# Contents

# Overview

Local network neighborhoods:

- Describe aggregation strategies
- Define computation graphs

Stacking multiple layers:

- Describe the model, parameters, training
- How to fit the model?
- Simple example for unsupervised and supervised training

- Basics of neural networks
  - Loss, Optimization, Gradient, SGD, non-linearity, MLP
- Idea for Deep Learning for Graphs
  - Multiple layers of embedding transformation
  - At every layer, use the embedding at previous layer as the input
  - Aggregation of neighbors and self-embeddings
- Graph Convolutional Network
  - Mean aggregation; can be expressed in matrix form
- GNN is a general architecture
  - CNN and Transformer can be viewed as a special GNN

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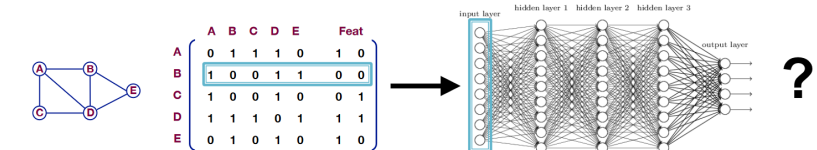
- **Assume we have a graph $G$:**
  - $V$ is the **vertex set**
  - $A$ is the **adjacency matrix** (assume binary)
  - $X \in \mathbb{R}^{m \times |V|}$ is a matrix of **node features**
  - $v$: a node in $V$; $N(v)$: the set of neighbors of $v$.
  - **Node features:**
    - Social networks: User profile,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

# Naive Approach: append node feaetures to adjacency matrix

## 3 problems: parameter size, graph size, node ordering

1. Parameter size
   a. One training example per node but for each node there are N + X (node features) number of features
   b. Training unstable / easy to overfit
2. Graph with different size
   o E.g. If 5 nodes, hard to fit in input size of 7
3. Node ordering
   o If the column order change, then the adjacency matrix changes
     ▪ Rows & cols permuted thou the info is the same
   o For images, the ordering can be top left pixel to bottom right
   o but for graphs, there is no fix node ordering i.e. unclear how to sort the graph to put them as input in the matrix
   o Has to be **invariant** to **node ordering**

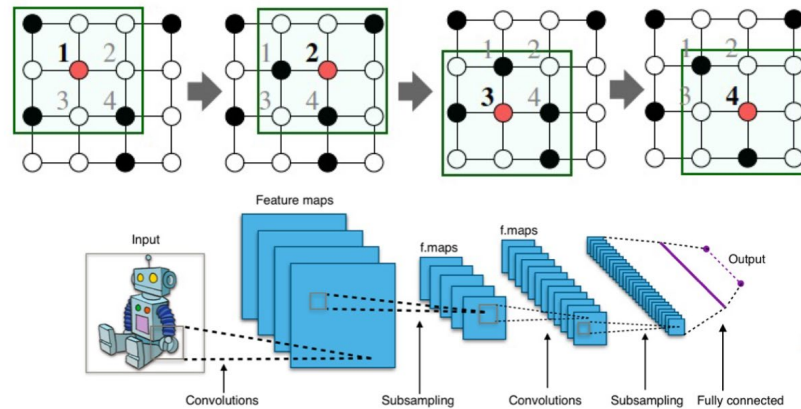- Join adjacency matrix and features
- Feed them into a deep neural net:



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

# Adopt CNN

**Goal.** generalize convolutions beyond simple lattices & Leverage node features/attributes (e.g., text, images)

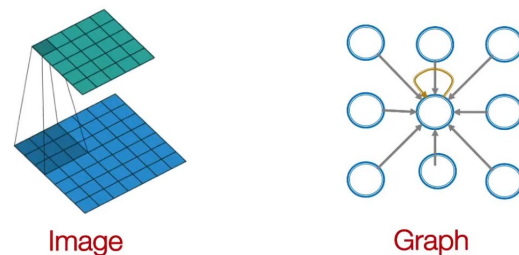**CNN = Sliding windows & locality.**

## CNN on an image:



### Problem

- no fixed notion of locality or sliding window on the graph
  - L: covers 3 nodes
  - R: covers more nodes
- Graph is permutation invariantss



or this:

### Solutions

- Aggregate information about a node based on its neighbouring nodes

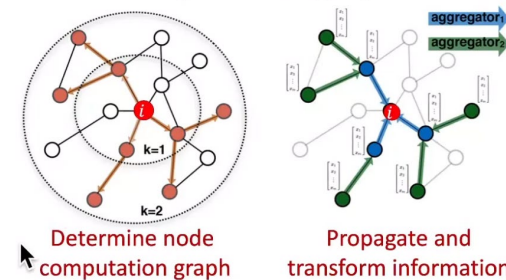Single Convolutional neural network (CNN) layer with 3x3 filter:



Image          Graph

**Idea:** transform information at the neighbors and combine it:
- Transform "messages" $h_i$ from neighbors: $W_i h_i$
- Add them up: $\sum_i W_i h_i$

## Kipf & Welling, ICLR 2017

- Neighbour nodes takes the message from the node and propagate.
  Steps
  1. Determine node computation graph
  2. Propagate & transform information

**Idea:** Node's neighborhood defines a computation graph



Determine node        Propagate and
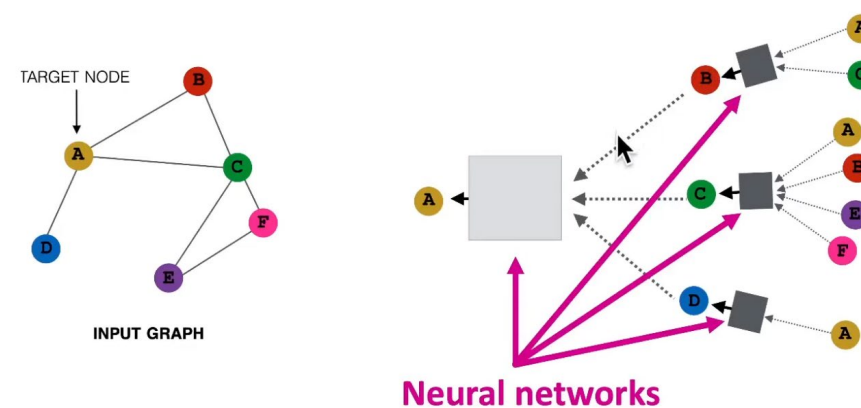computation graph     transform information

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

### Explain – transformation & aggregation

1. To decide node A informatino, we collect from its neighbours {B, C, D}, in which their info are based on their neighbours.
2. The message passing is then from the leaf to root
   a. First, transform the info from leaf
   b. Second, aggregate them in the parent node
   c. Repeat

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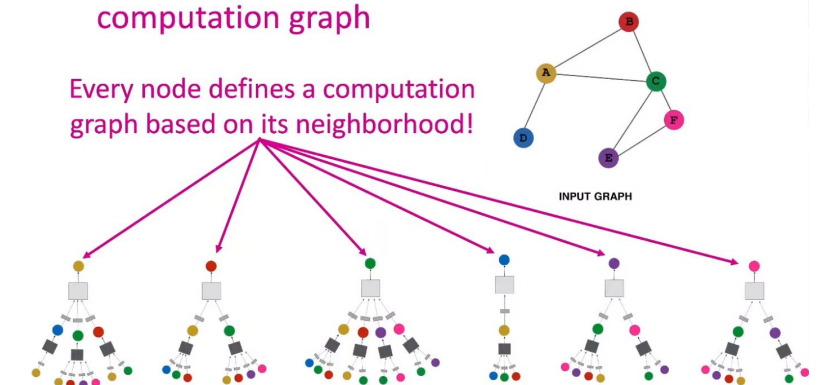


**Neural networks**

### Explain – each node has a computational graph

- Every node has its own computation graph / architecture
- The structure depends on other structure
- Different to classical DL

- **Intuition:** Network neighborhood defines a computation graph
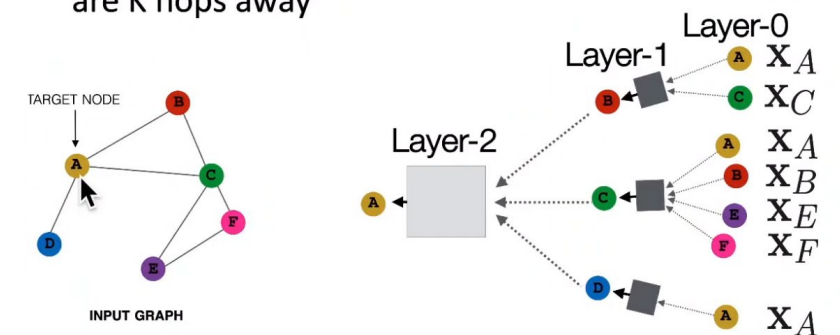
Every node defines a computation graph based on its neighborhood!



### Many layers – k-hops
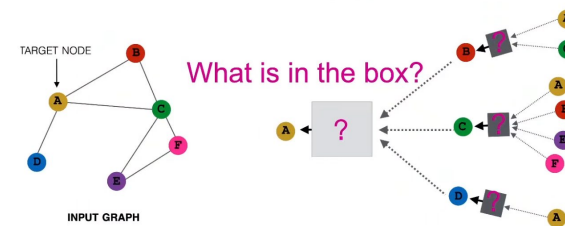
- **Model can be of arbitrary depth:**
  - Nodes have embeddings at each layer
  - Layer-0 embedding of node $u$ is its input feature, $x_u$
  - Layer-$k$ embedding gets information from nodes that are K hops away
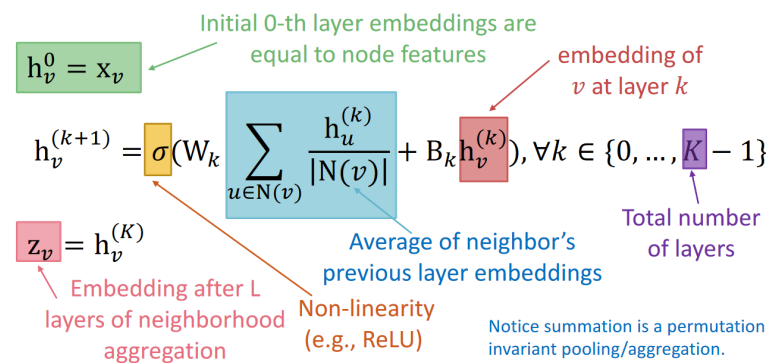
## Neighbourhood aggregation

- The aggregation result is the same regardless of the ordering

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



What is in the box?

## Transformation choice – average embedding

- Ordering invariant
- Examples
  - Average messages
  - Apply NN
    - Apply linear transformation
    - Follow by non-linearity

- Transform current layer features + aggregated previous child nodes messages



Initial 0-th layer embeddings are equal to node features

$h_v^0 = x_v$

$h_v^{(k+1)} = \sigma(W_k \sum_{u \in N(v)} \frac{h_u^{(k)}}{|N(v)|} + B_k h_v^{(k)}), \forall k \in \{0, \dots, K-1\}$

$z_v = h_v^{(K)}$

embedding of $v$ at layer $k$

Total number of layers

Average of neighbor's previous layer embeddings

Non-linearity (e.g., ReLU)

Embedding after L layers of neighborhood aggregation

Notice summation is a permutation invariant pooling/aggregation.

## Model Parameters: W, B; neighborhood aggregation, transformation

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

$h_v^{(0)} = x_v$

$h_v^{(k+1)} = \sigma(W_k \sum_{u \in N(v)} \frac{h_u^{(k)}}{|N(v)|} + B_k h_v^{(k)}), \forall k \in \{0..K-1\}$

$z_v = h_v^{(K)}$

Final node embedding

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

$h_v^k$: the hidden representation of node $v$ at layer $k$
- $W_k$: weight matrix for neighborhood aggregation
- $B_k$: weight matrix for transforming hidden vector of self

# Matrix Formulation - efficient sparse matrix

## Averging of neighbour embedding

- Many aggregations can be performed **efficiently** by (sparse) matrix operations
- Then
  - Node embedding is the avage of the neighbour embedding
  - i.e. the adjacency matrix * embedding spaces at a layer
- in short
  - averging of neighbour embedding
    - summing and averagign can be rewritten as
    - **matrix multiplication (dot prodct)**
  - the drawing is what correspond to what

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

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



$h_i^{(k-1)}$

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

## Re-writing update function in matrix form

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

- Re-writing update function in matrix form:

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

  $H^{(k)} = [h_1^{(k)} \dots h_{|V|}^{(k)}]^T$

  - Red: neighborhood aggregation
  - Blue: self transformation

# How to train GNN: supervised & unsupervised setting

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

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

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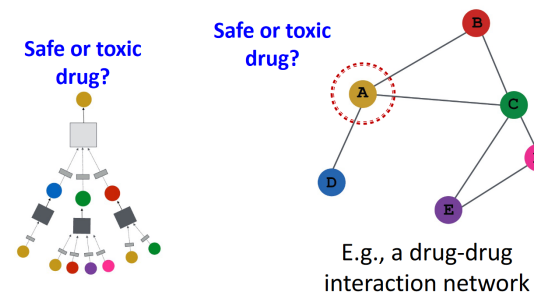
## Unsupervised Training

- No labels

  - **"Similar" nodes have similar embeddings**

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

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

## Supervised Training

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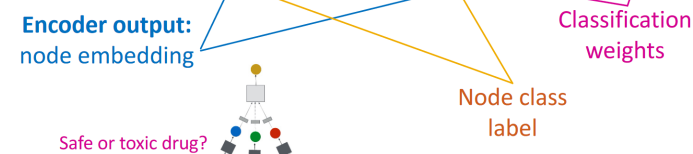


E.g., a drug-drug interaction network

Cross entropy loss
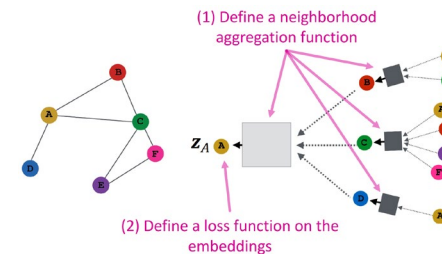
- If label is 1, want output to be 1
- If 0, want 0

- Use cross entropy loss (Slide 16)

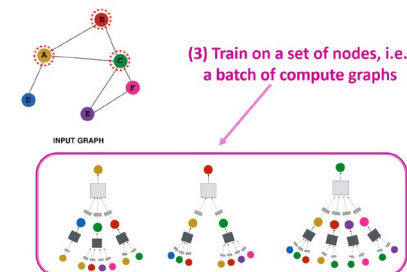$$\mathcal{L} = \sum_{v \in V} y_v \log(\sigma(z_v^T \theta)) + (1 - y_v)\log(1 - \sigma(z_v^T \theta))$$

**Encoder output:** node embedding

**Classification weights**

Node class label

Safe or toxic drug?

## Model design

### Step 1,2 - neighbour aggregation & loss function



(1) Define a neighborhood aggregation function

(2) Define a loss function on the embeddings

### Step 3 – train on a batch



(3) Train on a set of nodes, i.e., a batch of compute graphs
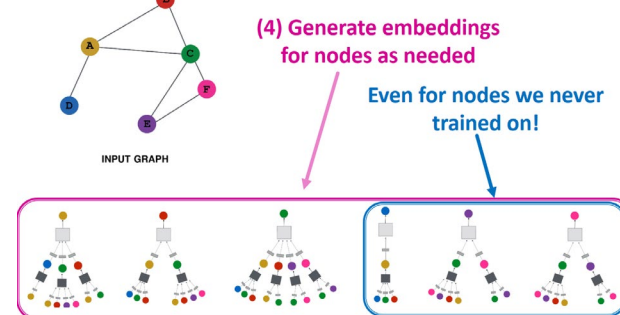
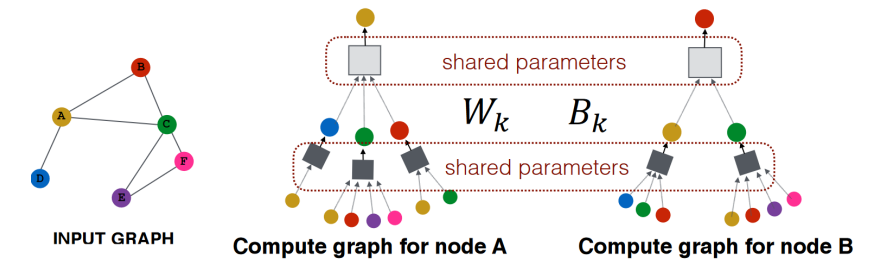### Step 4 - enerate embeddings for nodes as needed

- **Generalisability.** Train embedding for one graph and transfer to another graph



(4) Generate embeddings for nodes as needed
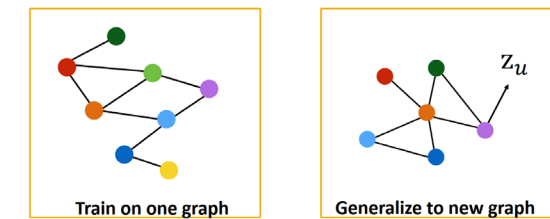
Even for nodes we never trained on!

# Inductive capcability

- The same aggregation parameters are **shared** for all nodes
- The number of model parameters is **sublinear** in $|V|$
- # model parameters **W & B**, depends on
  - #features / embedding dimensionality (since shared)
  - not the size of graph (#nodes)
- Thus, generalize to **unseen** nodes



INPUT GRAPH    Compute graph for node A    Compute graph for node B

shared parameters
$W_k$  $B_k$
shared parameters

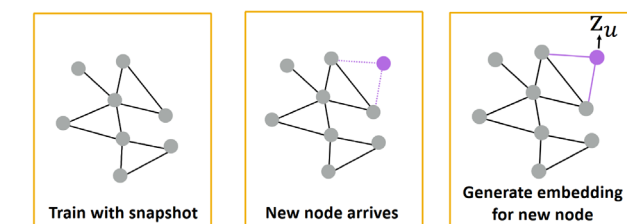## New graph



Train on one graph    Generalize to new graph

Inductive node embedding → Generalize to entirely unseen graphs

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

## New nodes

- One forward pass can generate new embedding for the new node as the graph evolving



Train with snapshot    New node arrives    Generate embedding for new node

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

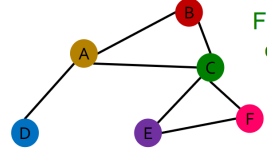# More

## Permutation Invariance

- Graph does not have a canonical order of the nodes
- For a graph with $m$ nodes, there are $m!$ different order plans.
- Are other neural network architectures, e.g., MLPs, permutation invariant / equivariant?
  - o No
  - o Switching the order of the input leads to different outputs!
  - o This explains why the **naïve MLP** approach **fails** for **graphs**

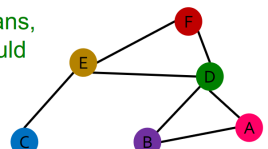    - Consider we learn a function $f$ that maps a graph $G = (A, X)$ to a vector $\mathbb{R}^d$ then
      $$f(A_1, X_1) = f(A_2, X_2)$$
      $A$ is the adjacency matrix
      $X$ is the node feature matrix

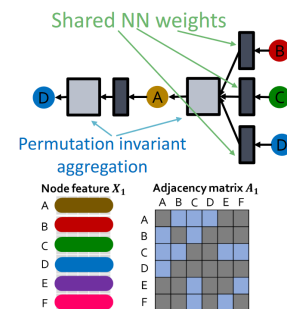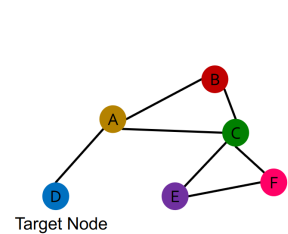**Order plan 1:** $A_1, X_1$     **Order plan 2:** $A_2, X_2$

For two order plans, output of $f$ should be the same!
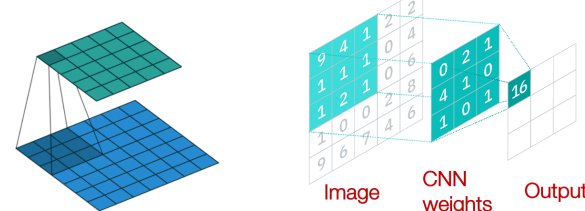
## Equivariant Property

- Message passing and neighbor aggregation in graph convolution networks is **permutation equivariant.**
- The target node (blue) has the **same computation graph** for **different order plans**

equivariant.

Shared NN weights

Permutation invariant aggregation

Node feature $X_1$

Adjacency matrix $A_1$

Target Node

# GNNs subsume CNNs and Transformers

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

Image    CNN weights    Output

CNN formulation: $h_v^{(l+1)} = \sigma(\sum_{u \in N(v) \cup \{v\}} W_l^u h_u^{(l)}), \quad \forall l \in \{0, \dots, L-1\}$

$N(v)$ represents the 8 neighbor pixels of $v$.

## GNN vs. CNN

- CNN can be seen as a special GNN with fixed neighbor size and ordering
  - o The size of the filter is pre-defined for a CNN
  - o The advantage of GNN is it processes arbitrary graphs with different degrees for each node
  - o CNN is not permutation equivariant
    - Switching the order of pixels will leads to different outputs.

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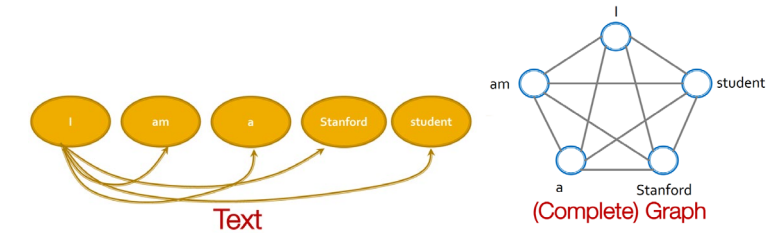
Image      Graph

- GNN formulation (previous slide): $h_v^{(l+1)} = \sigma(W_l \sum_{u \in N(v)} \frac{h_u^{(l)}}{|N(v)|} + B_l h_v^{(l)}), \forall l \in \{0, \dots, L-1\}$
- CNN formulation:   $h_v^{(l+1)} = \sigma(\sum_{u \in N(v) \cup \{v\}} W_l^u h_u^{(l)}), \forall l \in \{0, \dots, L-1\}$
  if we rewrite:   $h_v^{(l+1)} = \sigma(\sum_{u \in N(v)} W_l^u h_u^{(l)} + B_l h_v^{(l)}), \forall l \in \{0, \dots, L-1\}$

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

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

**Key difference**: We can learn different $W_l^u$ for different "neighbor" $u$ for pixel $v$ on the image. The reason is we can pick an order for the 9 neighbors using **relative position** to the center pixel: {(-1,-1). (-1,0), (-1, 1), …, (1, 1)}

# Transformer

- Key component: **self-attention**
  - o Every token/word attends to all the other tokens/words via matrix calculation

- Since each word attends to all the other words, the computation graph of a transformer layer is identical to that of a GNN on the fully-connected "word" graph.

Text

(Complete) Graph

# Takeaway

- What if the ordering plan does matter?