

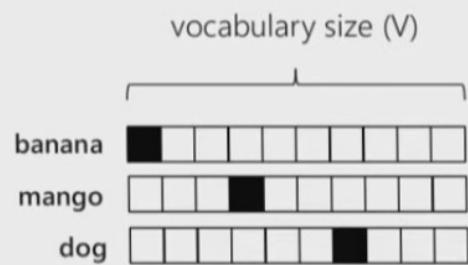
Graph Neural Networks

Representations

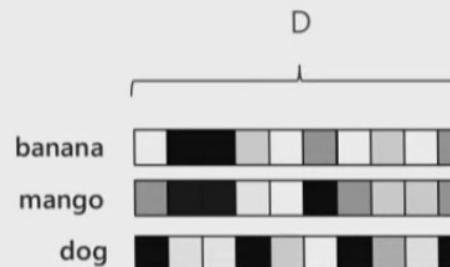
- Different ways of representing information
- Graph
- Matrix
- Vectors

- We can Transform from one representation to another.
- Encode it to obtain an embedded vector representation

Distributed Vector Representations



Local representation
(1-hot)

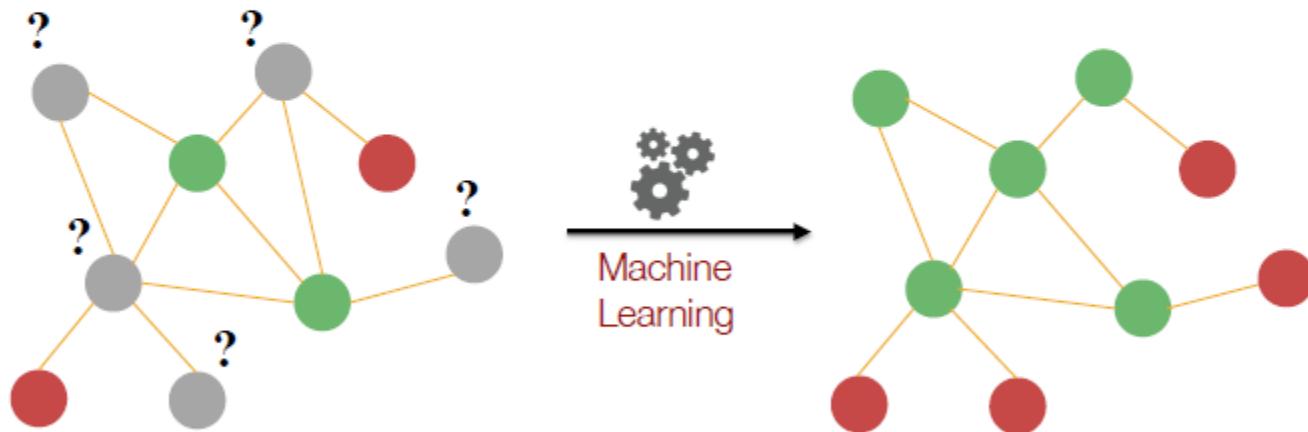


Distributed representation

Graph Representation Learning

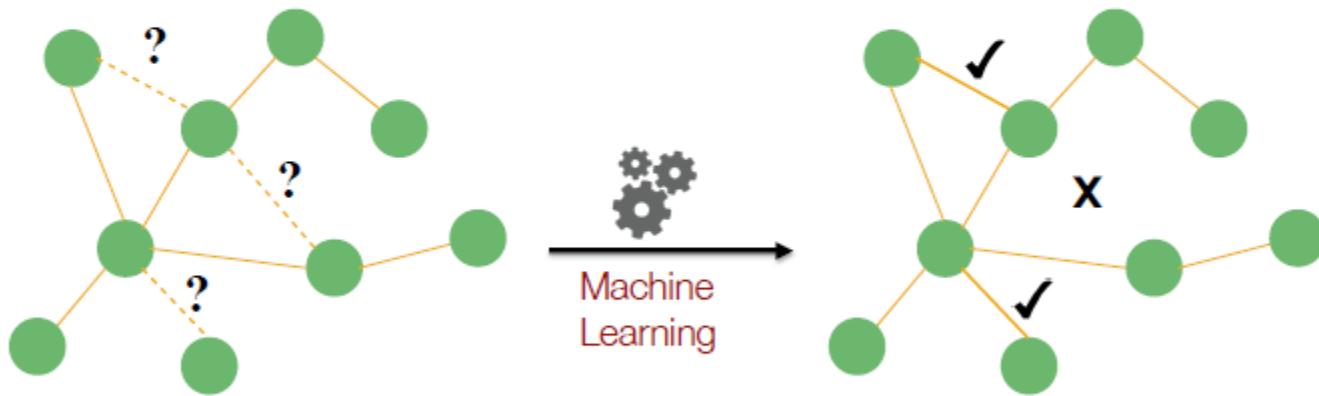


Machine Learning in Networks



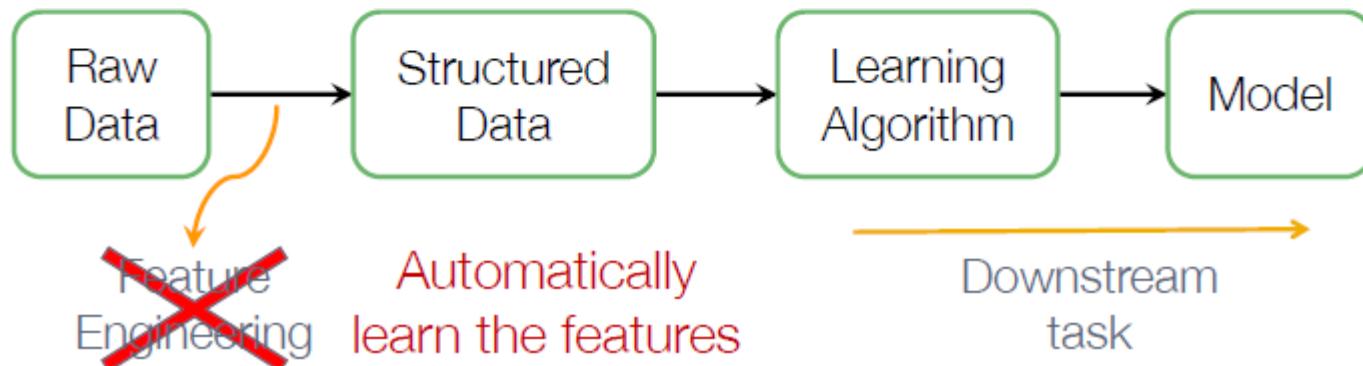
Node classification

Example: Link Prediction



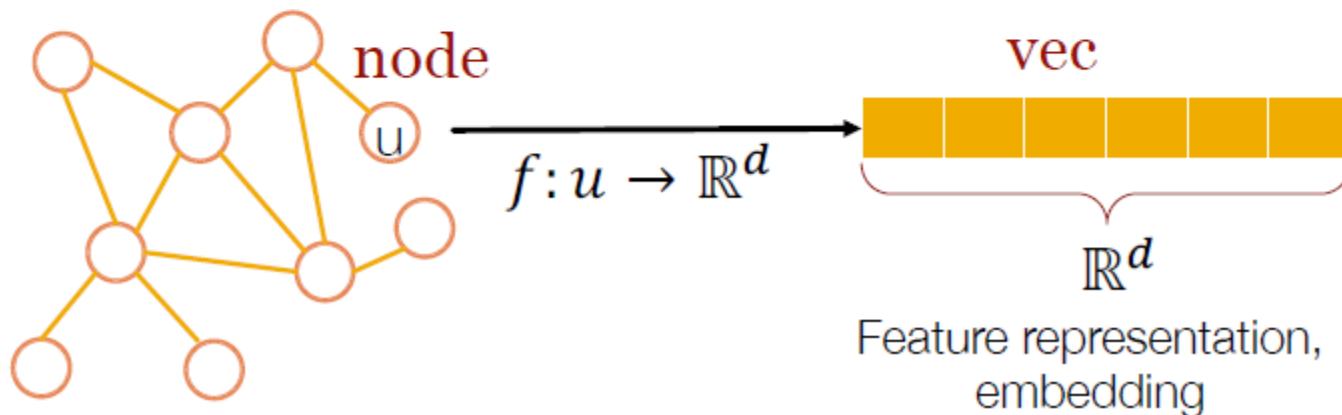
Machine Learning Lifecycle

- **(Supervised) Machine Learning Lifecycle requires feature engineering **every single time!****



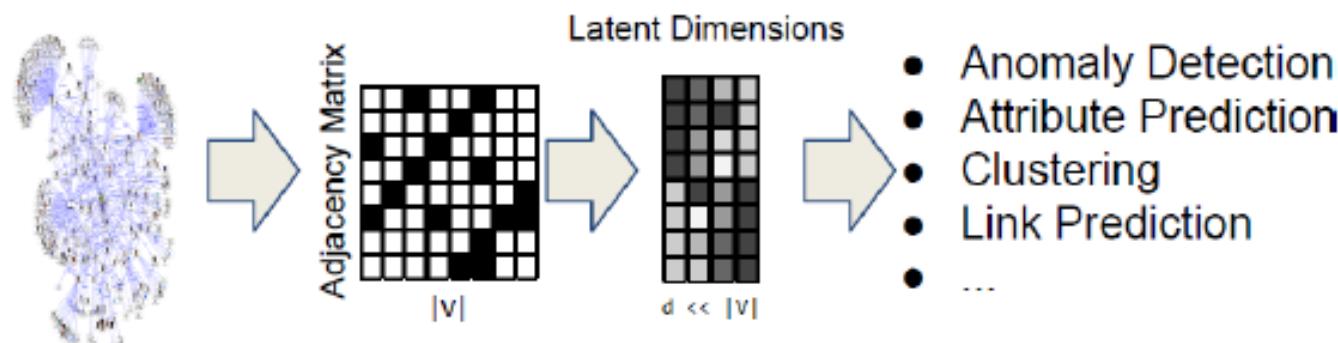
Feature Learning in Graphs

Goal: Efficient task-independent feature learning
for machine learning
in networks!



Why network embedding?

- Task: We map each node in a network into a low-dimensional space
 - Distributed representation for nodes
 - Similarity of embedding between nodes indicates their network similarity
 - Encode network information and generate node representation



Example Node Embedding

- 2D embedding of nodes of the Zachary's Karate Club network:

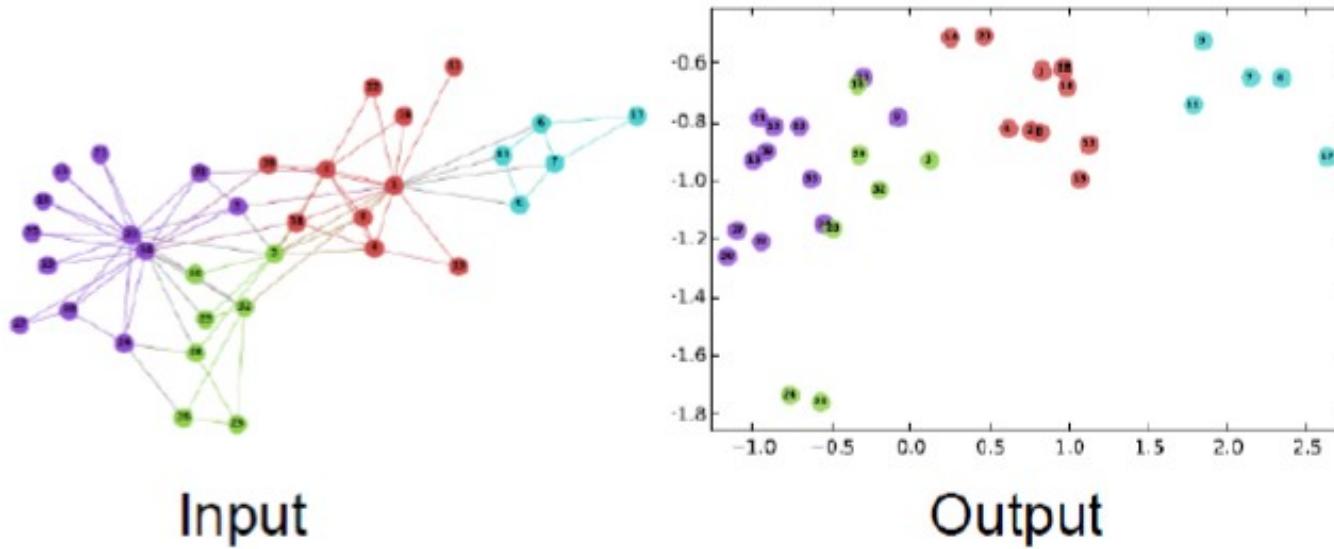
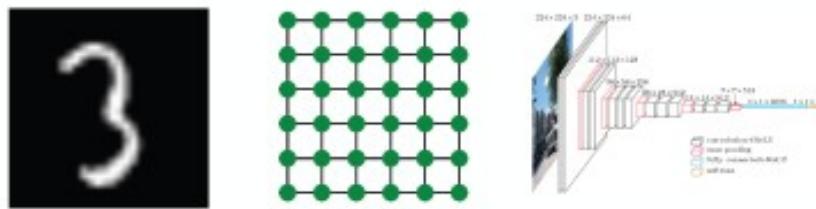


Image from: [Perozzi et al.](#). DeepWalk: Online Learning of Social Representations. *KDD 2014*.

Why Is It Hard?

- Modern deep learning toolbox is designed for simple sequences or grids.
 - CNNs for fixed-size images/grids....



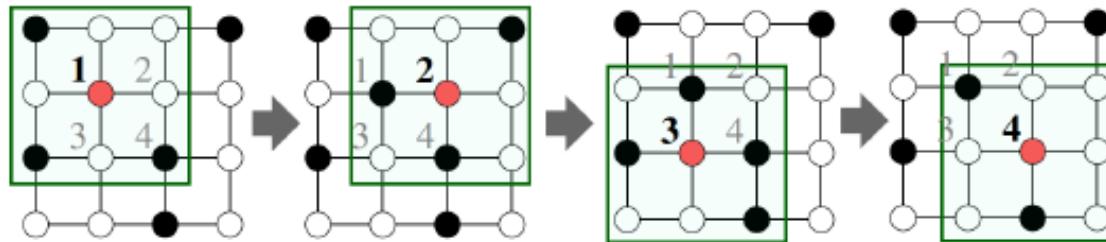
- RNNs or word2vec for text/sequences...



Why Is It Hard?

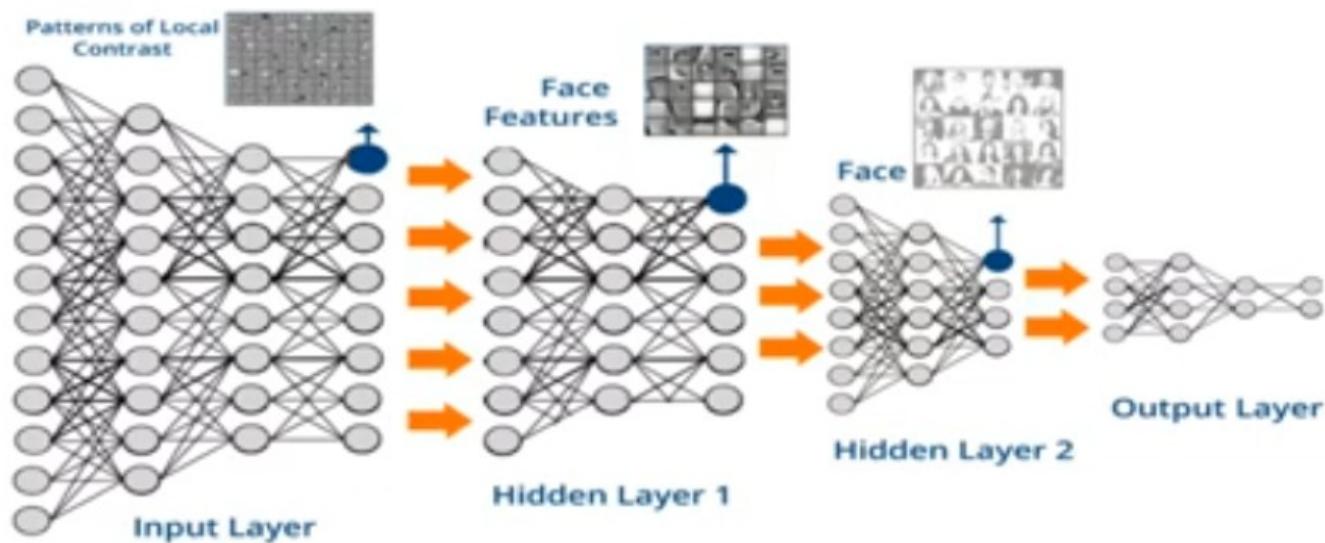
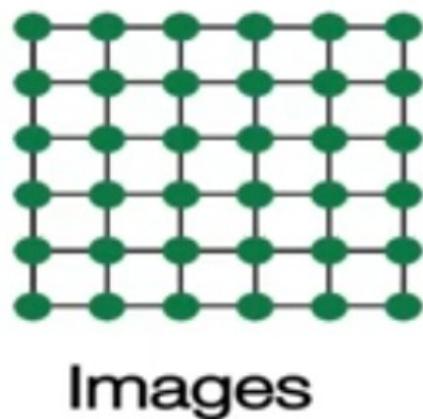
- But networks are far more complex!

- Complex topographical structure
(i.e., no spatial locality like grids)



- No fixed node ordering or reference point (i.e., the isomorphism problem)
 - Often dynamic and have multimodal features.

Why Convolution fails on Graphs?

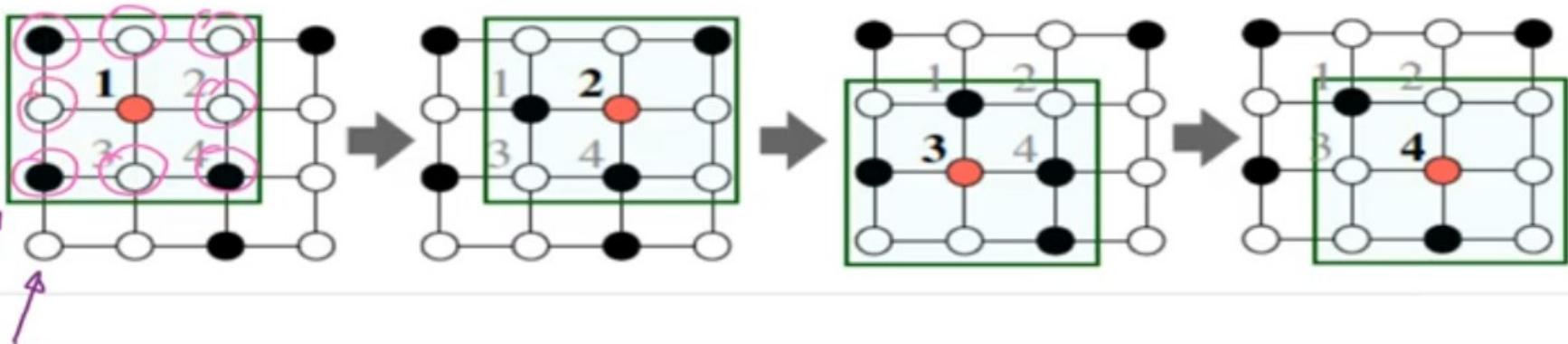


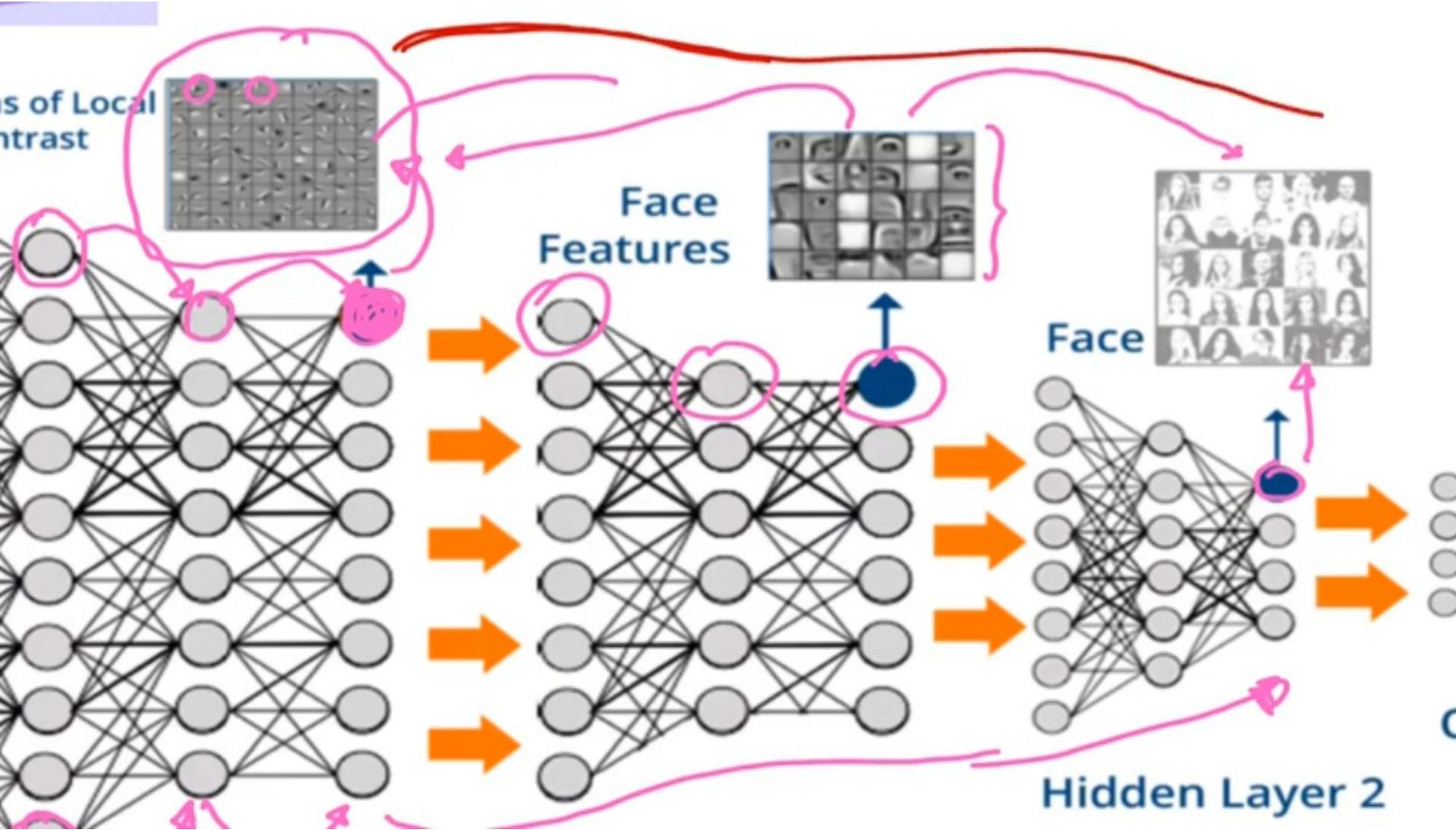
CNN

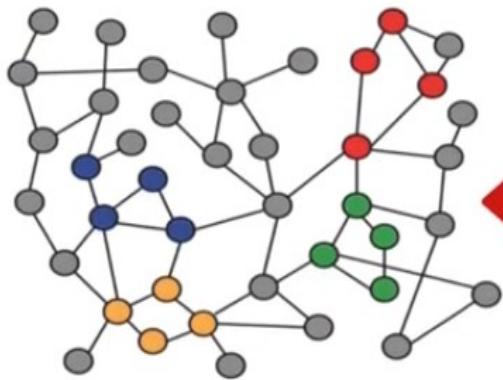
Locality } →

Aggregation }

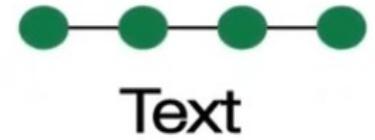
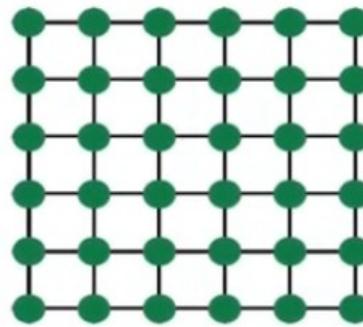
{ Composition (function of a function)







VS.



- Arbitrary size
- Complex topology
(no spatial locality)
- Graphs are Non-Euclidean
- No fixed node ordering

- Grid (spatial locality)
- Structured data
- Deep Learning is designed
for grids & sequences.

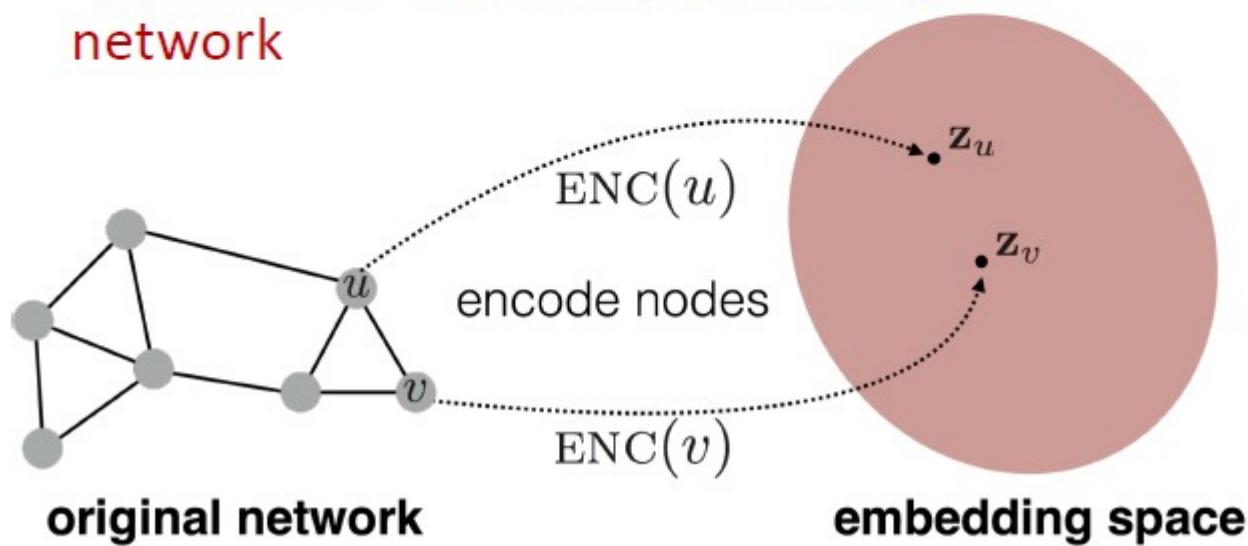
Embedding Nodes

Setup

- Assume we have a graph G :
 - V is the vertex set.
 - A is the adjacency matrix (assume binary).
 - No node features or extra information is used!

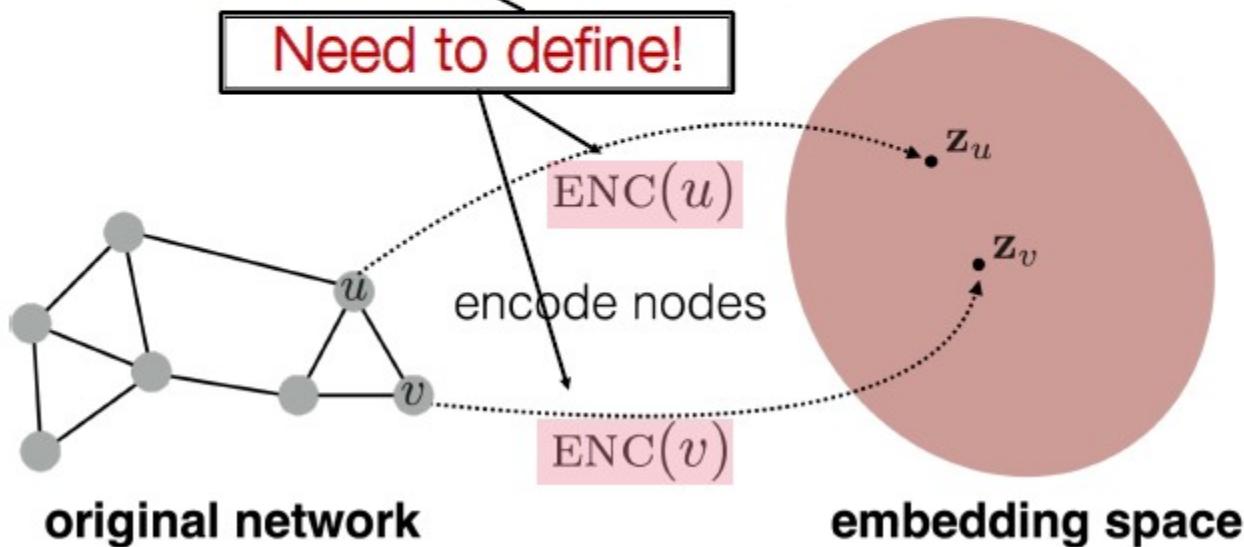
Embedding Nodes

- Goal is to encode nodes so that similarity in the embedding space (e.g., dot product) approximates similarity in the original network



Embedding Nodes

Goal: $\text{similarity}(u, v) \approx \mathbf{z}_v^\top \mathbf{z}_u$
in the original network Similarity of the embedding



Learning Node Embeddings

1. Define an encoder (i.e., a mapping from nodes to embeddings)
 2. Define a node similarity function (i.e., a measure of similarity in the original network).
 3. Optimize the parameters of the encoder so that:

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

in the original network
Similarity of the embedding

Two Key Components

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

$\text{ENC}(v) = \mathbf{z}_v$ embedding
node in the input graph

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

Similarity of u and v in the original network $\text{similarity}(u, v) \approx \mathbf{z}_v^\top \mathbf{z}_u$
dot product between node embeddings

“Shallow” Encoding

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

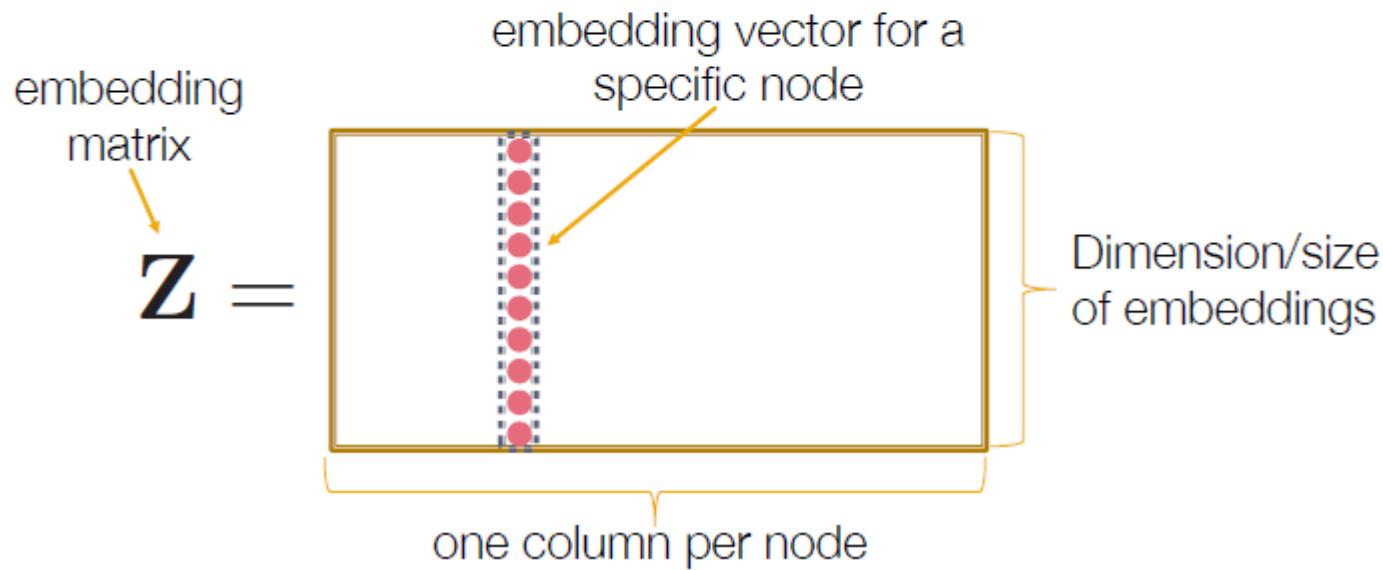
$$\text{ENC}(v) = \mathbf{Z}\mathbf{v}$$

$\mathbf{Z} \in \mathbb{R}^{d \times |\mathcal{V}|}$ matrix, each column is node embedding [what we learn!]

$\mathbf{v} \in \mathbb{I}^{|\mathcal{V}|}$ indicator vector, all zeroes except a one in column indicating node v

“Shallow” Encoding

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



“Shallow” Encoding

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

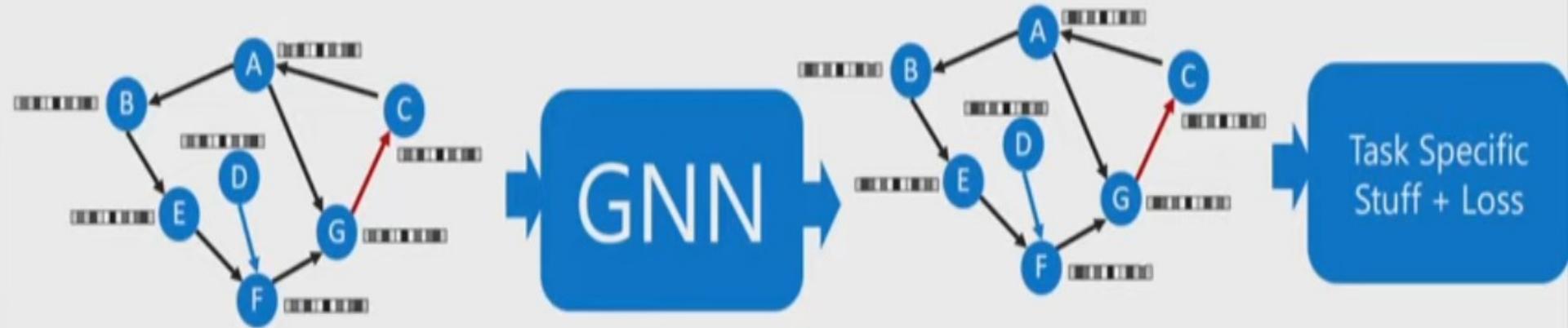
**Each node is assigned a unique
embedding vector**

Many methods: node2vec, DeepWalk, LINE

How to Define Node Similarity?

- Key choice of methods is **how they define node similarity**.
- E.g., should two nodes have similar embeddings if they....
 - are connected?
 - share neighbors?
 - have similar “structural roles”?
 - ...?

Graph Neural Networks

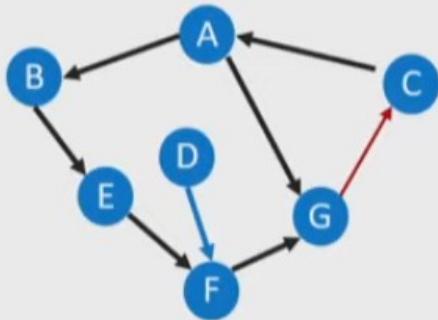


Initial Representation
of each node

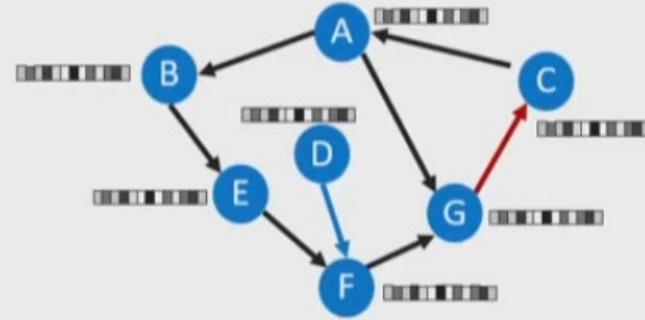
Output Representations
of each Node

Task Specific
Stuff + Loss

Graph Neural Networks

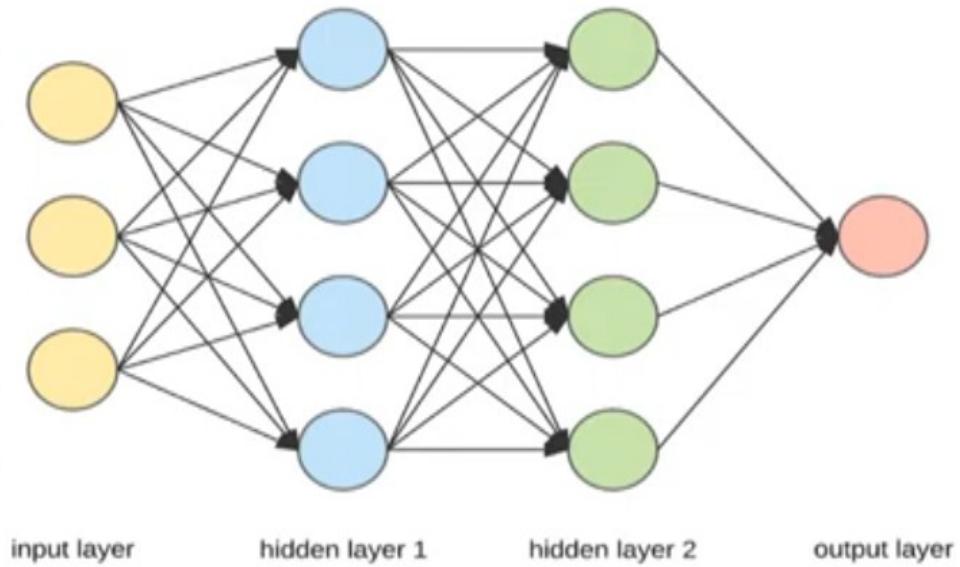
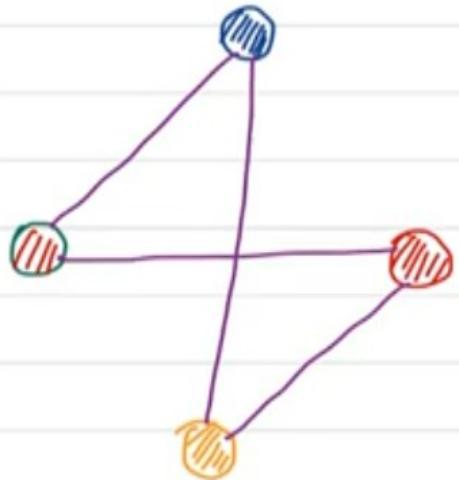


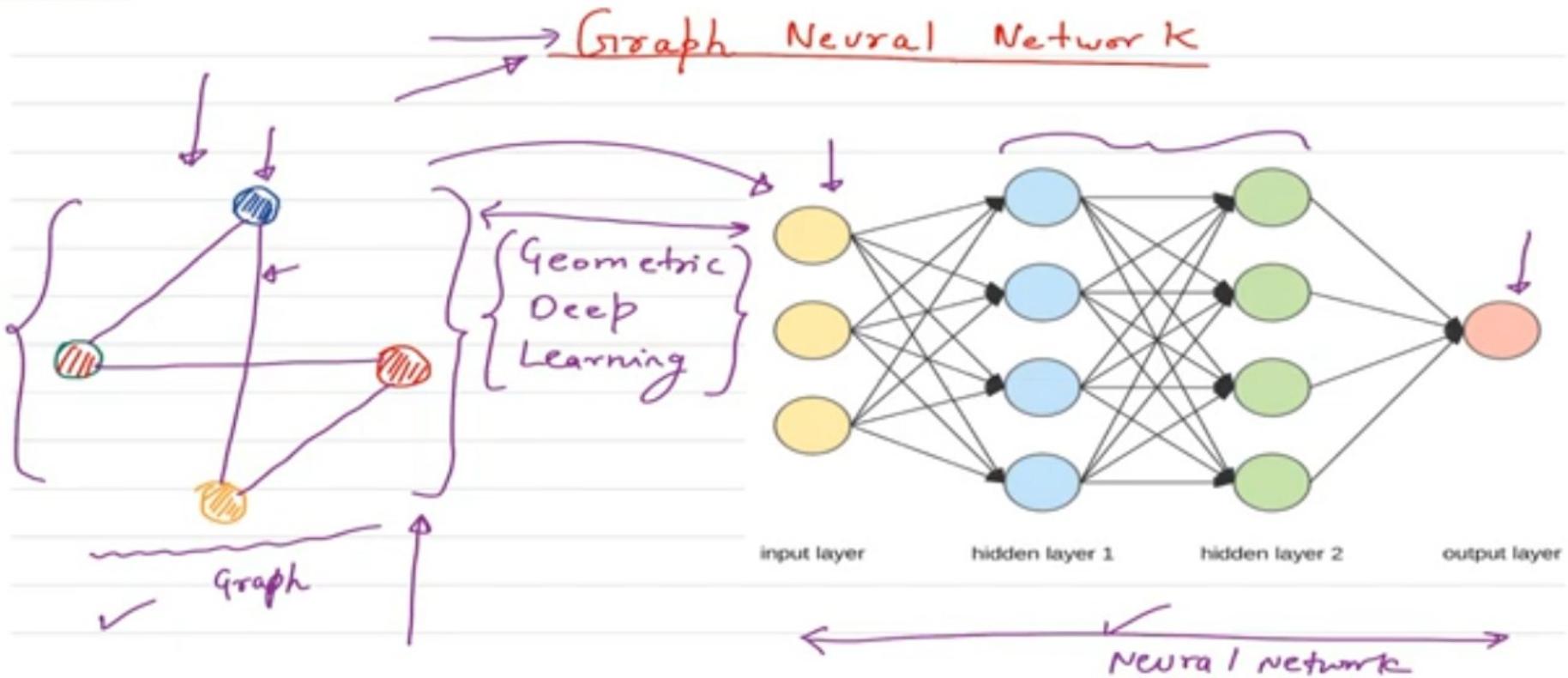
Graph Representation
of Problem



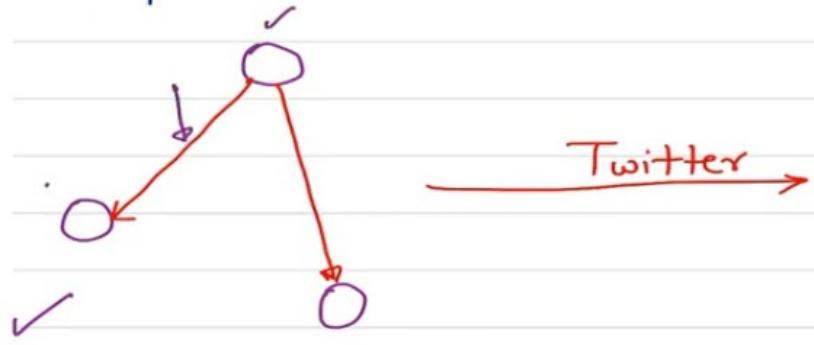
Initial Representation
of each node

→ Graph Neural Network

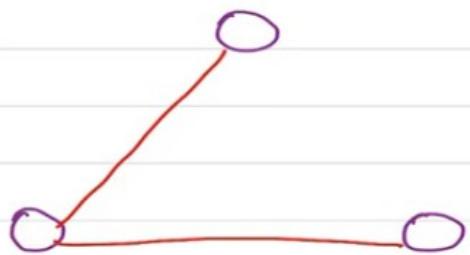
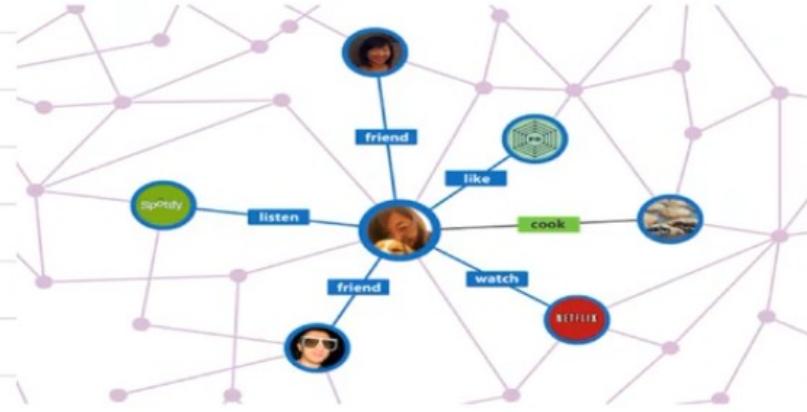




Graphs can be directed or undirected.



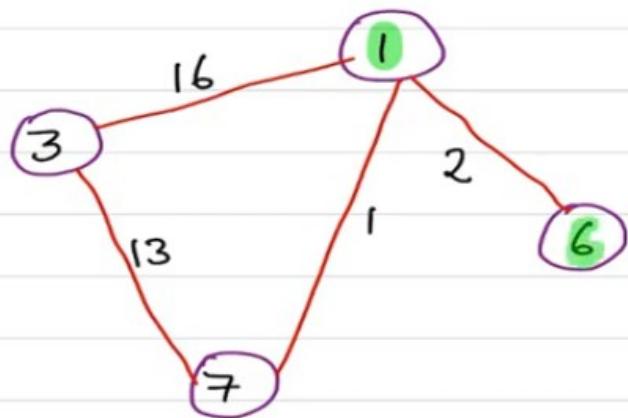
Twitter



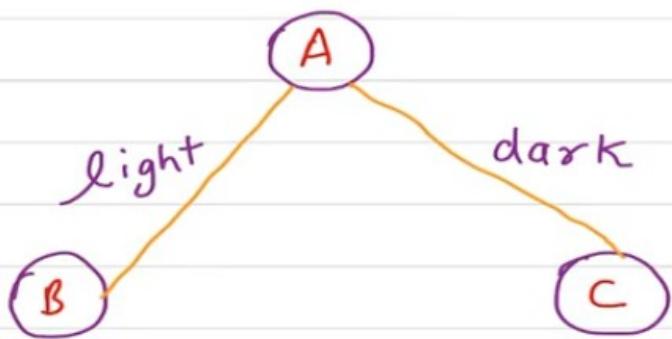
→

Facebook friend requests.

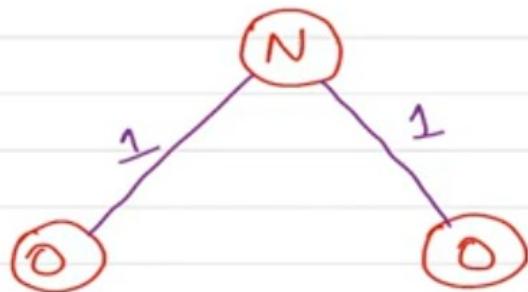
Graphs can have labels on their edges / nodes.



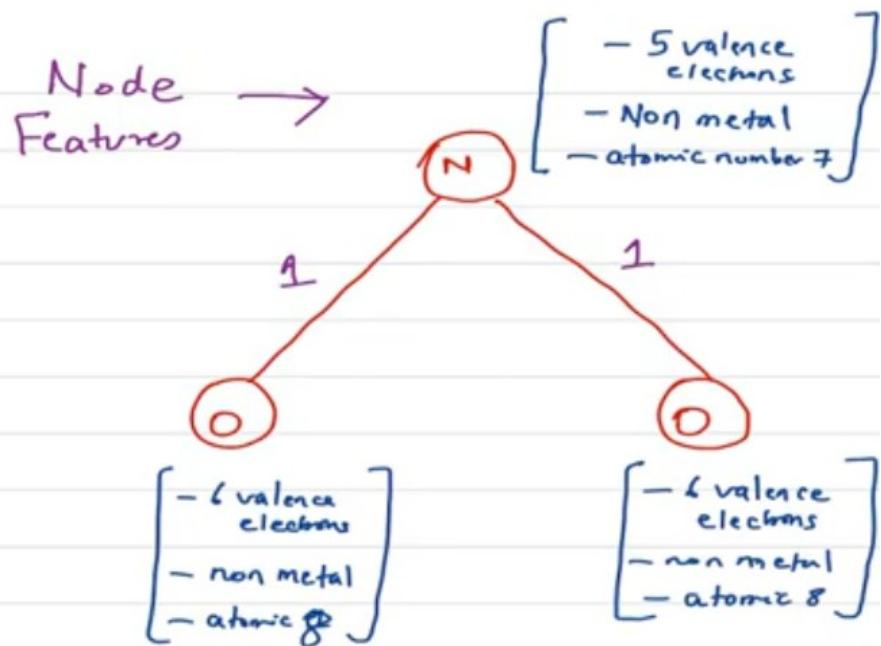
Numerical labels



Textual labels



NO₂: Unique labels



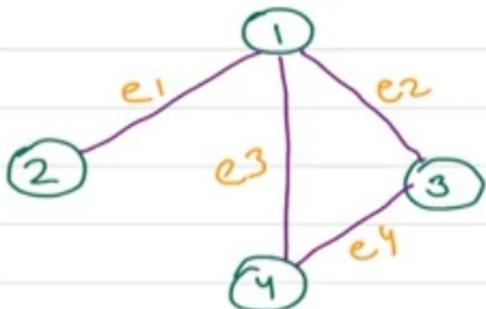
Note: Do not mix labels with features.

Analogy: Node is a person, node label is person name & node feature are person's characteristics.

How to convert graph in a format that computer can accept.

Incidence matrix $\xrightarrow{(I)}$ A matrix of size $n \times m$ where n are number of nodes & edges of a graph

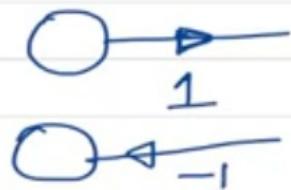
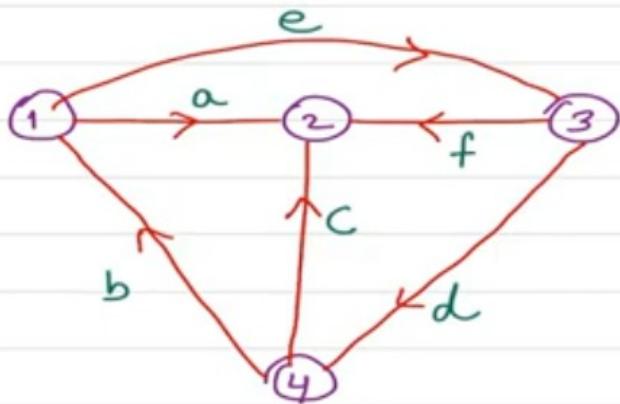
Undirected Graph



Entry is 1
if node &
edge are
incident
(related)
& 0 if not.

	e_1	e_2	e_3	e_4
1	1	1	1	0
2	1	0	0	0
3	0	1	0	1
4	0	0	1	1

Directed Graph.

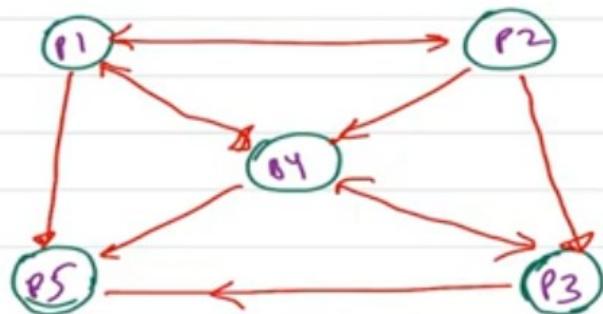


$$\begin{matrix} & \begin{matrix} a & b & c & d & e & f \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \end{matrix} & \left[\begin{matrix} 1 & -1 & 0 & 0 & 1 & 0 \\ -1 & 0 & -1 & 0 & 0 & -1 \\ 0 & 0 & 0 & 1 & -1 & 1 \\ 0 & 1 & 1 & -1 & 0 & 0 \end{matrix} \right] \end{matrix}$$

I

Adjacency Matrix (A)

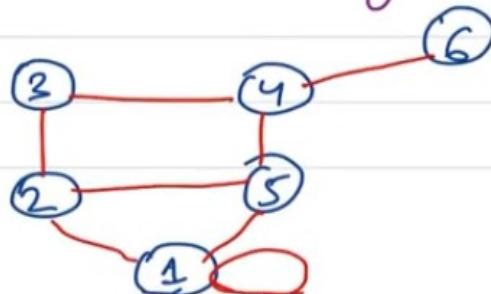
It is also called as connection matrix, is a square matrix. The component of the matrix is 0 or 1 depending upon whether the nodes of the graph are adjacent or not



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

An adjacency matrix can be weighted, which means each edge has an associated value. So instead of 1s, the value is put in the respective matrix co-ordinate.

Degree matrix (D) → A diagonal matrix which contains information about the degree of each node — that is the number of edges attached to each node.



$$\begin{matrix} & 1 & 2 & 3 & 4 & 5 & 6 \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{matrix} & \left[\begin{matrix} 4 & & & & & \\ & 3 & & & & \\ & & 2 & & & \\ & & & 3 & & \\ & & & & 3 & \\ & & & & & 1 \end{matrix} \right] \end{matrix}$$

Laplacian matrix (Graph Laplacian)

$$L = \{D - A\}$$

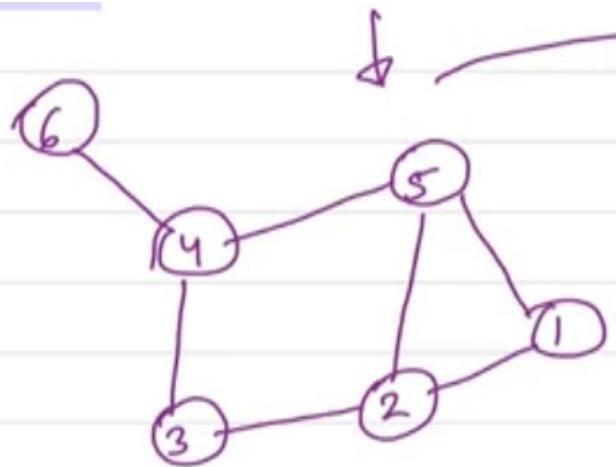
Also known as Laplace-Beltrami operator is a measure of smoothness of a vertex i.e how quickly it changes b/w adjacent vertex

Laplacian matrix (Graph Laplacian)

$$L = \{D - A\}$$

Also known as {Laplace-Beltrami operator} is a measure of smoothness of a vertex i.e. how quickly it changes b/w adjacent vertex.

$$\nabla f = \frac{f(x+1, y) - f(x, y)}{\text{smoothness of the image}}$$



Degree matrix

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

Adjacency

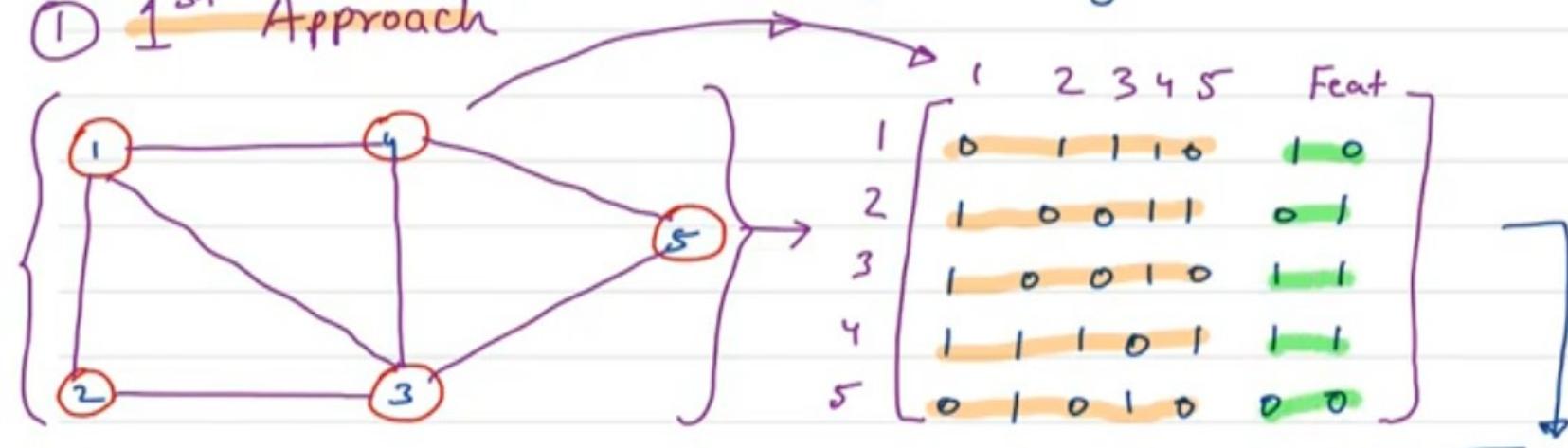
$$\begin{bmatrix} 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix}$$

Laplacian: $D - A$

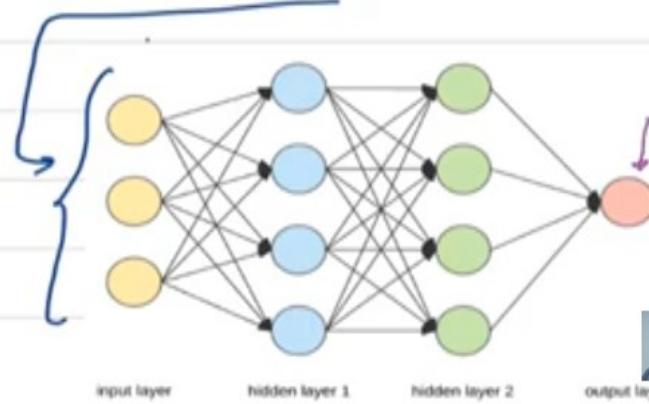
$$\begin{bmatrix} 2 & -1 & 0 & 0 & -1 & 0 \\ -1 & 3 & -1 & 0 & -1 & 0 \\ 0 & -1 & 2 & -1 & 0 & 0 \\ 0 & 0 & -1 & 3 & -1 & -1 \\ -1 & -1 & 0 & -1 & 3 & 0 \\ 0 & 0 & 0 & -1 & 0 & 1 \end{bmatrix}$$

How to train Graphs using Neural Network ?

① 1st Approach

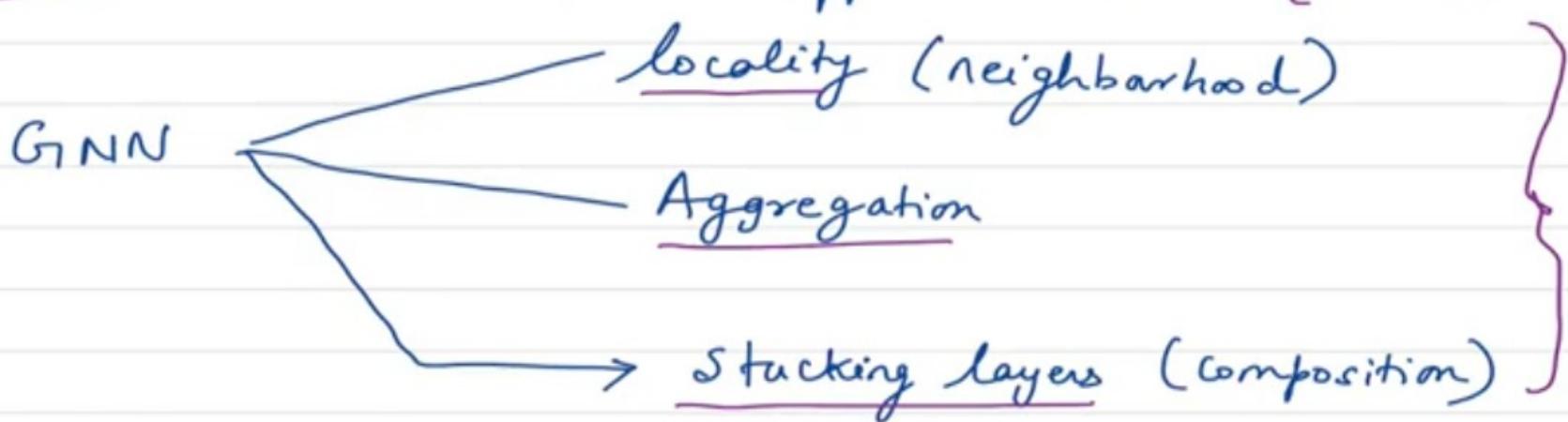


- ⊗ Not invariant to node ordering
- ⊗ Not applicable to graphs of different sizes.

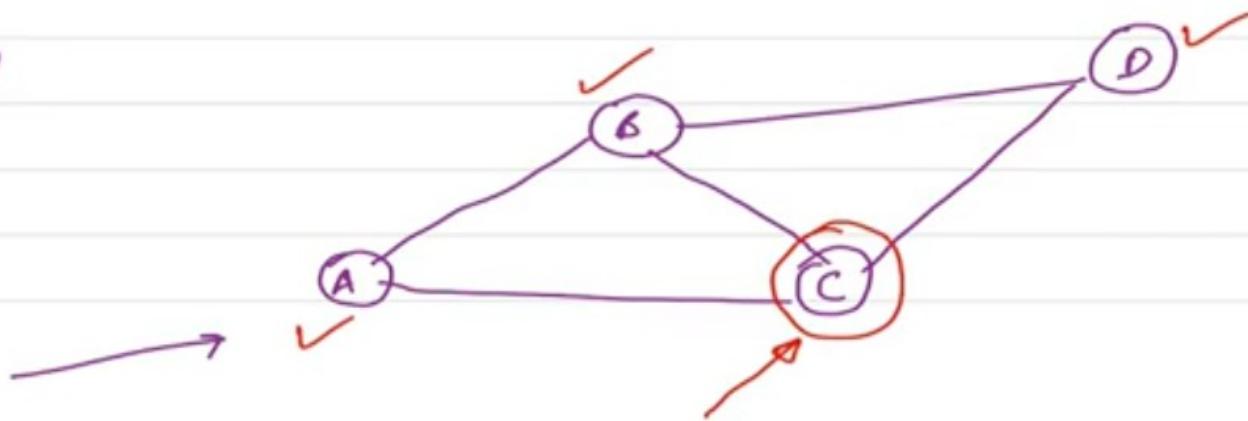


2nd Approach

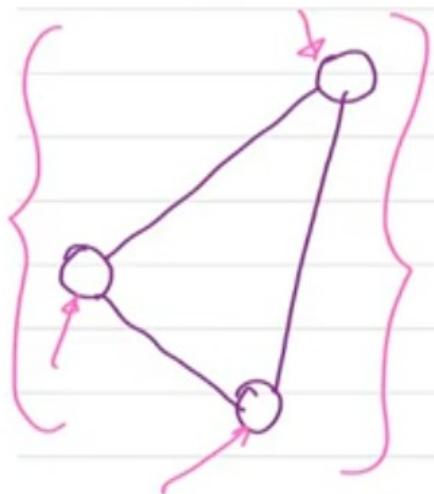
Use same approach as in {CNN}



(Ex)



For any graph



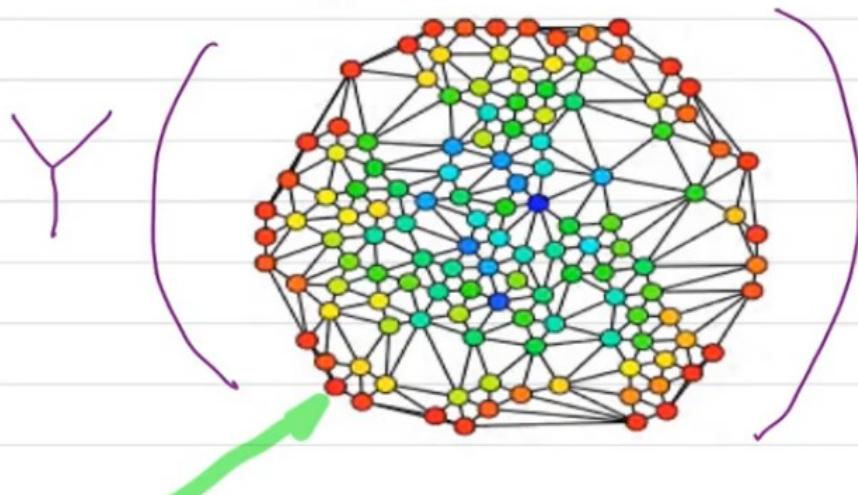
→ Graph $\underline{G = (V, E)}$

⊗ A is adjacency matrix of size
 $N \times N$

⊗ X is the node feature matrix of
size $(N \times M)$

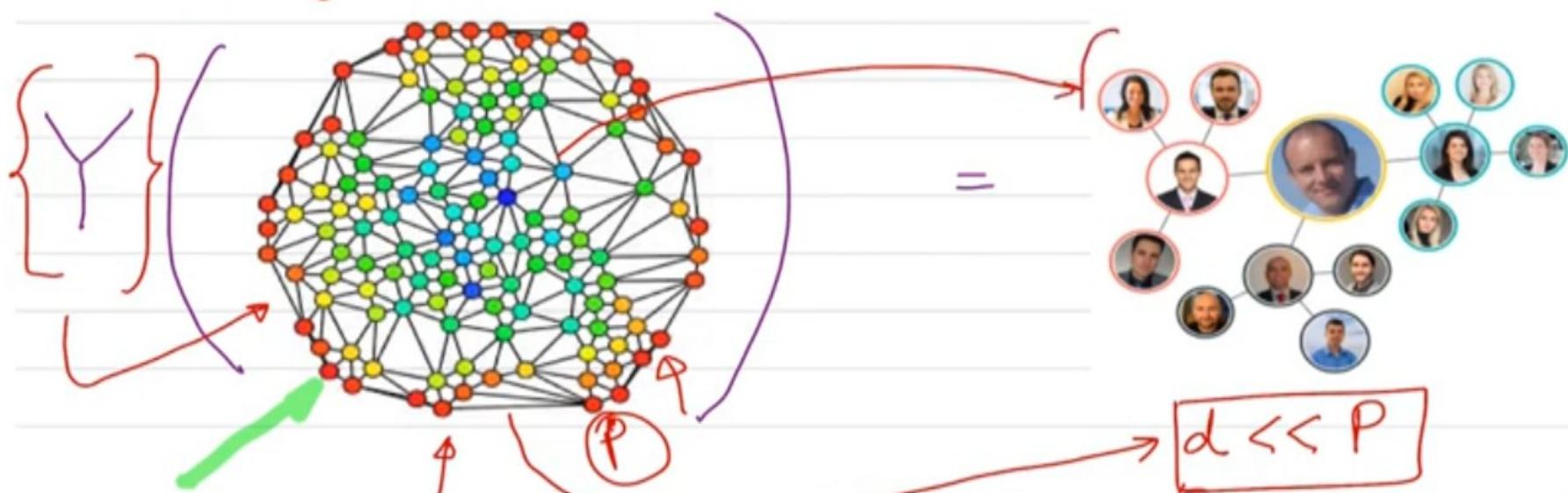
Node Embedding.

Target: Map nodes to d dimensional embeddings such that similar nodes in the graph are embedded close together

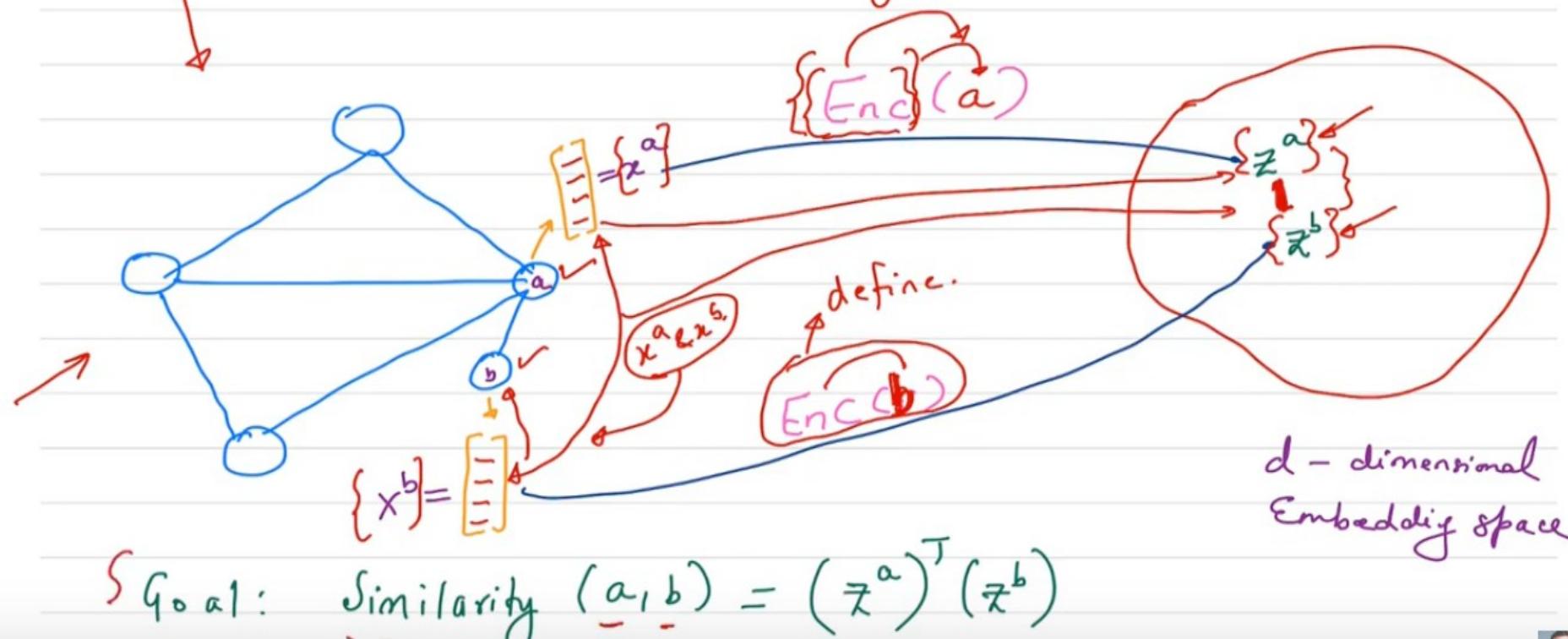


Node Embedding

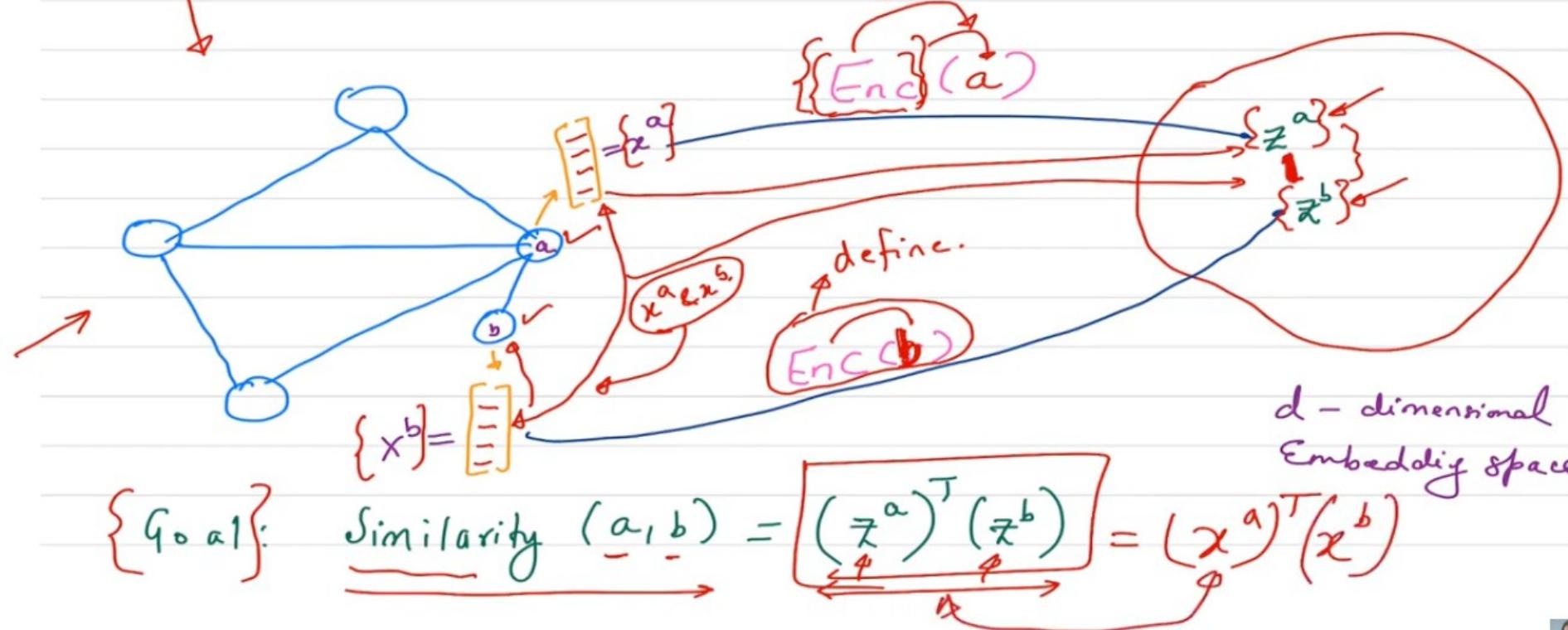
Target: Map nodes to d dimensional embeddings such that similar nodes in the graph are embedded close together



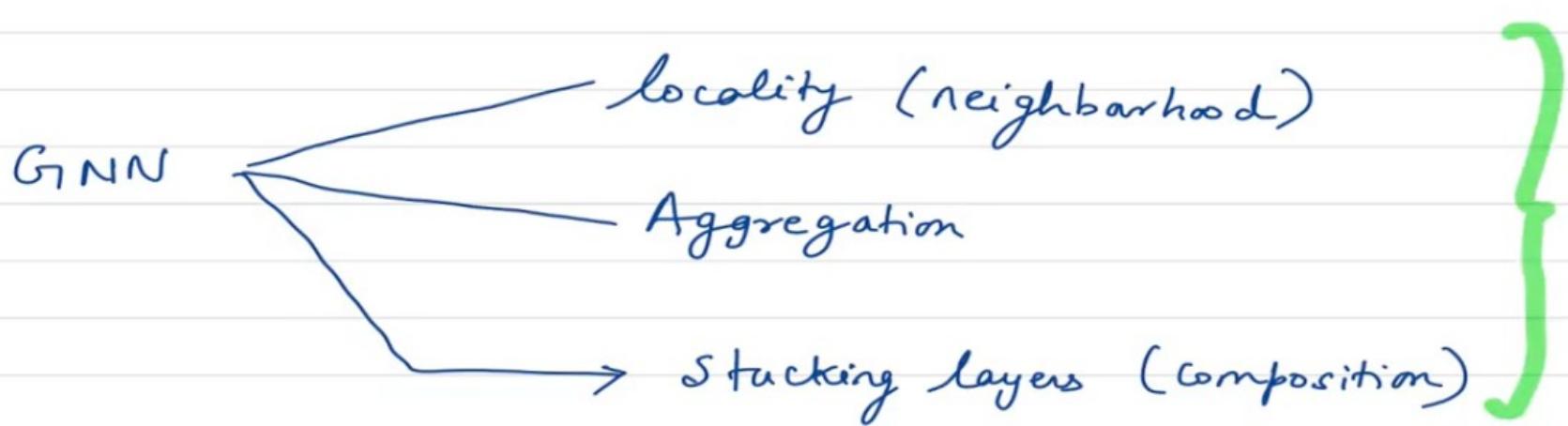
Node Embeddings



Node Embeddings



How to come up with this Encoder function.

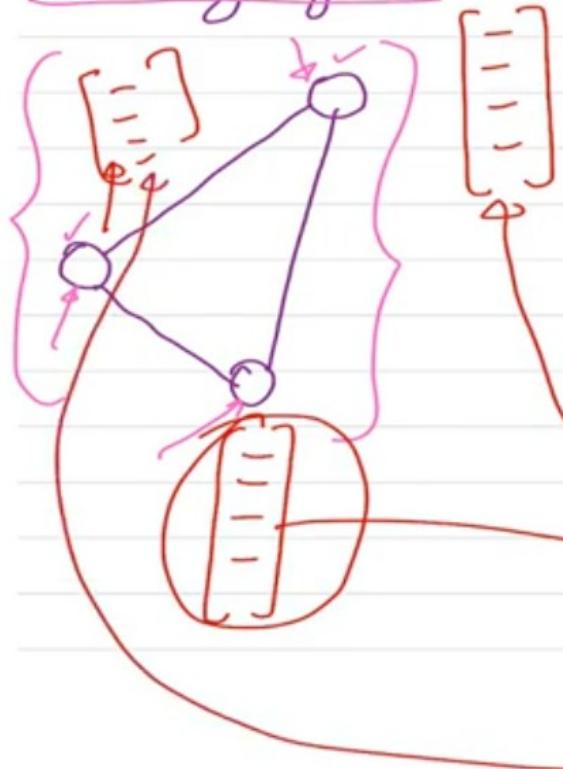


Graph Neural Networks

and Neural Message Passing



For any graph



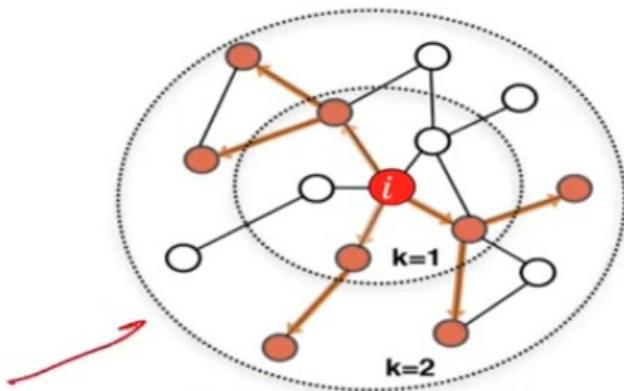
→ Graph $G = (V, E)$

⊗ A is adjacency matrix of size
 $\{N \times N\}$

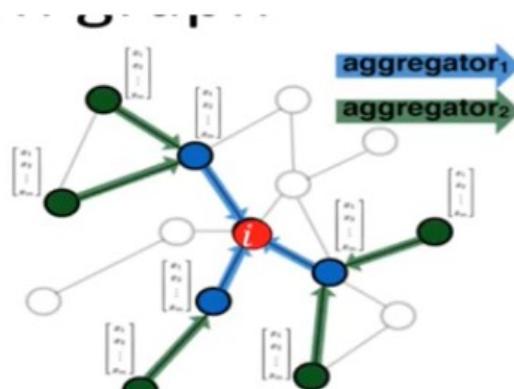
⊗ $\{X\}$ is the node feature matrix of
size $(N \times M)$



* Locality information using computational graph.

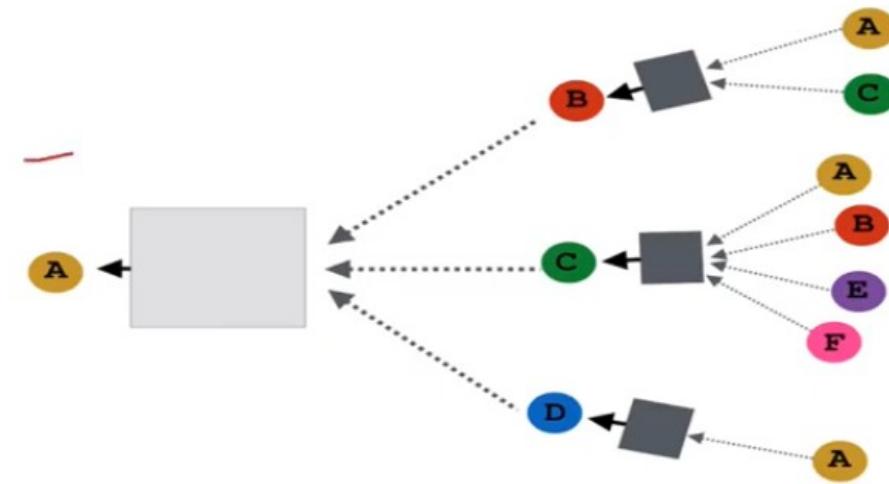
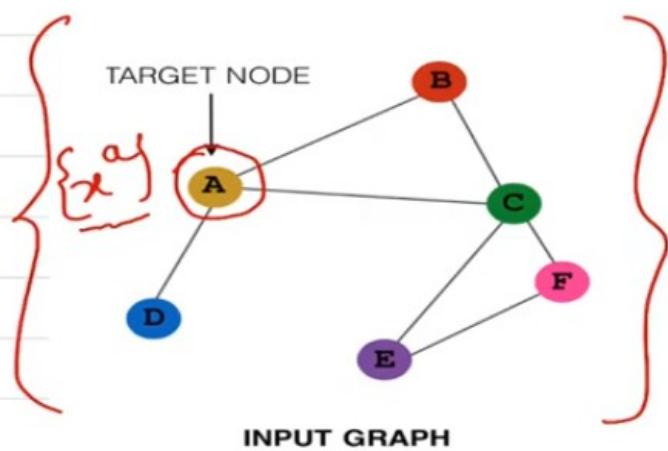


Determine node computation graph



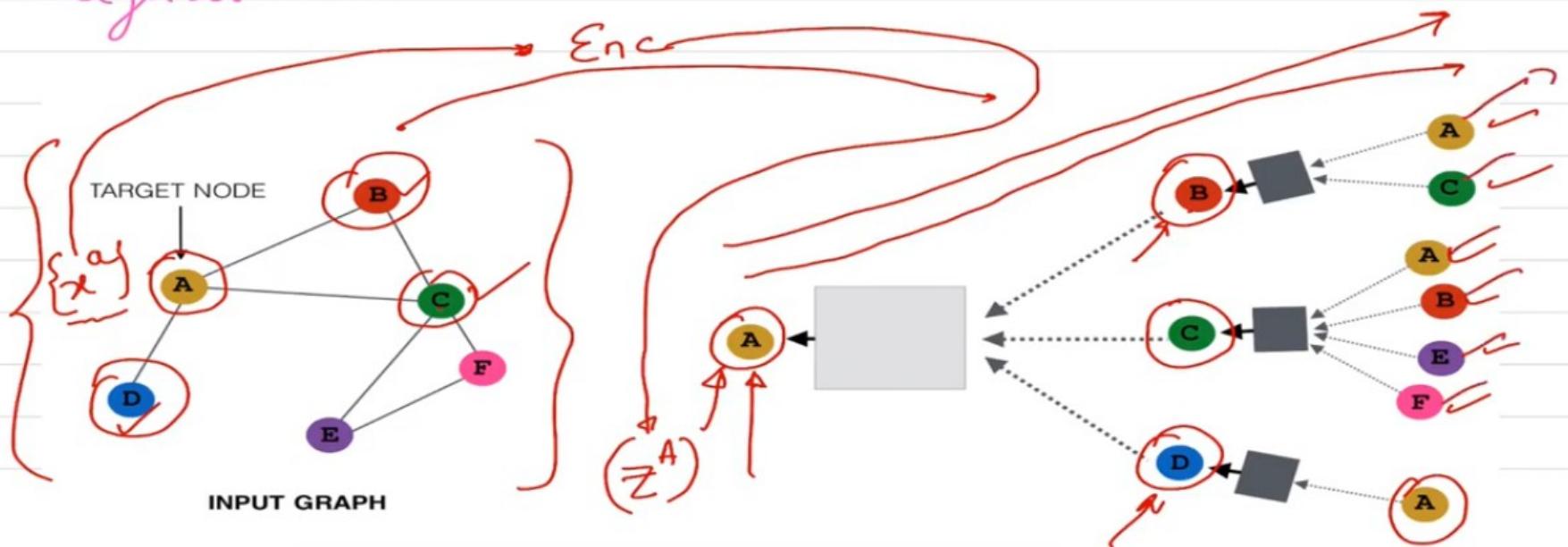
Propagate and transform information

→ Generate node embedding based on local neighborhood.



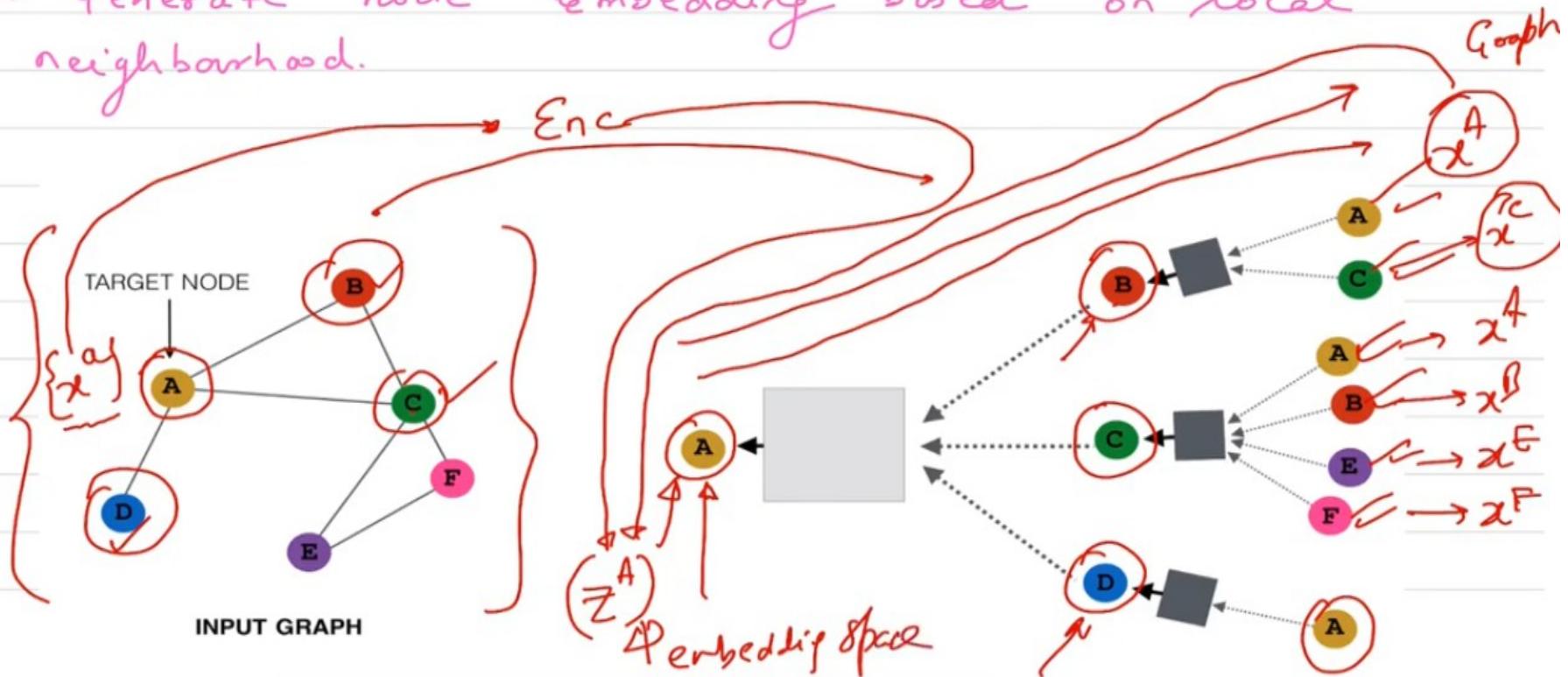
Aggregate (neighbors)

→ Generate node embedding based on local neighborhood.

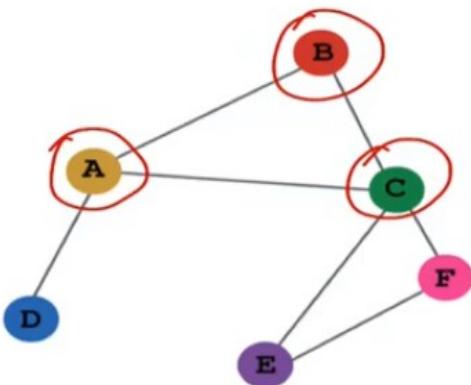


Aggregate (neighbors)

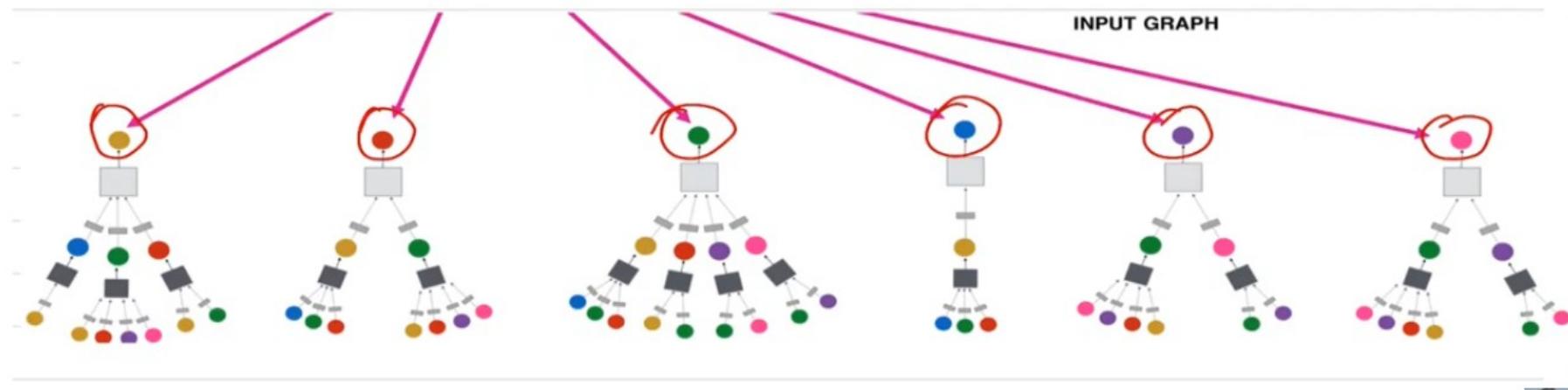
→ Generate node embedding based on local neighborhood.



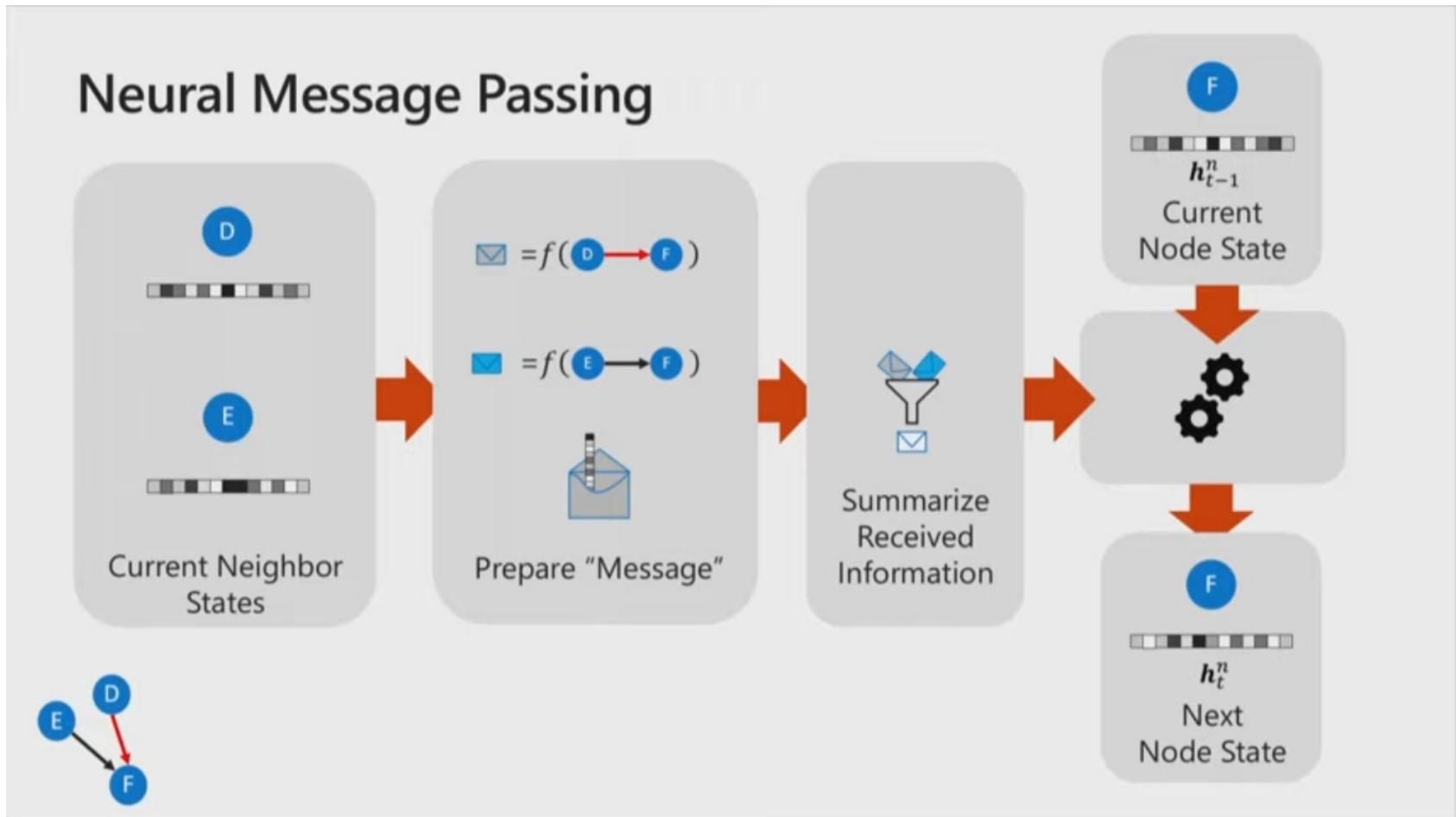
Aggregate (Neighbors)



Every node has
a computational graph



Neural Message Passing



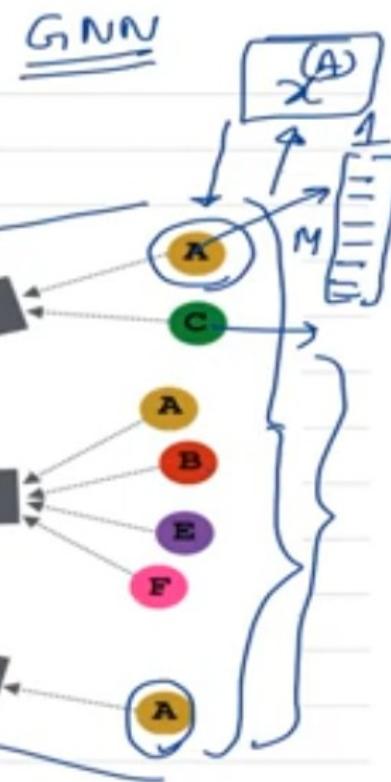
Forward - Propagation Rule

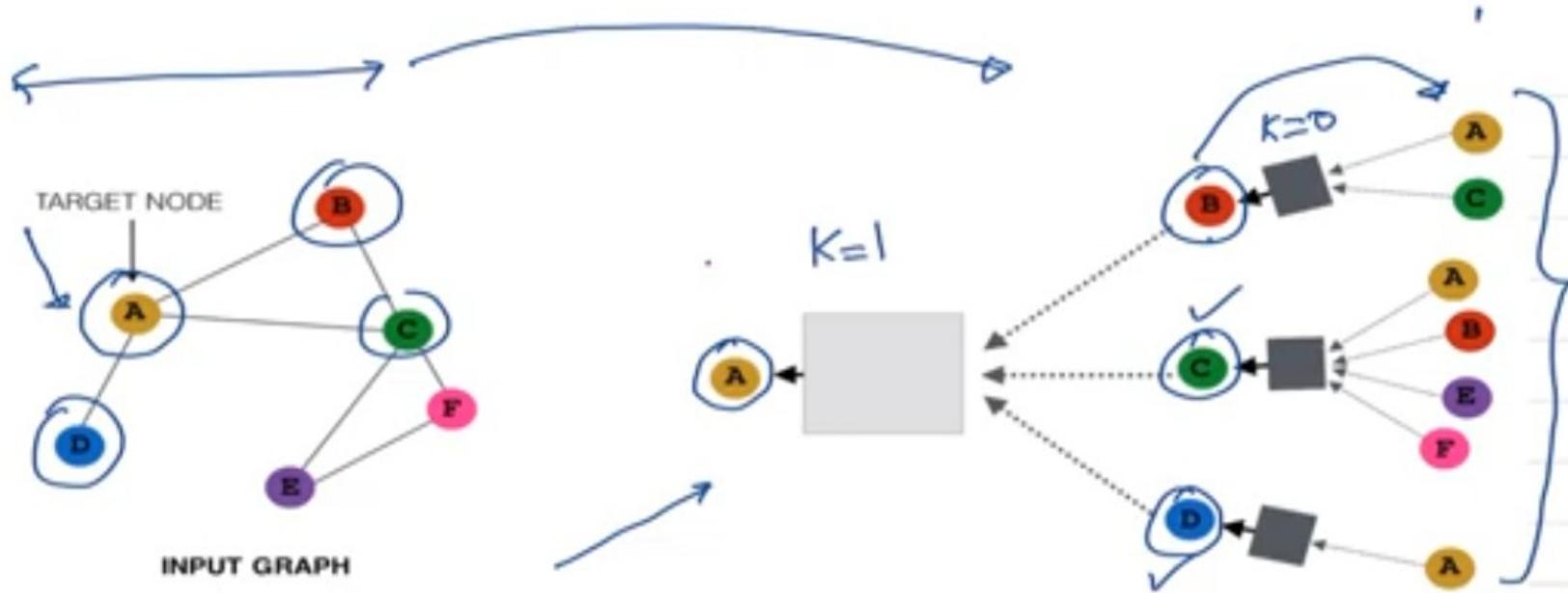


$$\begin{bmatrix} 1 \\ \vdots \\ N \end{bmatrix} = z^{(1)}$$

Embedding

$$(N \ll M)$$

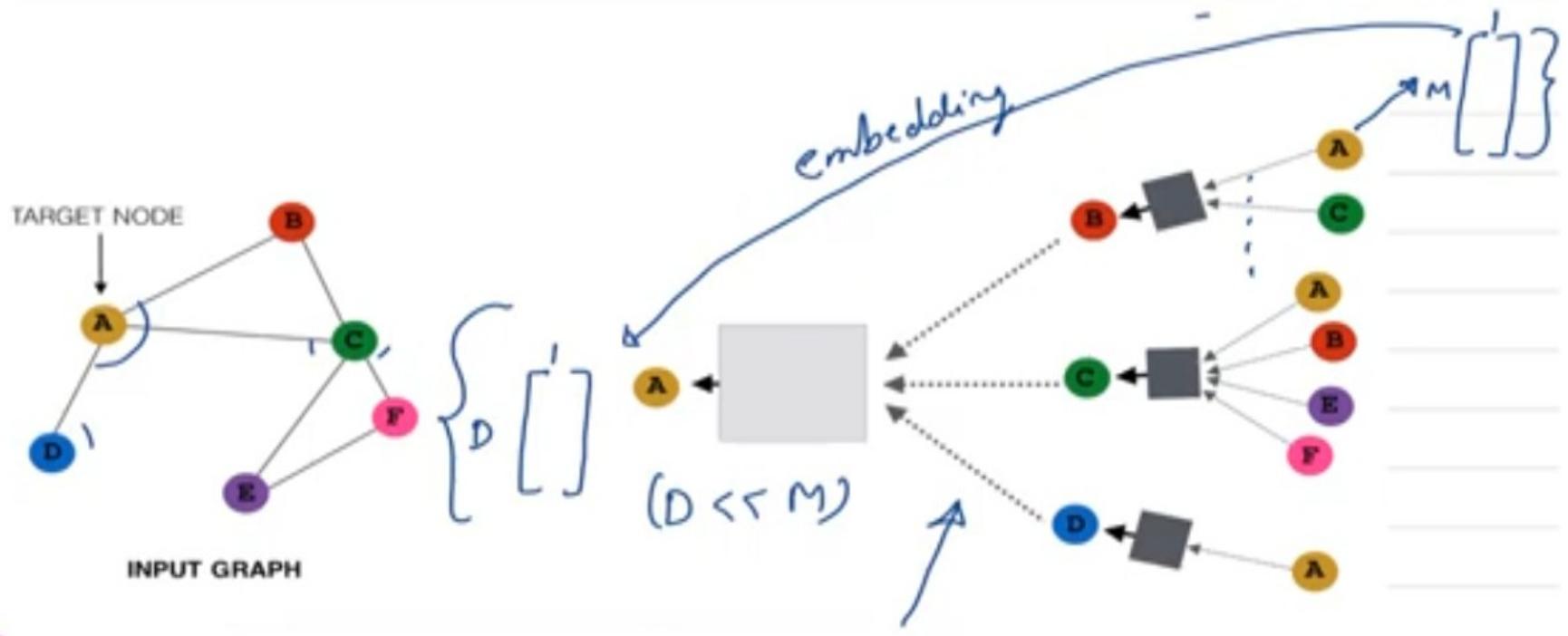




$$h_v^{(0)} = x_v$$

$$h_v^K = \sigma \left[w_K \sum_{u \in N(v)} \left[\frac{h_u^{(K-1)}}{N(v)} \right] + b_K h_v^{(K-1)} \right] \quad K = 1, \dots, K-1$$

$$z^{(A)} = h_v^{(K)}$$



$$h_v^0 = x_v$$

$$h_v^K = \sigma \left[w_K \sum_{u \in N(v)} \left[\frac{h_u^{(K-1)}}{N(v)} \right] + b_K h_v^{(K-1)} \right] \quad K=1, \dots, K-1$$

$$z^{(A)} = h_v^{(K)}$$

Graph $\underline{G} = \underline{(V, E)}$

* Input feature matrix $(N \times M) = X$

* Adjacency matrix $(N \times N) = A$

MATRIX NOTATION

Spectral analysis

Generalized Neural Network is represented as

$$H^{(l+1)} = f(H^{(l)}, A)$$

with $H^{(0)} = X$ & $H^{(L)} = Z$, L being number of layers.

Graph $G = (V, E)$

* Input feature matrix $(N \times M) = X$

* Adjacency matrix $(N \times N) = A$

MATRIX NOTATION

Spectral analysis

Generalized Neural network is represented as

$$H^{(l+1)} = f(H^{(l)}, A)$$

with $H^{(0)} = \{x\}$ & $H^{(L)} = \{z\}$, L being number of layers.

→ A simple layer wise propagation rule is

$$\left\{ H^{(l+1)} = \sigma \left(\underbrace{(A)H^{(l)}}_{\substack{\text{non-linear} \\ \text{activation} \\ \text{unit}}} \underbrace{W^{(l)}}_{\substack{N \times N \\ \text{for} \\ l=0}} \right) \right\}$$

where ($P \ll N$) for dim. reduction
needed for
Node embedding.

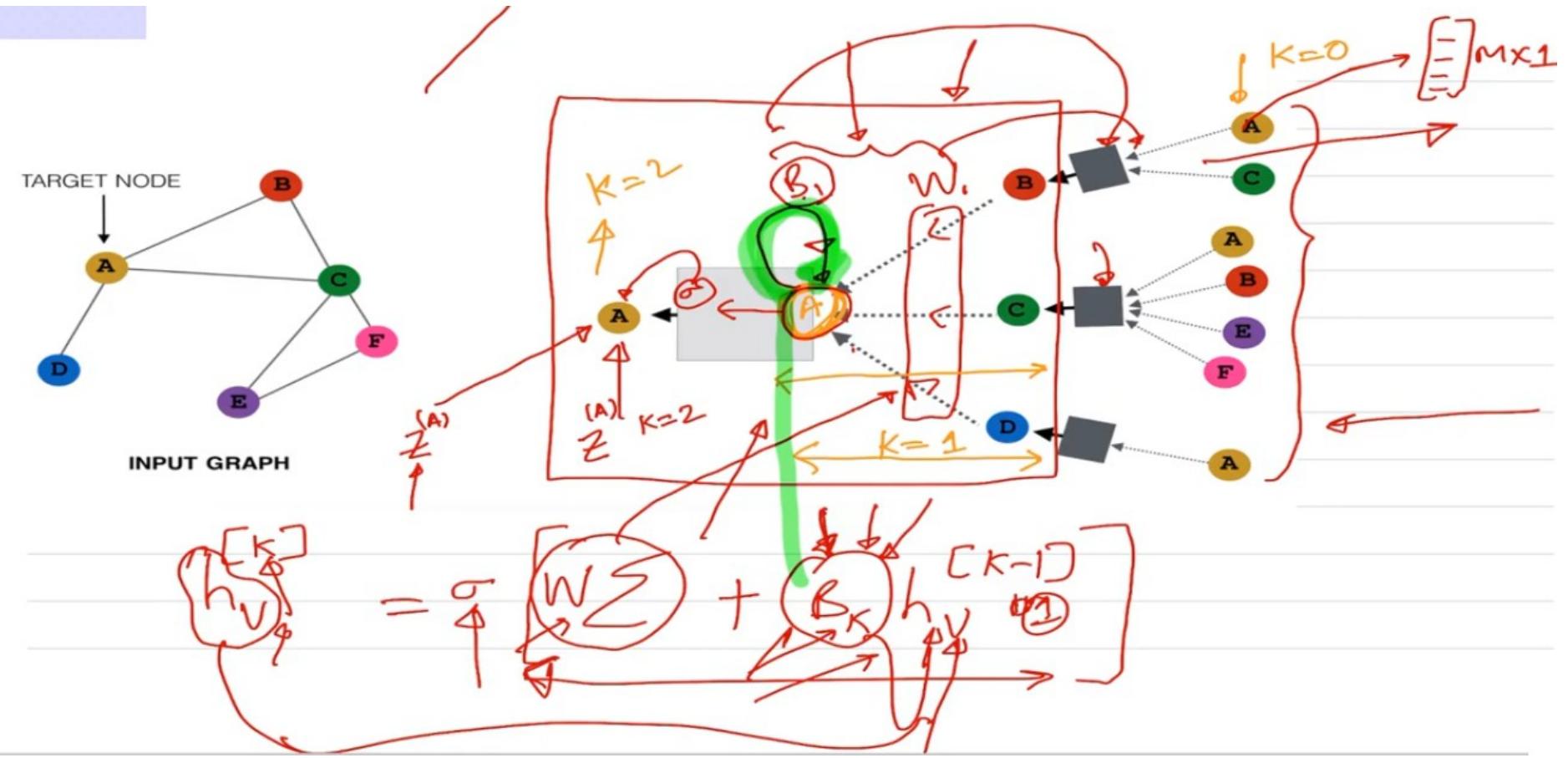
Limitations :

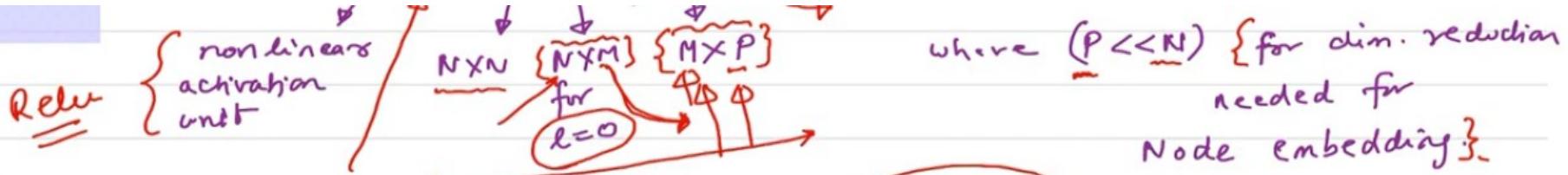
① No self loops in adjacency matrix. Fixed by adding I

② A is not normalized. $\xrightarrow{\text{results}}$ in change of feature vectors

Normalize A with Degree matrix $\Rightarrow D^{-1}A$

$\xrightarrow{\text{explodes gradient}}$
& optimizer sensitive to features





Limitations

$$(A + I)$$

$$(N \times M) \rightarrow (N \times P)$$

$$l=1 \quad l=0$$

① $\{ \text{No self loops} \}$ in adjacency Matrix. Fixed by adding (I)

② A is not normalized. results in change of feature vectors

Normalize A with Degree matrix $\Rightarrow D^{-1}A$

\Rightarrow explodes gradients
& optimizer sensitive
to feature scaling

$D^{-1}A \rightarrow$ corresponds to taking average of neighbouring node features.

① {No self loops} in adjacency matrix. Fixed by adding ($\frac{1}{n}I$)

② A is not normalized. results in change of feature vectors

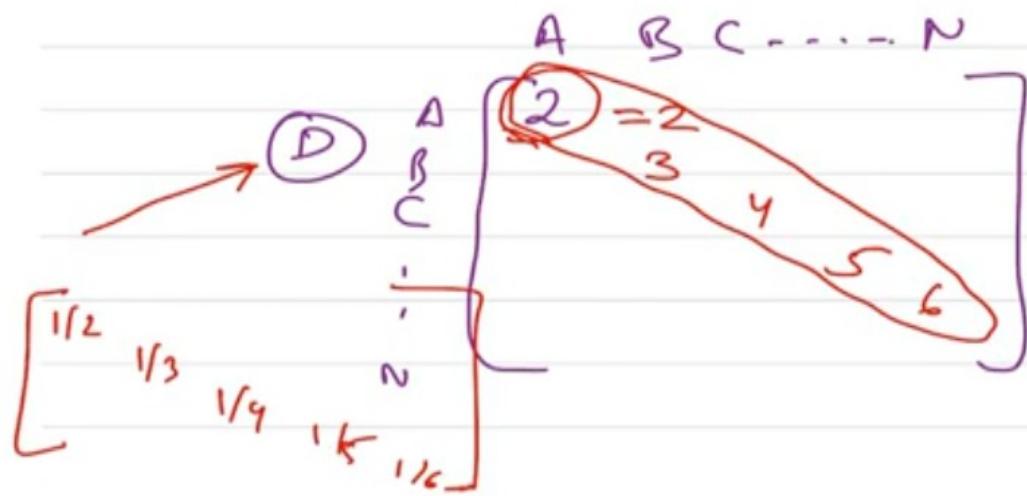
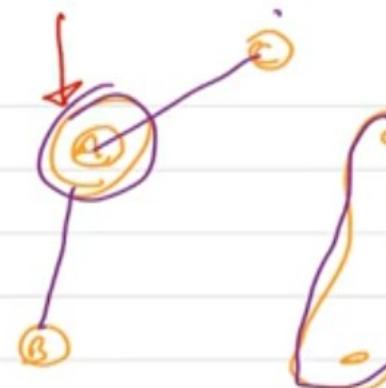
Normalize A with Degree matrix $\Rightarrow \{D^{-1}A\}$

explodes gradients
& optimizer sensitive
to feature scaling

$\{D^{-1}A\}$ → Corresponds to taking average of neighboring node features.

Paper use symmetric normalization $D^{-1/2} A D^{-1/2}$. With this new rule is

$$H^{(l+1)} = \sigma(D^{-1/2} \hat{A} D^{-1/2} H^{(l)} W^{(l)})$$



$\rightarrow \boxed{D^T A}$
 $\boxed{\begin{bmatrix} 0 & \left\{ \frac{1}{2} \right\} & \left\{ \frac{1}{2} \right\} & 0 \dots \dots \end{bmatrix}}$

$\begin{matrix} 0 \\ \frac{1}{2} \\ \frac{1}{2} \end{matrix}$
 (A)

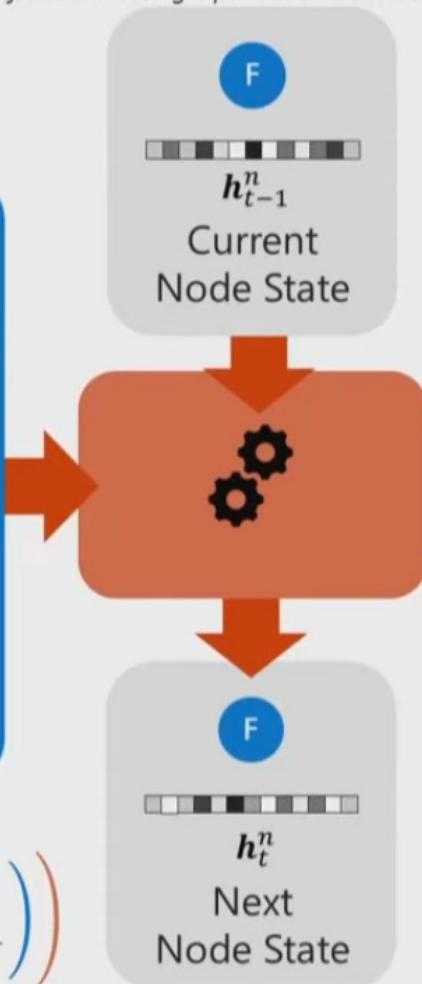
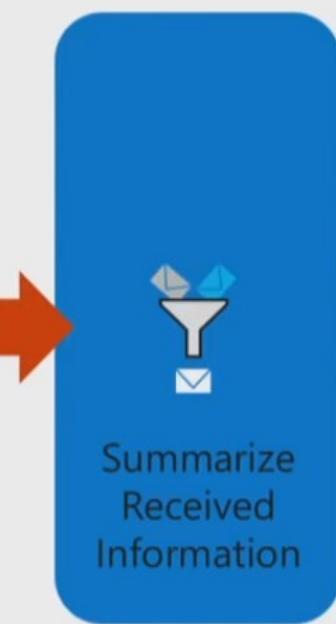
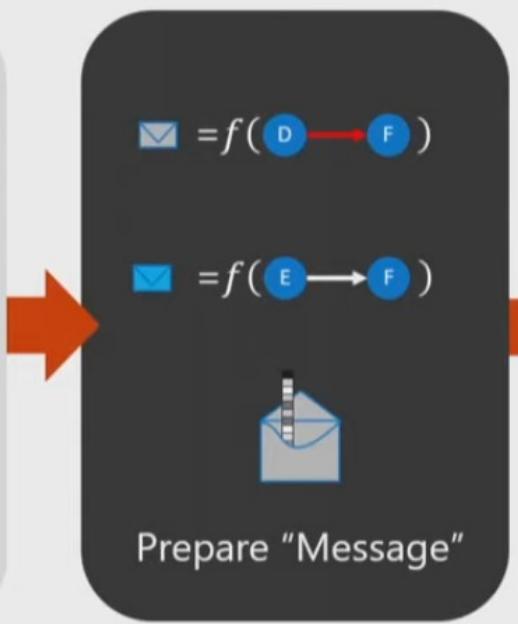
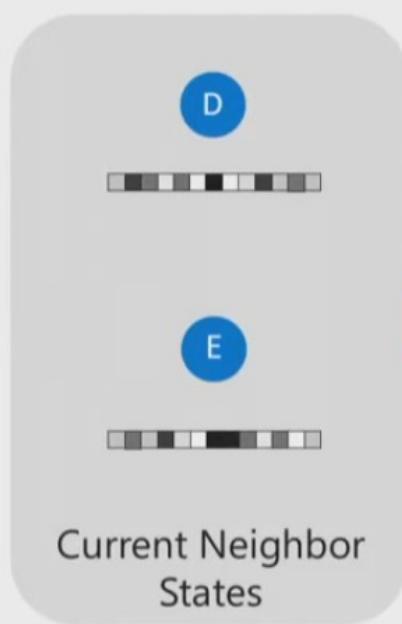
$\{D^{-1} A\}$ → Corresponds to taking average of neighbouring node features.

Paper use Symmetric normalization $D^{-1/2} A D^{-1/2}$. with this new rule is

$$H^{(l+1)} = \sigma(D^{-1/2} \hat{A} D^{-1/2} H^{(l)} W^{(l)})$$

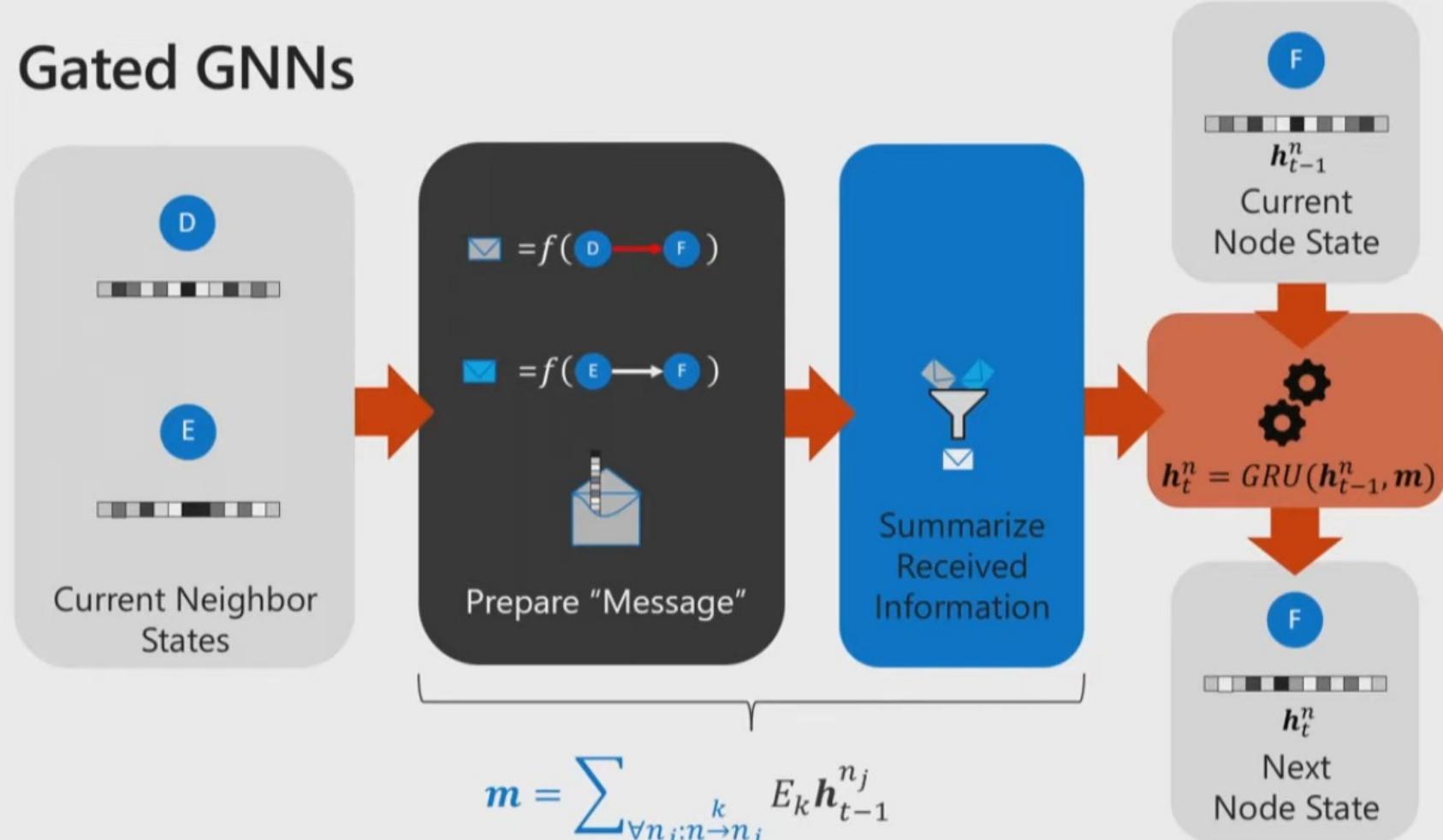
here $\hat{A} = A + I$, D is degree matrix of \hat{A}

GCNs



$$h_t^n = \sigma \left(\frac{1}{\text{numNeighbors} + 1} W_t \left(h_{t-1}^n + \sum_{\forall n_j: n \rightarrow n_j} h_{t-1}^{n'} \right) \right)$$

Gated GNNs



Li et al (2015). Gated graph sequence neural networks.

Expressing GGNNs as Matrix Operations



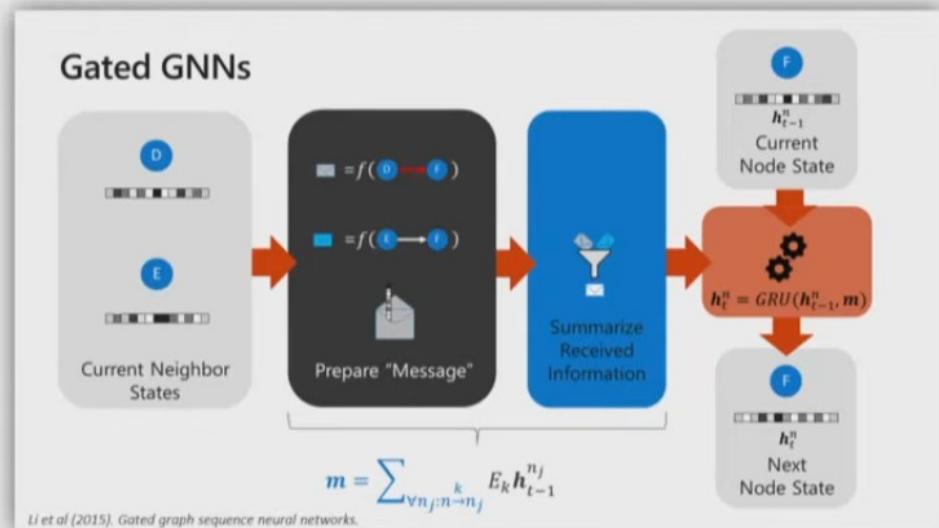
GGNN as Matrix Operation

Node States

$$H_t = \begin{bmatrix} h_t^{n_0} \\ \vdots \\ h_t^{n_K} \end{bmatrix} \quad (\text{num_nodes} \times D)$$

Messages to-be sent

$$M_t^k = E_k H_t \quad (\text{num_nodes} \times M)$$

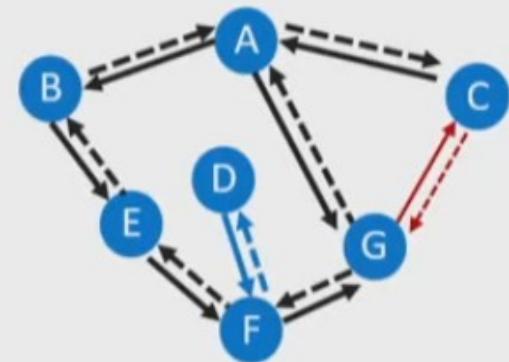
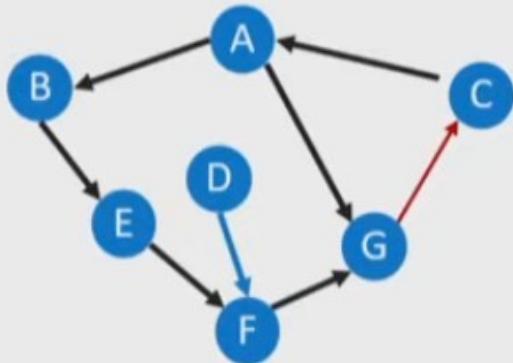


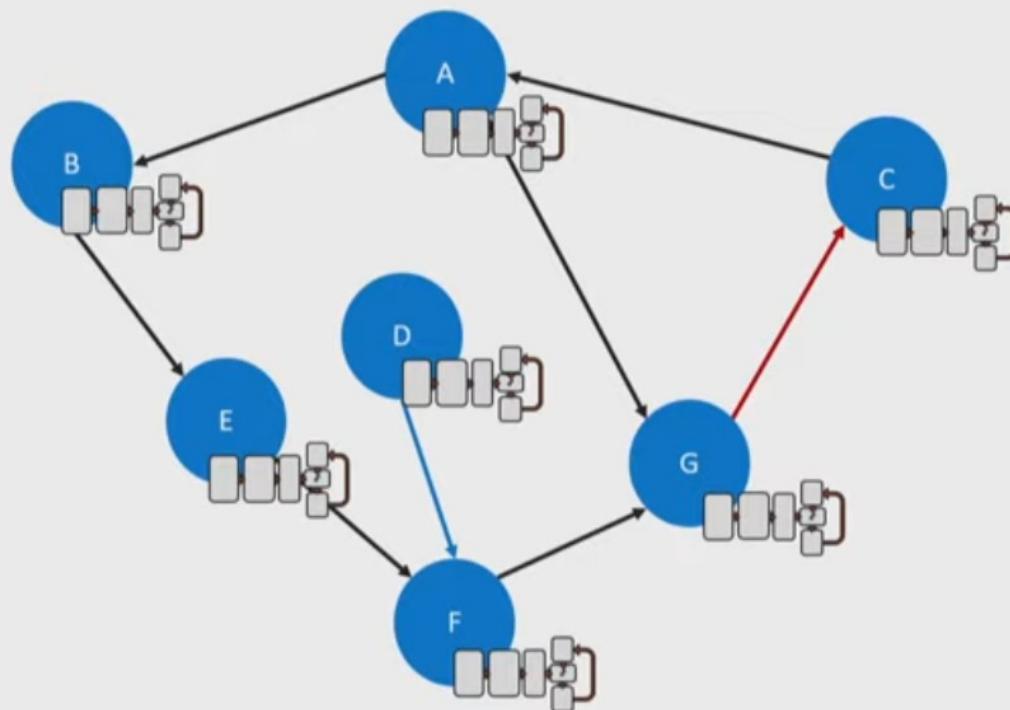
Received Messages

$$R_t = \sum_k A M_t^k \quad (\text{num_nodes} \times M)$$

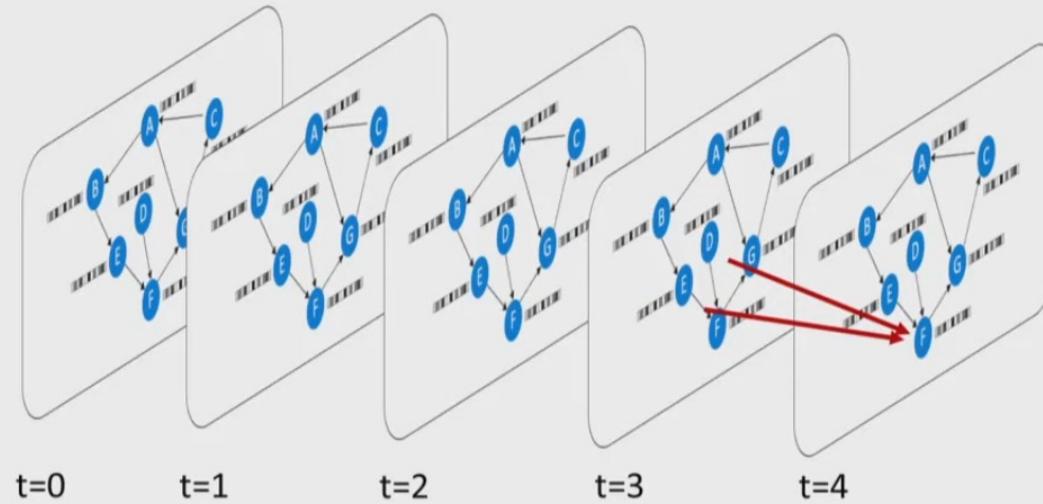
$$\text{Update } H_{t+1} = GRU(H_t, R_t)$$

Trick 1: Backwards Edges

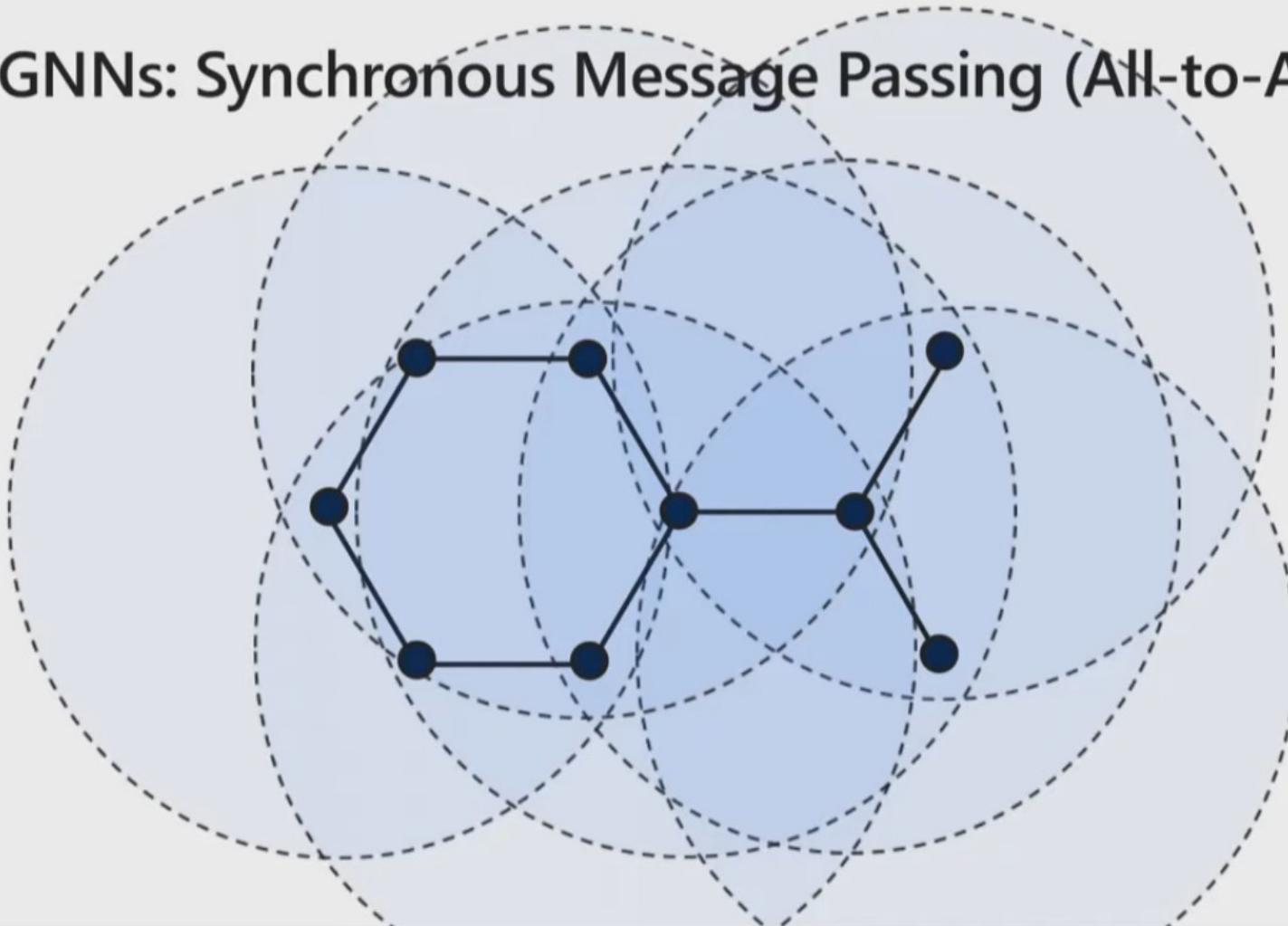




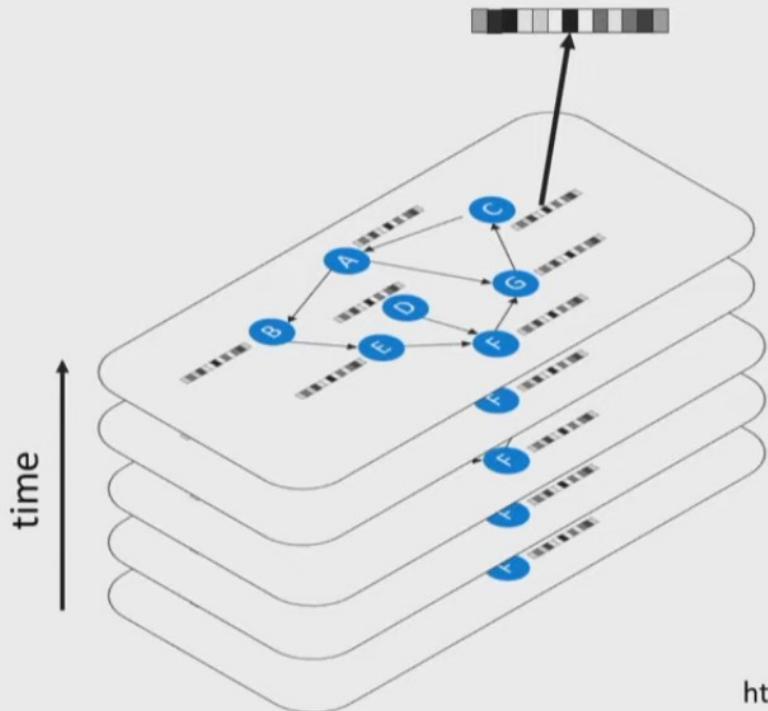
Graph Neural Networks: Message Passing



GNNs: Synchronous Message Passing (All-to-All)



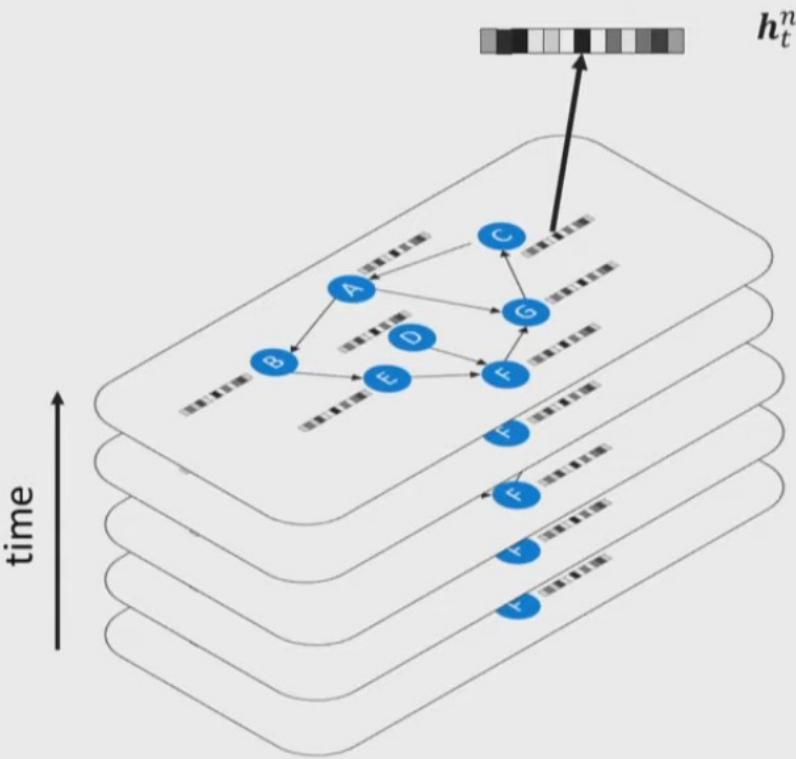
Graph Neural Networks: Output



- node selection
- node classification
- graph classification

<https://github.com/microsoft/tf-gnn-samples/>

Example: Node [Binary] Classification



$$\mathbf{h}_t^n$$

$$x_n = \sigma(\mathbf{w}^T \mathbf{h}_t^n + b)$$

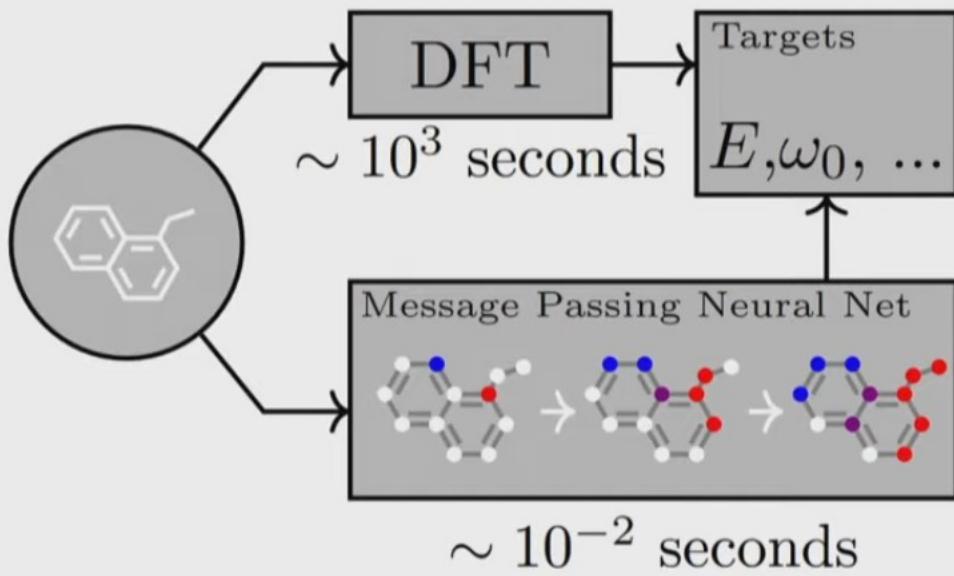
Binary cross entropy

$$\mathcal{L} = y_n \cdot \log x_n + (1 - y_n) \log(1 - x_n)$$

Two Sample Applications

...with graph neural networks.





Gilmer, Justin, et al. "Neural message passing for quantum chemistry." 2017.

Variable Misuse Task

```
var clazz=classTypes["Root"].Single() as JsonCodeGenerator.ClassType;  
Assert.NotNull(clazz);  
  
var first=classTypes["RecClass"].Single() as JsonCodeGenerator.ClassType;  
Assert.NotNull(first);  
  
Assert.Equal("string", first.Properties["Name"].Name);  
Assert.False(clazz.Properties["Name"].IsArray);
```

Possible type-correct options: `clazz`, `first`



Variable Misuse Task

```
var clazz=classTypes["Root"].Single() as JsonCodeGenerator.ClassType;  
Assert.NotNull(clazz);  
  
var first=classTypes["RecClass"].Single() as JsonCodeGenerator.ClassType;  
Assert.NotNull( clazz );  
  
Assert.Equal("string", first.Properties["Name"].Name);  
Assert.False(clazz.Properties["Name"].IsArray);
```

Possible type-correct options: `clazz`, `first`



Programs as Graphs

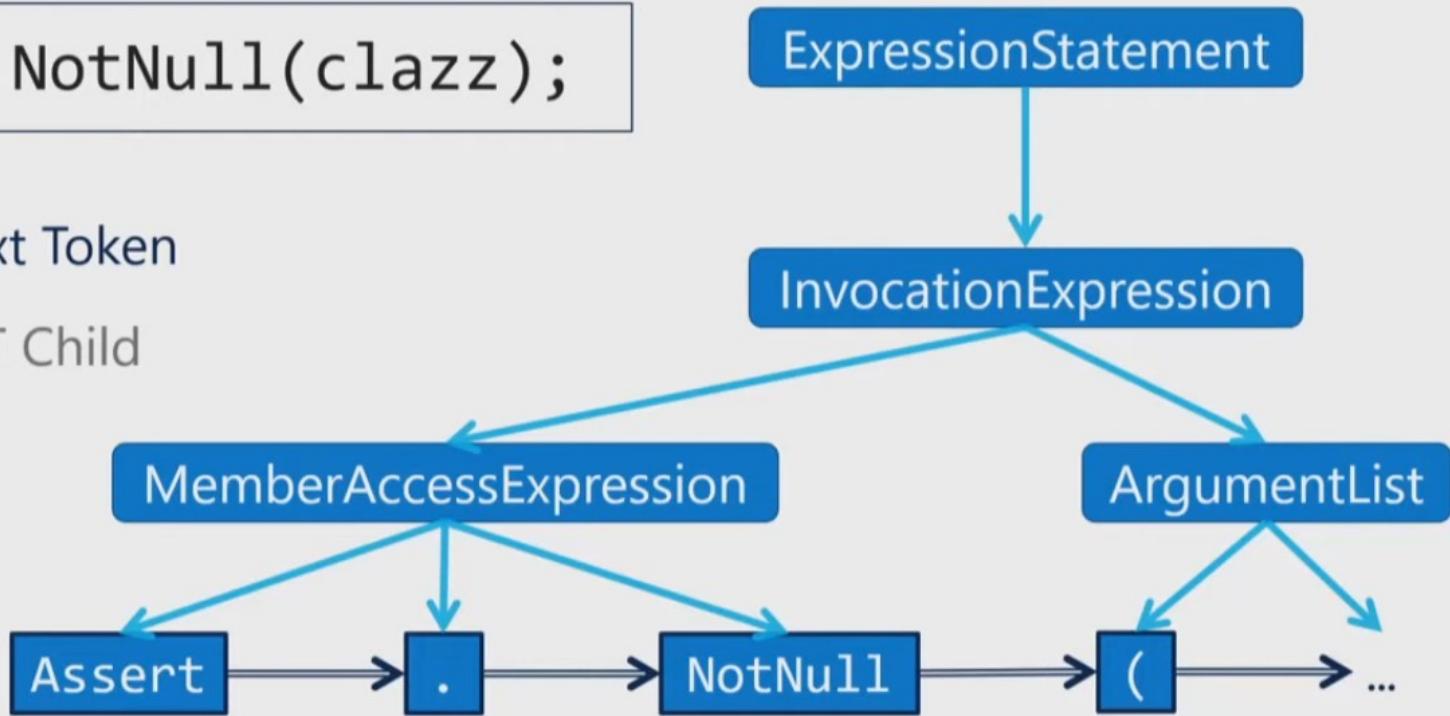
```
int SumPositive(int[] arr, int lim) {  
    int sum = 0;  
    for (int i=0; i < lim; i++)  
        if (arr[i] > 0)  
            sum += arr[i];  
    return sum;  
}
```

A control flow graph for the `SumPositive` function. The graph starts at the entry node (top left), which branches to the initialization of `sum` (black dot) and the start of the `for` loop (rectangle). The `for` loop has three nodes: the header (`i=0`), the condition (`i < lim`), and the increment (`i++`). The body of the loop contains an `if` statement (rectangle) with two outgoing edges: one to the `return` statement (black dot) if the condition is true, and one to the `for` loop header if it is false. The `if` statement also has two outgoing edges: one to the `return` statement if the array element is positive, and one to the `+=` assignment if it is not. The `+=` assignment has an outgoing edge to the `return` statement.

Programs as Graphs: Syntax

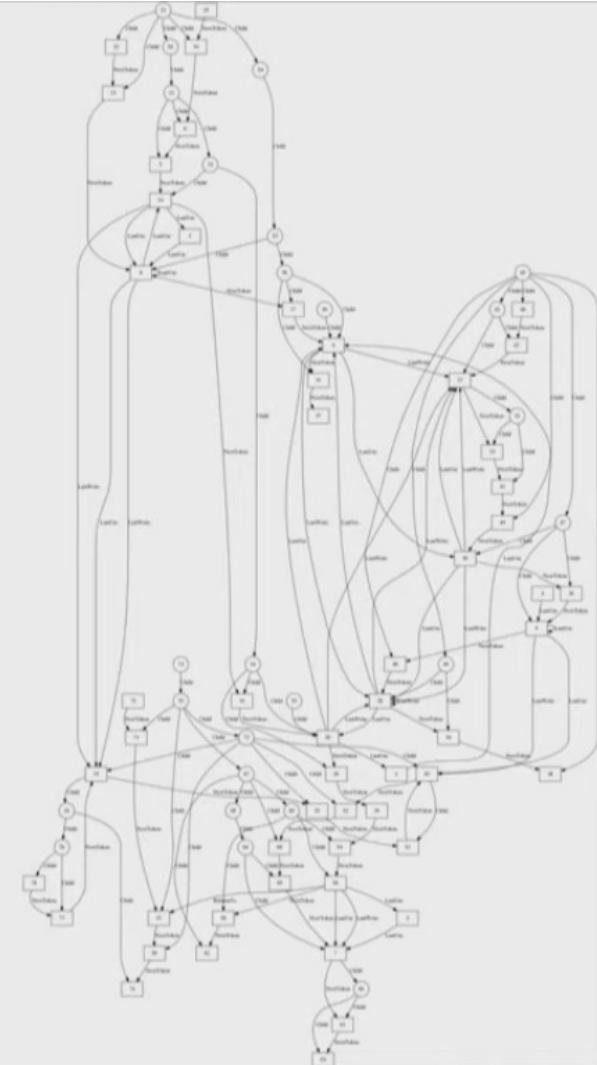
```
Assert.NotNull(clazz);
```

- Next Token
- AST Child



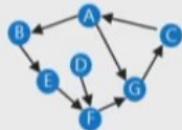
Programs as Graphs

```
int SumPositive(int[] arr, int lim) {  
    int sum = 0;  
    for (int i=0; i < lim; ++)  
        if (arr[i] > 0)  
            sum += arr[i];  
  
    return sum;  
}
```



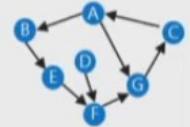
~900 nodes/graph ~8k edges/graph

Graph Representation for Variable Misuse



```
var clazz=classTypes["Root"].Single() as JsonCodeGenerator.ClassType;  
Assert.NotNull(clazz);  
  
var first=classTypes["RecClass"].Single() as JsonCodeGenerator.ClassType;  
Assert.NotNull( SLOT );      first      clazz  
  
Assert.Equal("string", first.Properties["Name"].Name);  
Assert.False(clazz.Properties["Name"].IsArray);
```

Graph Representation for Variable Misuse



```
var clazz=classTypes["Root"].Single() as JsonCodeGenerator.ClassType;  
Assert.NotNull(clazz);  
  
var first=classTypes["RecClass"].Single() as JsonCodeGenerator.ClassType;  
Assert.NotNull( SLOT );      first      clazz  
  
Assert.Equal("string", first.Properties["Name"].Name);  
Assert.False(clazz.Properties["Name"].IsArray);
```

Goal: make the representation of SLOT as close as possible to the representation of the correct candidate node

$$f(\mathbf{h}_T^{SLOT}, \mathbf{h}_T^{first}) \gg f(\mathbf{h}_T^{SLOT}, \mathbf{h}_T^{clazz})$$

```
bool TryFindGlobalDirectivesFile(string baseDirectory, string fullPath, out string path) {
    baseDirectory = baseDirectory.TrimEnd(Path.DirectorySeparatorChar);
    var directivesDirectory = Path.GetDirectoryName(fullPath)
        .TrimEnd(Path.DirectorySeparatorChar);
    while (directivesDirectory != null && directivesDirectory.Length >= baseDirectory.Length) {
        path = Path.Combine(directivesDirectory, GlobalDirectivesFileName);
        if (File.Exists(path)) return true;

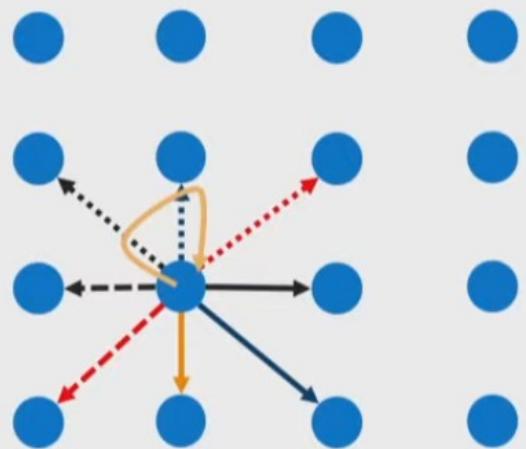
        directivesDirectory = Path.GetDirectoryName(directivesDirectory)
            .TrimEnd(Path.DirectorySeparatorChar);
    }
    path = null;
    return false;
}
```

What the model sees...

Other Models as Special Cases of GNNs



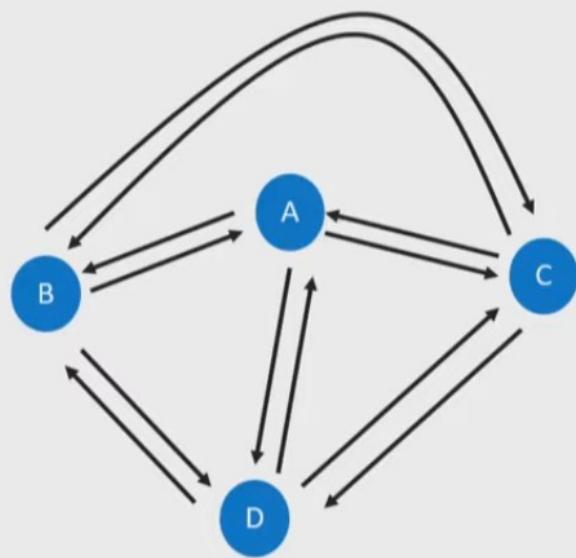
Special Case 1: Convolutions (CNN)



Special Case 2: “Deep Sets”



Special Case 2: “Deep Sets”



Stanford material

- <http://snap.stanford.edu/class/cs224w-2020/>