# Graph Convolution Networks

*Great things are not done by impulse, but by a series of small things brought together.*

*- Vincent Van Gogh*

Jay Urbain, Ph.D.

Electrical Engineering and Computer Science Department
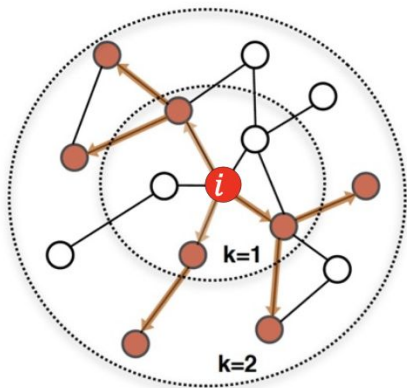
Milwaukee School of Engineering

# Credits

- **Graph Neural Networks: A review of Methods and Applications, Zhou, 2018.**
- **Graph Representation Learning, Hamilton, 2020.**
- **Graph Convolutional Networks, Kipf and Welling, 2016.**
- Multi-Layer Perceptron as Aggregator, Zaheer, 2017.
- Graph Attention Networks, Velickovic, 2017.
- Gated Graph Neural Networks, Li, 2015.
- DeepFindr
- Michael Bronstein, Oxford, Geometric Deep Learning Course
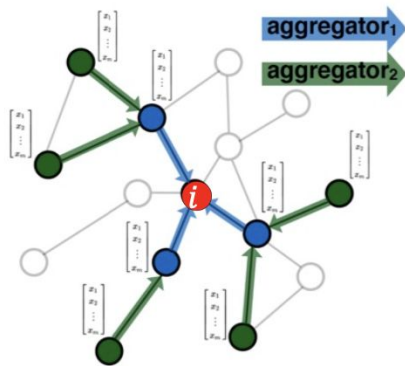- Jure Lesvovec, Stanford

# Review: Graph Neural Networks

Idea:

- Node's neighborhood defines a computation graph.
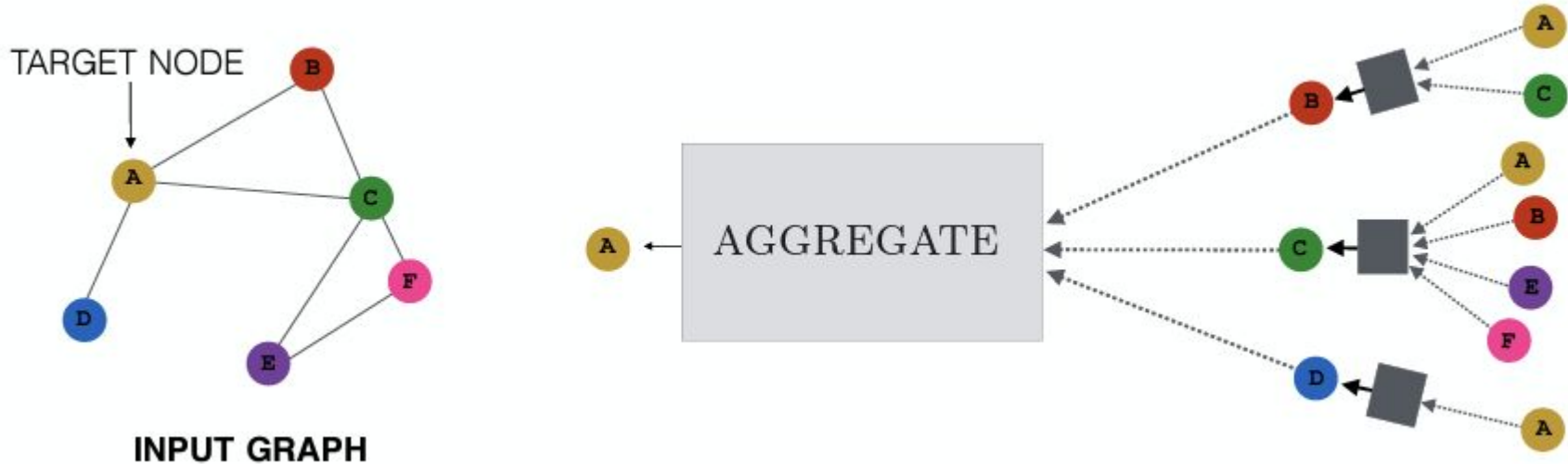- Propagate information across the graph to compute node features



Determine node
computation graph

Propagate and
transform information
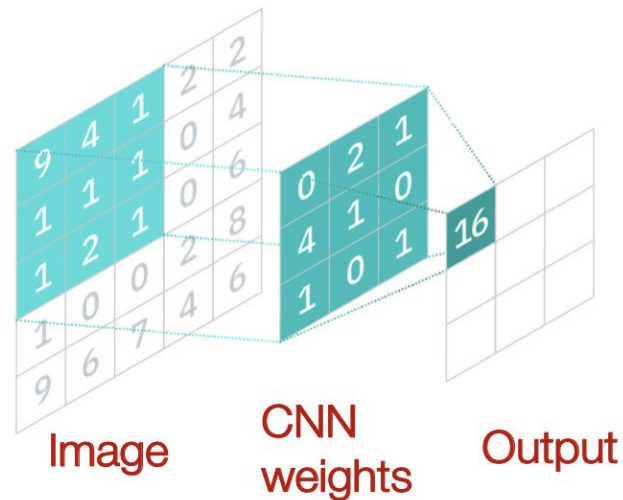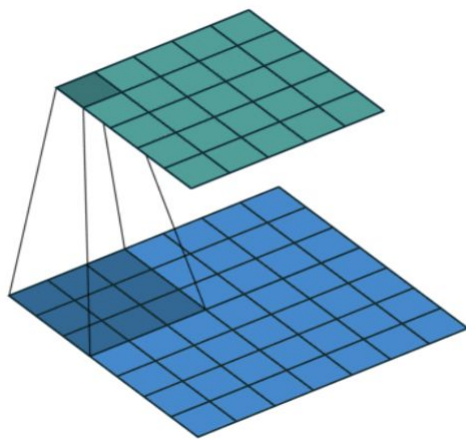
# Review: Aggregate from neighbors

- Nodes aggregate information from their neighbors using neural networks.
- Every node defines a computation graph based on its neighborhood!

# Why GNNs generalize other neural networks

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

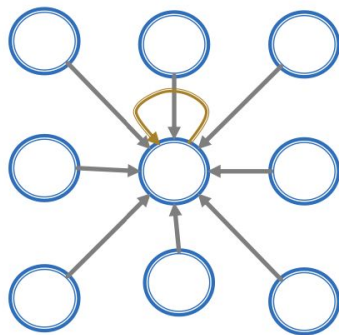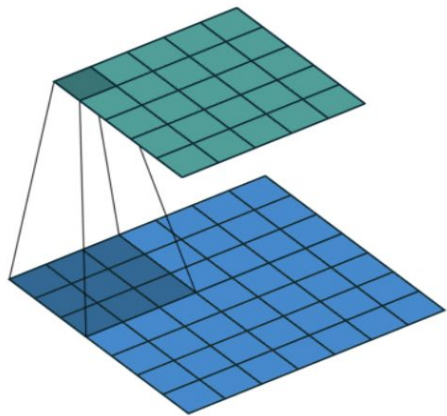Within the same layer, the same filter will be used throughout image, this is referred to as **weight sharing**.



Image

CNN weights

Output

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

# GNN versus CNN

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



Graph

Node order equivariance: graphs often have no inherent ordering present amongst the nodes.

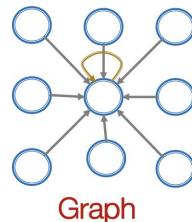Compare this to images, where every pixel is uniquely determined by its absolute position within the image!

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

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

# GNN versus CNN

- Convolutional neural network (CNN) layer with 3x3 filter.
- Key difference: We can learn different $\mathbf{W}_l^u$ for each "neighbor" $u$ for pixel $v$ on the image. The reason: we can pick an ordering for the 9 neighbors using relative position to the center pixel: $\{(-1,-1). (-1,0), (-1, 1), ..., (1, 1)\}$



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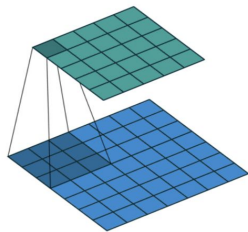
GNN formulation: $h_v^{(l+1)} = \sigma(\mathbf{W}_l \sum_{u \in N(v)} \frac{h_u^{(l)}}{|N(v)|} + B_l h_v^{(l)}), \forall l \in \{0, ..., L-1\}$

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

# GNN versus CNN

- CNN can be seen as a special GNN with fixed neighbor size and ordering.
- The size of the filter is pre-defined for a CNN.
- The advantage of GNN is it processes arbitrary graphs with different degrees for each node.
- CNN is not permutation invariant:
  - Switching the pixels changes the results.

# Transformer

- Transformer is one of the most popular architectures that achieves great performance in many sequence modeling tasks.
- Key component: *self-attention*.
  - Every word "attends" to every other word via matrix calculation.

# Transformer

- Transformer layer can be seen as a special GNN that runs on a fully connected word graph.
- 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.



(Complete) Graph

# A GNN Layer

- GNN Layer = Message + Aggregation
- Different instantiations under this perspective: GCN, GraphSage, GAT, etc.



TARGET NODE

INPUT GRAPH

GNN Layer 1

(2) Aggregation

(1) Message

# A GNN Layer

Idea of a GNN Layer:

- Compress a set of vectors into a single vector
- Two-step process: **Message** and **Aggregation**



Node $v$

(2) Aggregation

(1) Message

Output node embedding $\mathbf{h}_v^{(l)}$

$l$-th GNN Layer

Input node embedding $\mathbf{h}_v^{(l-1)}$ , $\mathbf{h}_{u \in N(v)}^{(l-1)}$
(from node itself + neighboring nodes)

# Message Computation

- Each node will create a message, which will be sent to other nodes later

$$\mathbf{m}_u^{(l)} = \text{MSG}^{(l)}\left(\mathbf{h}_u^{(l-1)}\right)$$

- Example: Linear layer multiply node features with weight matrix.

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



TARGET NODE

INPUT GRAPH

Node $v$

(2) Aggregation

(1) Message

# Message Aggregation

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

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

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

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



TARGET NODE

INPUT GRAPH

Node $v$

(2) Aggregation

(1) Message

# Message Aggregation Issues

Information from node $v$ itself could get lost.

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

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

- Message: compute message from node $v$ itself
- Aggregation: After aggregating from neighbors, we can aggregate the message from node $v$ itself

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

# Single GNN Layer Summary

Message: each node computes a message

$$\mathbf{m}_u^{(l)} = \text{MSG}^{(l)}\left(\mathbf{h}_u^{(l-1)}\right)\ \mathbf{h}_v^{(l-1)}$$

Aggregation: aggregate messages from neighbors

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

Nonlinearity (activation): Adds expressiveness/non-linear learning

- Can be added to message or aggregation
- $\sigma(\cdot)$: ReLU($\cdot$), Sigmoid($\cdot$) , ...



(2) Aggregation

(1) Message

# Classic GCN Layer - 1

Graph Convolutional Networks (GCN)

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

Aggregation: aggregate messages from neighbors



**Message**

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

(2) Aggregation

(1) Message

# Classic GCN Layer - how to write as message

Graph Convolutional Networks (GCN)

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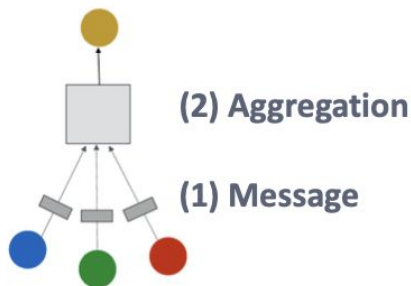
**Message**

**Aggregation**

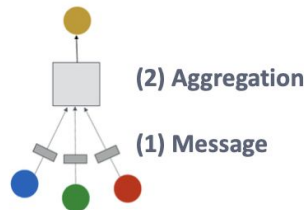Message: each neighbor normalizes by node degree

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

Aggregation: sum over messages from neighbors, then apply activation.

GCN graph assumes self-edges that are included in calculation.

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



(2) Aggregation

(1) Message

# GraphSAGE

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

Message is computed within the AGG()

Two-stage aggregation

- Stage 1: Aggregate from node neighbors

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

- Stage 2: Further aggregate over the node itself

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

# GraphSAGE Neighbor Aggregation

Mean: take weighted average of neighbors

$$AGG = \sum_{u \in N(v)} \frac{\mathbf{h}_u^{(l-1)}}{|N(v)|}$$

**Aggregation**          **Message computation**

Pool: Transform neighbor vectors and apply symmetric vector function Mean($\cdot$) or Max($\cdot$)

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

**Aggregation**          **Message computation**

LSTM: Apply LSTM to reshuffled neighbors

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

**Aggregation**

# GraphSAGE L2 Normalization

- Without L2 normalization, the embedding vectors have different scales.
- In some cases (not always), normalization of embedding results in performance improvement.
- After L2 normalization, all vectors will have the same L2-norm.

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

# Adding Attention

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

**Attention weights**

Weighting factor (importance) of node $u$'s message to node $v$. Node degree.

$$\alpha_{vu} = \frac{1}{|N(v)|}$$

# Graph Attention Network

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

**Attention weights**

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 Network

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

**Attention weights**

Can we do better than simple neighborhood aggregation?

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

Goal: Learn importance to different neighbors of each node in the graph based on context.

Compute embedding $h$ 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

Compute attention coefficients **e** across pairs of nodes **u, v** based on their messages.

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

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

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