# Graph Neural Networks (cont)

COMP9312_23T2

# GraphSage & GAT

# Classical GNN Layers: 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 $\text{AGG}(\cdot)$

  - **Two-stage aggregation**

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

      $$\mathbf{h}_{N(v)}^{(l)} \leftarrow \text{AGG}\left(\left\{\mathbf{h}_u^{(l-1)}, \forall u \in N(v)\right\}\right)$$
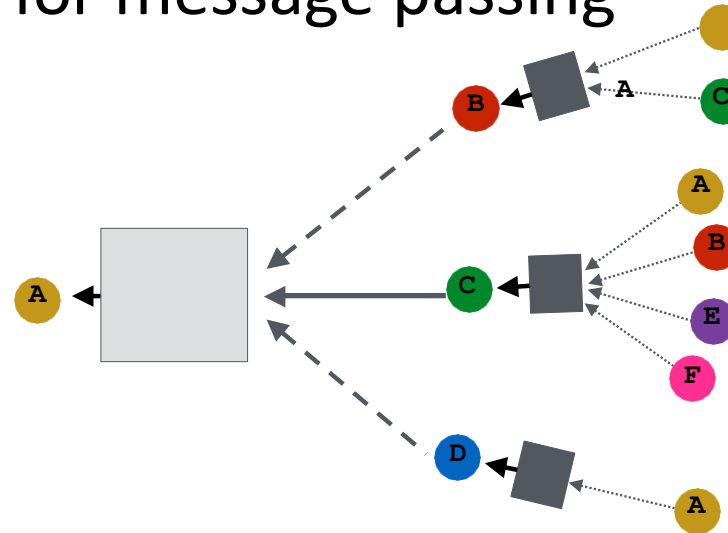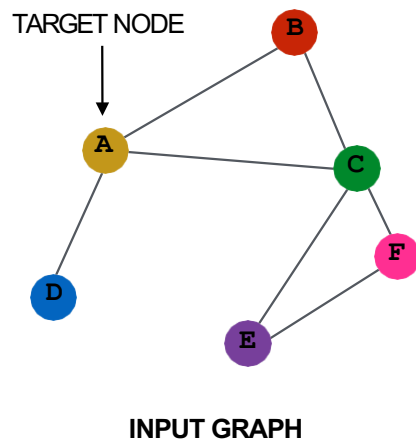
    - **Stage 2:** Further aggregate over the node itself

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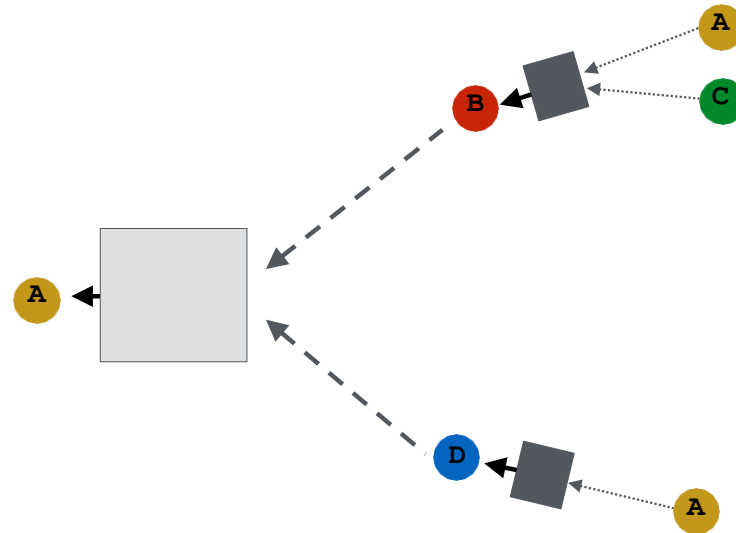
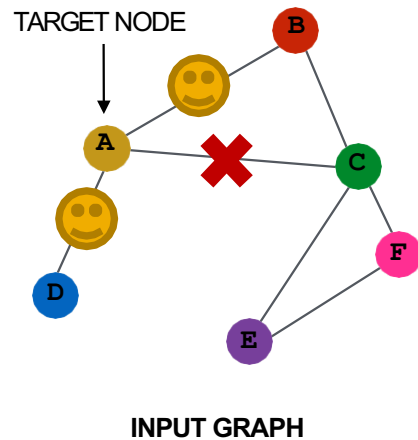# Node Neighborhood Sampling

- **Previously:**
  - All the nodes are used for message passing

TARGET NODE

INPUT GRAPH

- **New idea: (Randomly)** sample a node's neighborhood for message passing
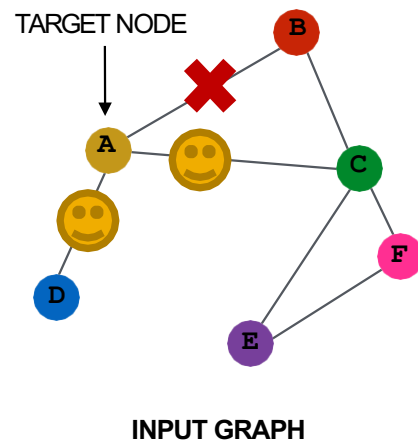
**UNSW COMP9312_23T2**

# Neighborhood Sampling Example

- **For example, we can randomly choose 2 neighbors to pass messages in a given layer**
  - Only nodes $B$ and $D$ will pass messages to $A$



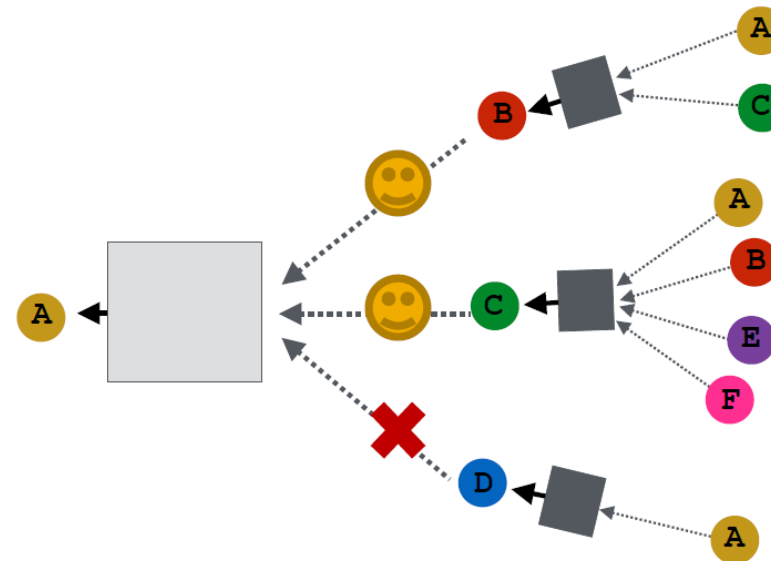TARGET NODE

INPUT GRAPH

# Neighborhood Sampling Example

- **In the next layer when we compute the embeddings, we can sample different neighbors**
  - Only nodes $C$ and $D$ will pass messages to $A$

TARGET NODE

INPUT GRAPH

UNSW COMP9312_23T2

# Neighborhood Sampling Example

- **In expectation, we get embeddings similar to the case where all the neighbors are used**
  - **Benefits:** Greatly **reduces** computational cost
  - And in practice it works great!

TARGET NODE

**INPUT GRAPH**

# GraphSAGE: L$_2$ Normalization

- ## $\ell_2$ **Normalization:**

  - **Optional:** Apply $\ell_2$ normalization to $\mathbf{h}_v^{(l)}$ at every layer

  - $\mathbf{h}_v^{(l)} \leftarrow \dfrac{\mathbf{h}_v^{(l)}}{\left\| \mathbf{h}_v^{(l)} \right\|_2} \ \forall v \in V$ where $\|u\|_2 = \sqrt{\sum_i u_i^2}$ ($\ell_2$-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

# Classical GNN Layers: GAT(1)

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

**Attention weights**

- **In GCN / GraphSAGE**

  - $\alpha_{vu} = \frac{1}{|N(v)|}$ is the **weighting factor (importance)** of node $u$'s message to node $v$

  - $\Longrightarrow \alpha_{vu}$ is defined **explicitly** based on the structural properties of the graph (node degree)

  - $\Longrightarrow$ All neighbors $u \in N(v)$ are equally important to node $v$

# Classical GNN Layers: GAT(2)

- **Graph Attention Networks**

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

**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?

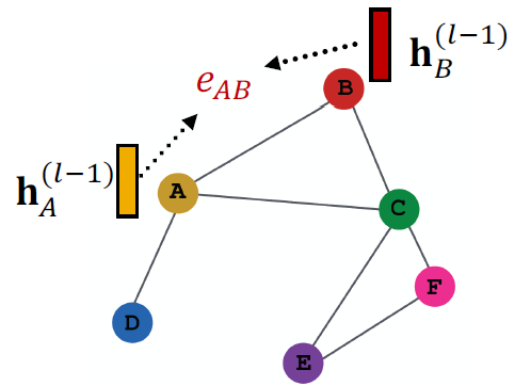**Can we let weighting factors $\alpha_{vu}$ to be learned?**

- **Goal:** Specify **arbitrary importance** to different neighbors of each node in the graph
- **Idea:** Compute embedding $\boldsymbol{h}_v^{(l)}$ 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 (1)

- 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:

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

  - $e_{vu}$ **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 (2)

- **Normalize** $e_{vu}$ into the **final attention weight** $\alpha_{vu}$
  - Use the **softmax** function, so that $\sum_{u \in N(v)} \alpha_{vu} = 1$:

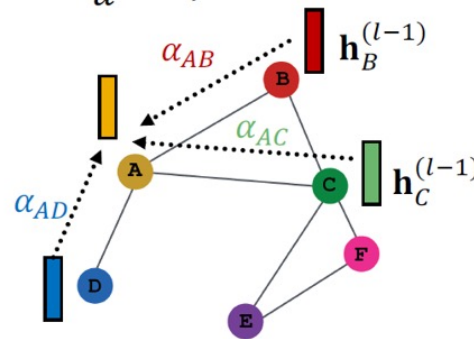$$\alpha_{vu} = \frac{\exp(e_{vu})}{\sum_{k \in N(v)} \exp(e_{vk})}$$

- **Weighted sum** based on the **final attention weight** $\boldsymbol{\alpha_{vu}}$

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

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

- 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

# Attention Mechanism (3)

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

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

  $$\mathbf{h}_v^{(l)}[1] = \sigma(\sum_{u \in N(v)} \alpha_{vu}^1 \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)})$$
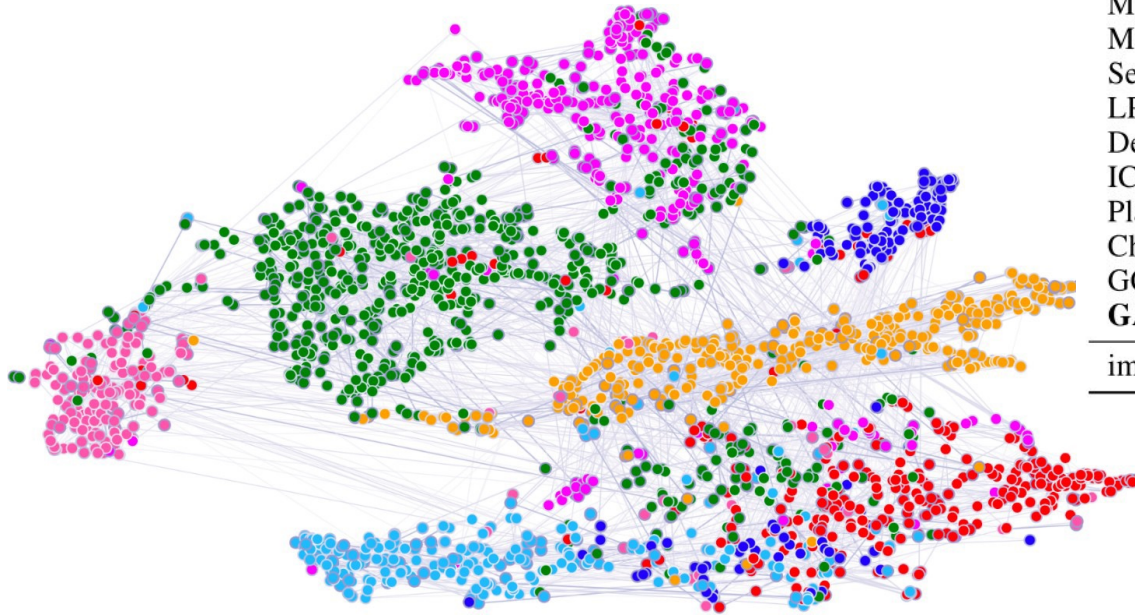  $$\mathbf{h}_v^{(l)}[2] = \sigma(\sum_{u \in N(v)} \alpha_{vu}^2 \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)})$$
  $$\mathbf{h}_v^{(l)}[3] = \sigma(\sum_{u \in N(v)} \alpha_{vu}^3 \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)})$$

  - **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 Attention Mechanism

- **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: 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% |
| **GAT** | **83.3%** |
| 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 Layer in Practice

- **In practice, these classic GNN layers are a great starting point**
  - We can often get better performance by considering a general GNN layer design
  - Concretely, we can include modern deep learning modules that proved to be useful in many domains
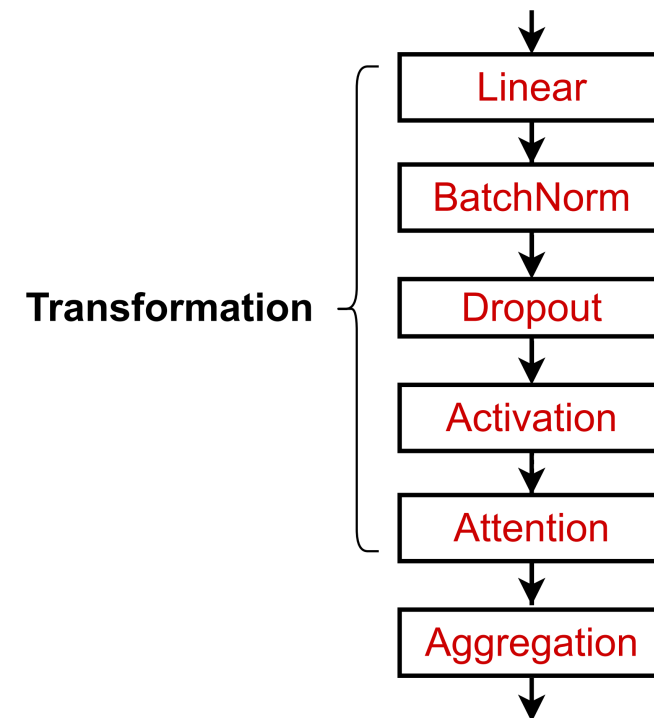
**A suggested GNN Layer**

**Transformation**

- Linear
- BatchNorm
- Dropout
- Activation
- Attention
- Aggregation

# GNN Layer in Practice

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

**A suggested 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

**Transformation**

Linear

BatchNorm

Dropout

Activation

Attention

Aggregation

**UNSW COMP9312_23T2**

# Dropout

- **Goal**: Regularize a neural net to prevent overfitting.
- **Idea**:
  - **During training**: with some probability $p$, randomly set neurons to zero (turn off)
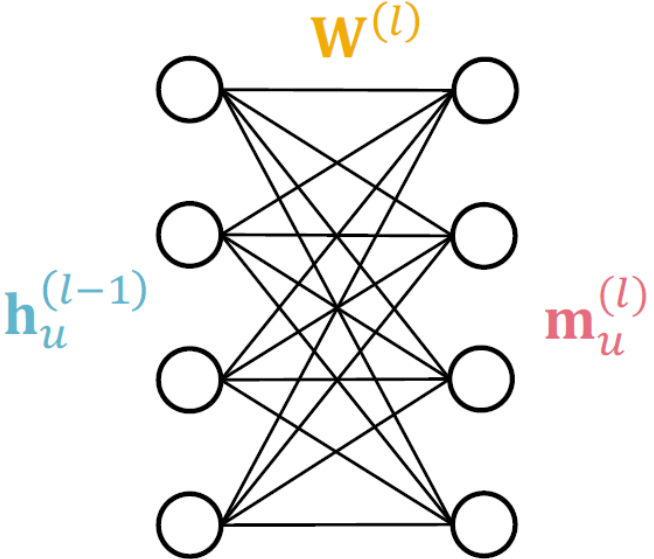  - **During testing:** Use all the neurons for computation

**Dropout**

**Removed neurons**

UNSW COMP9312_23T2

# Dropout for GNNs

- In GNN, Dropout is applied to **the <u>linear layer</u> in the message function**

  - **A simple message function with linear layer:** $\mathbf{m}_u^{(l)} = \mathbf{W}^{(l)}\mathbf{h}_u^{(l-1)}$



(2) Aggregation

(1) Message

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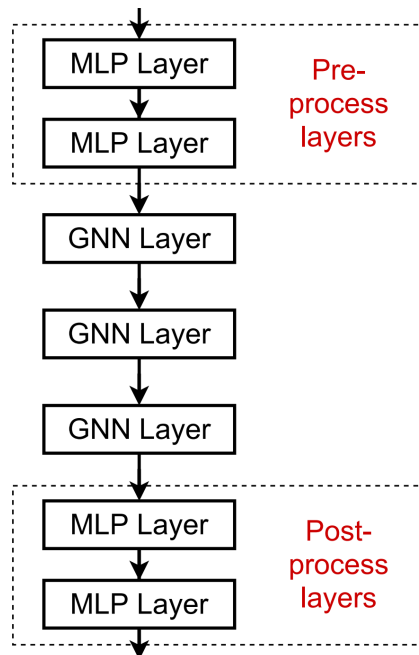
$\mathbf{h}_u^{(l-1)}$     $\mathbf{m}_u^{(l)}$

**Dropout**

**Visualization of a linear layer**

# Expressive Power for Shallow GNNS

- **How to make a shallow GNN more expressive?**
- **Solution:** 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**

MLP Layer
MLP Layer
Pre-process layers

GNN Layer

GNN Layer

GNN Layer

MLP Layer
MLP Layer
Post-process layers

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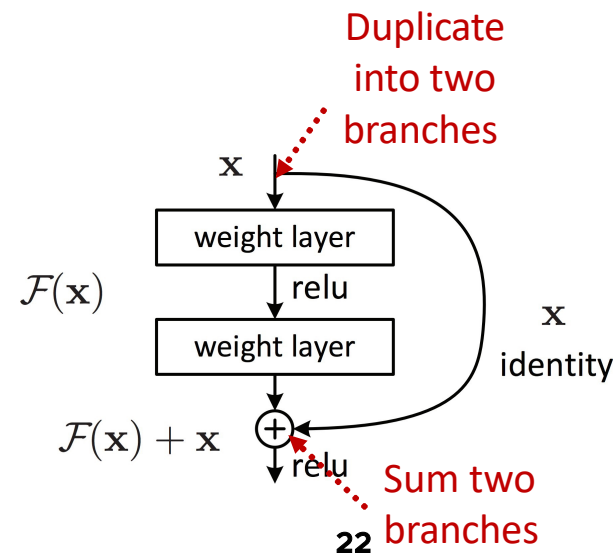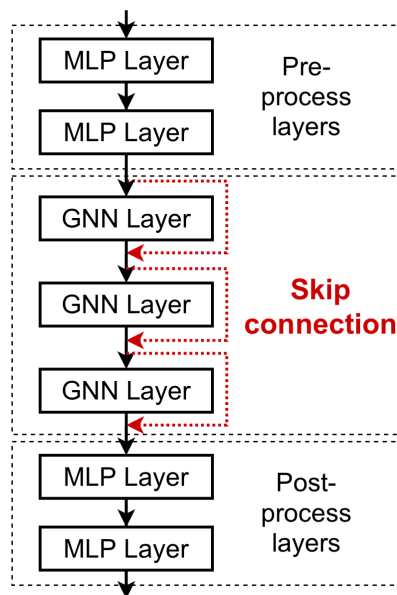E.g., when nodes represent images/text

**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!**

21

# Design GNN Layer Connectivity

- **What if my problem still requires many GNN layers?**

- **Lesson: 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**



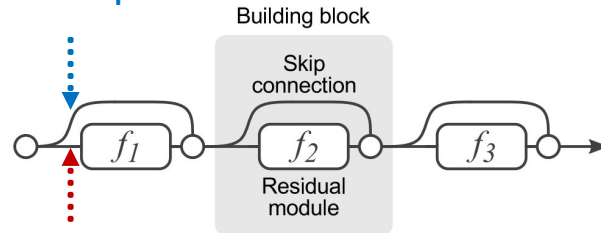**Idea of skip connections:**
Before adding shortcuts:
$$F(\mathbf{x})$$
After adding shortcuts:
$$F(\mathbf{x}) + \mathbf{x}$$

# Skip Connections

- ## Why do skip connections work?

  - **Intuition:** Skip connections create **a mixture of models**

  - $N$ skip connections → $2^N$ possible paths

  - Each path could have up to $N$ modules

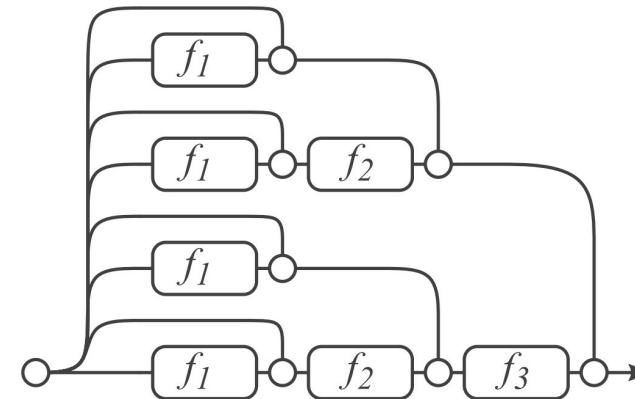  - We automatically get **a mixture of shallow GNNs and deep GNNs**

**All the possible paths:**
$$2 * 2 * 2 = 2^3 = 8$$

**Path 2:** skip this module

Building block

Skip connection

$f_1$ $f_2$ $f_3$

Residual module

**Path 1:** include this module

(a) Conventional 3-block residual network

$=$

$f_1$

$f_1$ $f_2$

$f_1$

$f_1$ $f_2$ $f_3$

(b) Unraveled view of (a)

# GCN with Skip Connections
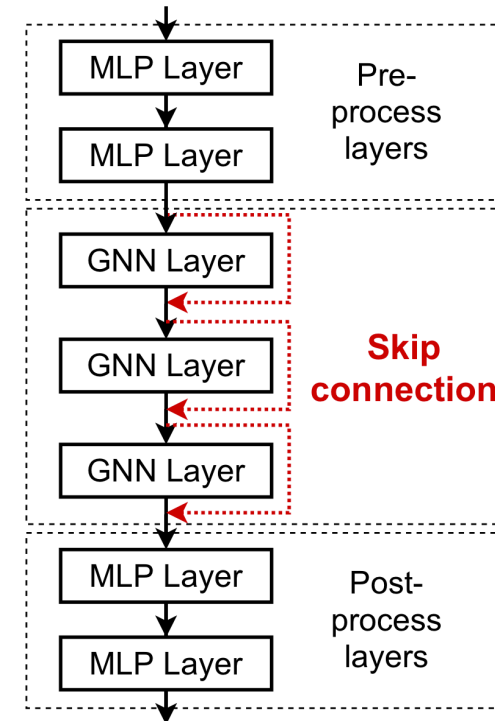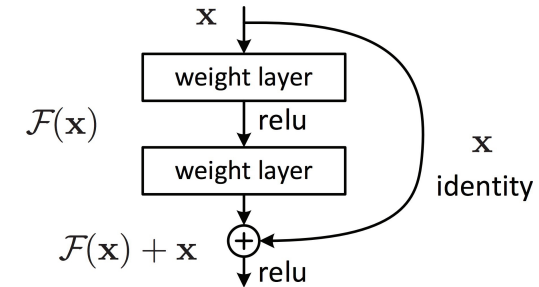
- **A standard GCN layer**

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

  **This is our $F(\mathbf{x})$**
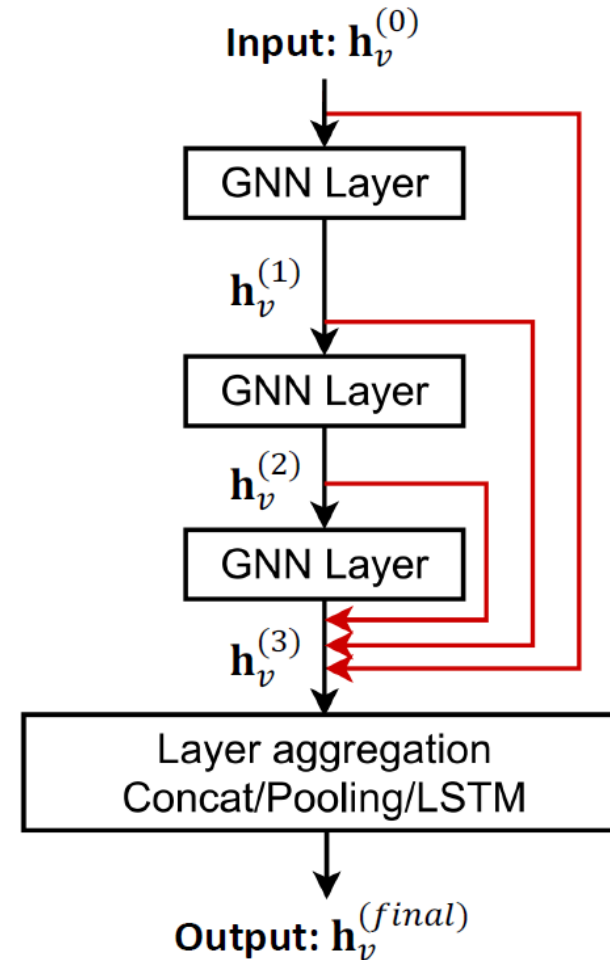
- **A GCN layer with skip connection**

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

  $F(\mathbf{x})$  +  $\mathbf{x}$



**24**

# Other Skip Connections

- **Other options:** Directly skip to the last layer
  - The final layer directly **aggregates from the all the node embeddings** in the previous layers



Input: $\mathbf{h}_v^{(0)}$

GNN Layer

$\mathbf{h}_v^{(1)}$

GNN Layer

$\mathbf{h}_v^{(2)}$

GNN Layer

$\mathbf{h}_v^{(3)}$

Layer aggregation
Concat/Pooling/LSTM

Output: $\mathbf{h}_v^{(final)}$

UNSW COMP9312_23T2

# Graph Augmentation (Optional)

# Why Augment Graphs

Problems in training a GNN

- **Features:**
  - The input graph **lacks features**
- **Graph structure:**
  - The graph is **too sparse** → inefficient message passing
  - The graph is **too dense** → message passing is too costly
  - The graph is **too large** → cannot fit the computational graph into a GPU

# Graph Augmentation Approaches

- **Graph Feature augmentation**
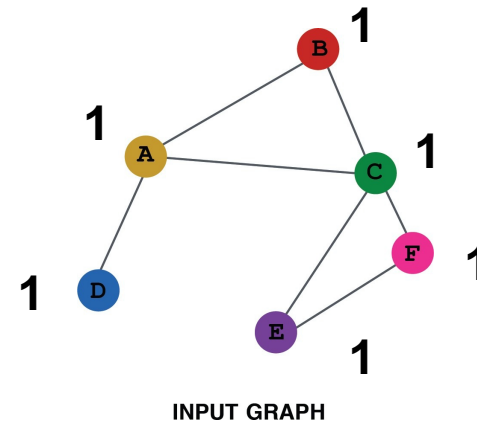
  - The input graph **lacks features** → **feature augmentation**

- **Graph Structure augmentation**

  - The graph is **too sparse** → **Add virtual nodes / edges**

  - The graph is **too dense** → **Sample neighbors when doing message passing**

  - The graph is **too large** → **Sample subgraphs to compute embeddings**

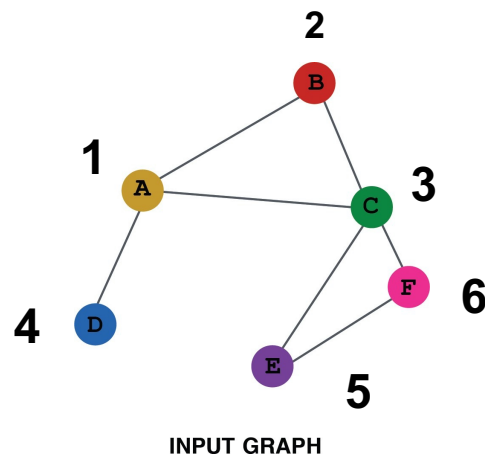# Features Augmentation on Graphs

When might we need feature augmentation?

- **(1) Input graph does not have node features**

  - This is common when we only have the adj. matrix

- **Standard approaches:**

- **a) Assign constant values to nodes**



INPUT GRAPH

# Features Augmentation on Graphs

When might we need feature augmentation?

- **(1) Input graph does not have node features**

  - This is common when we only have the adj. matrix

- **Standard approaches:**

- **b) Assign unique IDs to nodes**

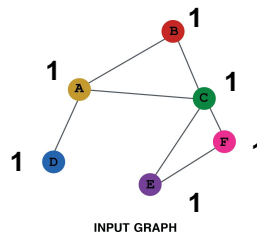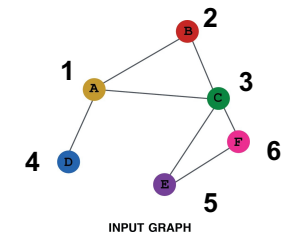  - These IDs are converted into **one-hot vectors**

**One-hot vector for node with ID=5**

**ID = 5**

↓

**[0, 0, 0, 0, 1, 0]**

**Total number of IDs = 6**

INPUT GRAPH

# Features Augmentation on Graphs

| | Constant node feature | One-hot node feature |
|---|---|---|
| | INPUT GRAPH | INPUT GRAPH |
| Expressive power | Medium. All the nodes are identical, but GNN can still learn from the graph structure | High. Each node has a unique ID, so node-specific information can be stored |
| Inductive learning (Generalize to unseen nodes) | High. Simple to generalize to new nodes: we assign constant feature to them, then apply our GNN | Low. Cannot generalize to new nodes: new nodes introduce new IDs, GNN doesn't know how to embed unseen IDs |
| Computational cost | Low. Only 1 dimensional feature | High. $O(|V|)$ dimensional feature, cannot apply to large graphs |
| Use cases | Any graph, inductive settings (generalize to new nodes) | Small graph, transductive settings (no new nodes) |

# Features Augmentation on Graphs
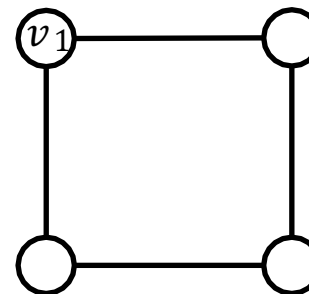
When might we need feature augmentation?

- **(2) Certain structures are hard to learn by GNN**
- **Example:** Cycle count feature:
  - Can GNN learn the length of a cycle that $v_1$ resides in?
  - **Unfortunately, no**

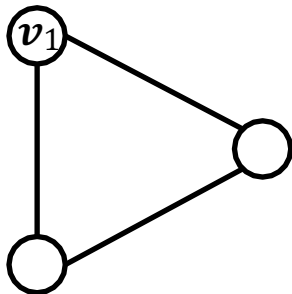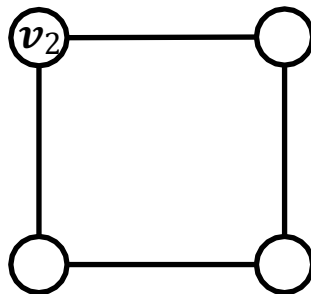$v_1$ resides in a cycle with length 3      $v_1$ resides in a cycle with length 4

# Features Augmentation on Graphs

- $v_1$ **cannot differentiate which graph it resides in**

  - Because all the nodes in the graph have degree of 2

  - The computational graphs will be the same binary tree
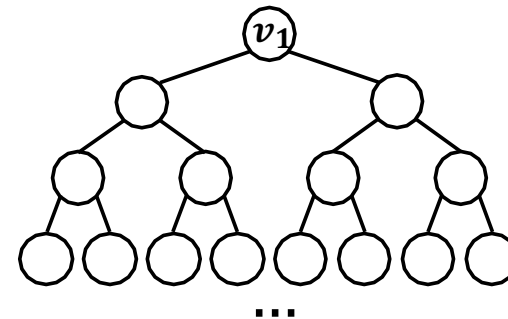
$v_!$ resides in a cycle with length 3

$v_!$ resides in a cycle with length 4

The computational graphs for node $v_1$ are always the same



$v_!$ resides in a cycle with infinite length

# Features Augmentation on Graphs

When might we need feature augmentation?

- **(2) Certain structures are hard to learn by GNN**
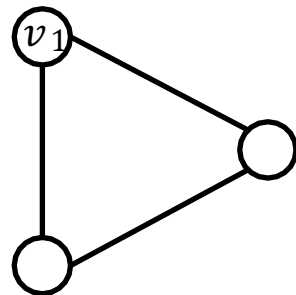- Solution: We can use cycle count as augmented node features

We start from cycle with length 0

**Augmented node feature for $v_1$**

$$[0, 0, 0, 1, 0, 0]$$

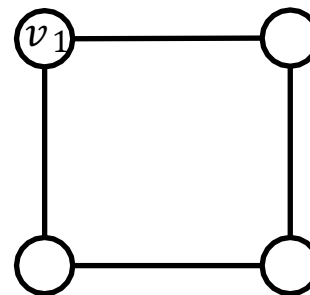↑

$v_1$ resides in a cycle with length 3

**Augmented node feature for $v_1$**

$$[0, 0, 0, 0, 1, 0]$$

↑

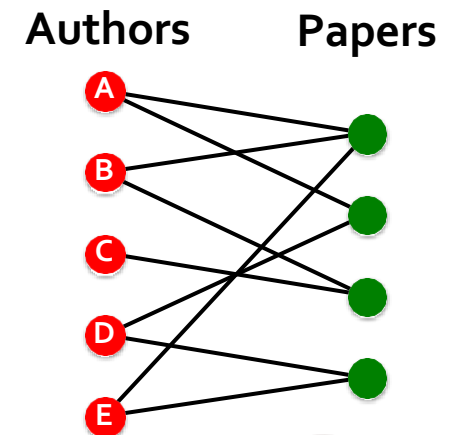$v_1$ resides in a cycle with length 4

# Features Augmentation on Graphs

When might we need feature augmentation?

- **(2) Certain structures are hard to learn by GNN**
- Other commonly used augmented features:
  - **Node degree**
  - **Clustering coefficient**
  - **Centrality**
  - **…**

# Add Virtual Nodes/ Edges

- **Motivation:** Augment sparse graphs
- **(1) Add virtual edges**

  - **Common approach:** Connect 2-hop neighbors via virtual edges

  - **Intuition:** Instead of using adj. matrix $A$ for GNN computation, use $A + A^2$

  - **Use cases:** Bipartite graphs
    - Author-to-papers (they authored)
    - 2-hop virtual edges make an author-author collaboration graph

**Authors**   **Papers**

# Add Virtual Nodes/Edges

**The virtual node**



INPUT GRAPH

- **Motivation:** Augment sparse graphs

- **(2) Add virtual nodes**

  - The virtual node will connect to all the nodes in the graph
    - Suppose in a sparse graph, two nodes have shortest path distance of 10
    - After adding the virtual node, **all the nodes will have a distance of two**
      - **Node A – Virtual node – Node B**

  - **Benefits:** Greatly **improves message passing in sparse graphs**