phi(x_v, x_w, e_{vw}) \right)$$

Diagram illustrating the computation of node feature  $h_v$ :

- UPDATE:** An upward arrow from the node feature  $x_v$  to the first term  $x_v$  in the equation.
- AGGREGATE:** A downward arrow from the summation symbol ( $\bigoplus$ ) to the second term  $w \in N(v)$ .
- TRANSFORM:** A downward arrow from the summation symbol ( $\phi$ ) to the third term  $e_{vw}$ .

The result is enclosed in a large bracket labeled  $\gamma$ . To the right, a yellow box highlights the edge feature  $e_{vw}$ , which is shown as a 5x5 matrix of portraits. Below the matrix is the text:  $Shape = [\#N, \#N, F_{edge}]$ .



## MULTIDIMENSIONAL EDGE FEATURES: MP-GNN

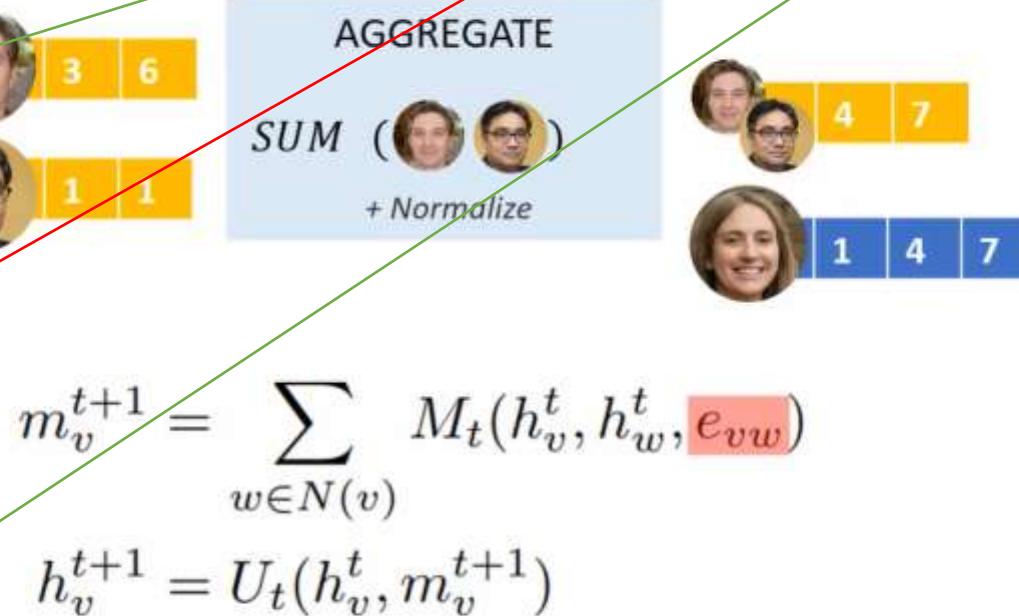
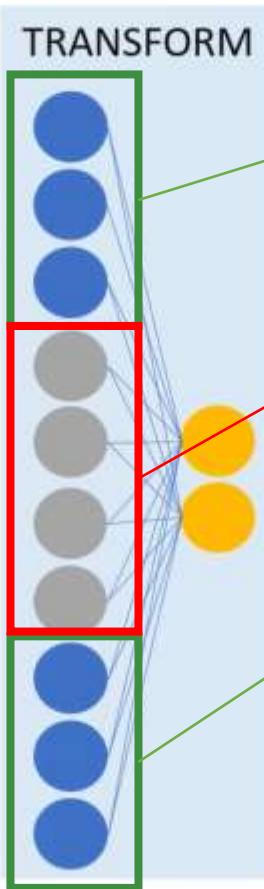
$$m_v^{t+1} = \sum_{w \in N(v)} M_t(h_v^t, h_w^t, e_{vw})$$

$$h_v^{t+1} = U_t(h_v^t, m_v^{t+1})$$

### MP-GNN

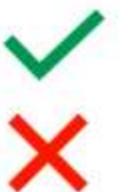
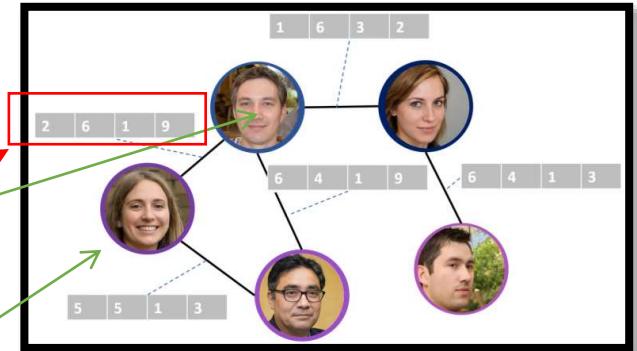
Neural Message Passing for Quantum Chemistry  
Gilmer et al.

# MP-GNN



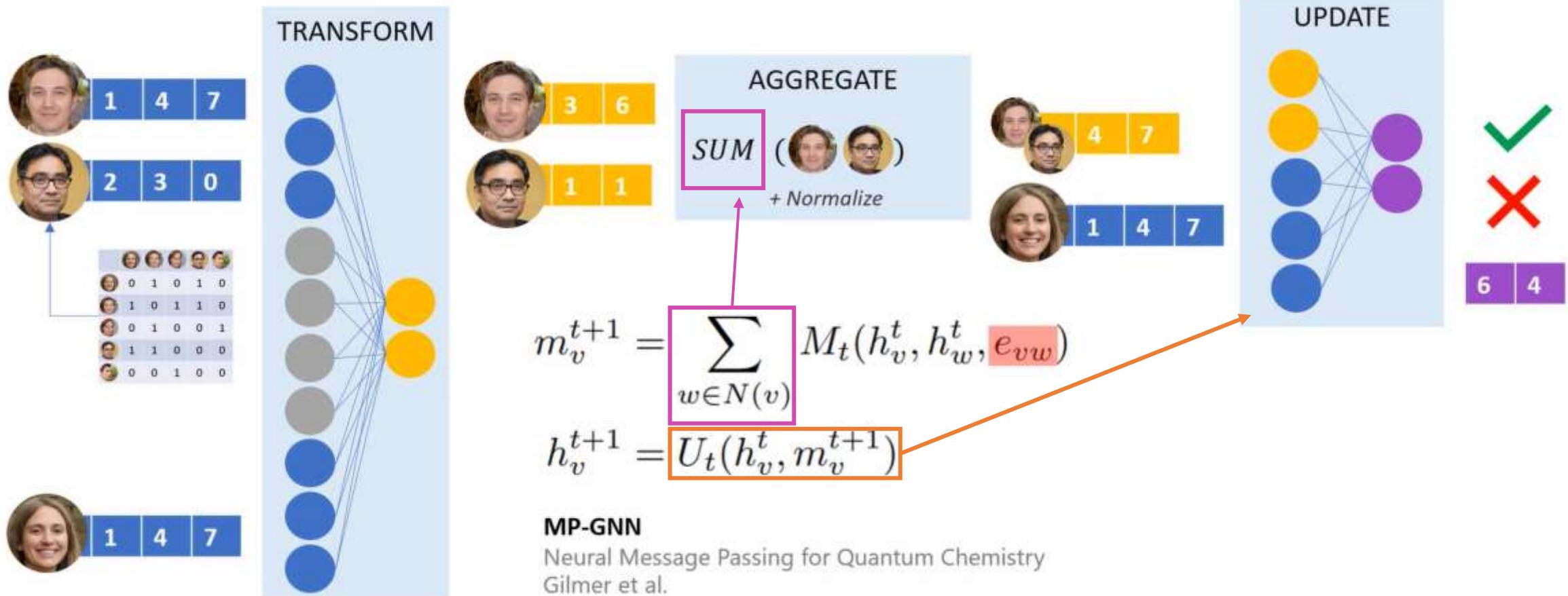
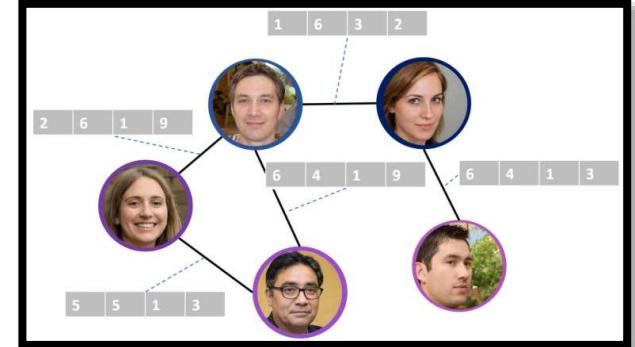
**MP-GNN**

Neural Message Passing for Quantum Chemistry  
Gilmer et al.

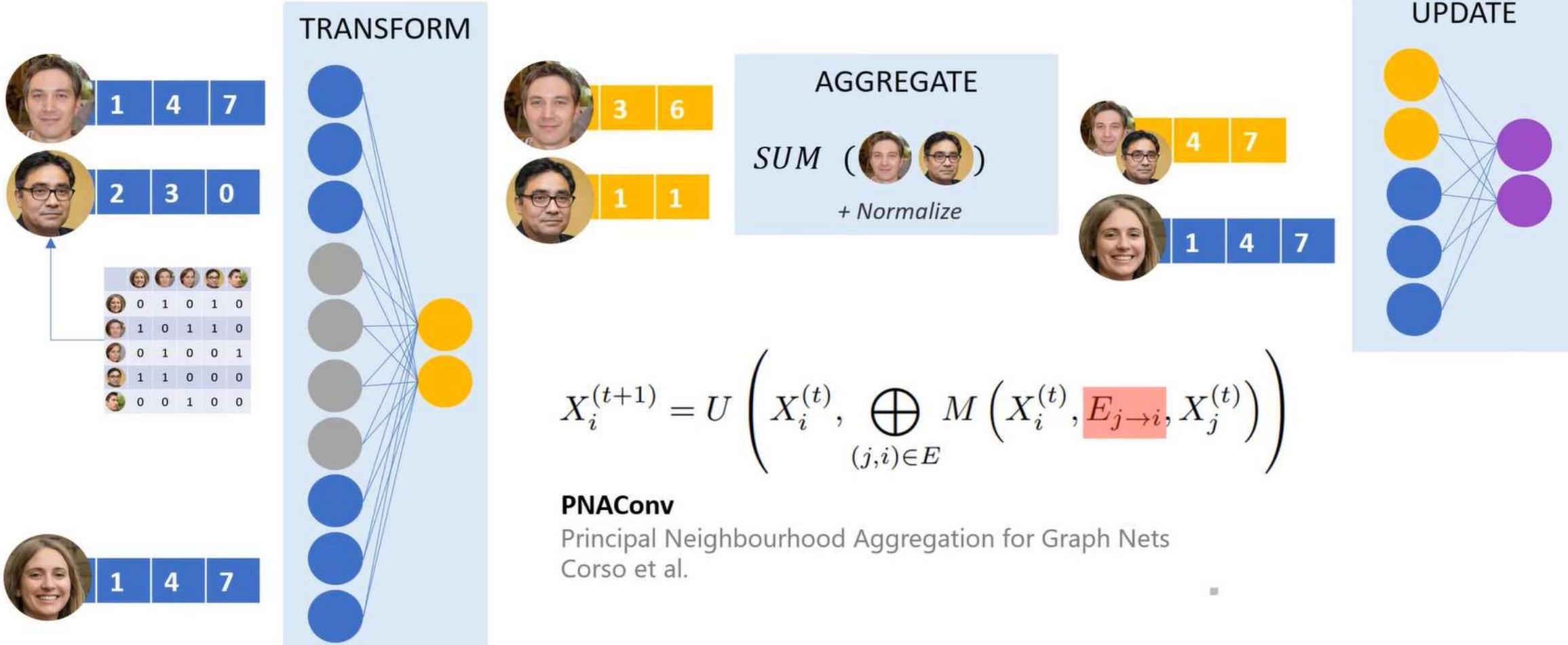


6 | 4

# MP-GNN



# MULTIDIMENSIONAL EDGE FEATURES: PNACONV



# MULTIDIMENSIONAL EDGE FEATURES

## OTHER EXAMPLES

A collage of academic papers from the Graph Neural Network (GNN) field, including titles like 'Graph Convolutional Networks for Graphs with Multi-Dimensionally Weighted Edges' and 'SplineCNN: Fast Geometric Deep Learning with Continuous B-Spline Kernels'. The collage includes various paper covers with titles, authors, and abstracts.

# USING EDGE FEATURES IN PYTORCH GEOMETRIC



PyTorch  
geometric

The screenshot shows the PyTorch Geometric documentation for the `torch_geometric.nn` module. The page title is `torch_geometric.nn`. The contents section lists various layers and components:

- Convolutional Layers
- Aggregation Operators
- Normalization Layers
- Pooling Layers
- Unpooling Layers
- Models
- KGE Models
- Encodings
- Functional

The sidebar includes links for `INSTALL PYG`, `Installation`, `GET STARTED`, and `Read the Docs`.

# USING EDGE FEATURES IN PYTORCH GEOMETRIC



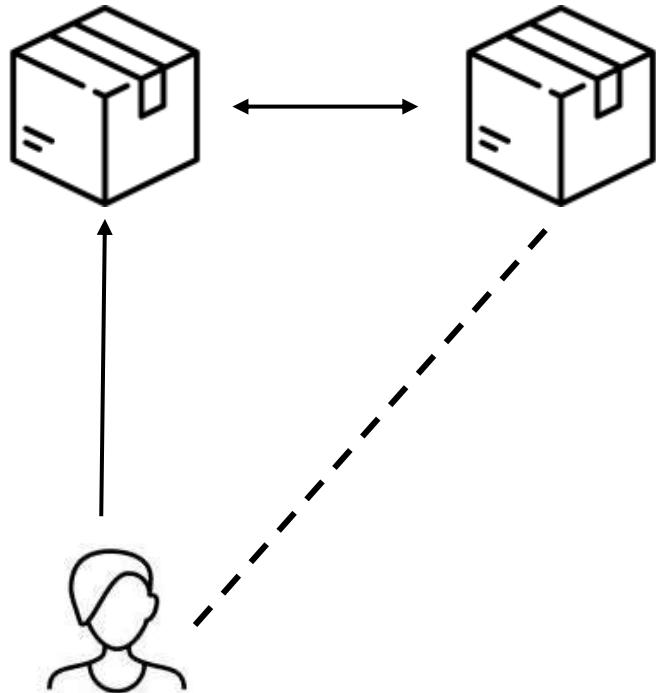
- **edge\_weight** → GNN Layer can use weight values on the adjacency matrix
- **edge\_type** → GNN Layer can use different edge types / relations
- **edge\_attr** → GNN Layer can use edge features

```
forward ( x: Union[Tensor, None, Tuple[Optional[Tensor], Tensor]], edge_index:  
Union[Tensor, SparseTensor], edge_type: Optional[Tensor] = None ) [source]
```

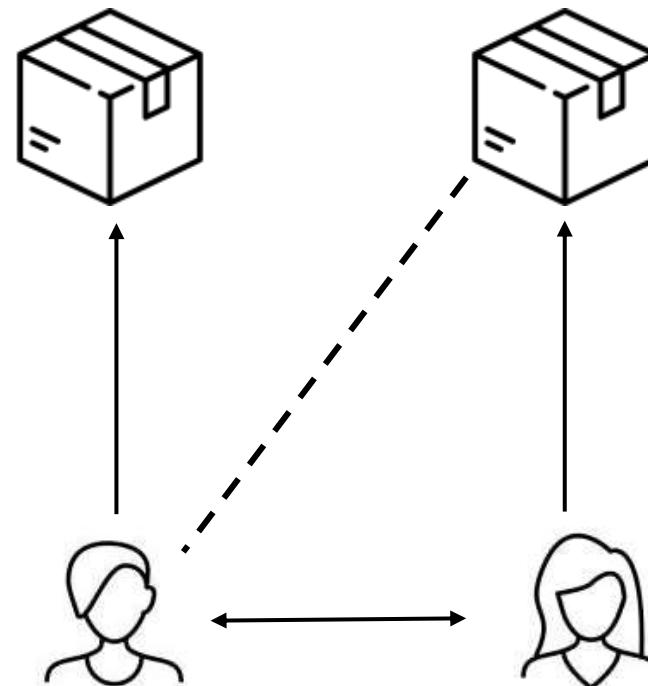


# **LINK PREDICTION AND GRAPH AUTOENCODER**

# WHAT IS A RECOMMENDER SYSTEM?

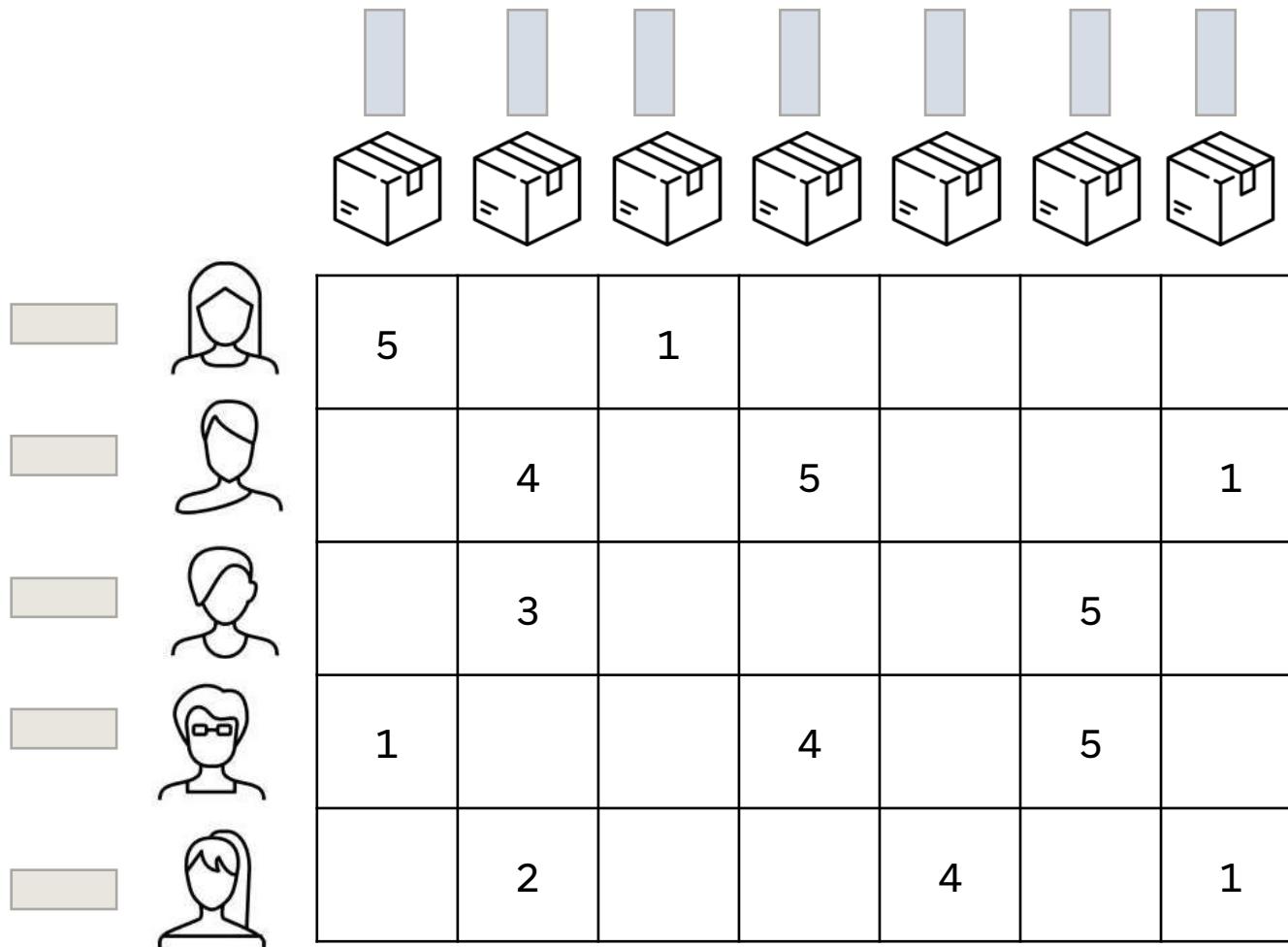


Content-based filtering

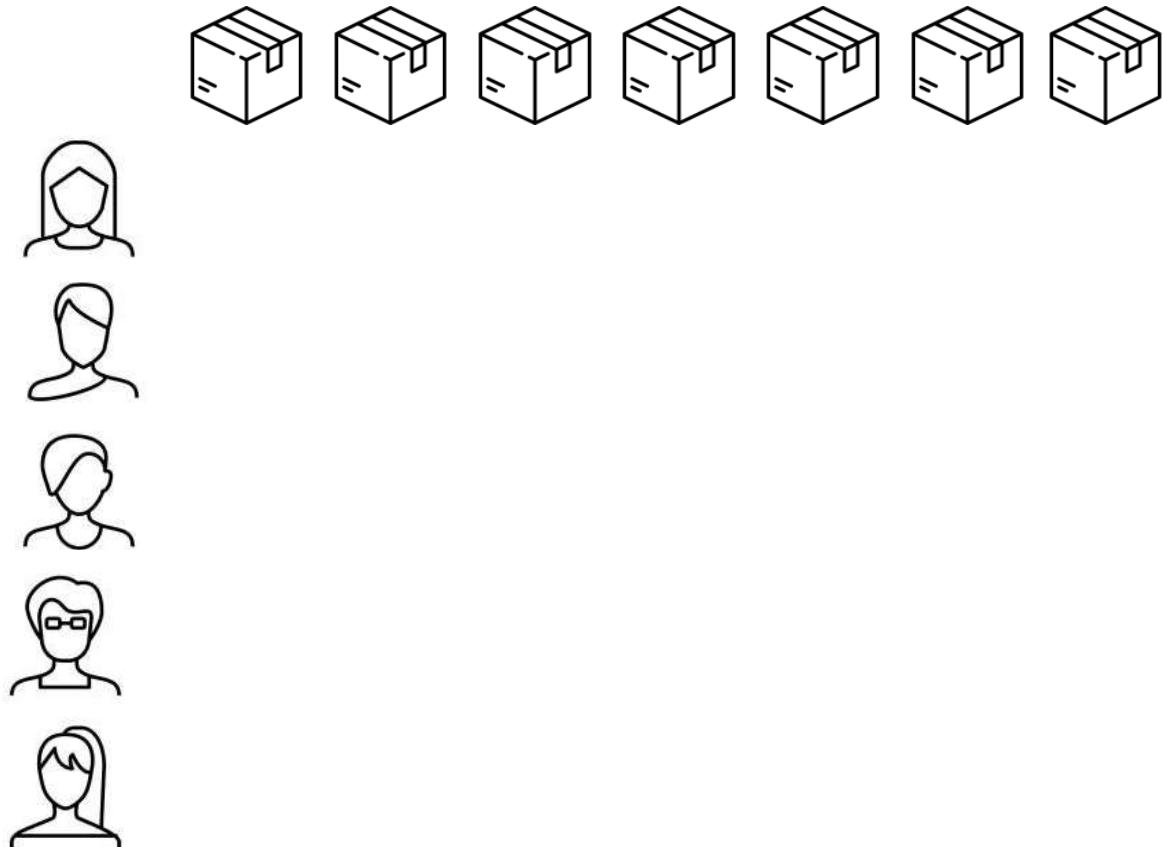


Collaborative filtering

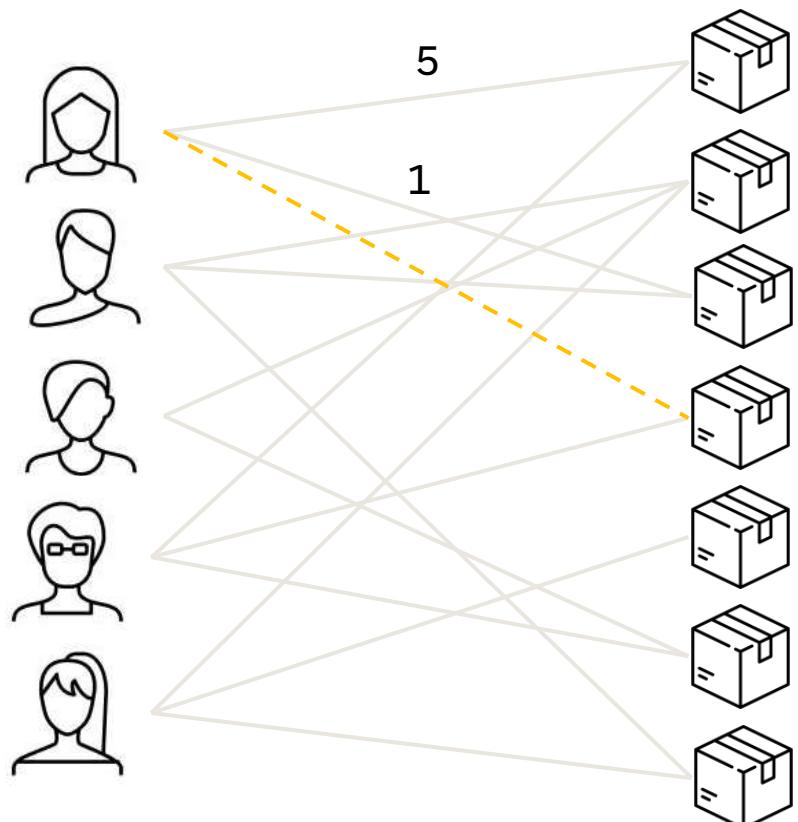
# COLLABORATIVE FILTERING



# BIPARTITE GRAPH



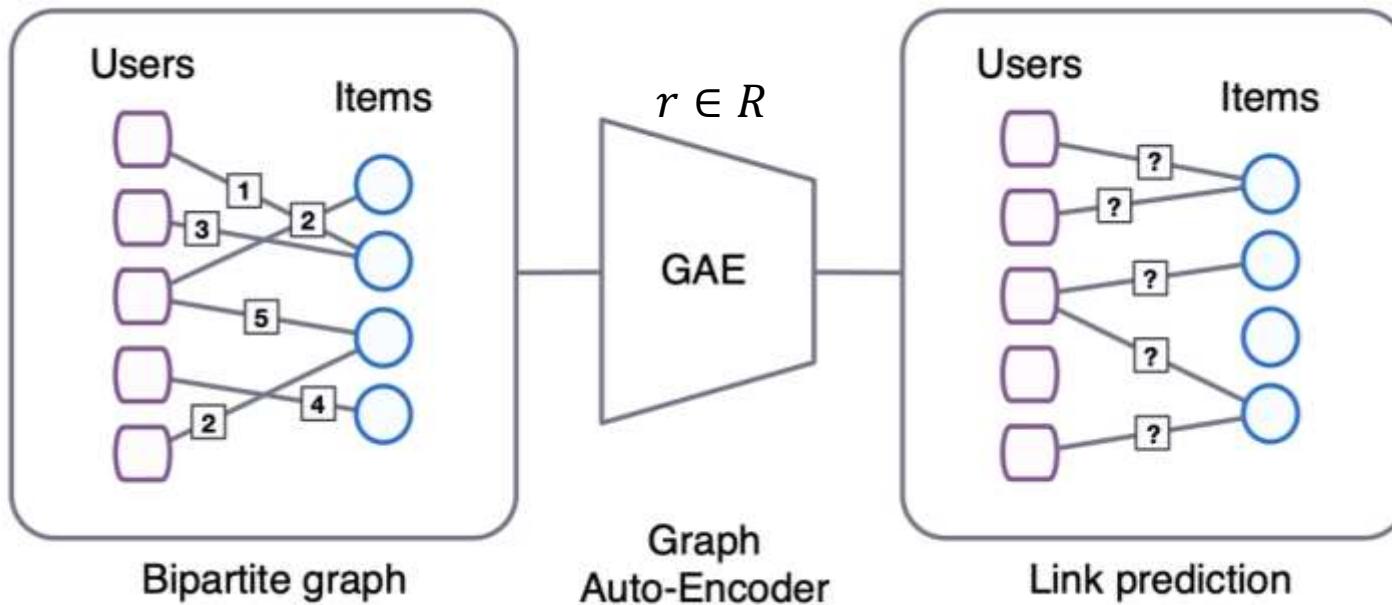
# BIPARTITE GRAPH



# GRAPH CONVOLUTIONAL MATRIX COMPLETION

		Items			
		5	1	0	0
Users		0	3	0	0
		0	0	5	0
		0	0	0	4
		0	0	2	0

Rating matrix  $M$

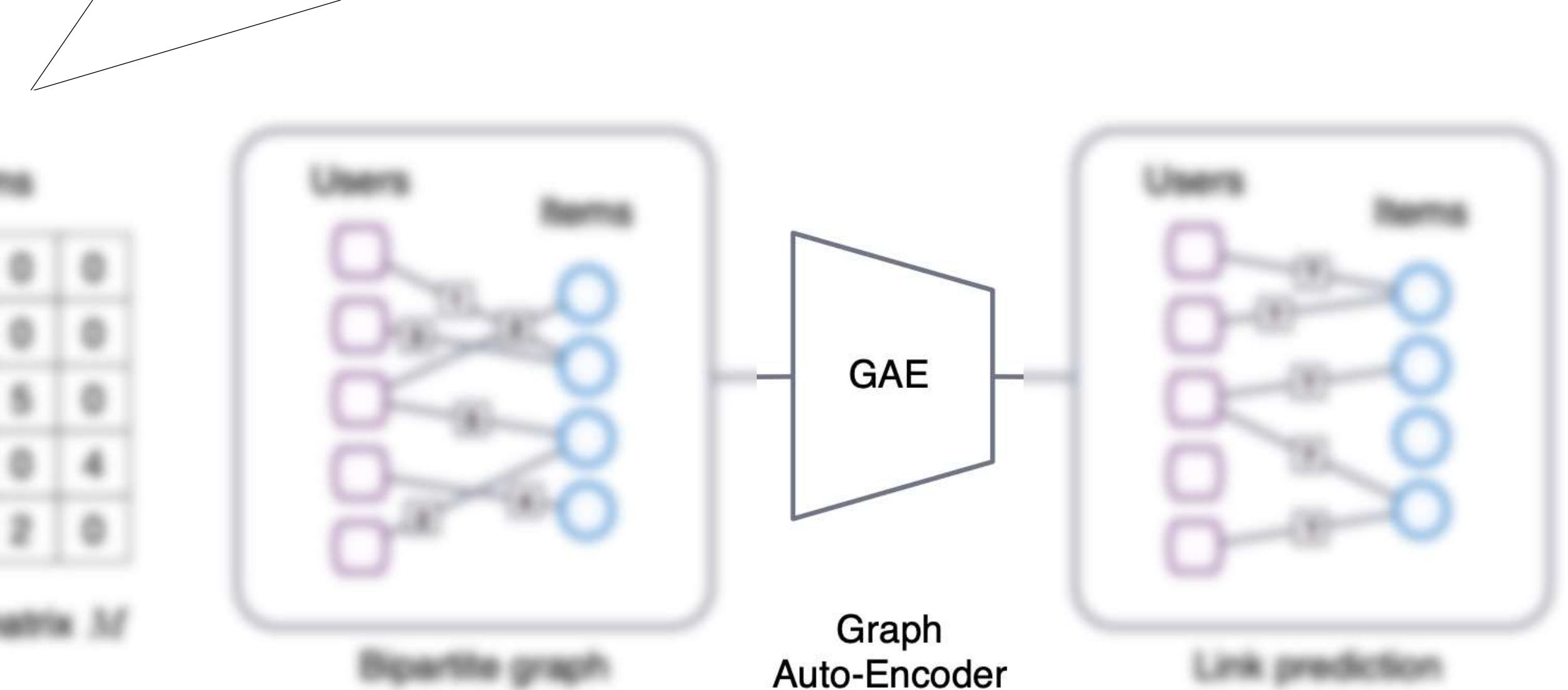


SoftMax  
Which edge type

$$p(\check{M}_{ij} = r) = \frac{e^{u_i^T Q_r v_j}}{\sum_{s \in R} e^{u_i^T Q_s v_j}}$$

Graph Convolutional Matrix Completion

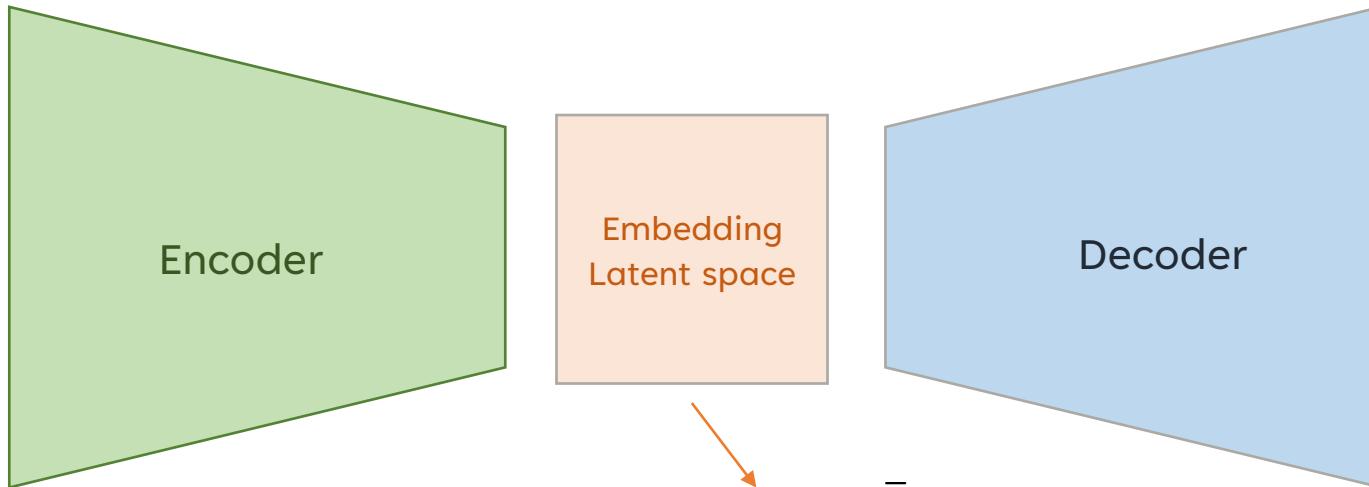
Rianne van den Berg, Thomas N. Kipf, Max Welling 2017



## Graph Convolutional Matrix Completion

Rianne van den Berg, Thomas N. Kipf, Max Welling 2017

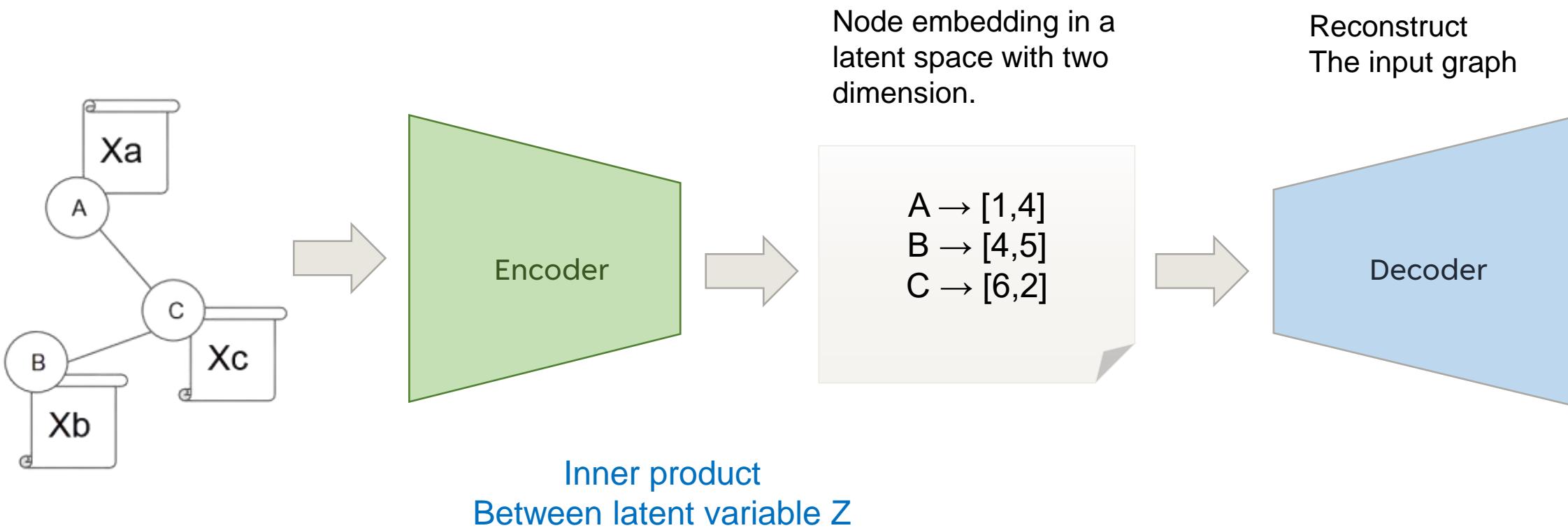
# GRAPH AUTOENCODERS (GAE)



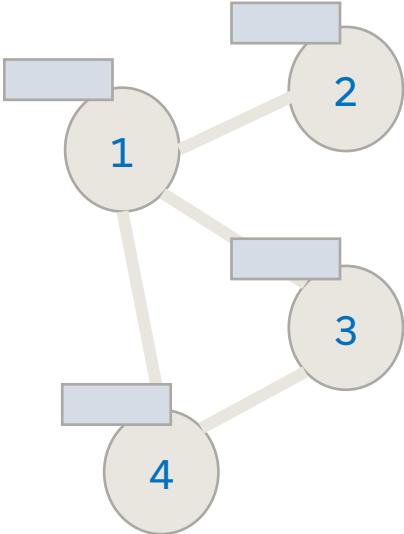
A graph convolutional Neural Network  
produces a low dimensional embedding representation

$$\begin{aligned}\bar{X} &= GCN(A, X) = \text{ReLU}(\tilde{A}XW_0) \\ \text{With } \tilde{A} &= D^{-1/2} A D^{-1/2}\end{aligned}$$

# GRAPH AUTOENCODERS (GAE)



# WHY INNER PRODUCT?



1	2.4	8.1	0.3
2	0.7	0.6	0.2
3	0.3	9.2	1.2
4	2.1	1.8	0.8

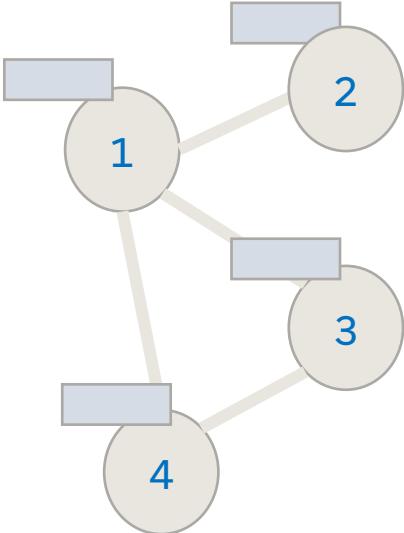
$Z$

1	2	3	4
2.4	0.7	0.3	2.1
8.1	0.6	9.2	1.8
0.3	0.2	1.2	0.8

$Z^T$

$$\hat{\mathbf{A}} = \sigma(\mathbf{Z}\mathbf{Z}^\top), \text{ with } \mathbf{Z} = \text{GCN}(\mathbf{X}, \mathbf{A})$$

# WHY INNER PRODUCT?



1	2.4	8.1	0.3
2	0.7	0.6	0.2
3	0.3	9.2	1.2
4	2.1	1.8	0.8

$Z \quad 4 \times 3$

1	2	3	4
2.4	0.7	0.3	2.1
8.1	0.6	9.2	1.8
0.3	0.2	1.2	0.8

$Z^T \quad 3 \times 4$

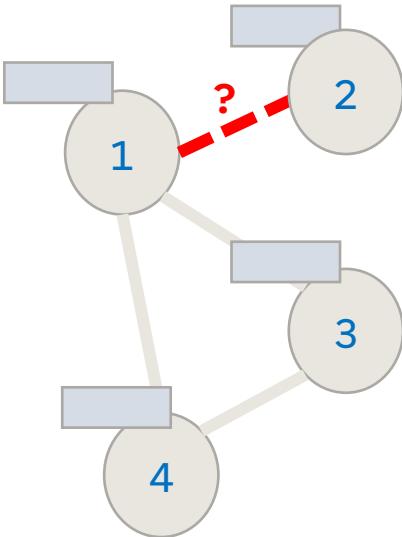
1	2	3	4
?	?	?	?
?	?	?	?
?	?	?	?
?	?	?	?

Adjacency

$4 \times 4$

$$\hat{\mathbf{A}} = \sigma(\mathbf{Z}\mathbf{Z}^\top), \text{ with } \mathbf{Z} = \text{GCN}(\mathbf{X}, \mathbf{A})$$

# WHY INNER PRODUCT?



1	2.4	8.1	0.3
2	0.7	0.6	0.2
3	0.3	9.2	1.2
4	2.1	1.8	0.8

$Z \quad 4 \times 3$

1	2	3	4
2.4	0.7	0.3	2.1
8.1	0.6	9.2	1.8
0.3	0.2	1.2	0.8

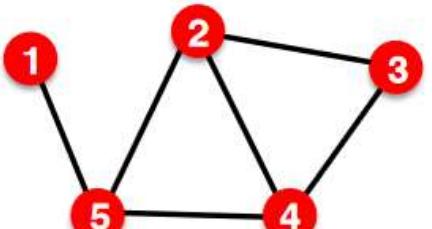
$Z^T \quad 3 \times 4$

1	2	3	4
1	?	?	?
2	?	?	?
3	?	?	?
4	?	?	?

Adjacency  
 $4 \times 4$

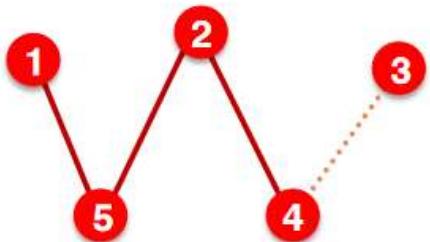
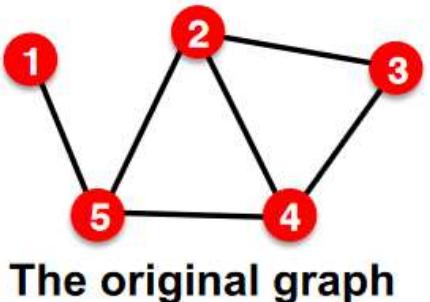
$$\hat{\mathbf{A}} = \sigma(\mathbf{Z}\mathbf{Z}^\top), \text{ with } \mathbf{Z} = \text{GCN}(\mathbf{X}, \mathbf{A})$$

# SETTING UP LINK PREDICTION



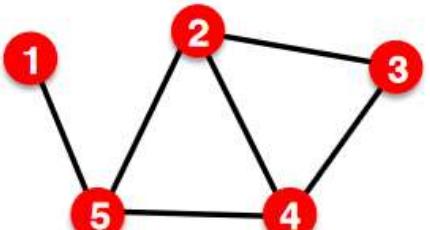
The original graph

# SETTING UP LINK PREDICTION

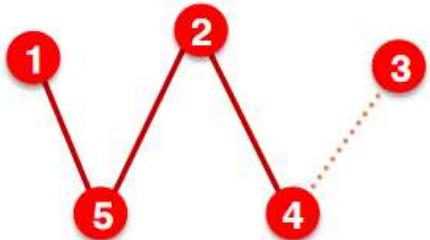


(1) At training time:  
Use **training message edges** to predict **training supervision edges**

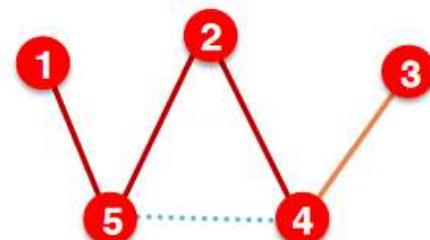
# SETTING UP LINK PREDICTION



The original graph

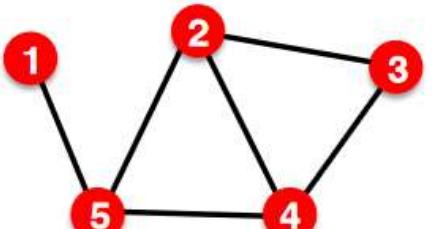


(1) At training time:  
Use **training message edges** to predict **training supervision edges**

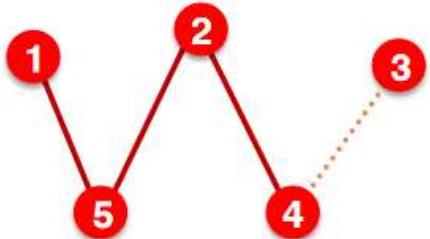


(2) At validation time:  
Use **training message edges & training supervision edges** to predict **validation edges**

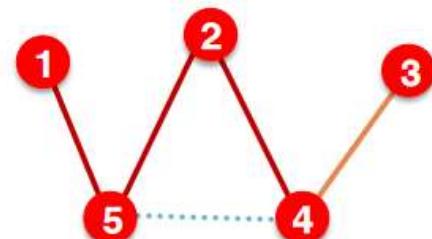
# SETTING UP LINK PREDICTION



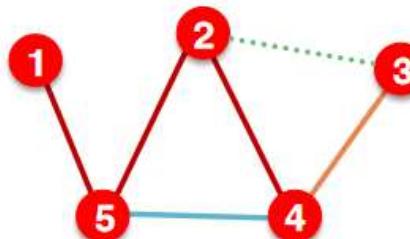
The original graph



(1) At training time:  
Use **training message edges** to predict **training supervision edges**



(2) At validation time:  
Use **training message edges & training supervision edges** to predict **validation edges**



(3) At test time:  
Use **training message edges & training supervision edges & validation edges** to predict **test edges**

# **HETEROGENEOUS & KNOWLEDGE GRAPH EMBEDDING**