

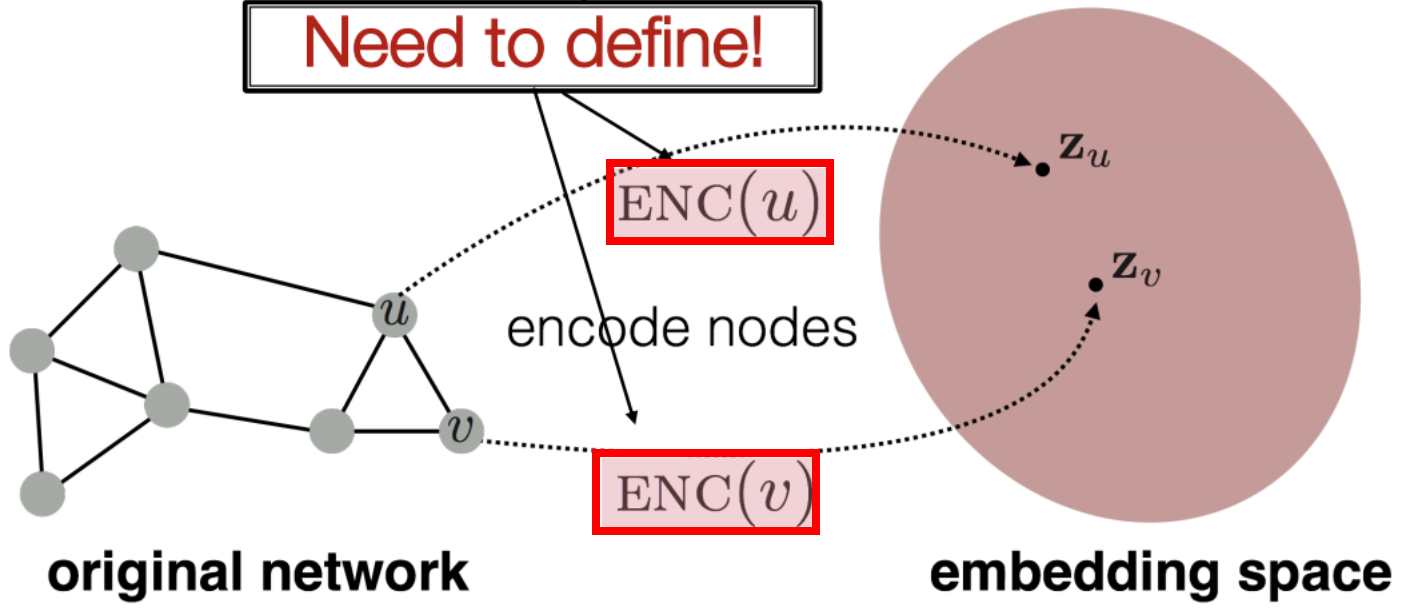
# Graph Neural Networks

Neural Networks Design And Application

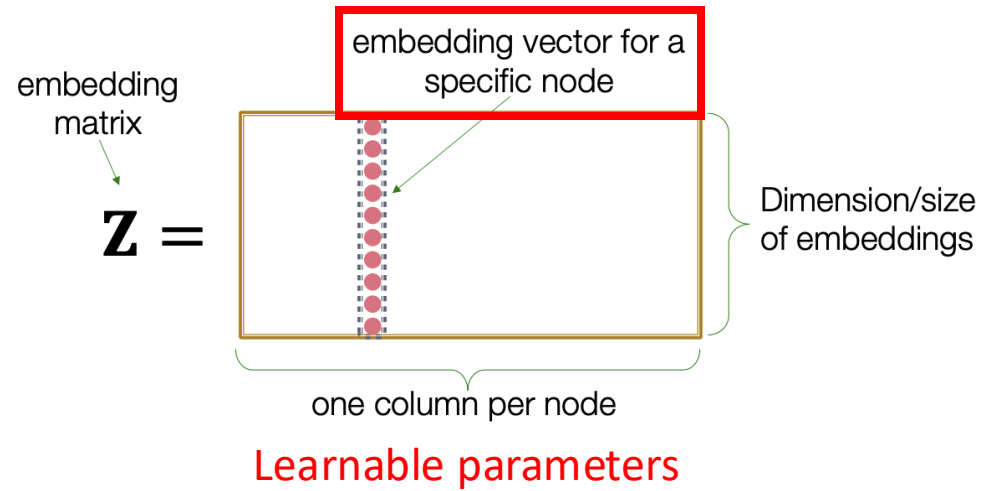
# Encoder-decoder for graph data

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

**Need to define!**

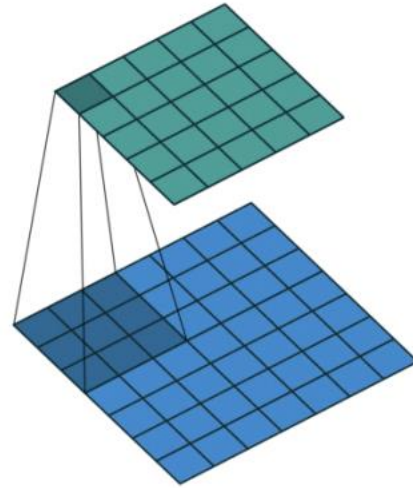


Q: how to learn ENC?



Q: how can we replace linear node embedding with nonlinear functions?

# Graph neural networks

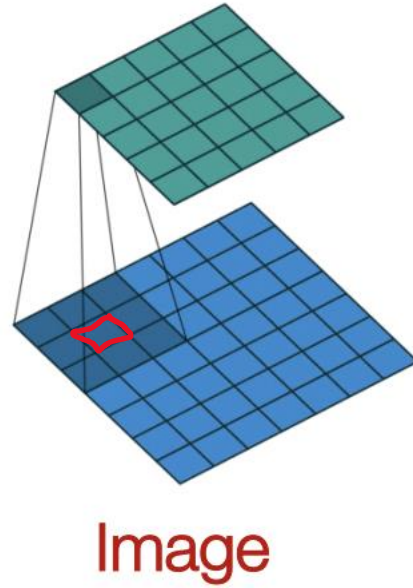


Image

Q: can we use convolution operation on graph?

Graph convolutional neural networks

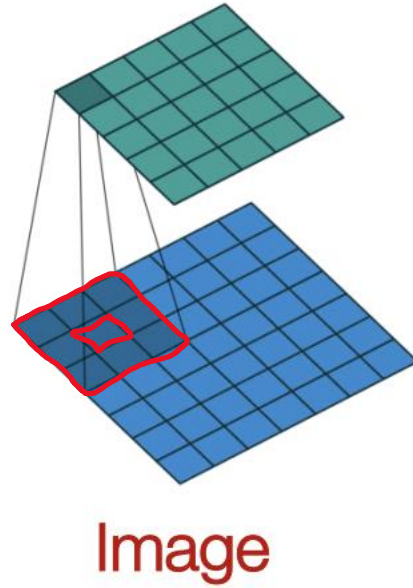
# Graph neural networks



Q: can we use convolution operation on graph?

Graph convolutional neural networks

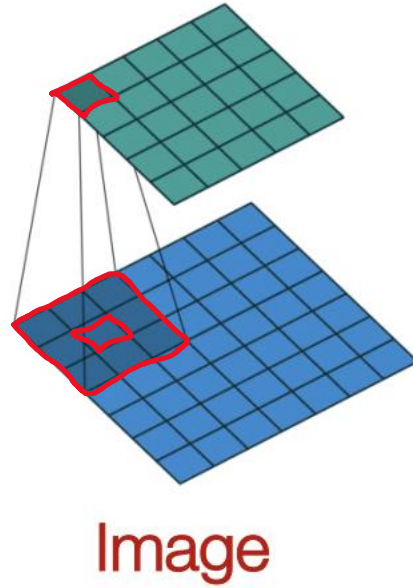
# Graph neural networks



Q: can we use convolution operation on graph?

Graph convolutional neural networks

# Graph neural networks

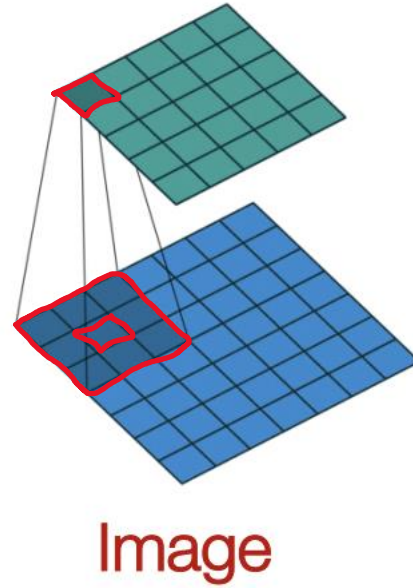


Q: can we use convolution operation on graph?

Graph convolutional neural networks

# Graph neural networks

Transform information at neighbors

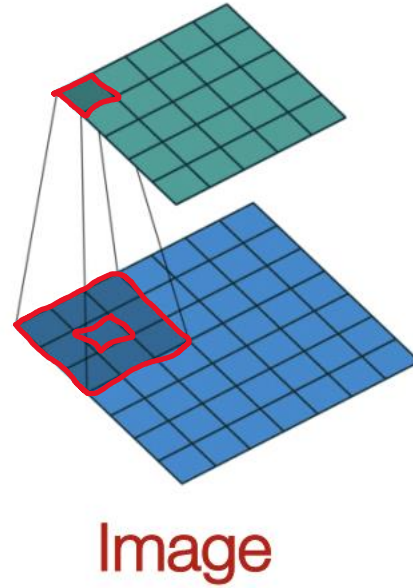


Q: can we use convolution operation on graph?

Graph convolutional neural networks

# Graph neural networks

Transform information at neighbors  
Combine them



Q: can we use convolution operation on graph?

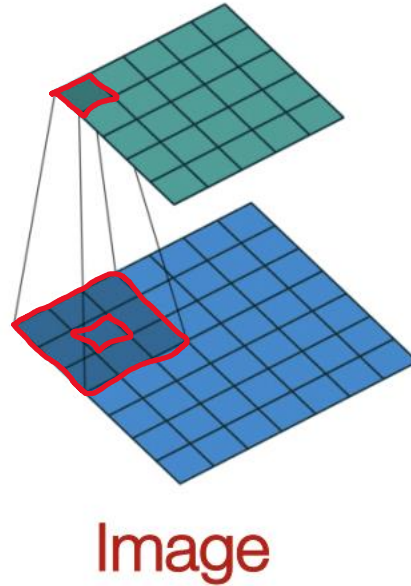
Graph convolutional neural networks



# Graph neural networks

Q: What is **information**?

Transform **information** at neighbors  
Combine them



Q: can we use convolution operation on graph?

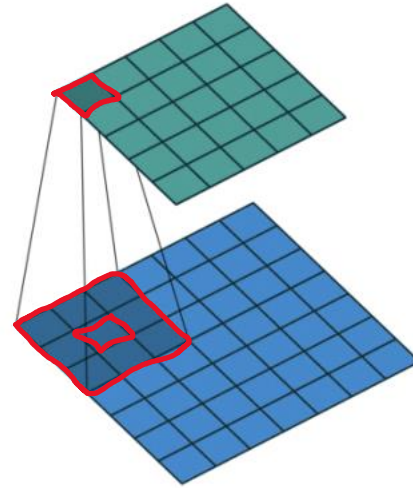
Graph convolutional neural networks

# Graph neural networks

Q: What is **information**?

Raw features (pixel values) or previous feature maps

Transform **information** at neighbors  
Combine them



Image

Q: can we use convolution operation on graph?

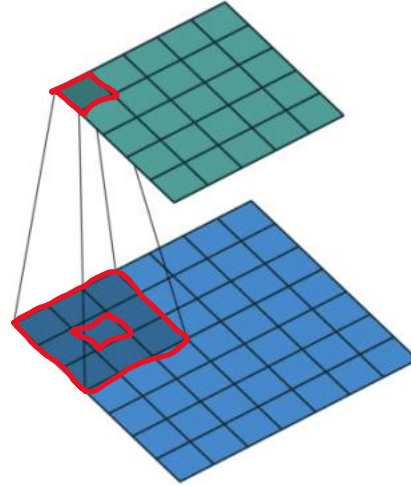
Graph convolutional neural networks

# Graph neural networks

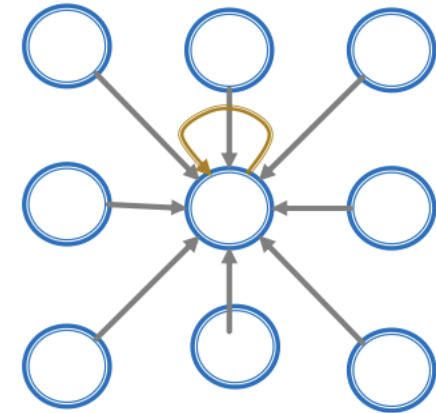
Q: What is **information**?

Raw features (pixel values) or previous feature maps

Transform information at neighbors  
Combine them



Image

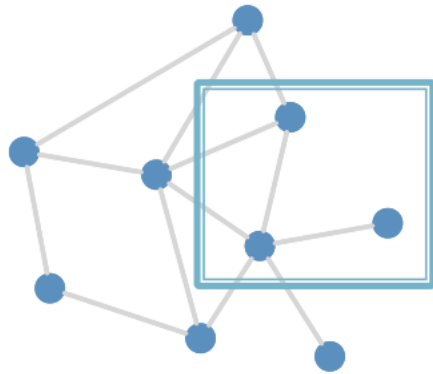


Graph

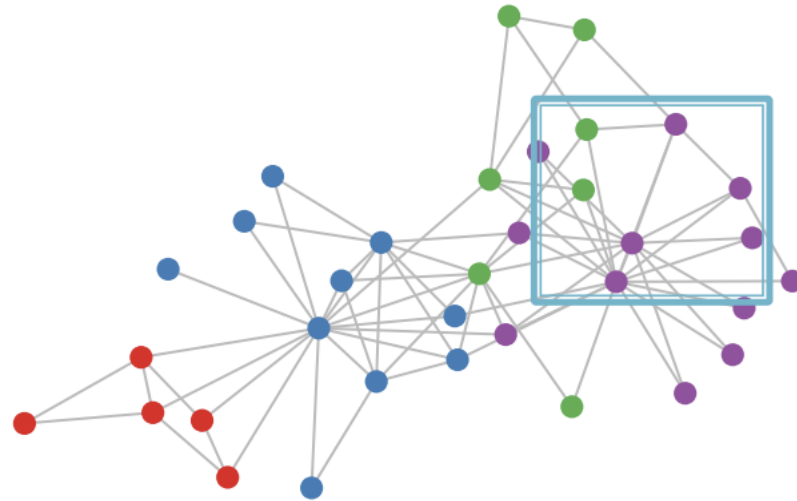
Q: can we use convolution operation on graph?

Graph convolutional neural networks

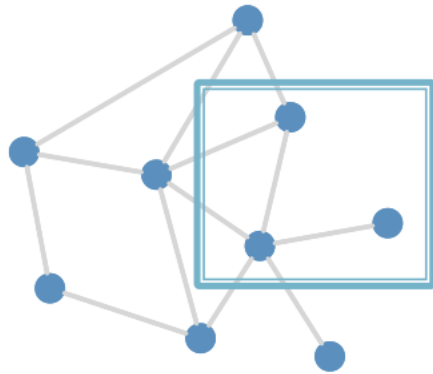
# General graph data



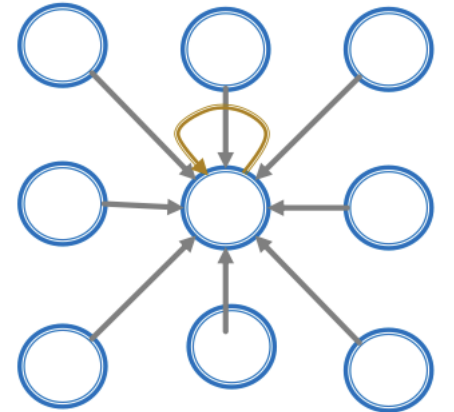
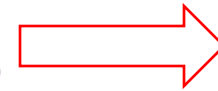
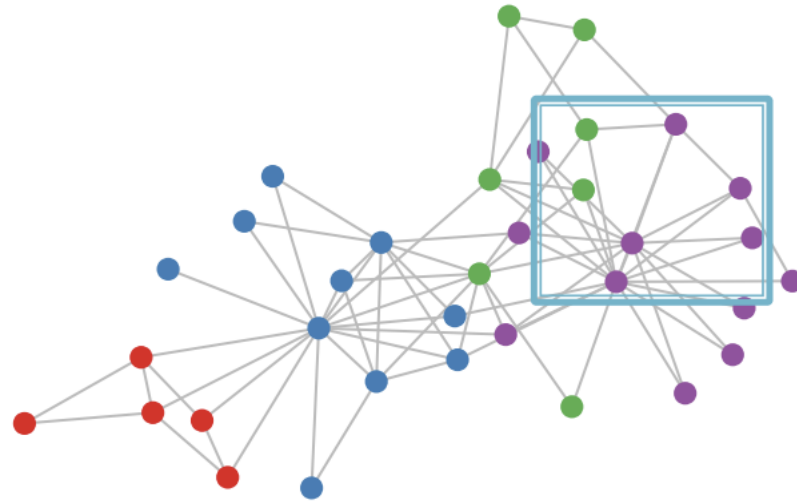
or this:



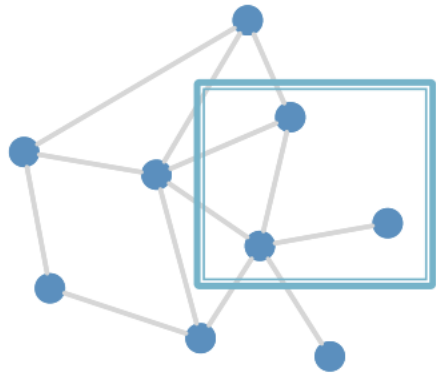
# General graph data



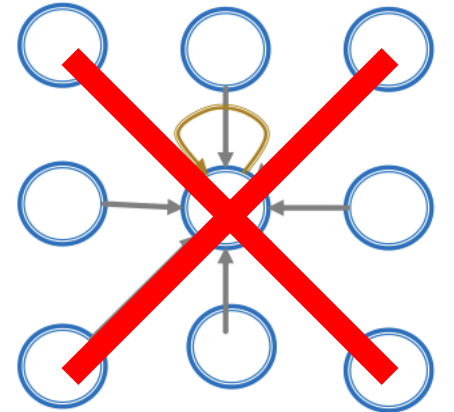
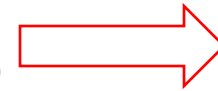
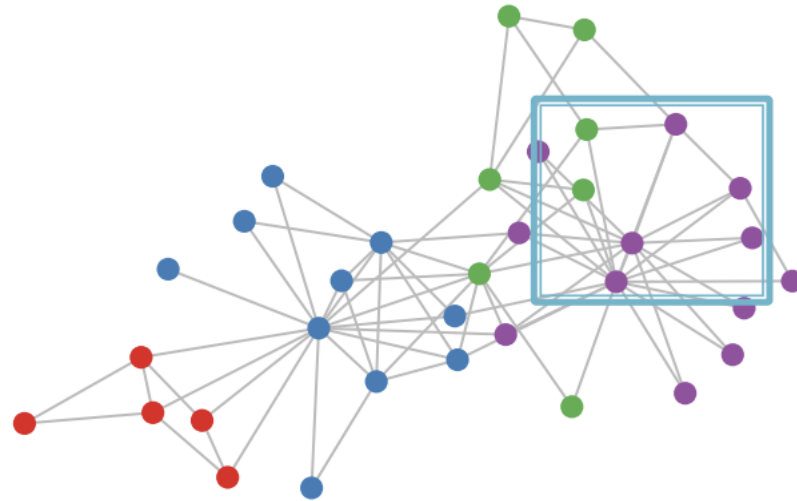
or this:



# General graph data

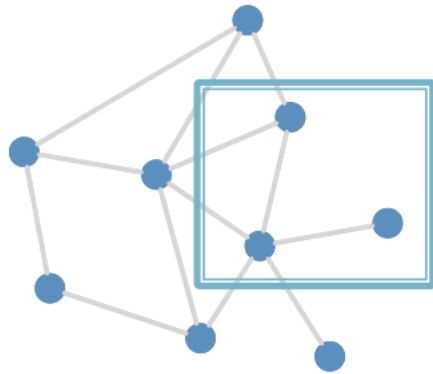


or this:

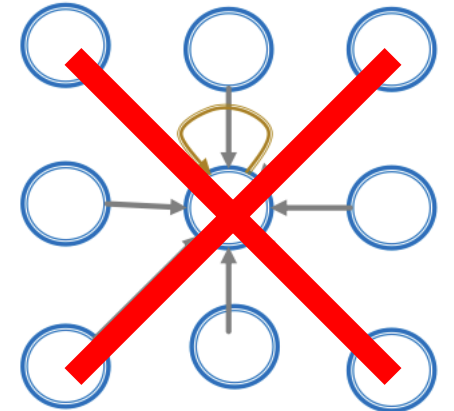
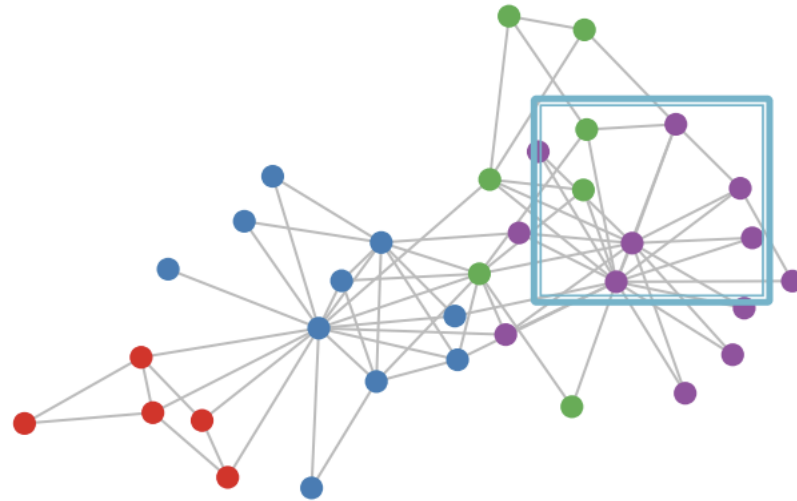


Q: can we extend similar operation to general graph?

# General graph data



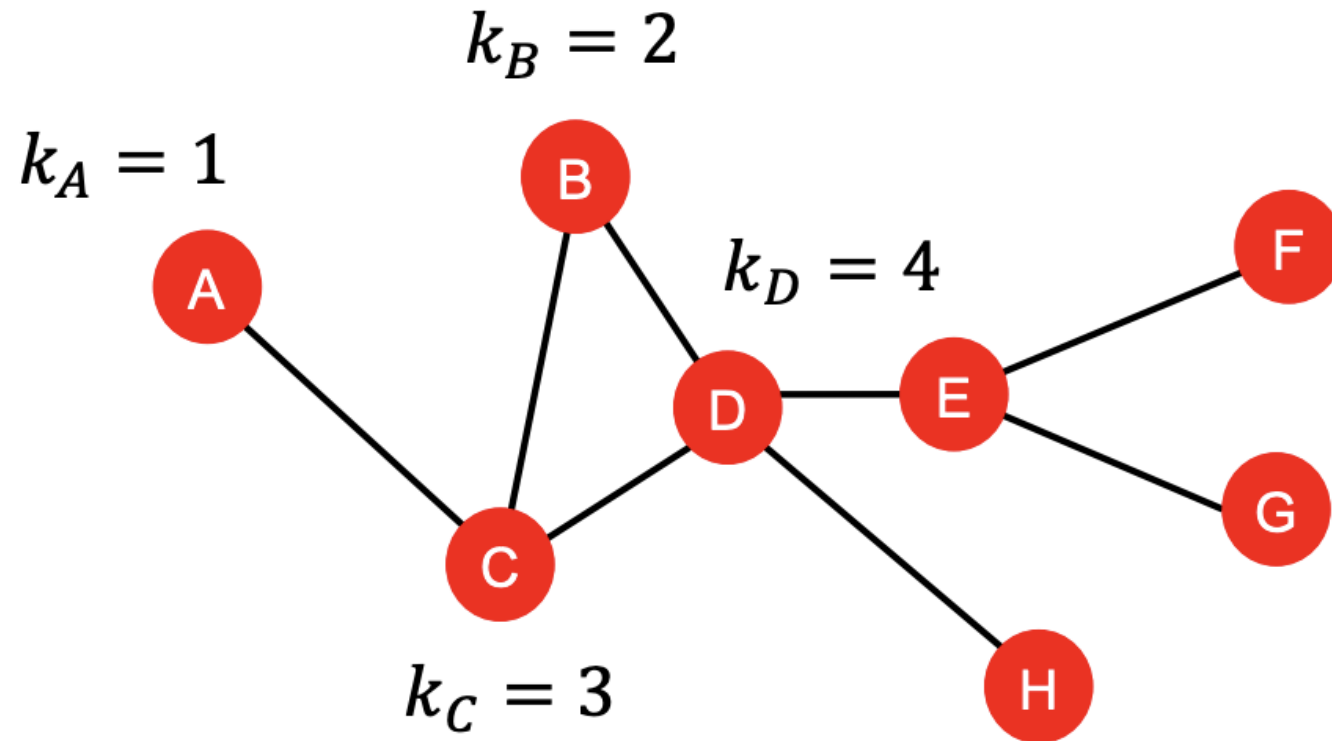
or this:



**Q:** can we extend similar operation to general graph?

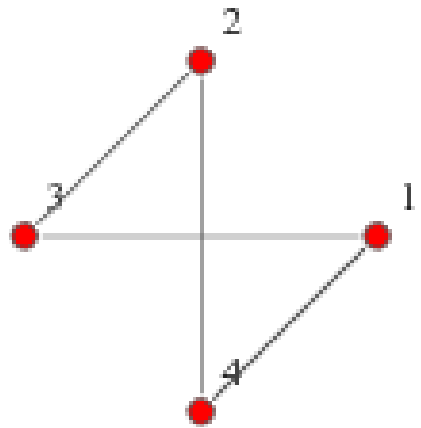
**Key:** aggregate information from **neighbors**

# Node degree

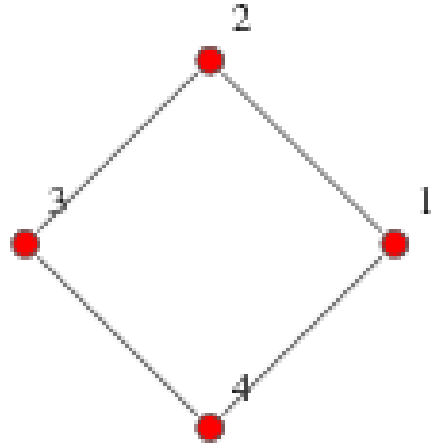




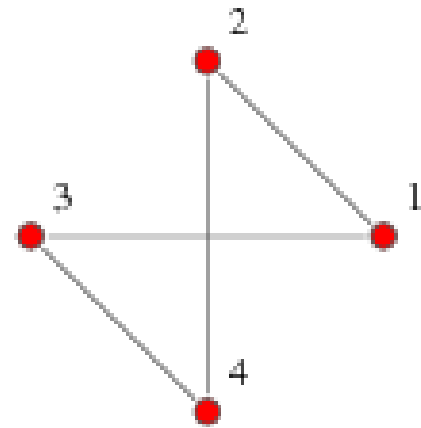
# Adjacency matrix



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

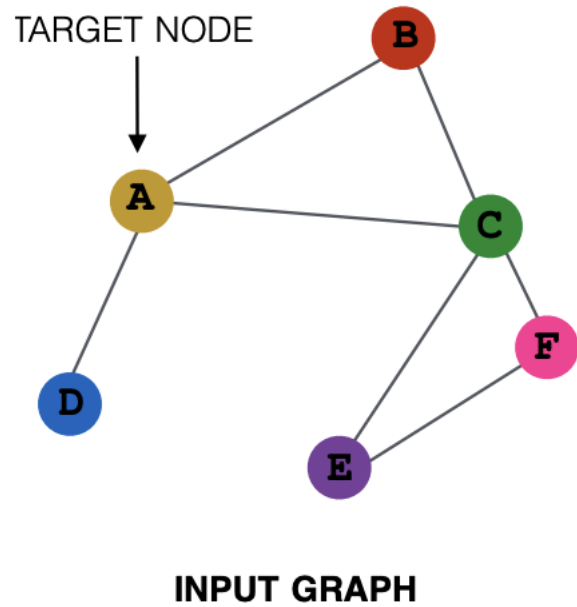


$$\begin{pmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{pmatrix}$$



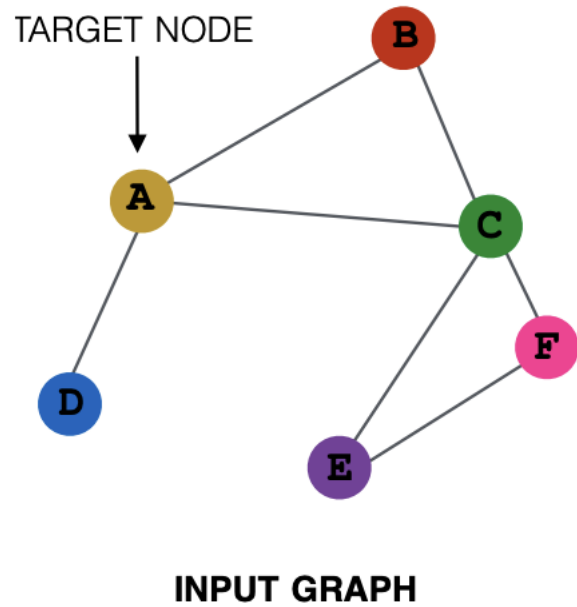
$$\begin{pmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{pmatrix}$$

# Graph networks: aggregate neighbors



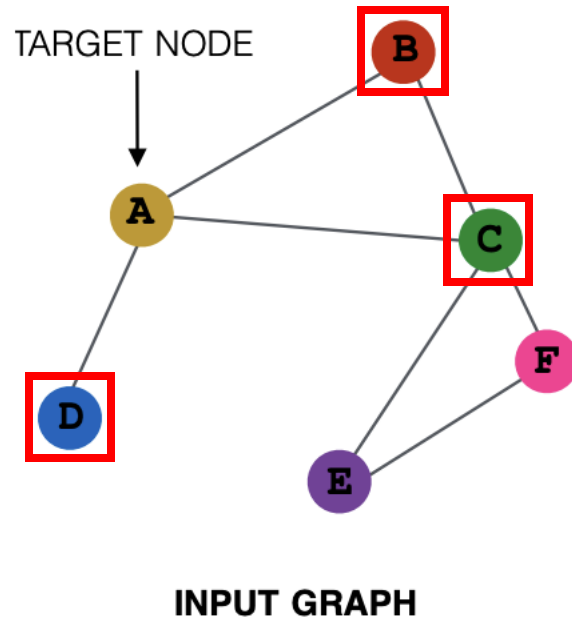
# Graph networks: aggregate neighbors

Nearest neighbors of A?



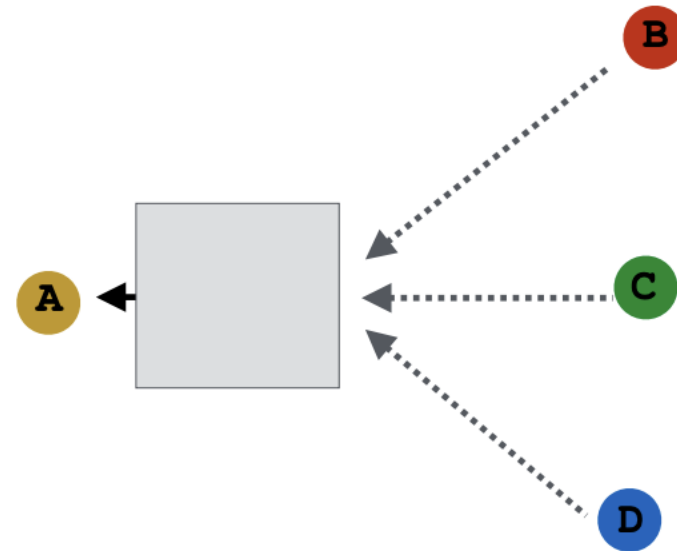
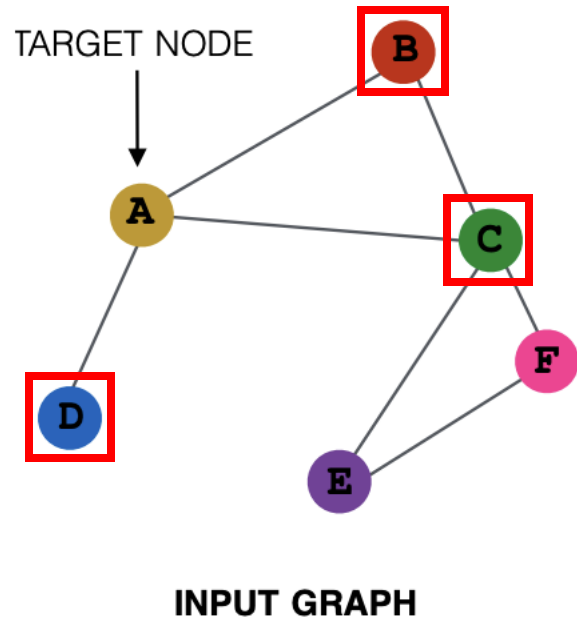
# Graph networks: aggregate neighbors

Nearest neighbors of A:  
B, C, D



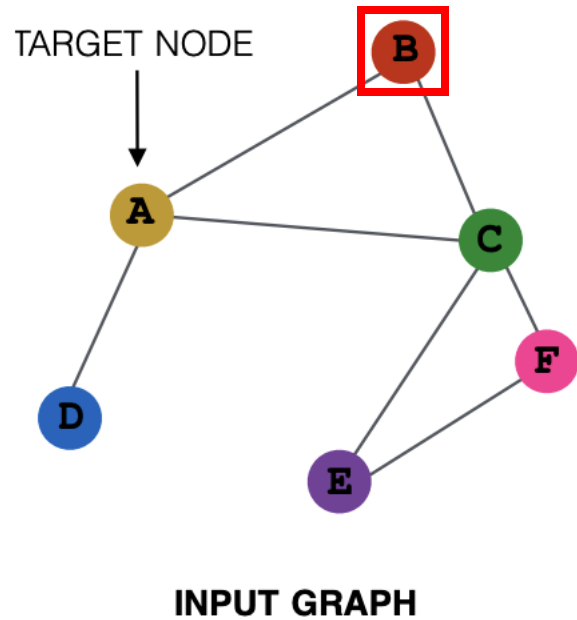
# Graph networks: aggregate neighbors

Nearest neighbors of A:  
B, C, D



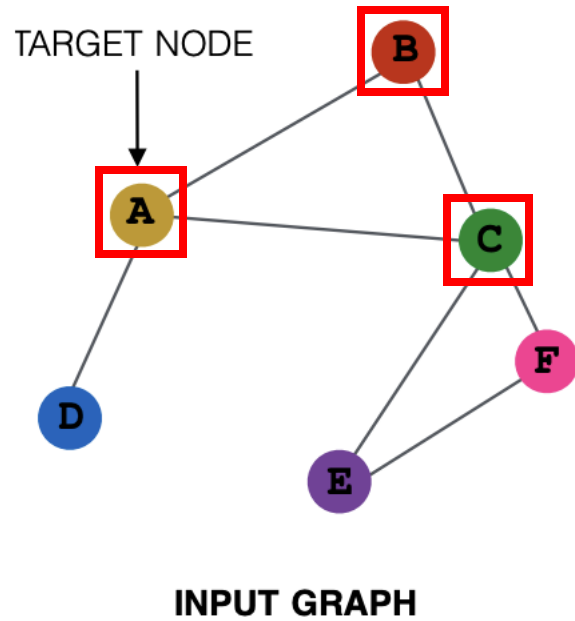
# Graph networks: aggregate neighbors

Nearest neighbors of **B**?



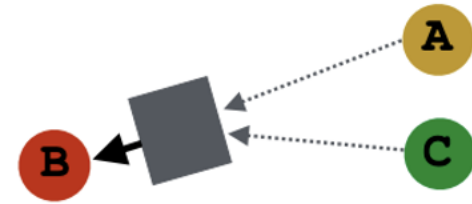
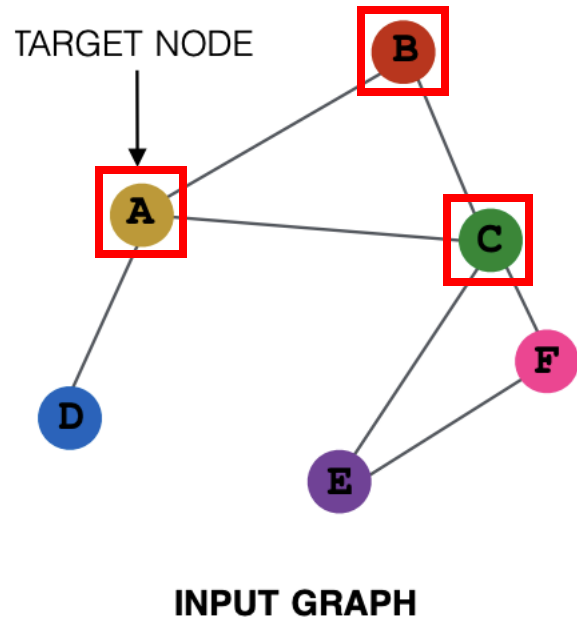
# Graph networks: aggregate neighbors

Nearest neighbors of **B**: **A**, **C**



# Graph networks: aggregate neighbors

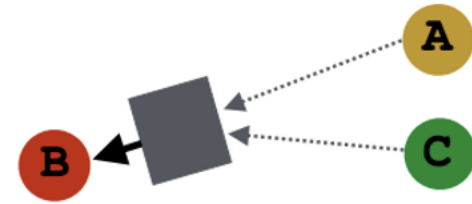
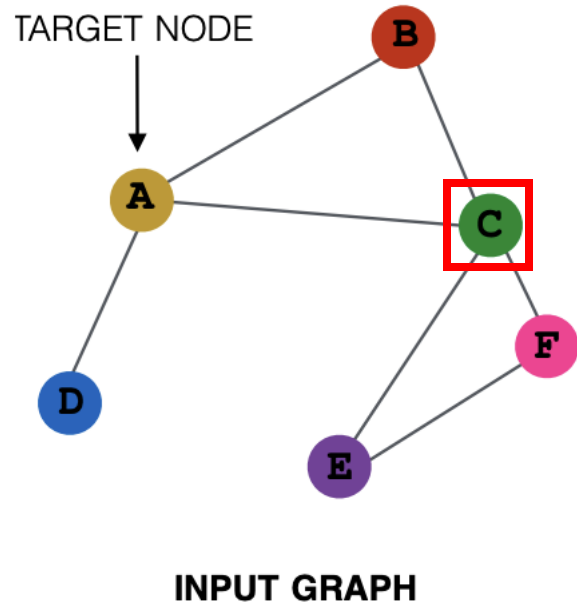
Nearest neighbors of **B**: **A**, **C**





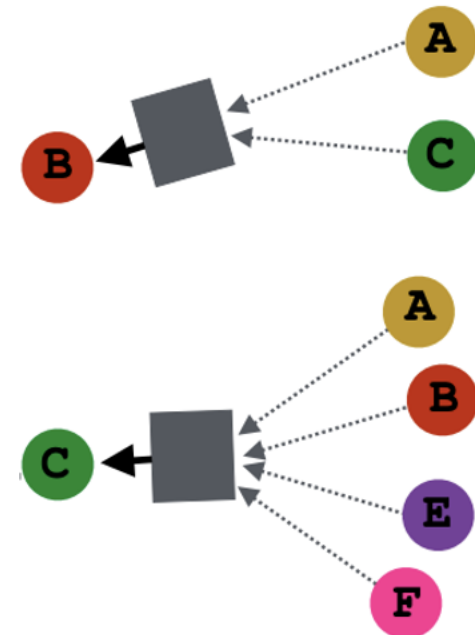
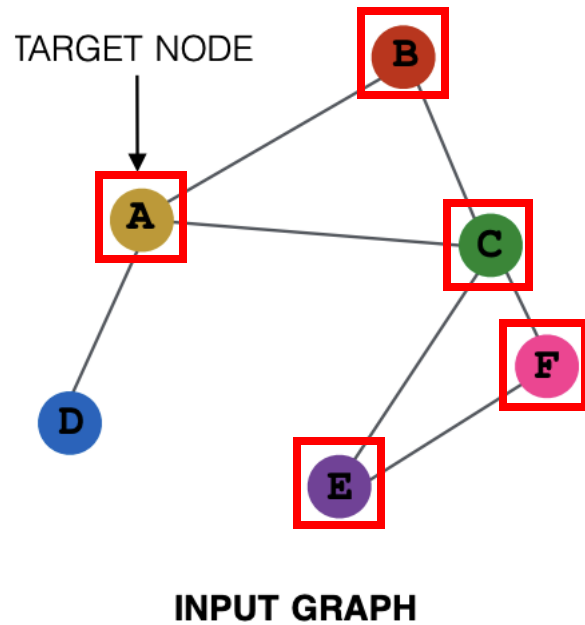
# Graph networks: aggregate neighbors

Nearest neighbors of **C**?



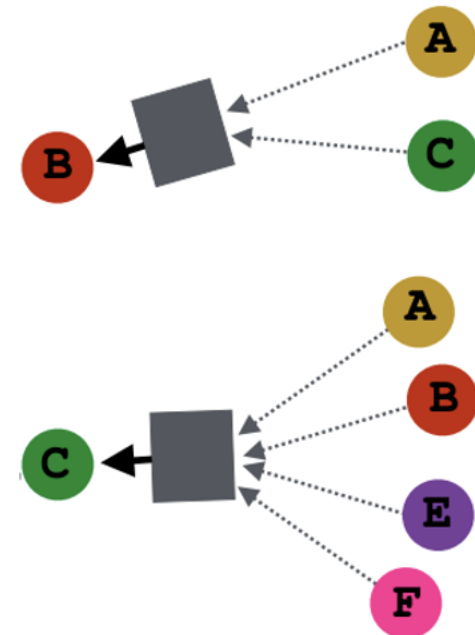
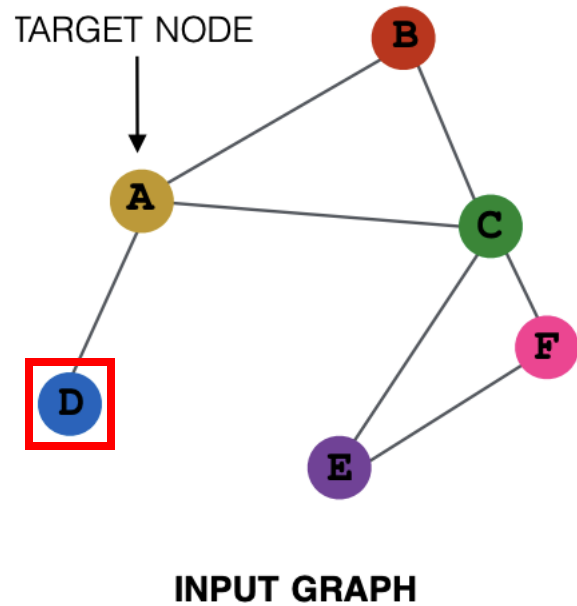
# Graph networks: aggregate neighbors

Nearest neighbors of **C**: **A**, **B**, **E**, **F**



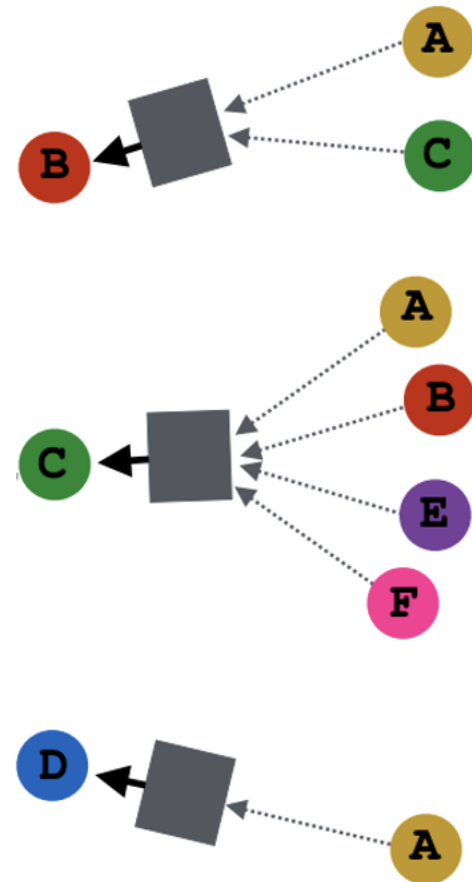
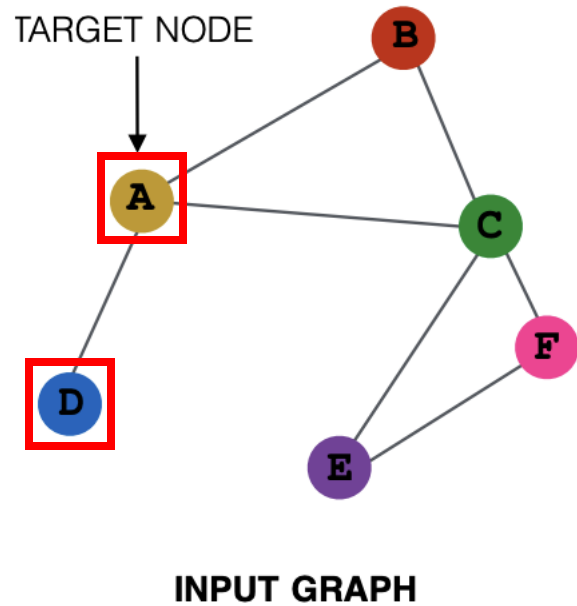
# Graph networks: aggregate neighbors

Nearest neighbors of **D**?



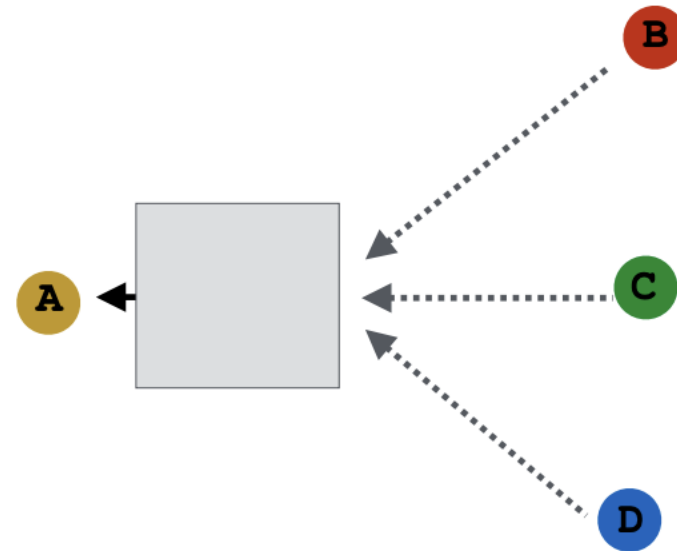
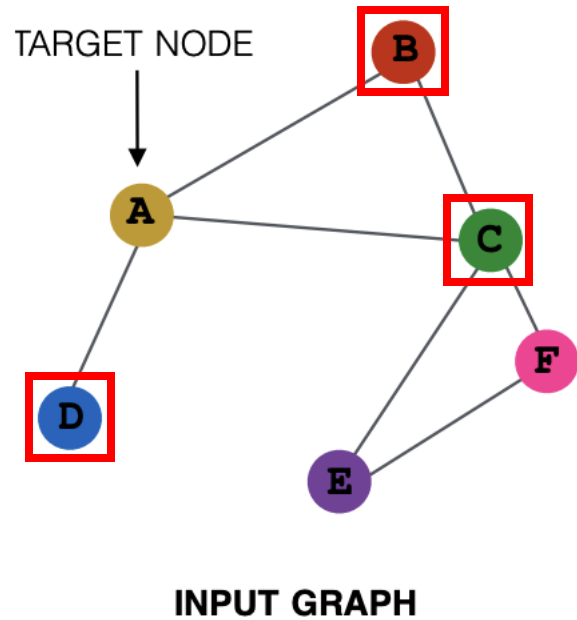
# Graph networks: aggregate neighbors

Nearest neighbors of **D**: **A**



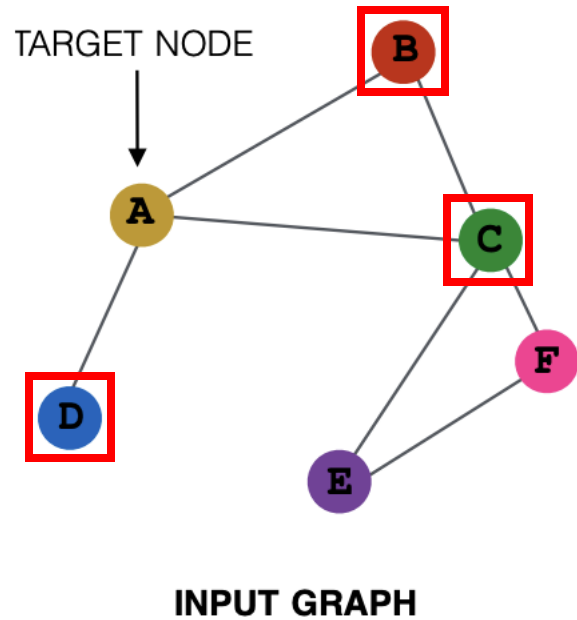
# Graph networks: aggregate neighbors

Nearest neighbors of A:  
B, C, D

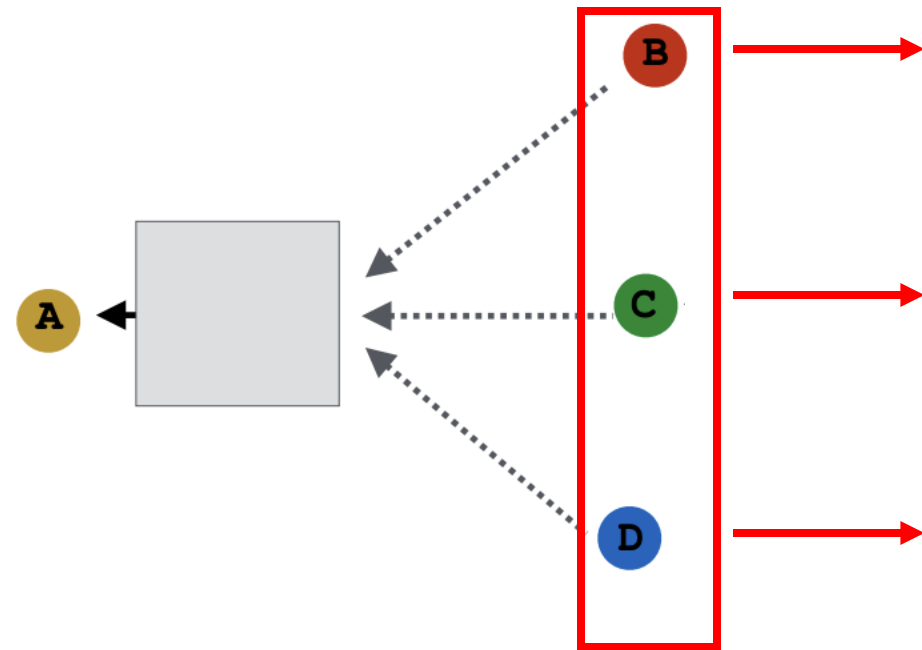


# Graph networks: aggregate neighbors

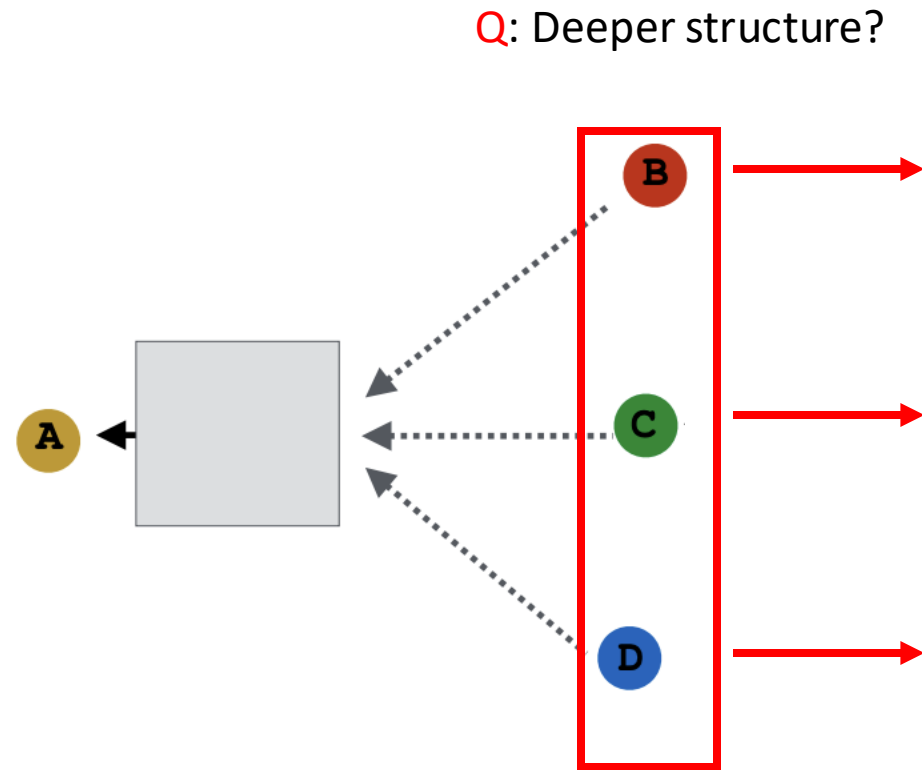
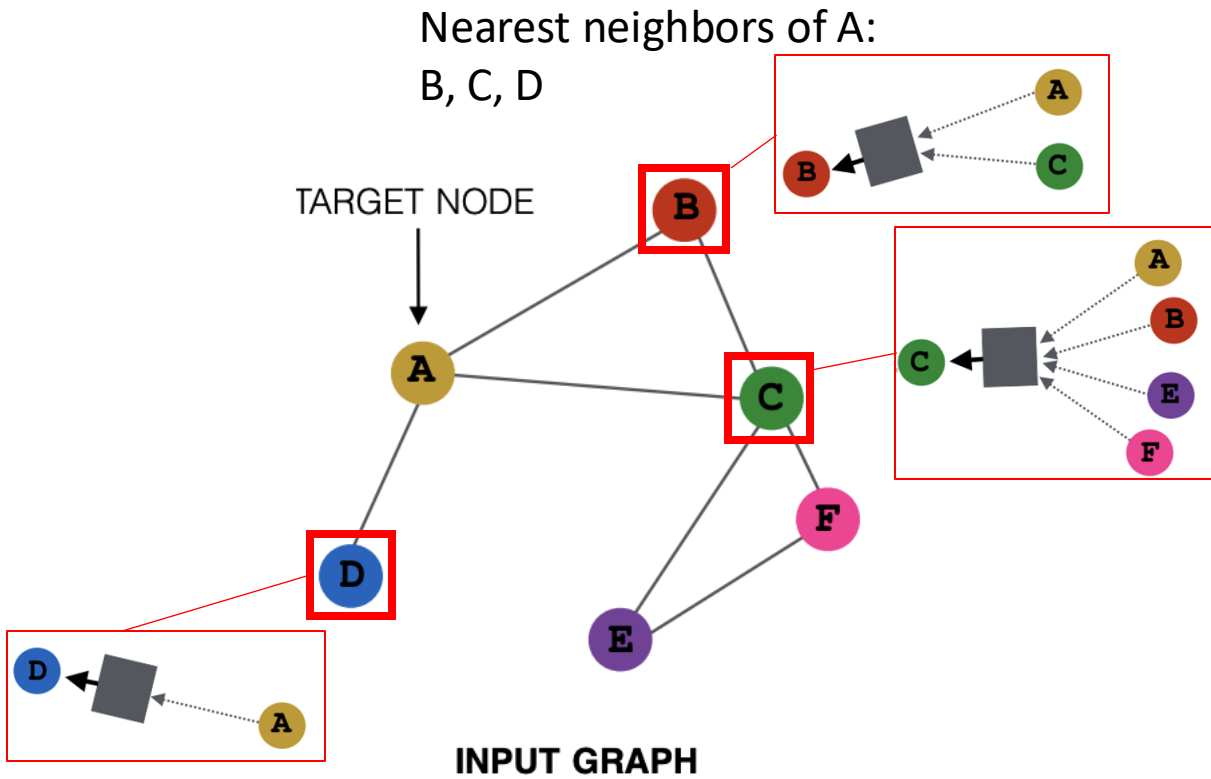
Nearest neighbors of A:  
B, C, D



Q: Deeper structure?



# Graph networks: aggregate neighbors



# Graph networks: aggregate neighbors

Nearest neighbors of A:  
B, C, D

TARGET NODE

A

B

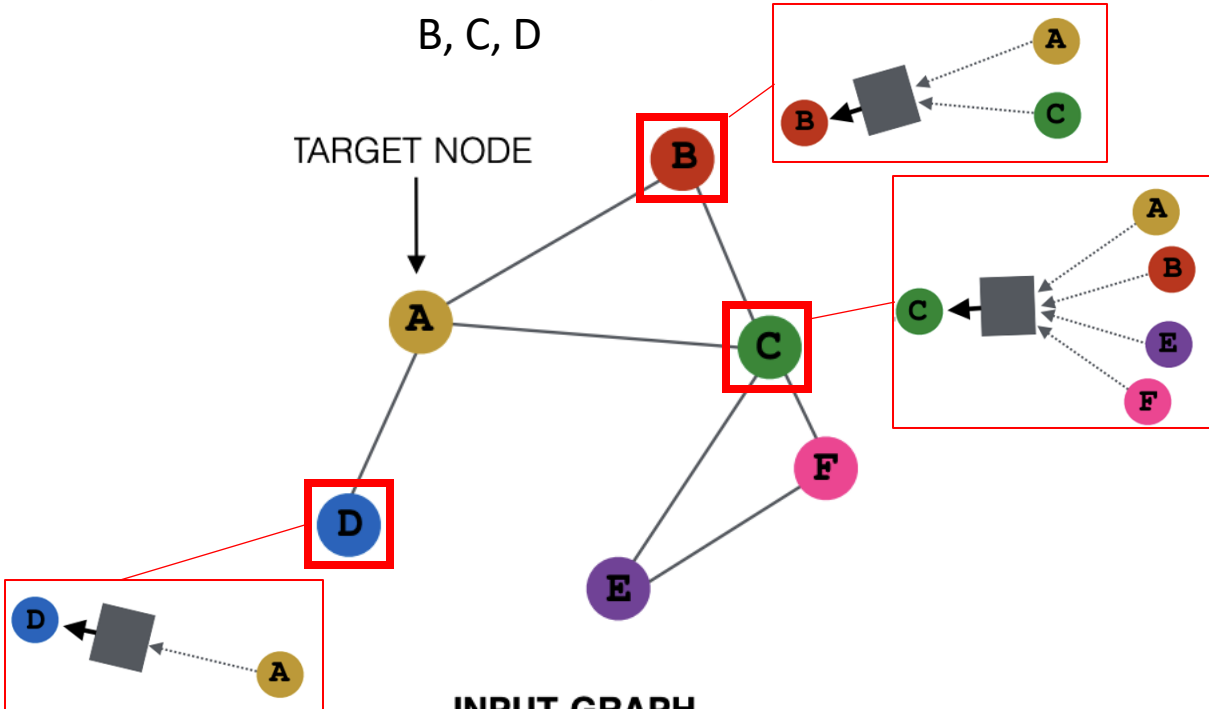
C

D

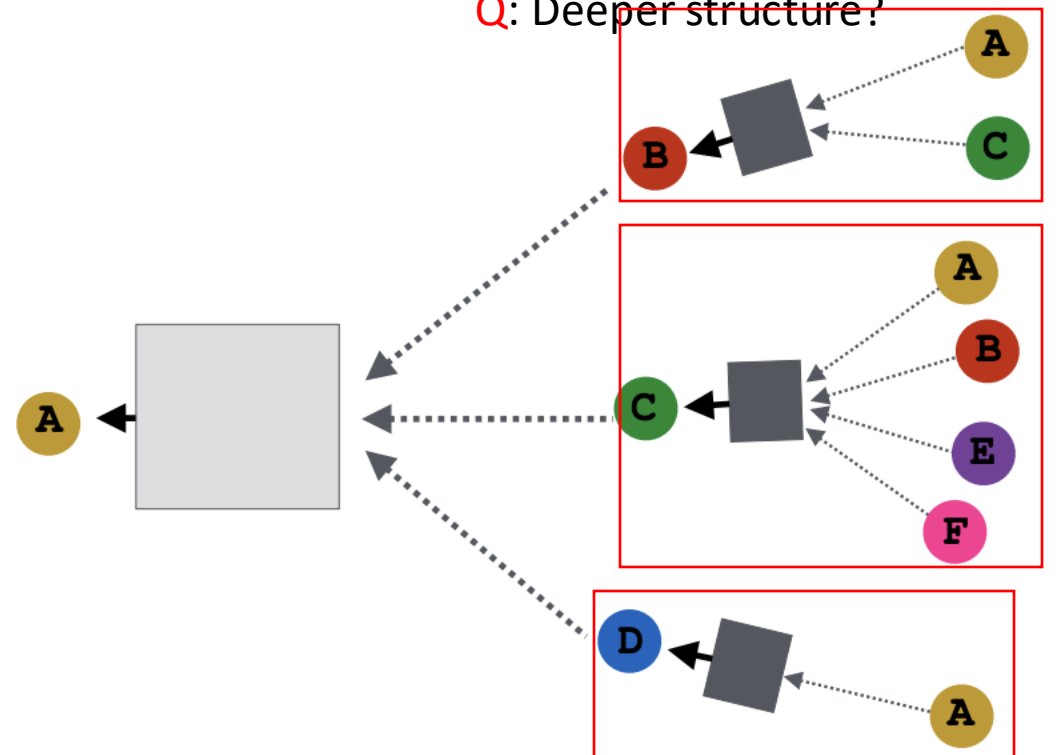
E

F

INPUT GRAPH

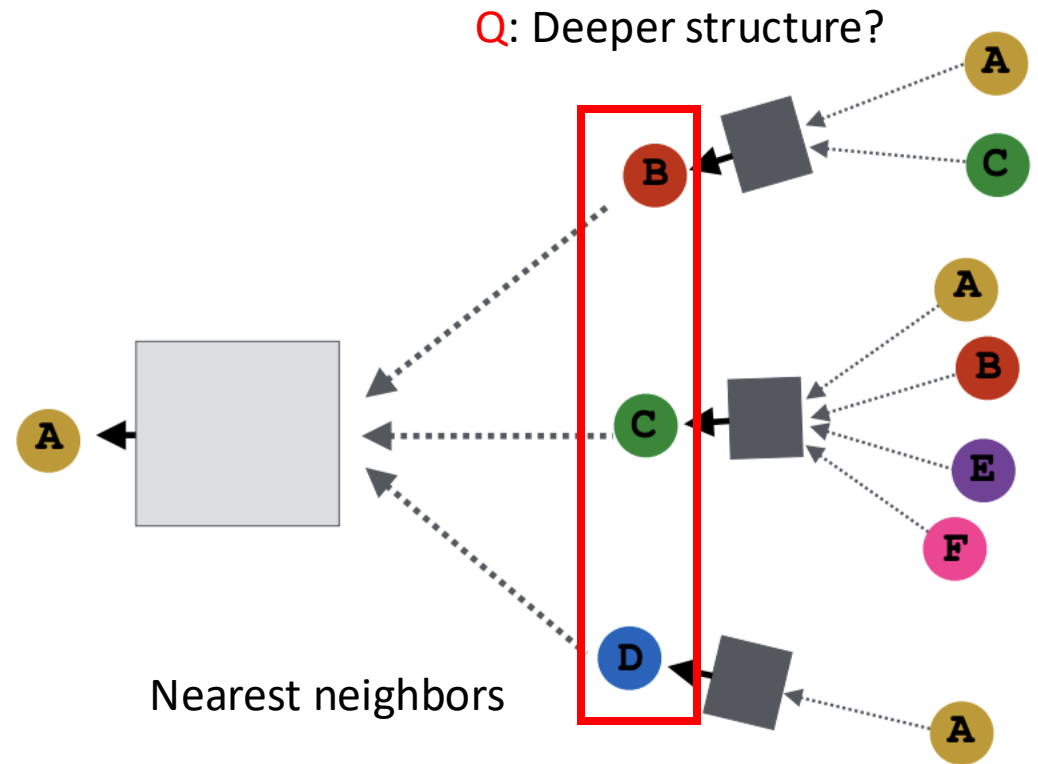
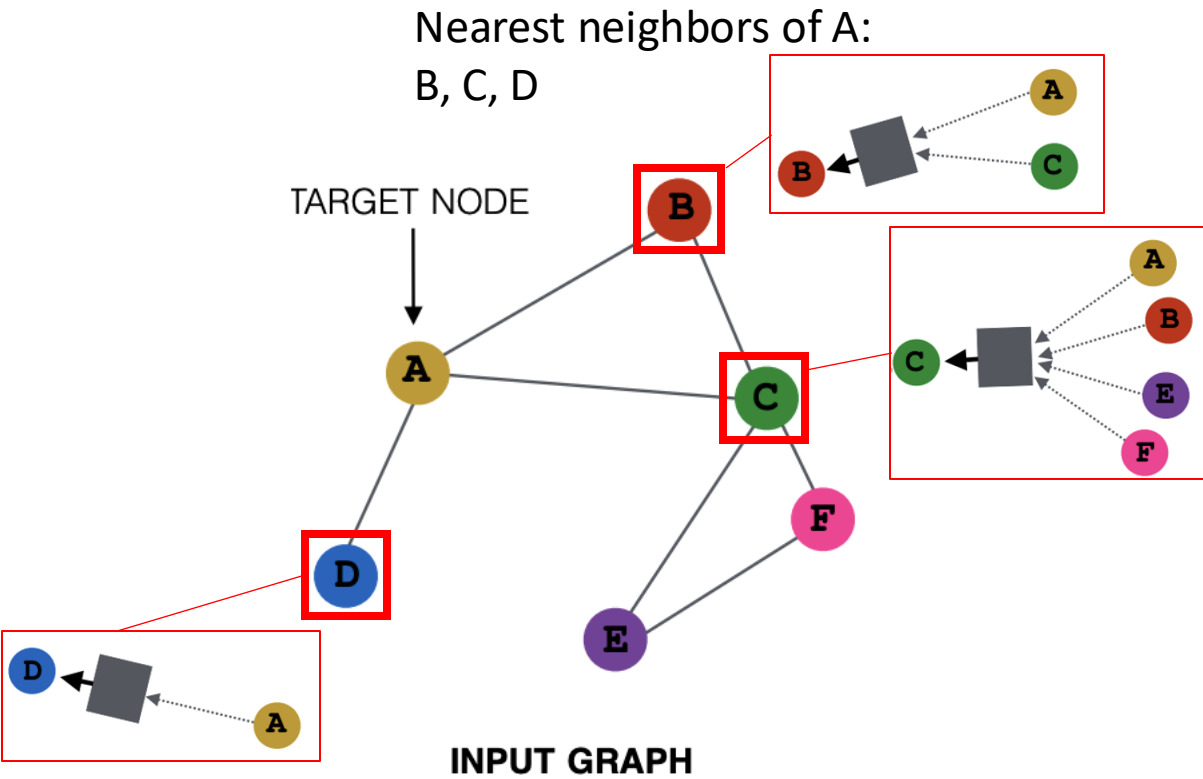


Q: Deeper structure?

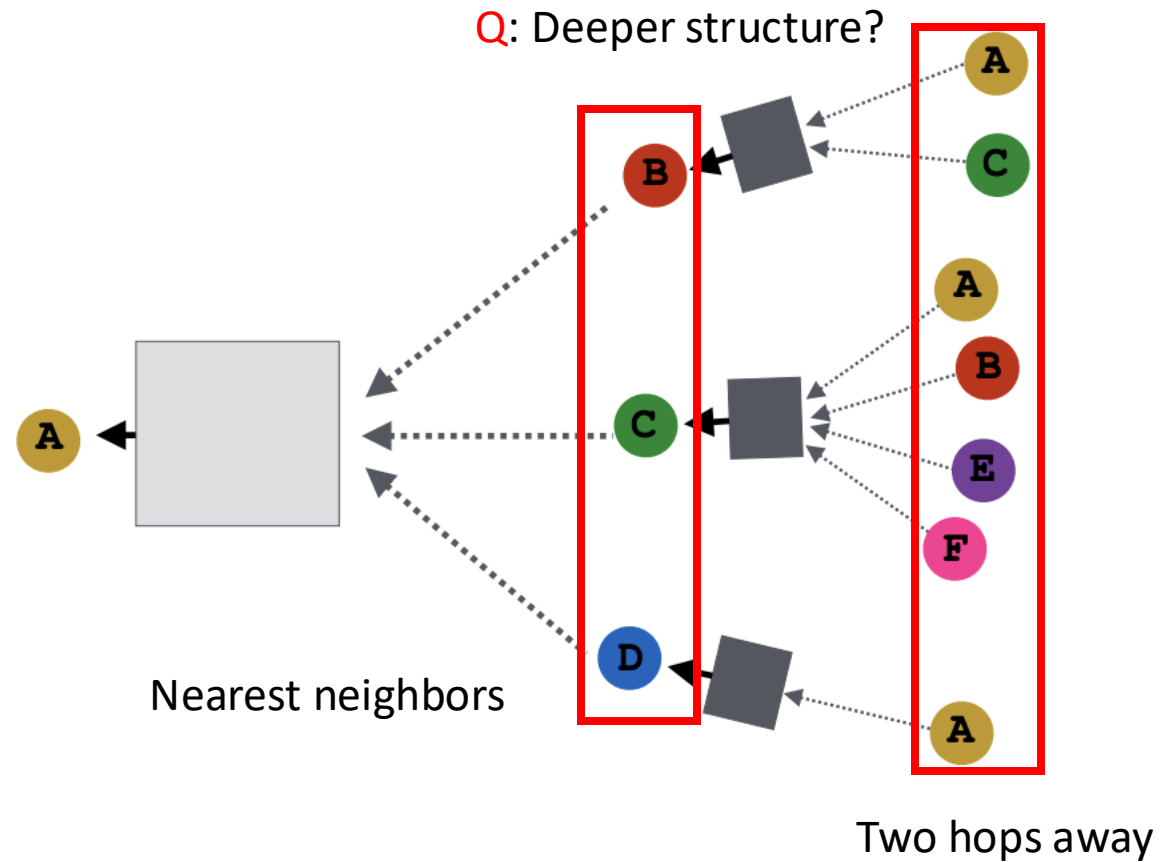
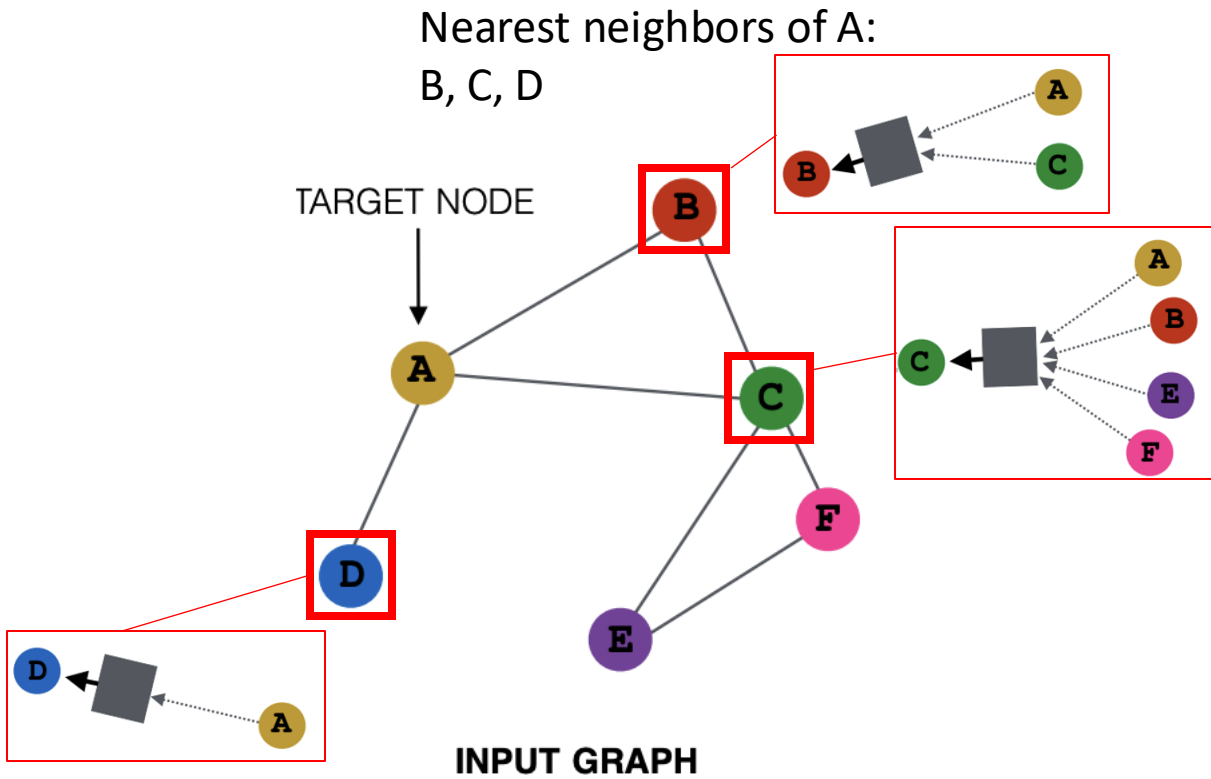




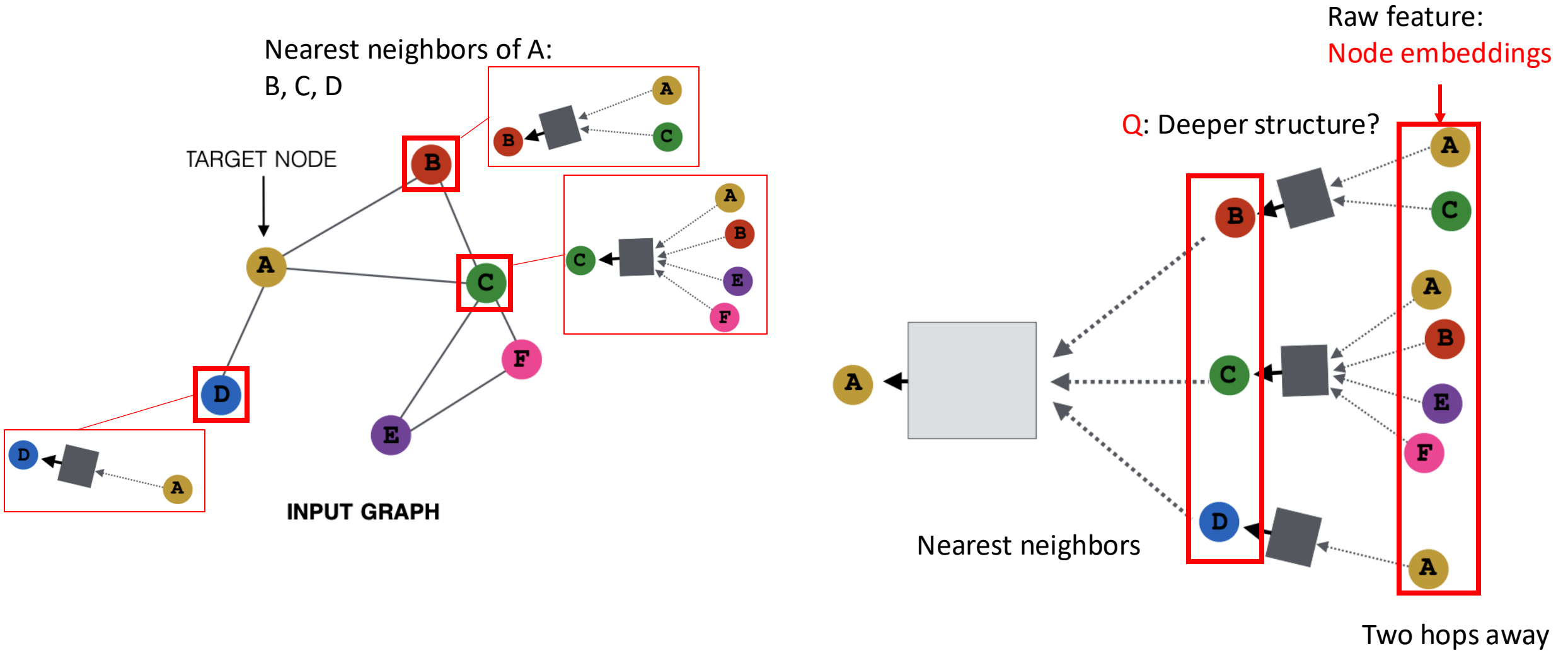
# Graph networks: aggregate neighbors



# Graph networks: aggregate neighbors



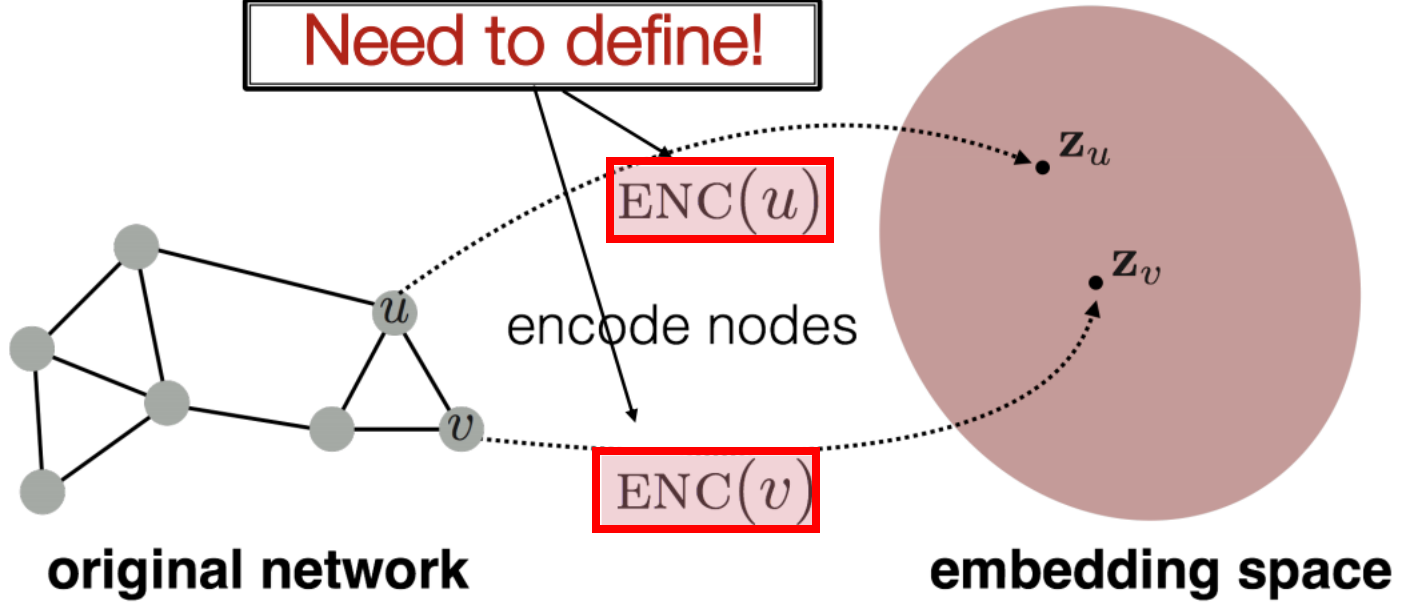
# Graph networks: aggregate neighbors



# Encoder-decoder for graph data

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

Need to define!



Q: how to learn ENC?  
Linear transformation

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

$\mathbf{Z} \in \mathbb{R}^{d \times |\mathcal{V}|}$

$v \in \mathbb{I}^{|\mathcal{V}|}$

Each node

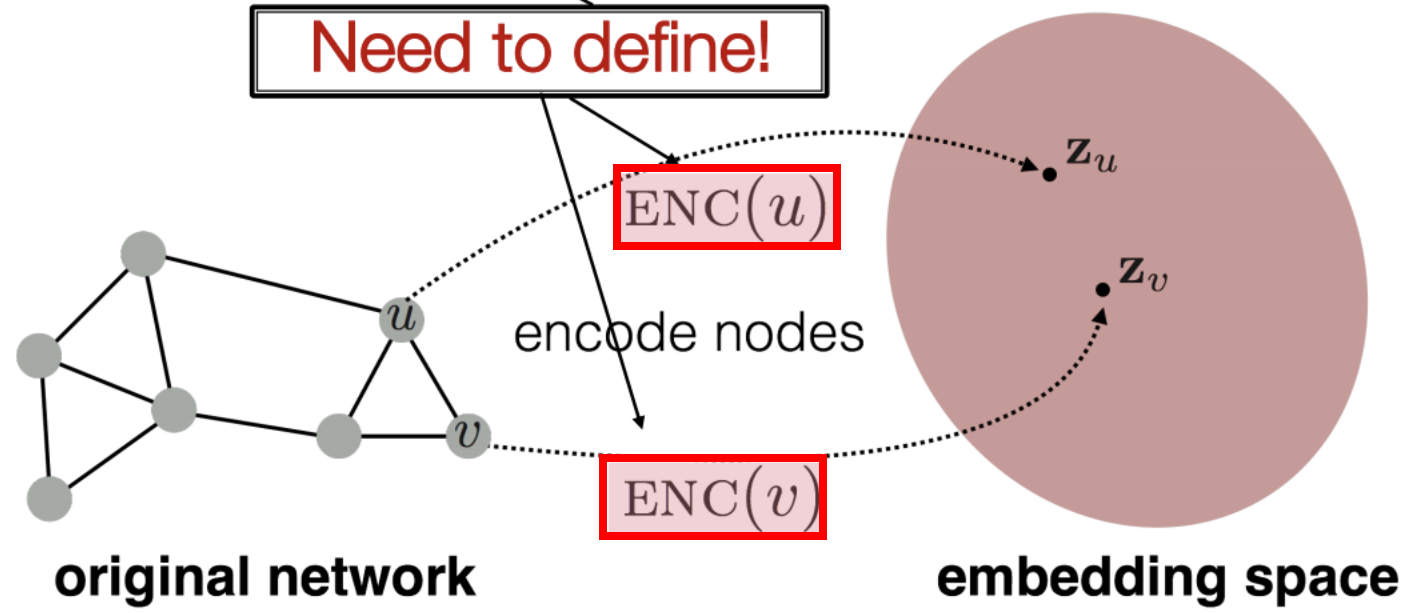
Binary	Gray code	One-hot
000	000	00000001
001	001	00000010
010	011	00000100
011	010	00001000
100	110	00010000
101	111	00100000
110	101	01000000
111	100	10000000

# Encoder-decoder for graph data

Other embedding methods:  
random walk embedding  
node2vec

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

Need to define!



Q: how to learn ENC?  
Linear transformation

Each node

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

$\mathbf{Z} \in \mathbb{R}^{d \times |V|}$

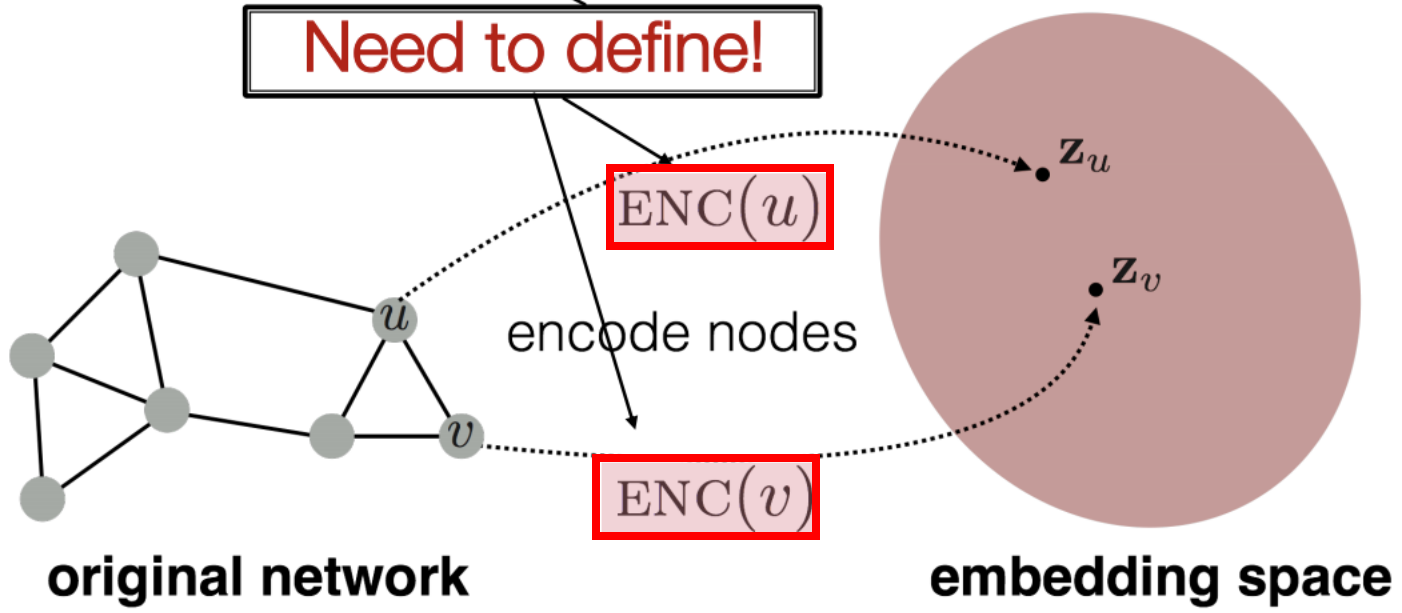
$v \in \mathbb{I}^{|V|}$

Binary	Gray code	One-hot
000	000	00000001
001	001	00000010
010	011	00000100
011	010	00001000
100	110	00010000
101	111	00100000
110	101	01000000
111	100	10000000

# Encoder-decoder for graph data

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

Need to define!



Other embedding methods:  
random walk embedding  
node2vec

Q: how to learn ENC?  
Linear transformation

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

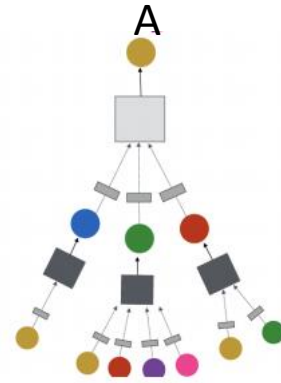
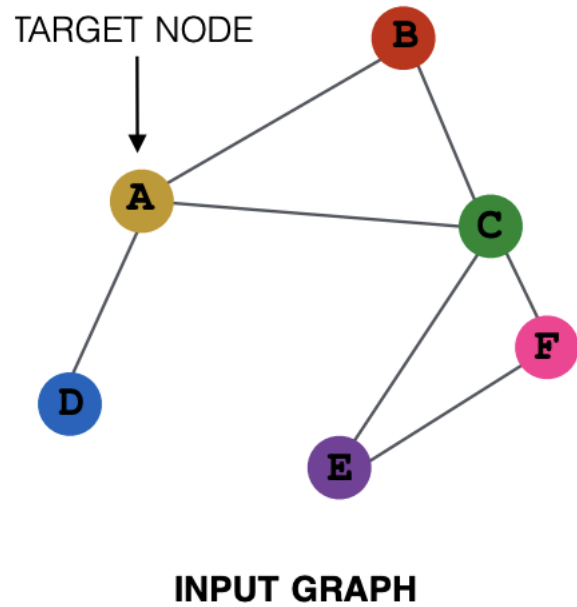
$\mathbf{Z} \in \mathbb{R}^{d \times |V|}$

$v \in \mathbb{I}^{|V|}$

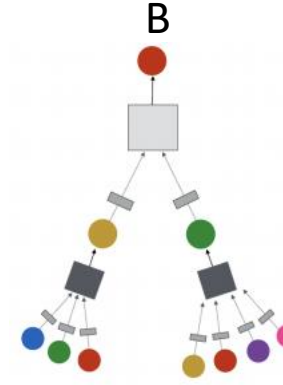
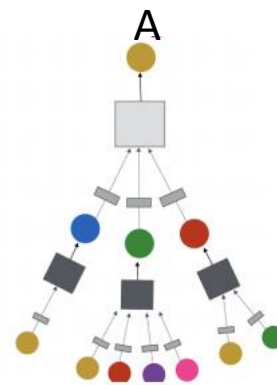
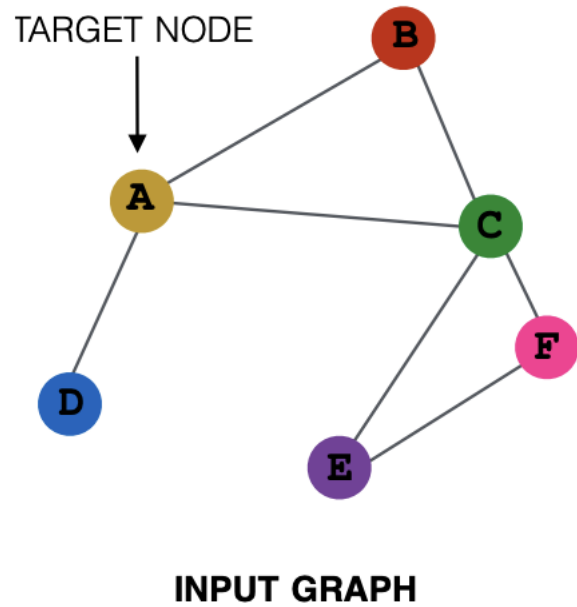
Each node

Binary	Gray code	One-hot
000	000	00000001
001	001	00000010
010	011	00000100
011	010	00001000
100	110	00010000
101	111	00100000
110	101	01000000
111	100	10000000

# Graph networks: aggregate neighbors

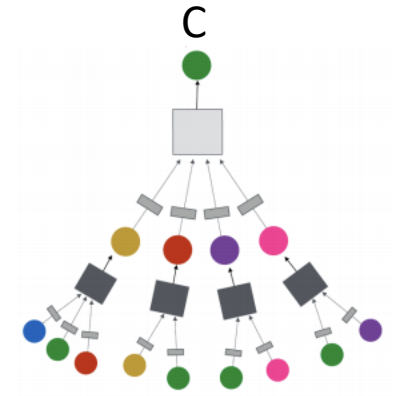
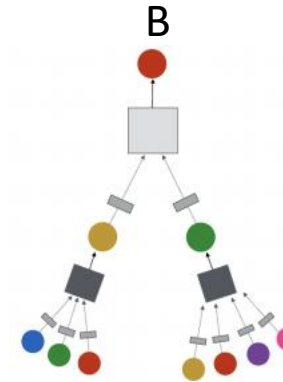
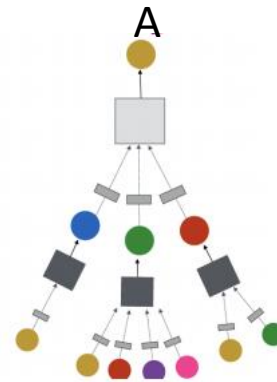
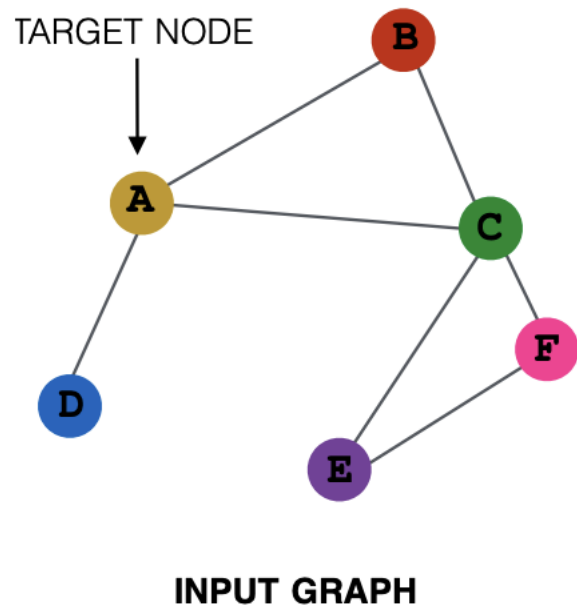


# Graph networks: aggregate neighbors

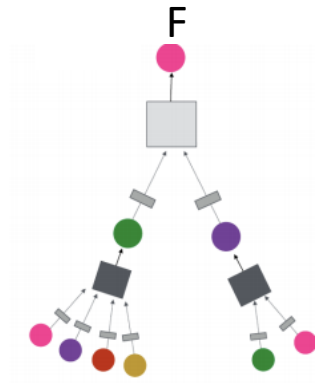
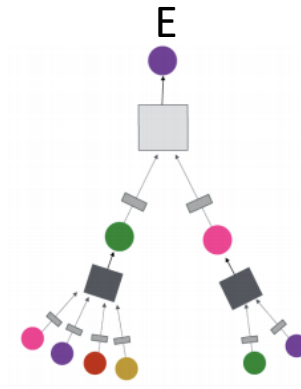
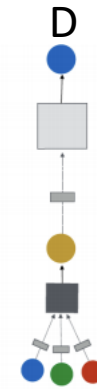
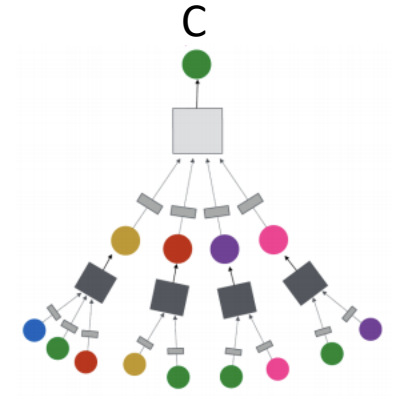
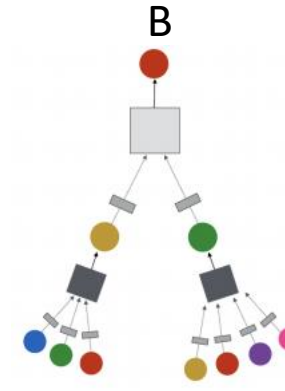
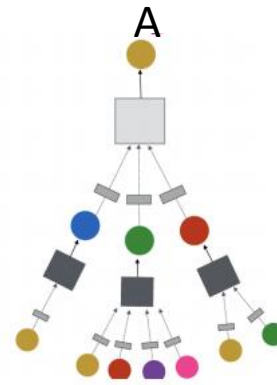
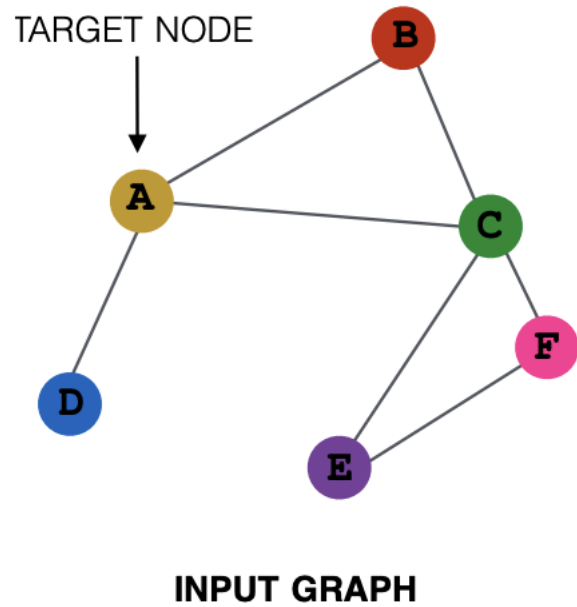




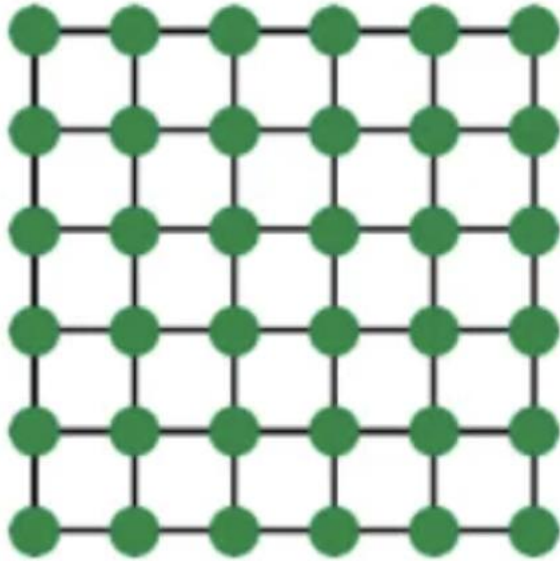
# Graph networks: aggregate neighbors



# Graph networks: aggregate neighbors

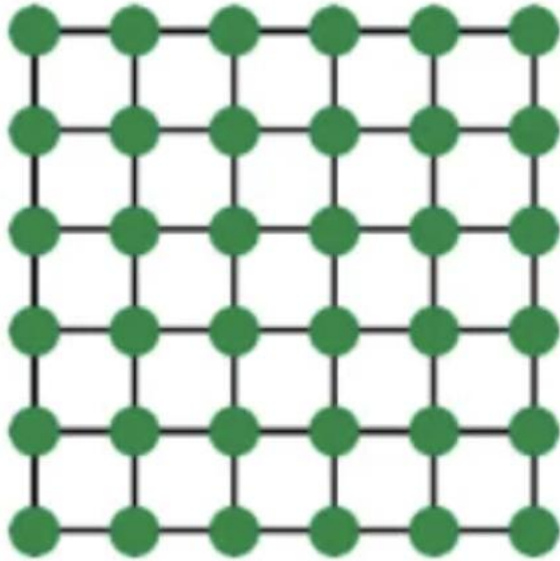


# Graph networks: aggregate neighbors

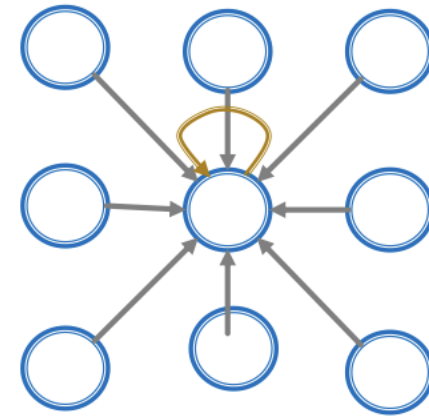


Grid graph

# Graph networks: aggregate neighbors

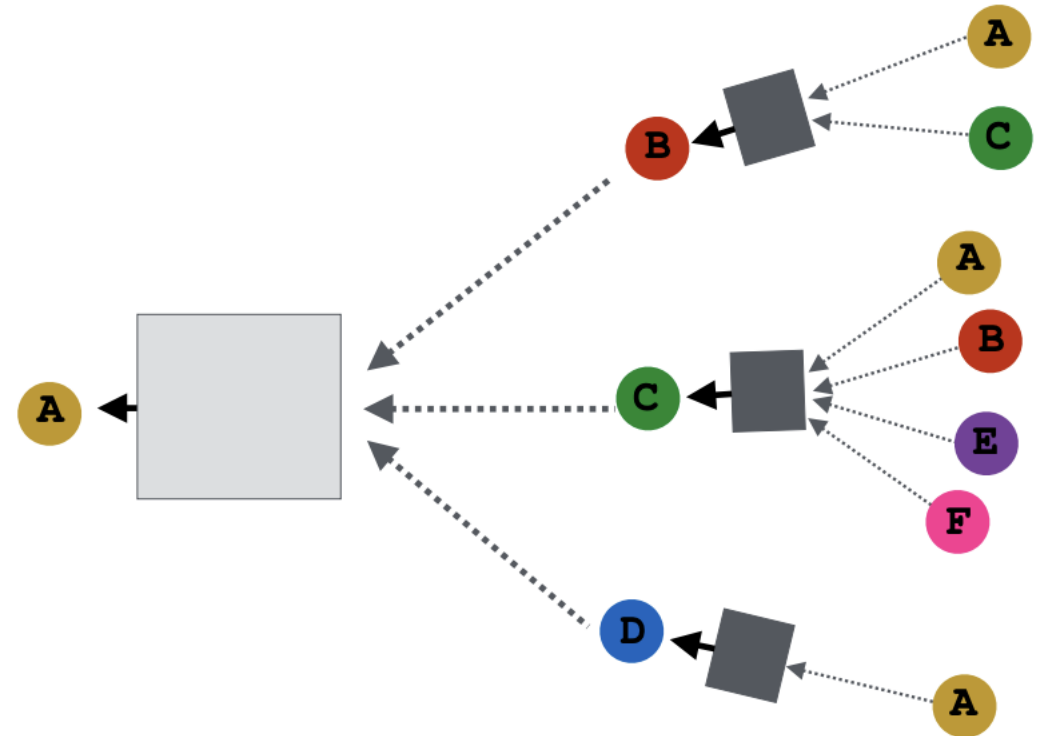
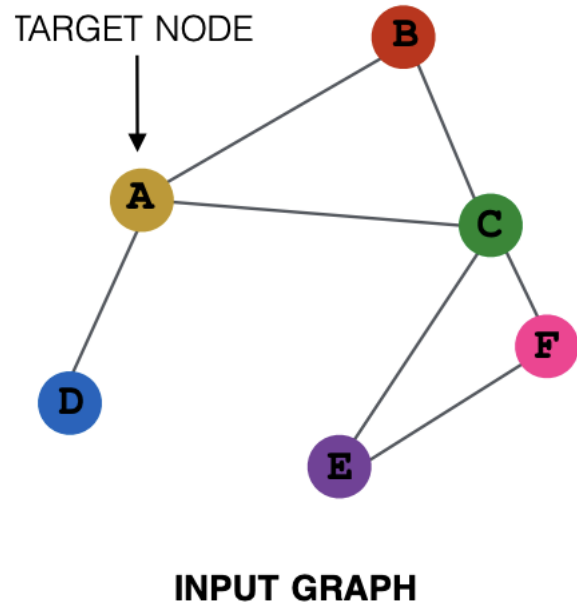


Grid graph

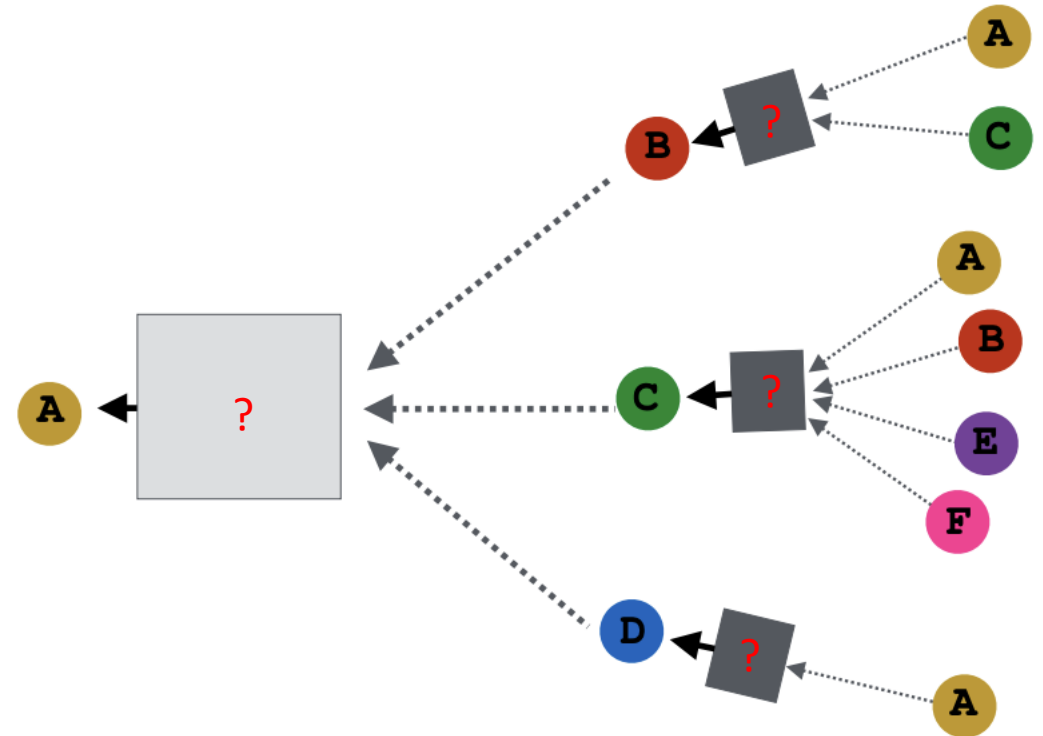
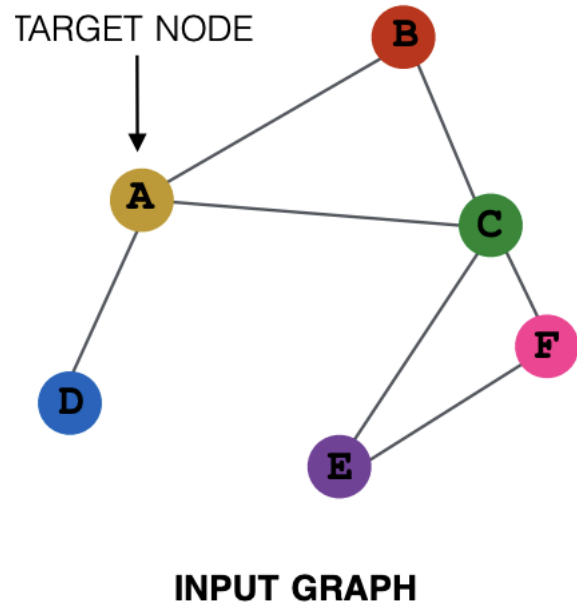


General graph

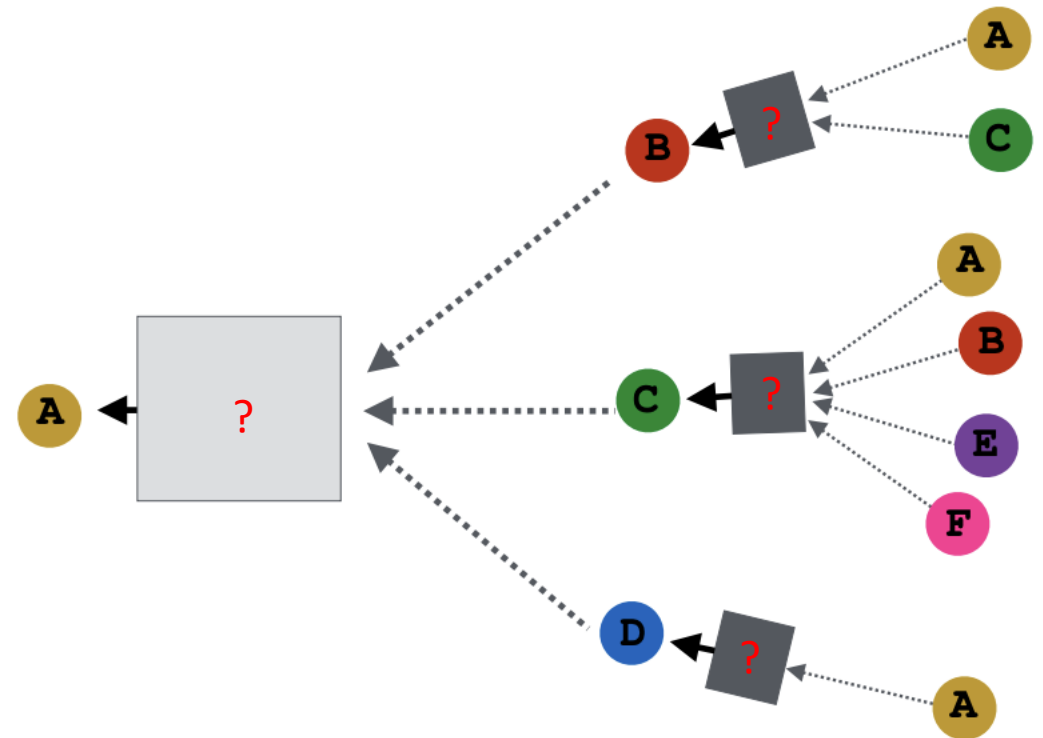
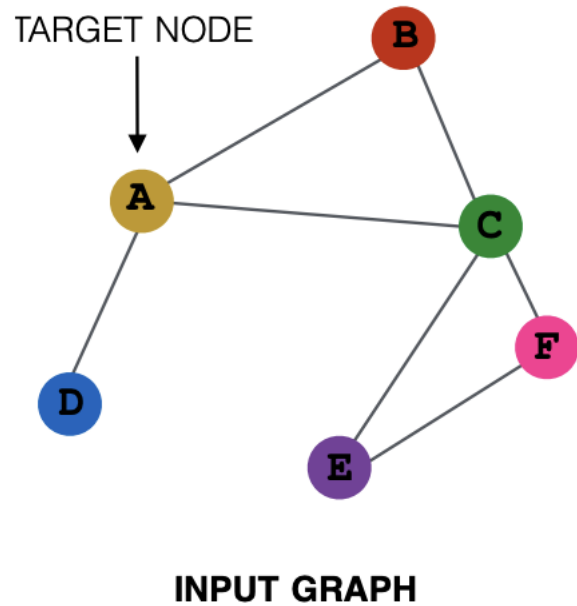
# Graph networks: aggregate neighbors



# Graph networks: aggregate neighbors

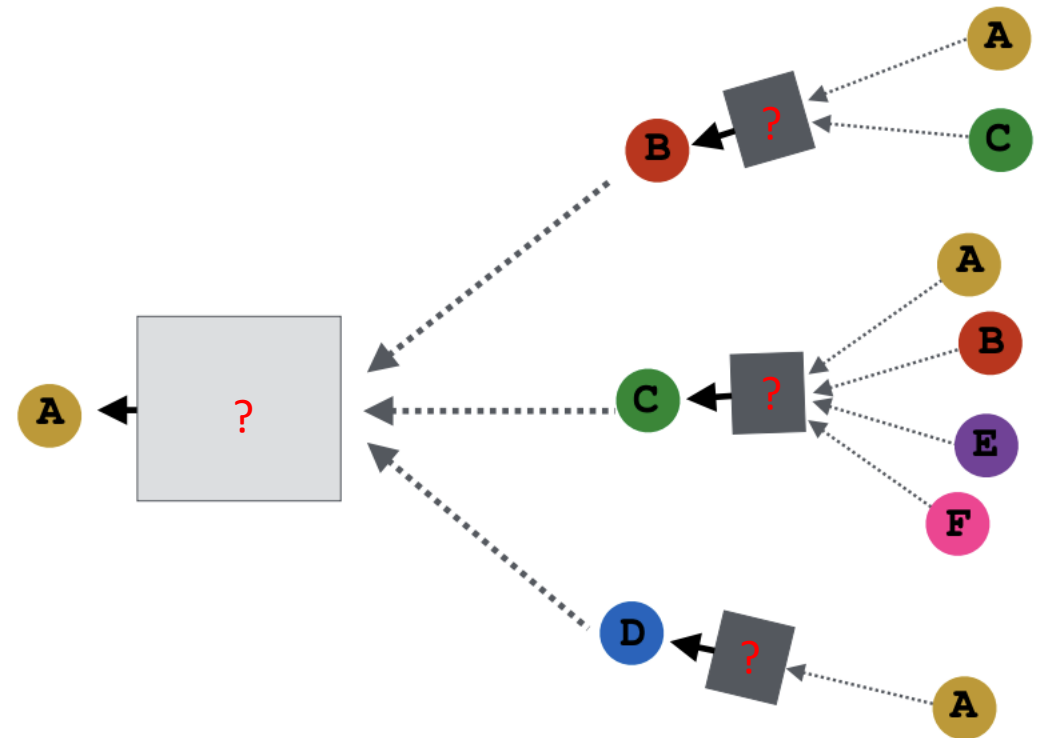
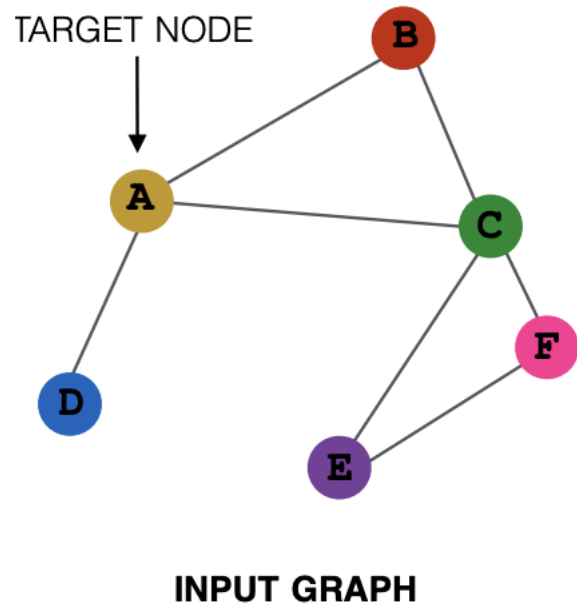


# Graph networks: aggregate neighbors



Q: what can we do in the box to aggregate information?

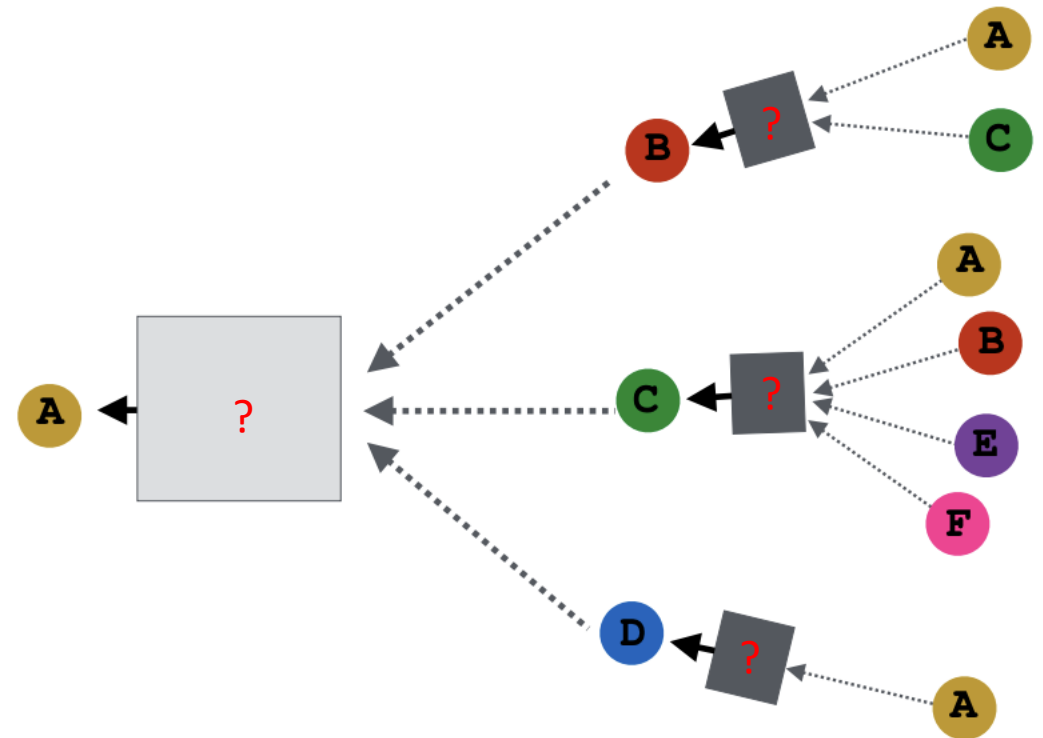
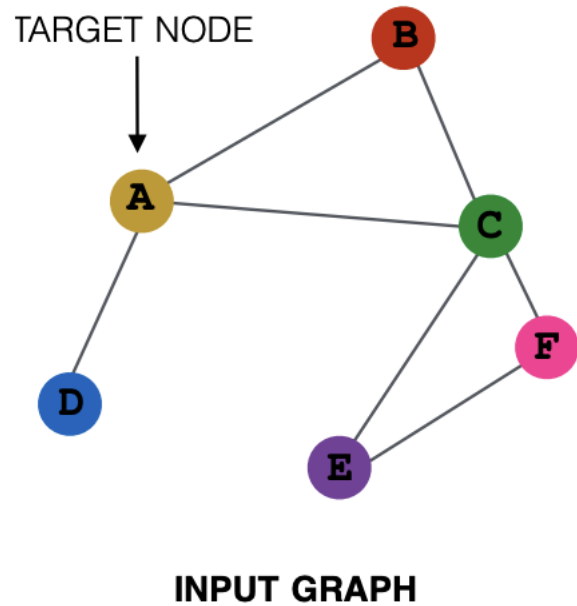
# Graph networks: aggregate neighbors



Q: what can we do in the box to aggregate information?  
Average/summation?

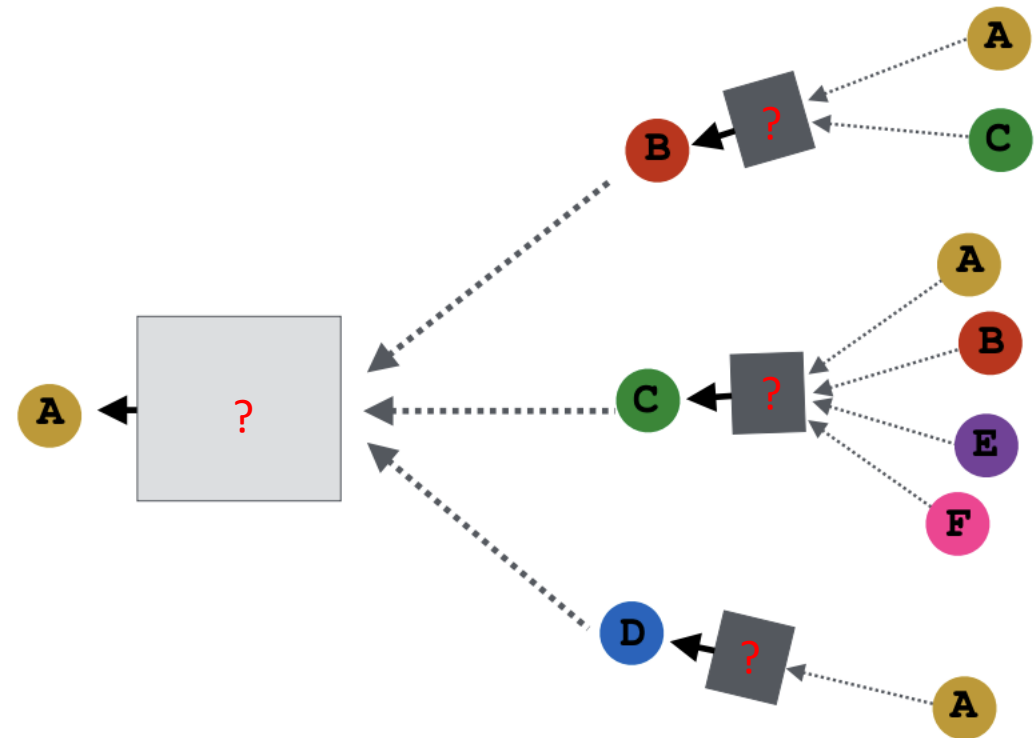
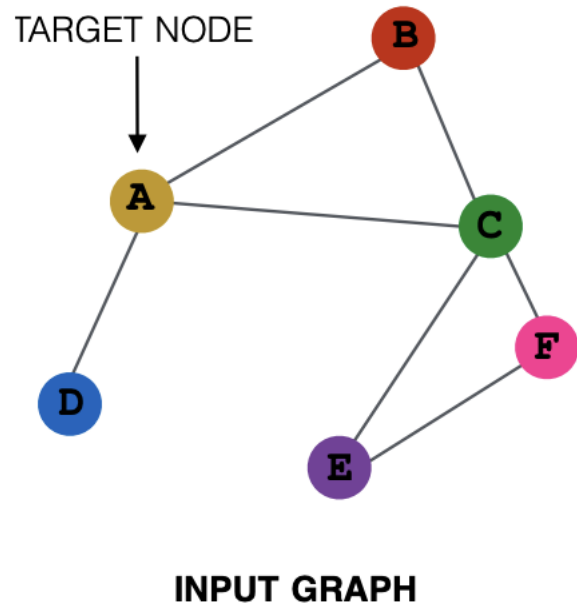


# Graph networks: aggregate neighbors



Q: what can we do in the box to aggregate information?  
Average/summation → linear model

# Graph networks: aggregate neighbors

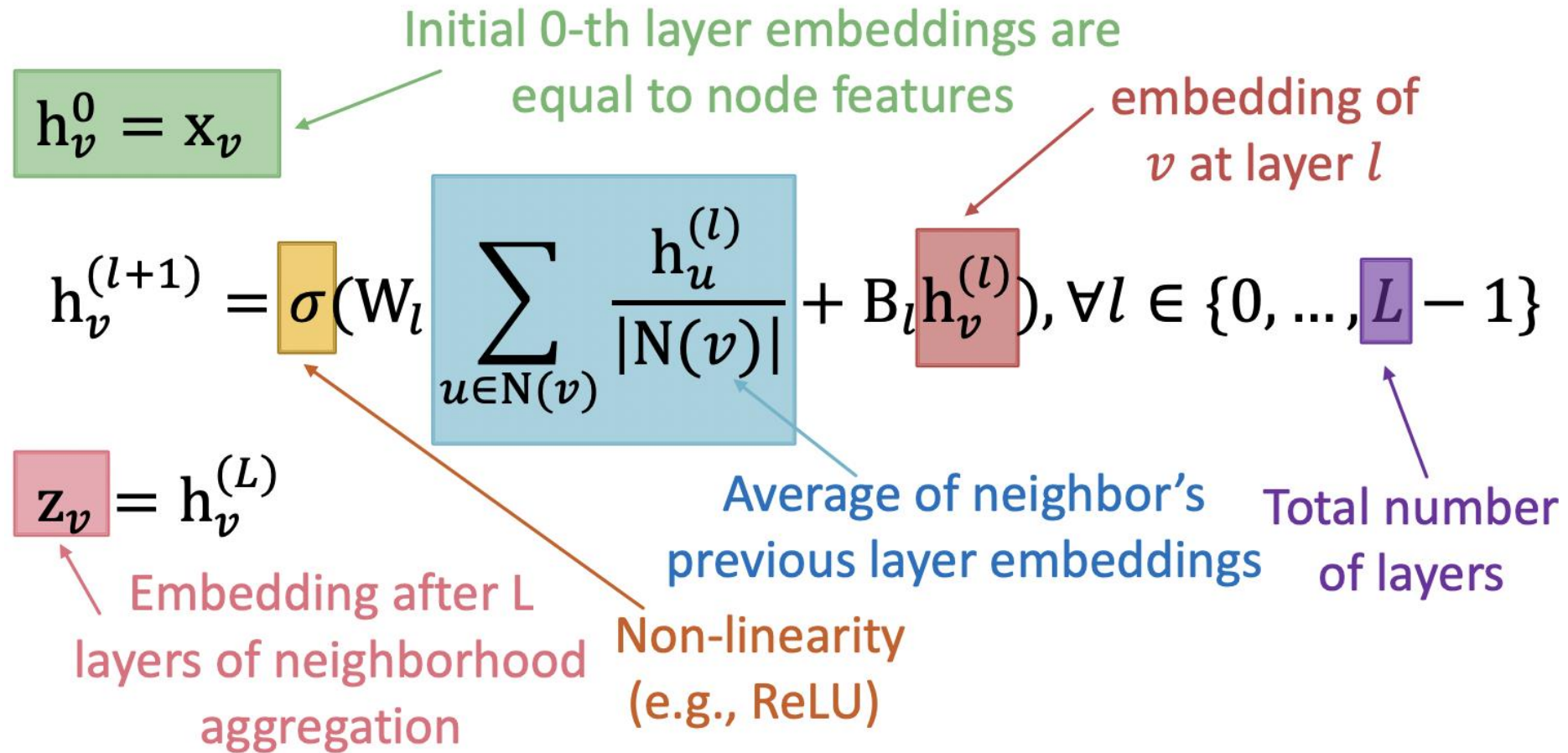


Q: what can we do in the box to aggregate information?

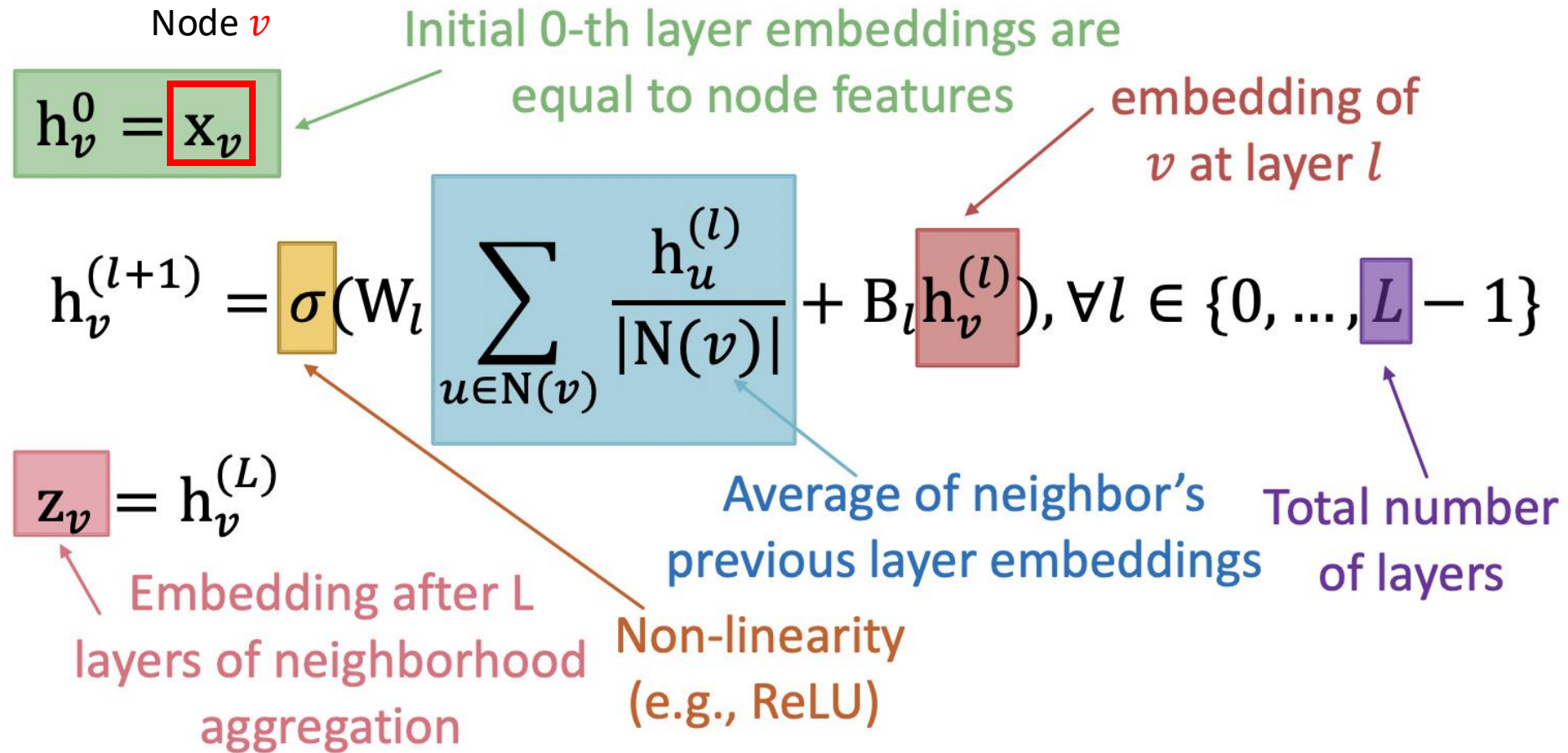
Average/summation → linear model

A neural network nonlinear layer?

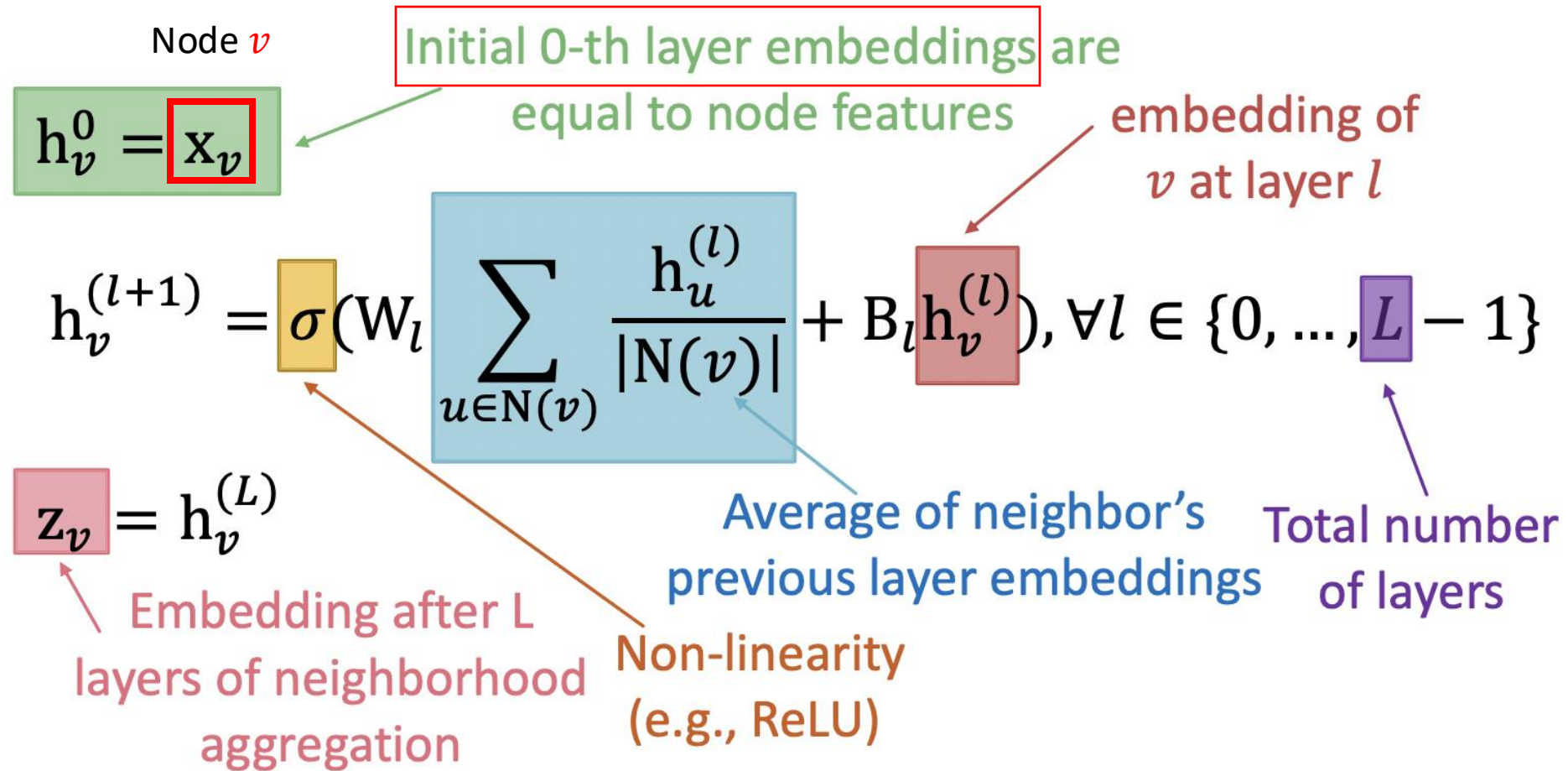
# Graph networks: aggregate neighbors



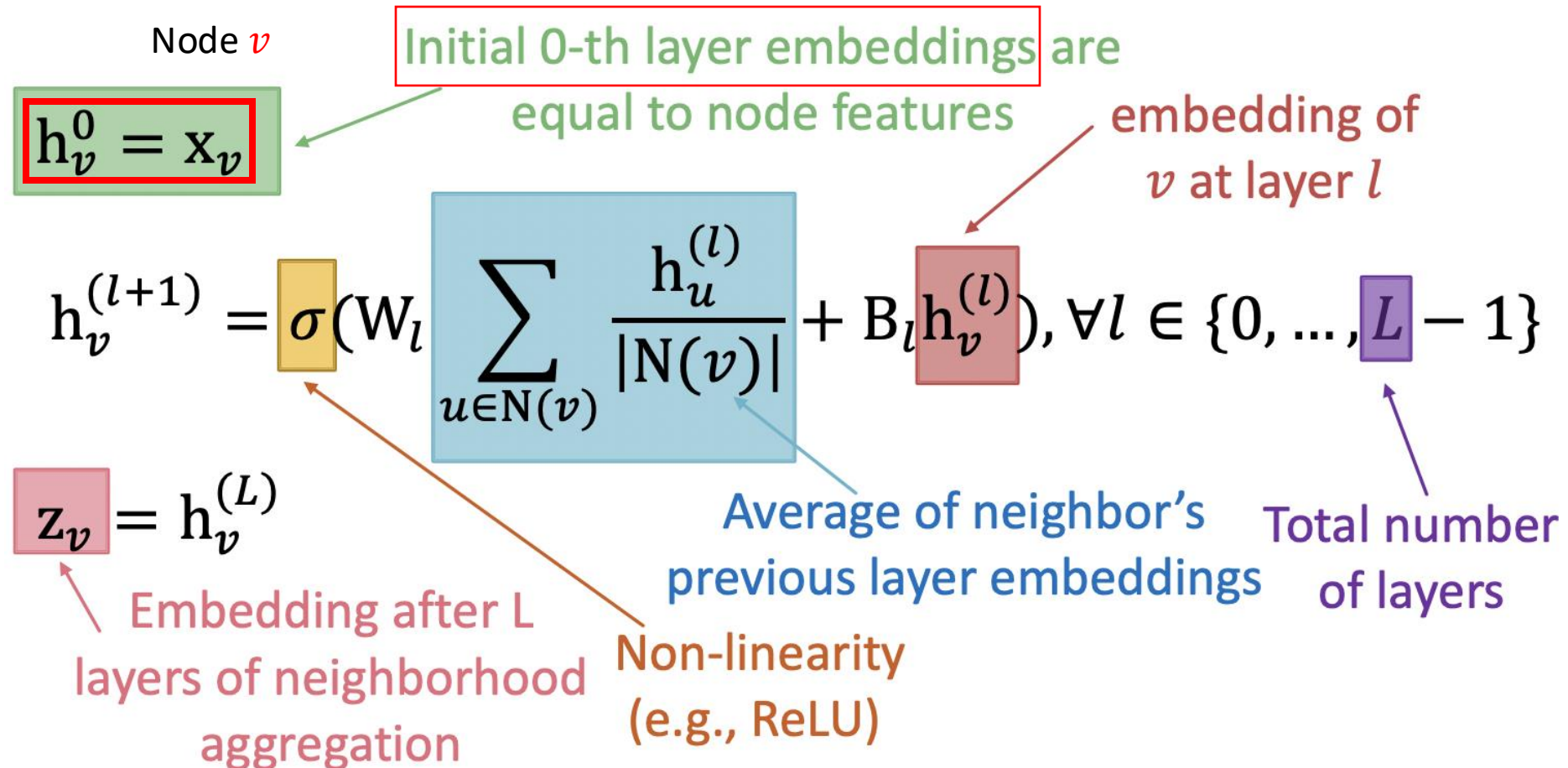
# Graph networks: aggregate neighbors



# Graph networks: aggregate neighbors

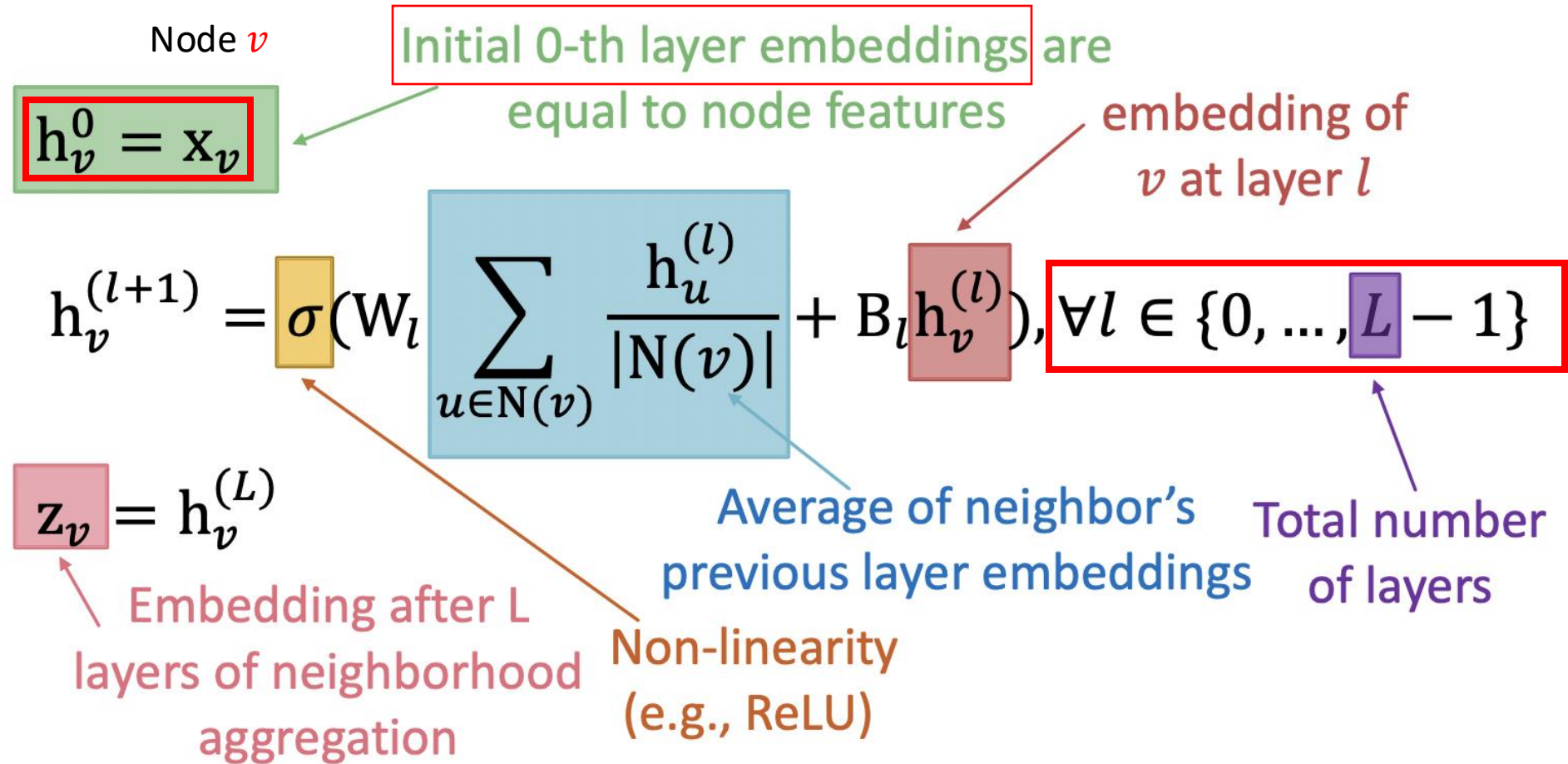


# Graph networks: aggregate neighbors

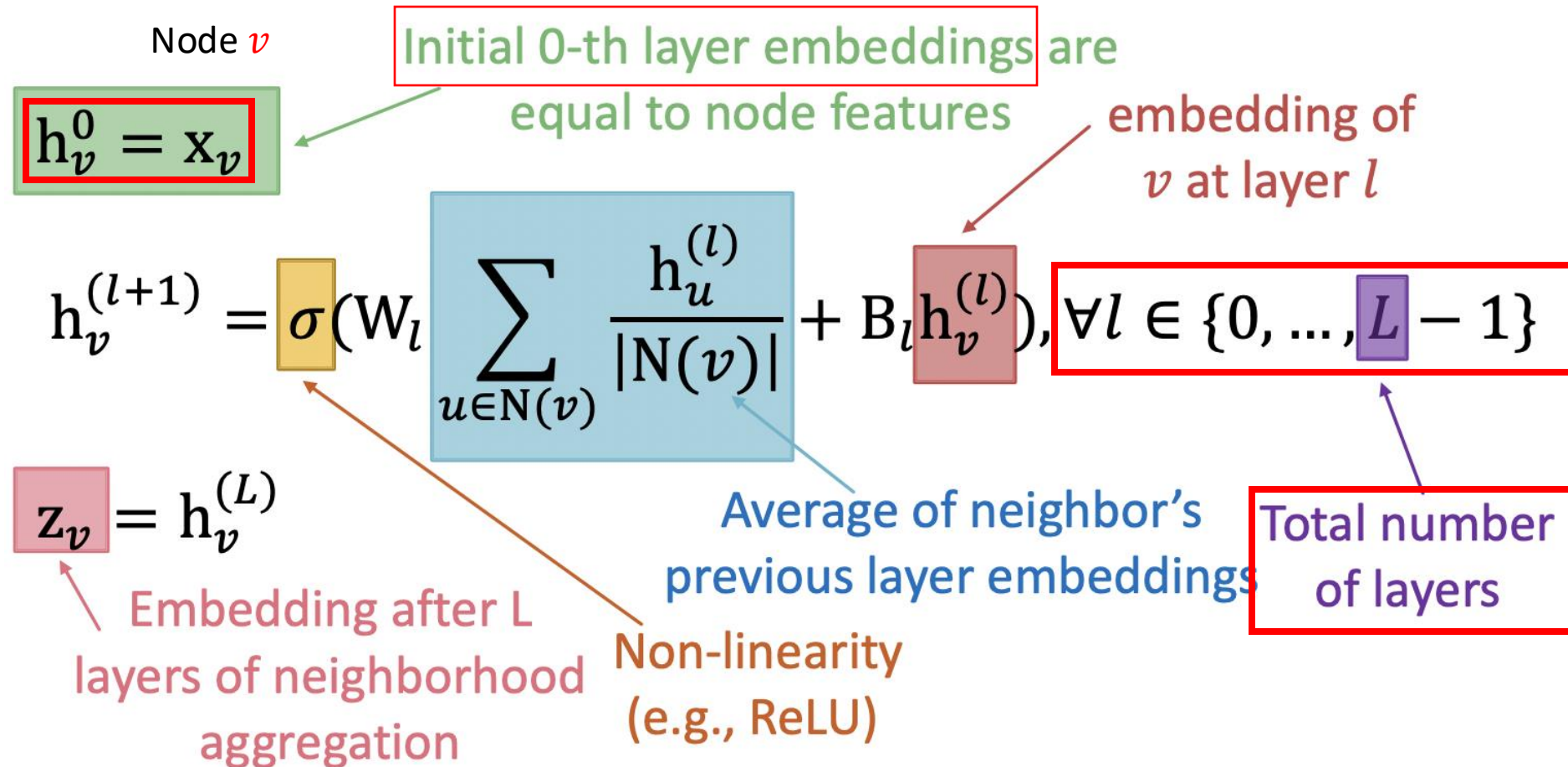




# Graph networks: aggregate neighbors

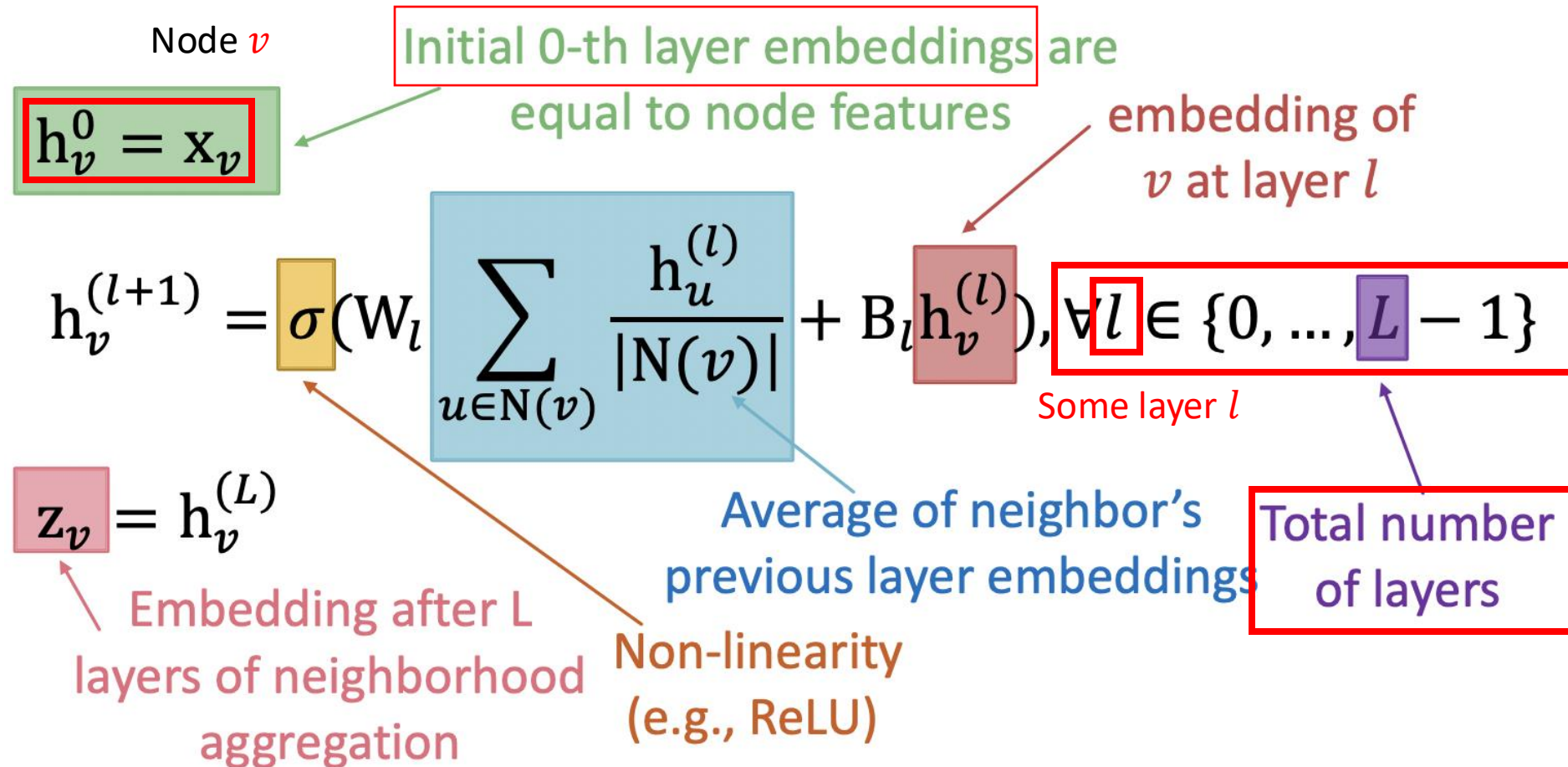


# Graph networks: aggregate neighbors

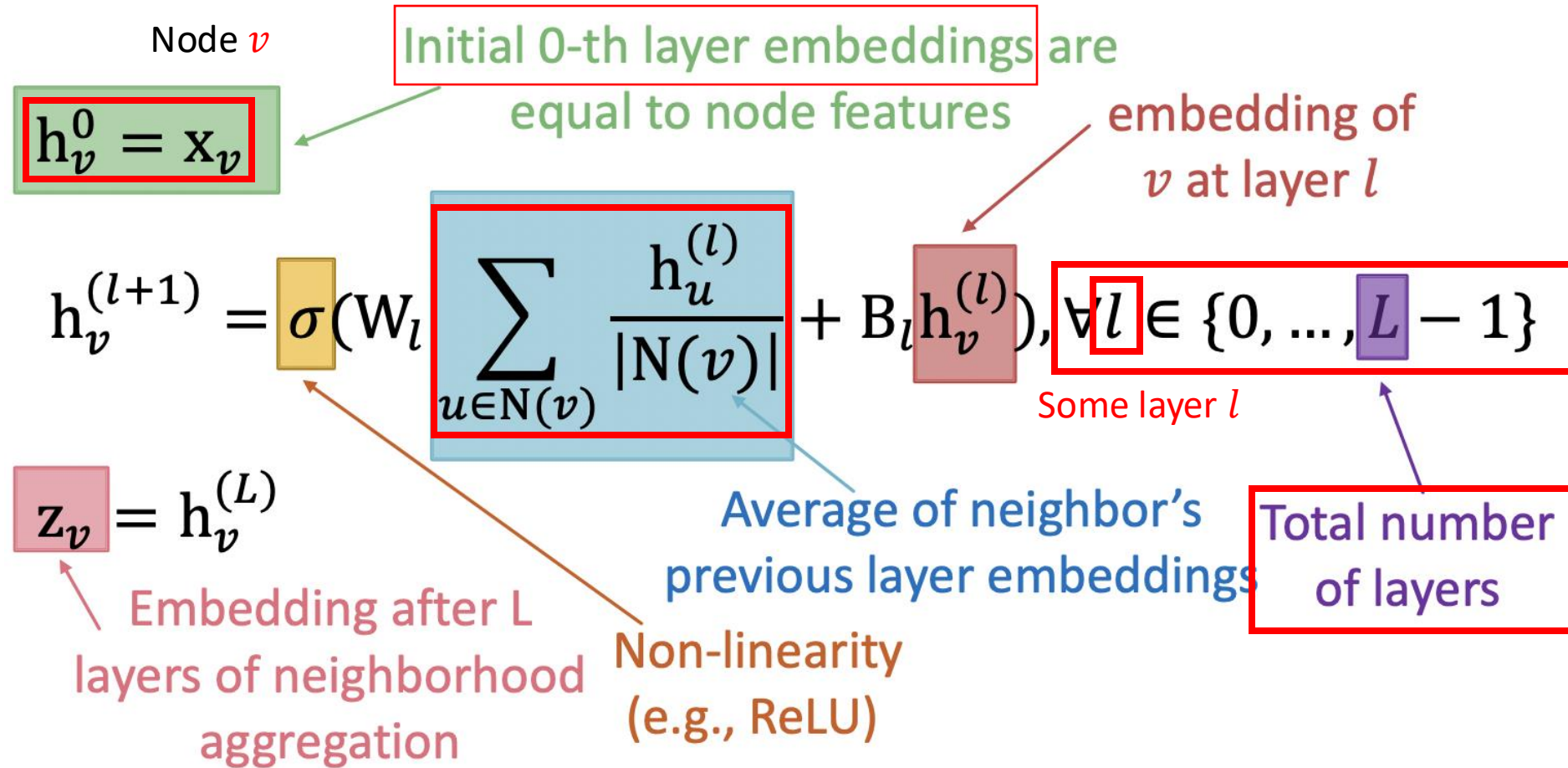




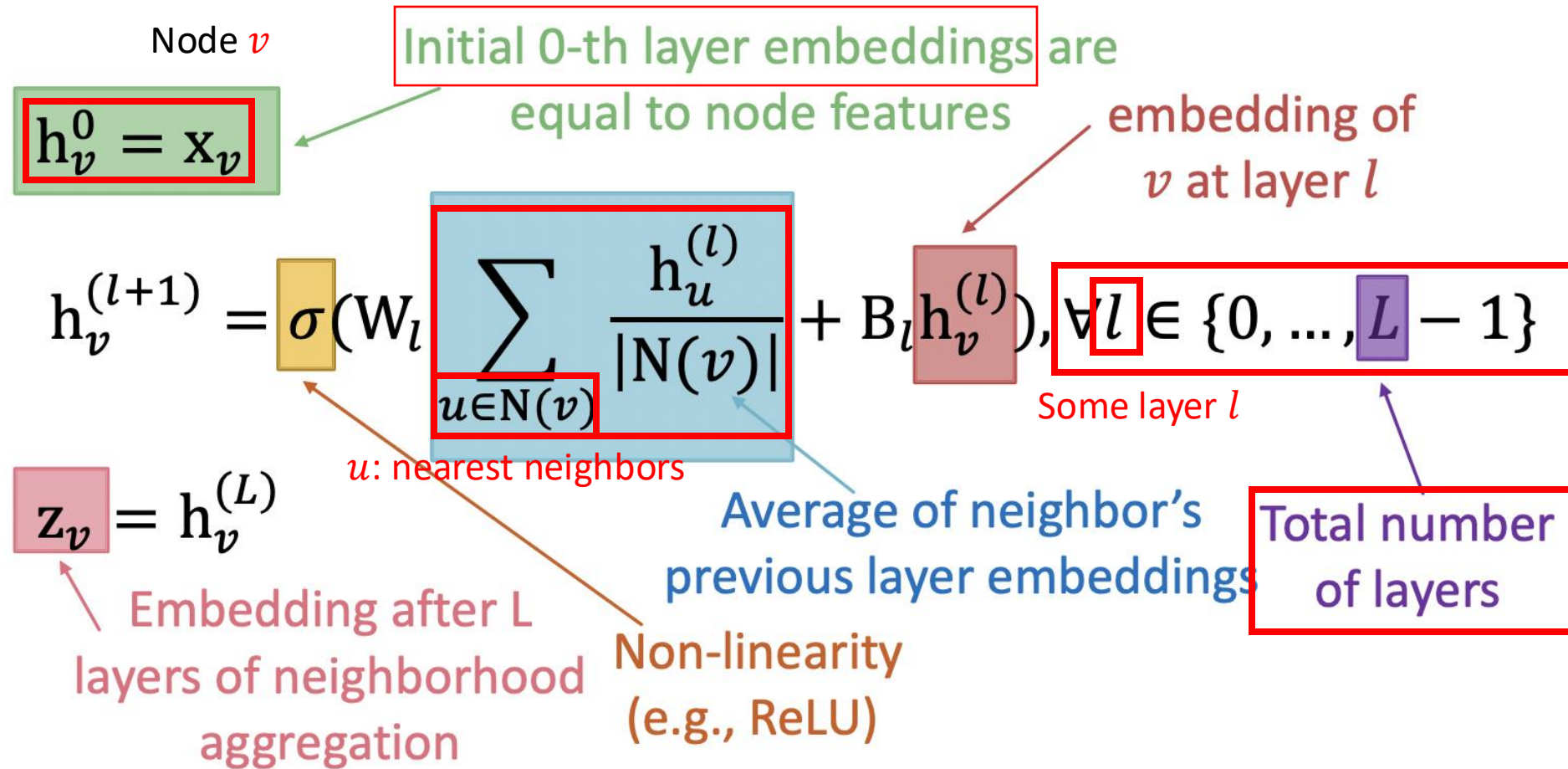
# Graph networks: aggregate neighbors



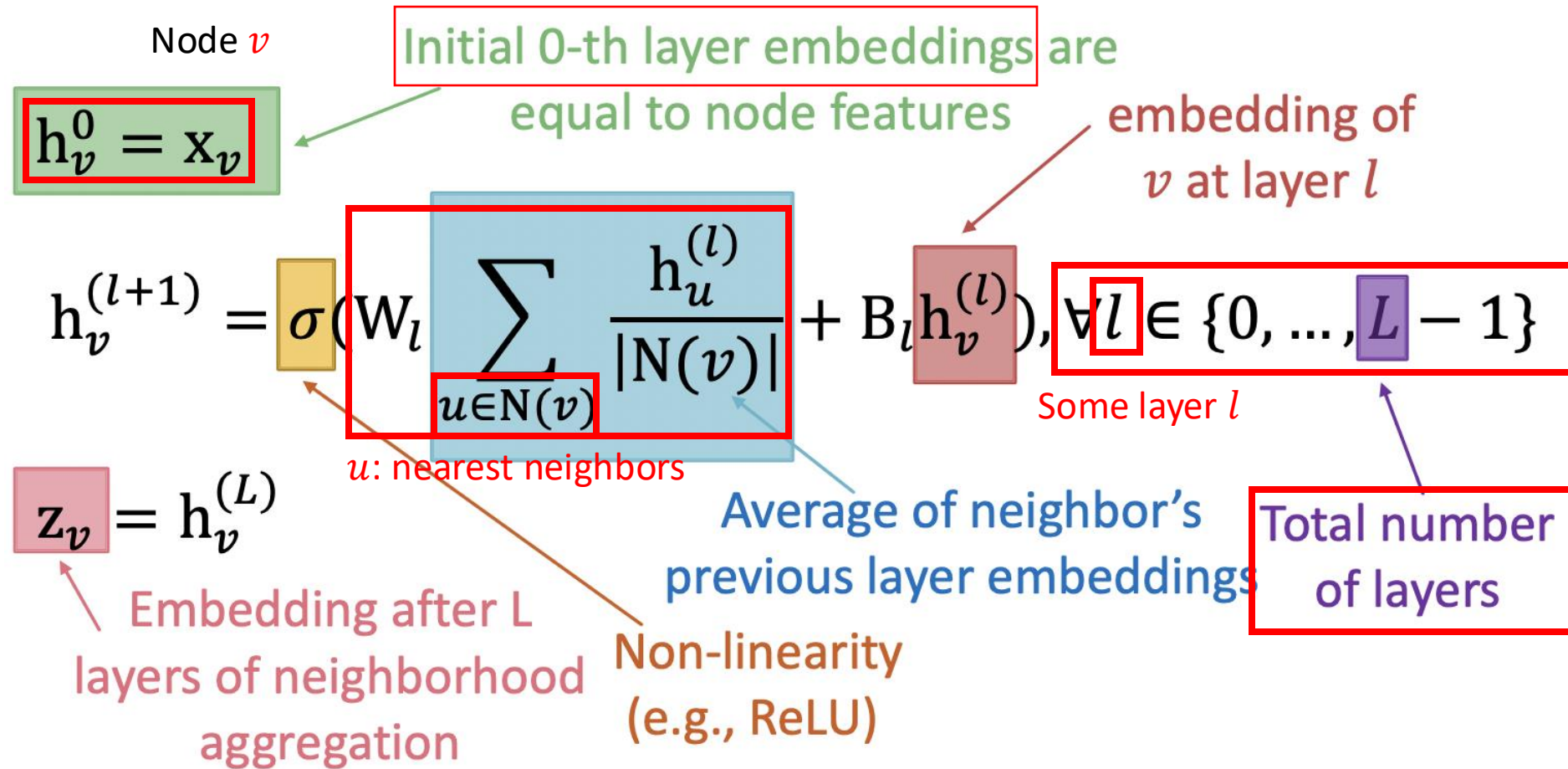
# Graph networks: aggregate neighbors



# Graph networks: aggregate neighbors

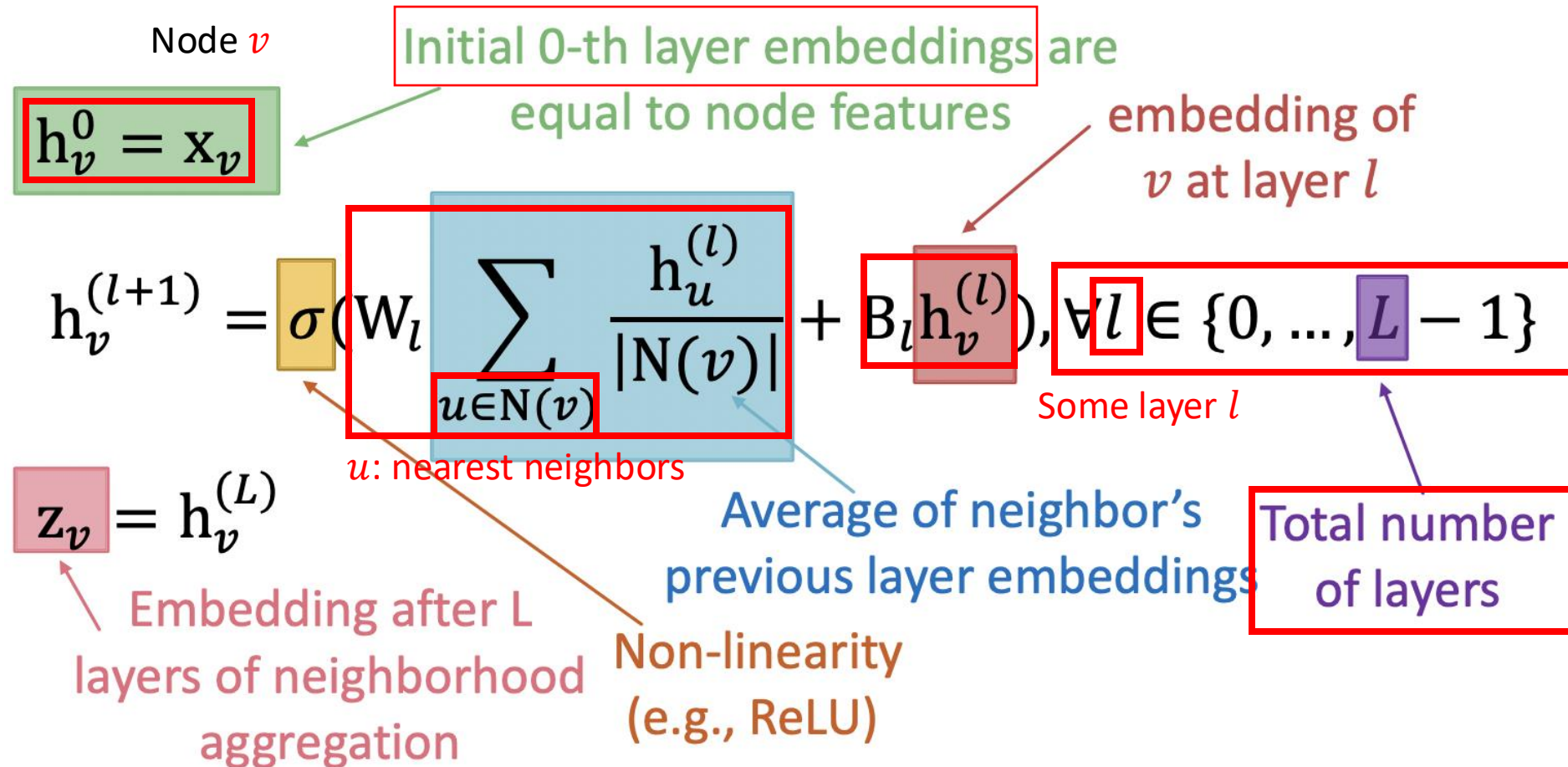


# Graph networks: aggregate neighbors

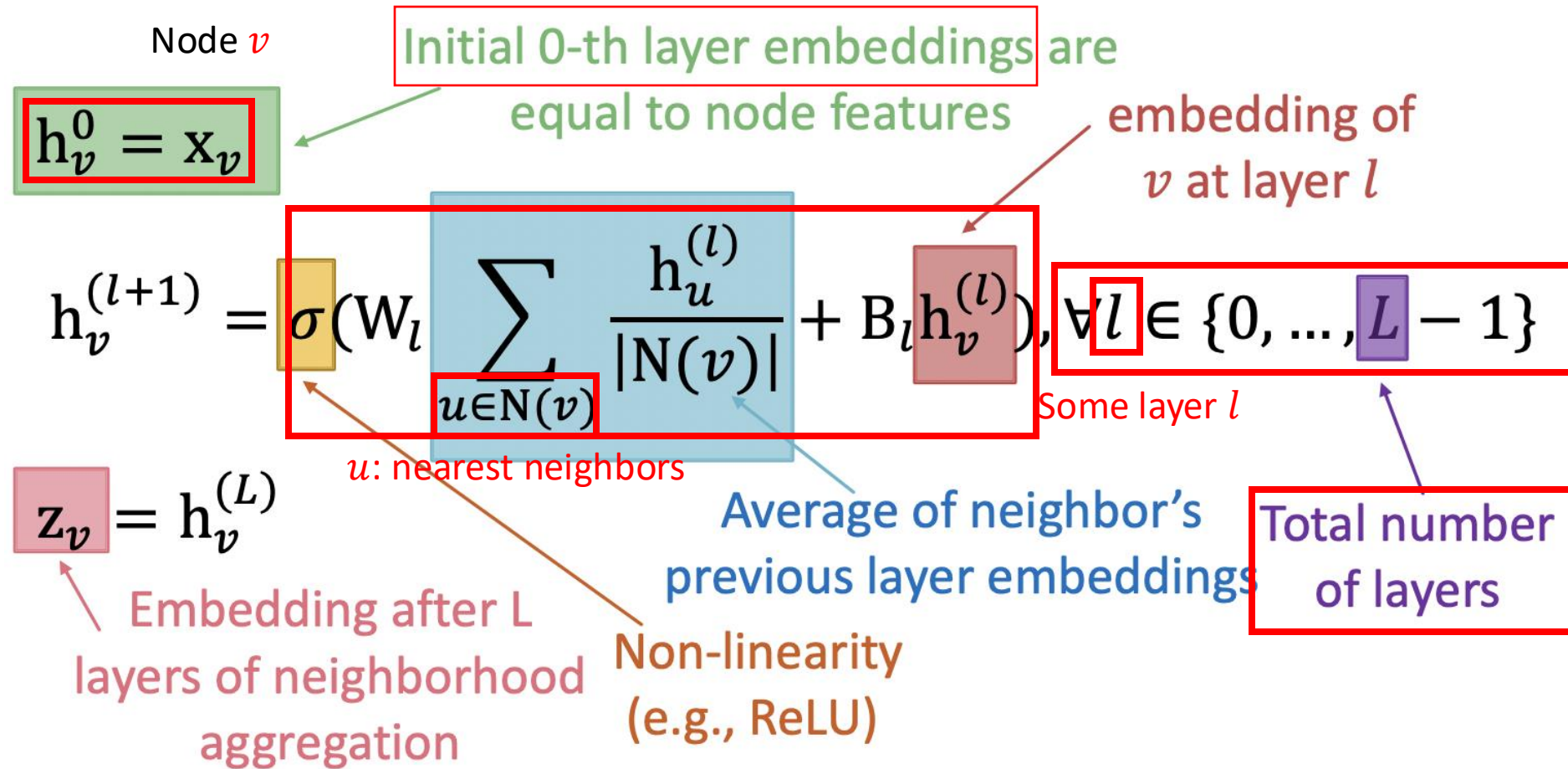




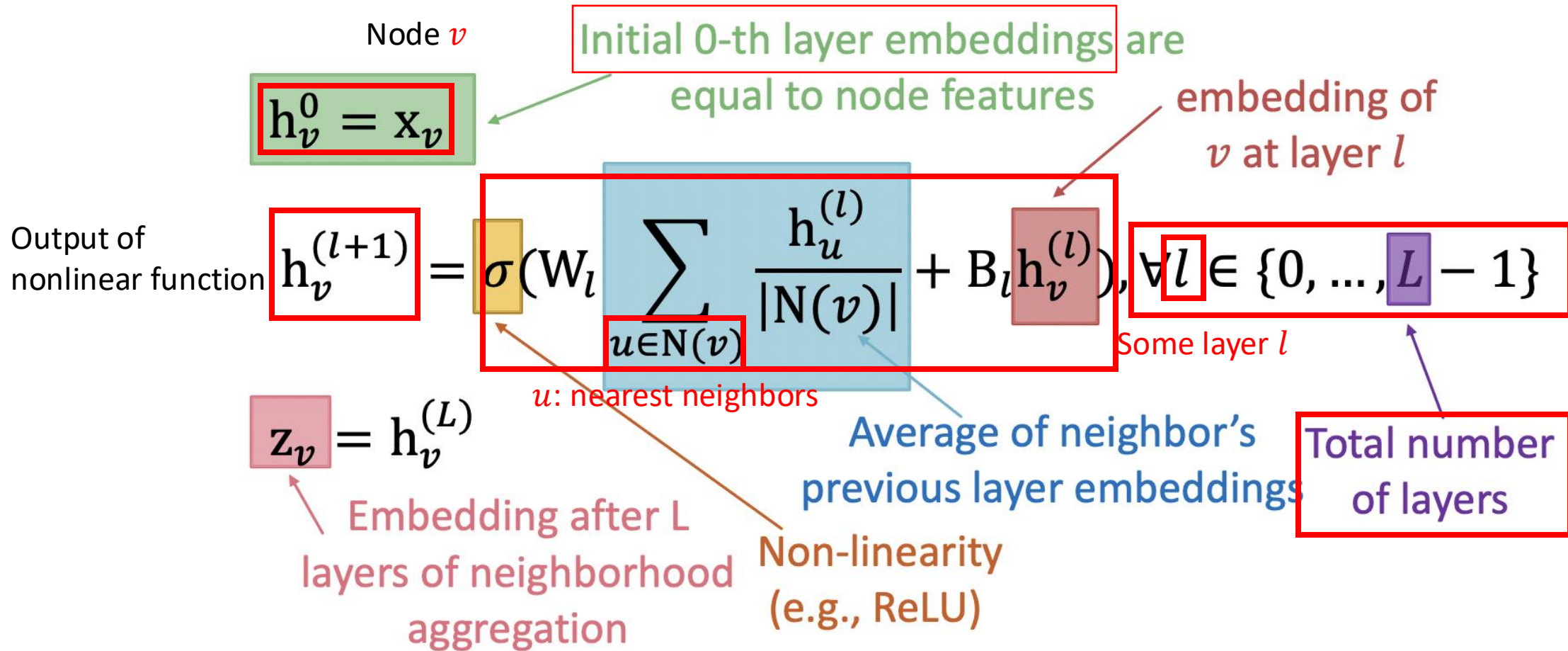
# Graph networks: aggregate neighbors



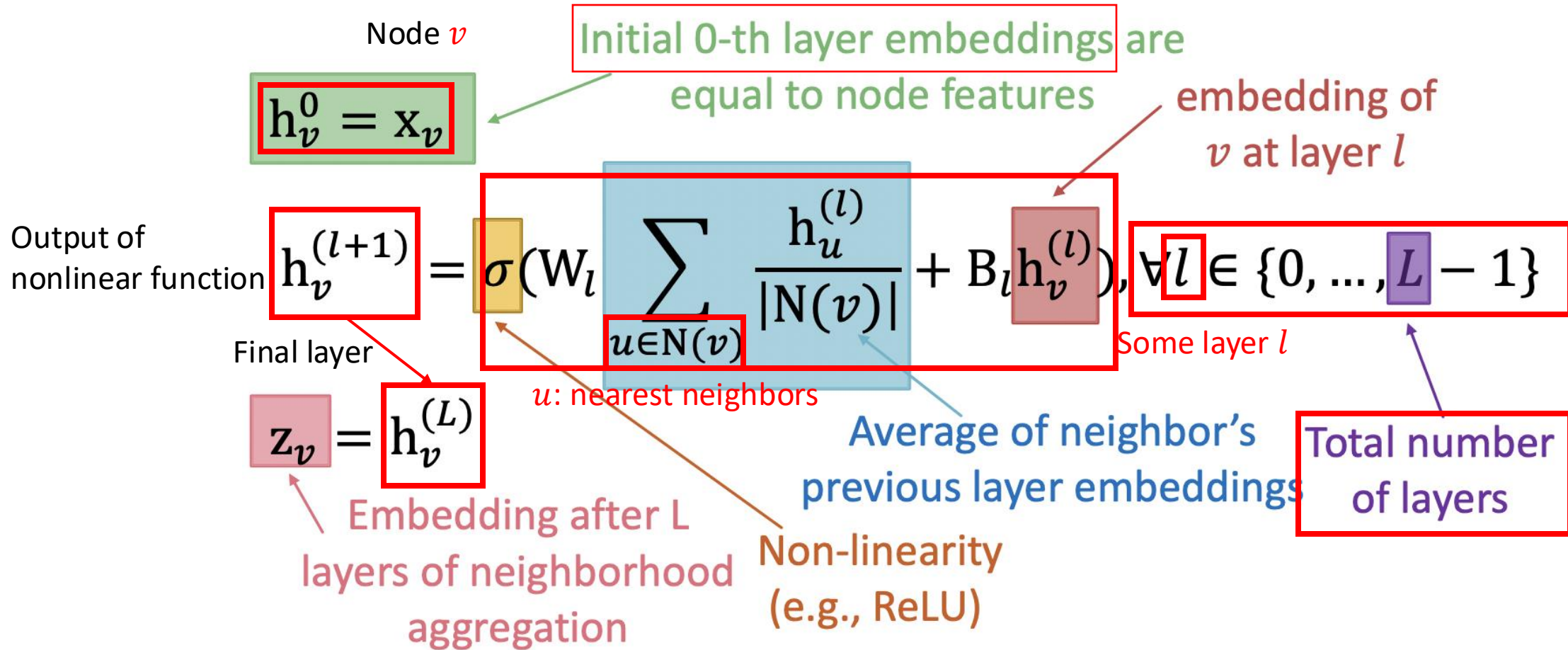
# Graph networks: aggregate neighbors



# Graph networks: aggregate neighbors



# Graph networks: aggregate neighbors





# Graph networks: aggregate neighbors

Initial 0-th layer embeddings are equal to node features

$$h_v^0 = x_v$$

embedding of  $v$  at layer  $l$

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

GraphSAGE

$$h_v^{(l+1)} = \sigma \left( [W_l \cdot \text{AGG} \left( \{h_u^{(l)}, \forall u \in N(v)\} \right), B_l h_v^{(l)}] \right)$$

total number of layers

layers of neighborhood aggregation

Non-linearity (e.g., ReLU)

# Graph networks: aggregate neighbors

Initial 0-th layer embeddings are equal to node features

$$h_v^0 = x_v$$

embedding of  $v$  at layer  $l$

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

GraphSAGE

$$h_v^{(l+1)} = \sigma \left( [W_l \cdot \text{AGG} \left( \{h_u^{(l)}, \forall u \in N(v)\} \right), B_l h_v^{(l)}] \right)$$

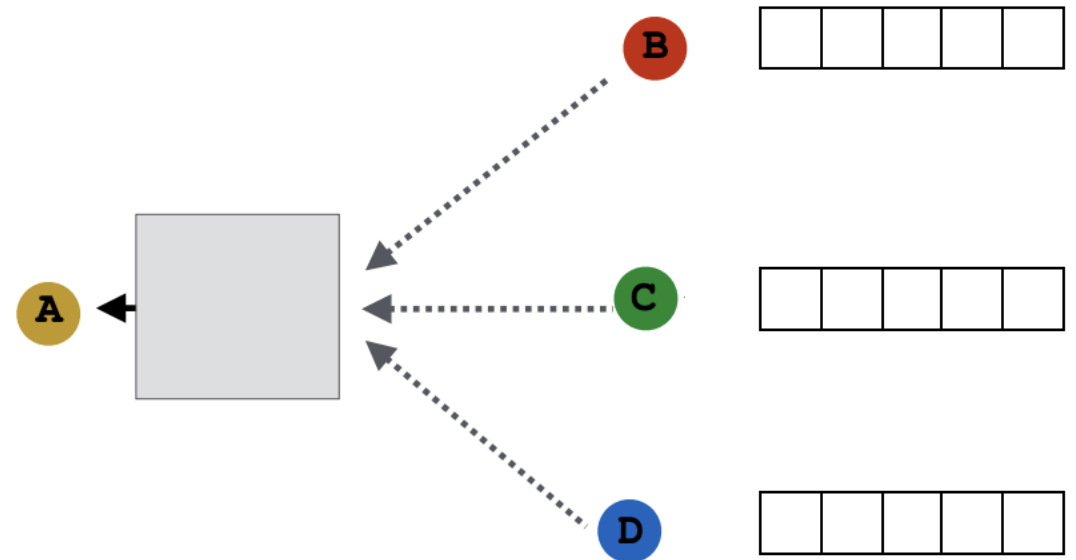
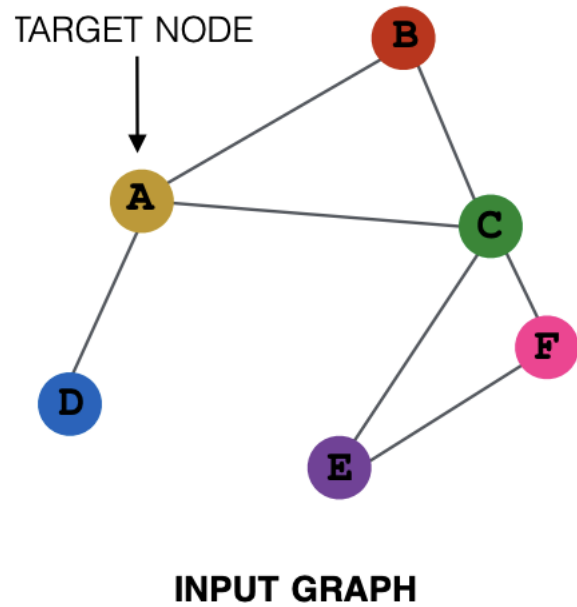
total number of layers

layers of neighborhood aggregation

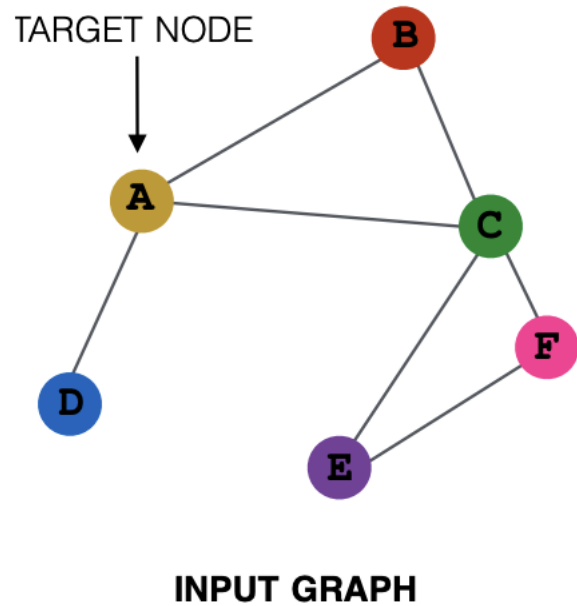
Non-linearity (e.g., ReLU)

Q: why GraphSAGE can be more general?

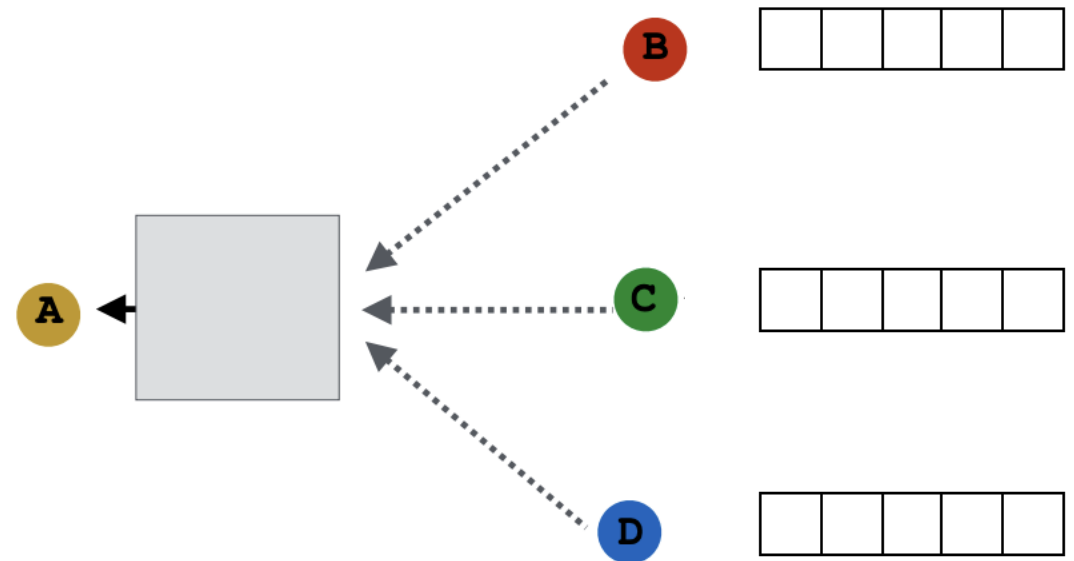
# Graph networks: aggregate neighbors



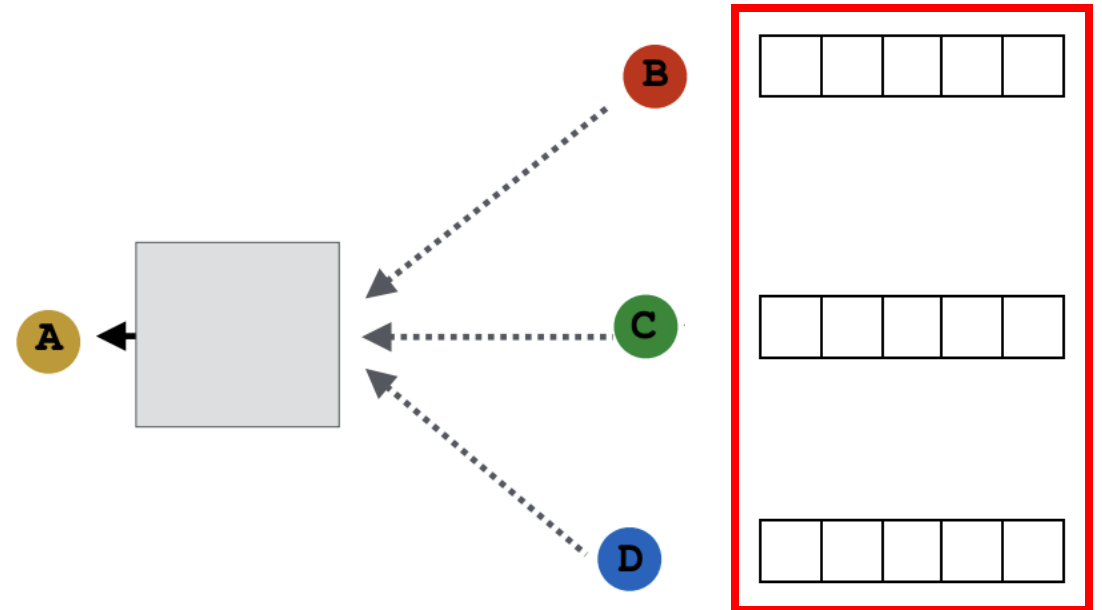
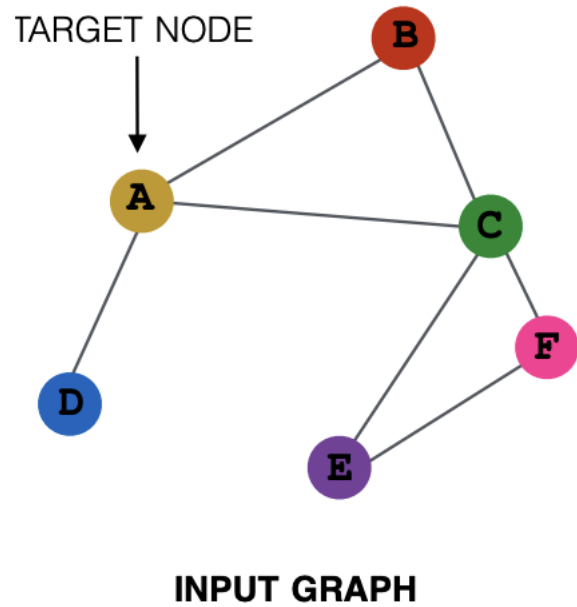
# Graph networks: aggregate neighbors



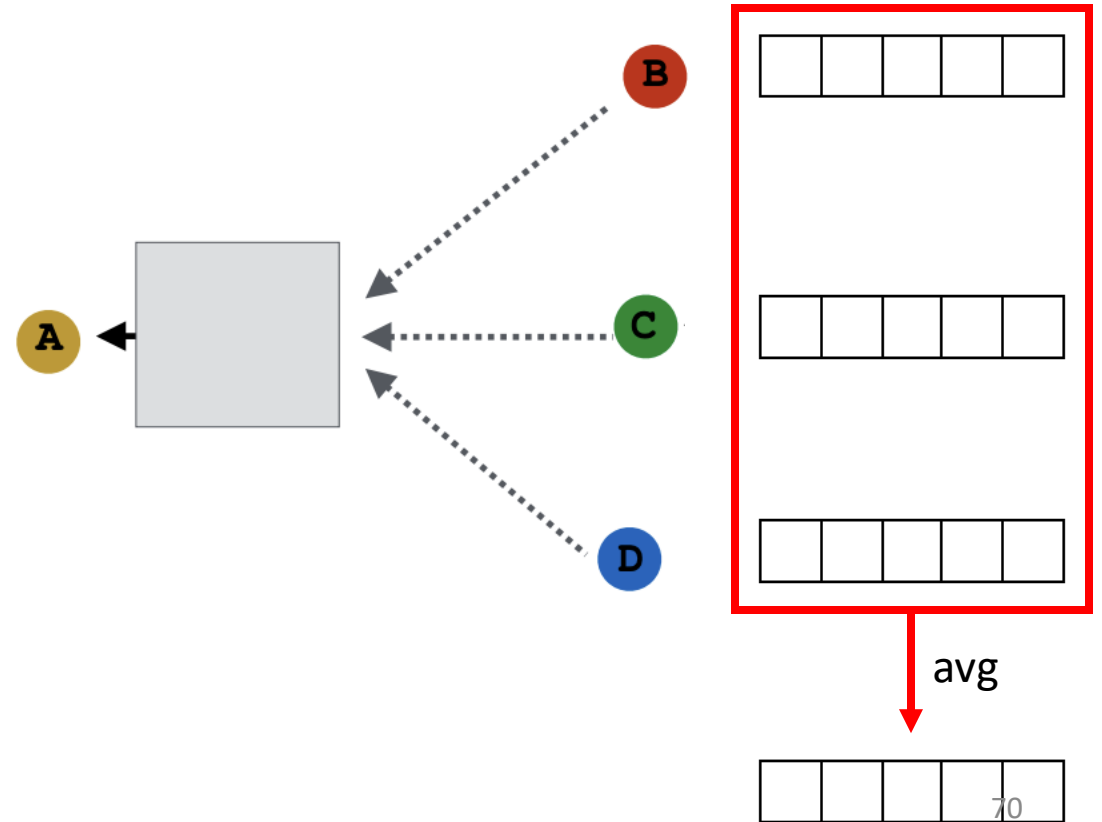
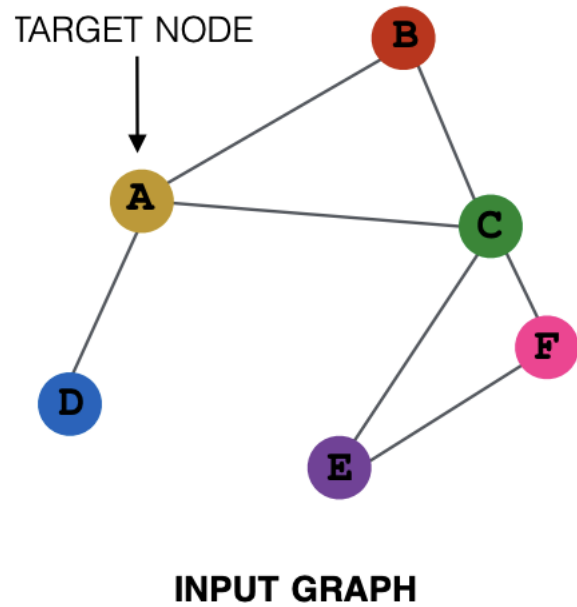
Suppose: we have 5d node features



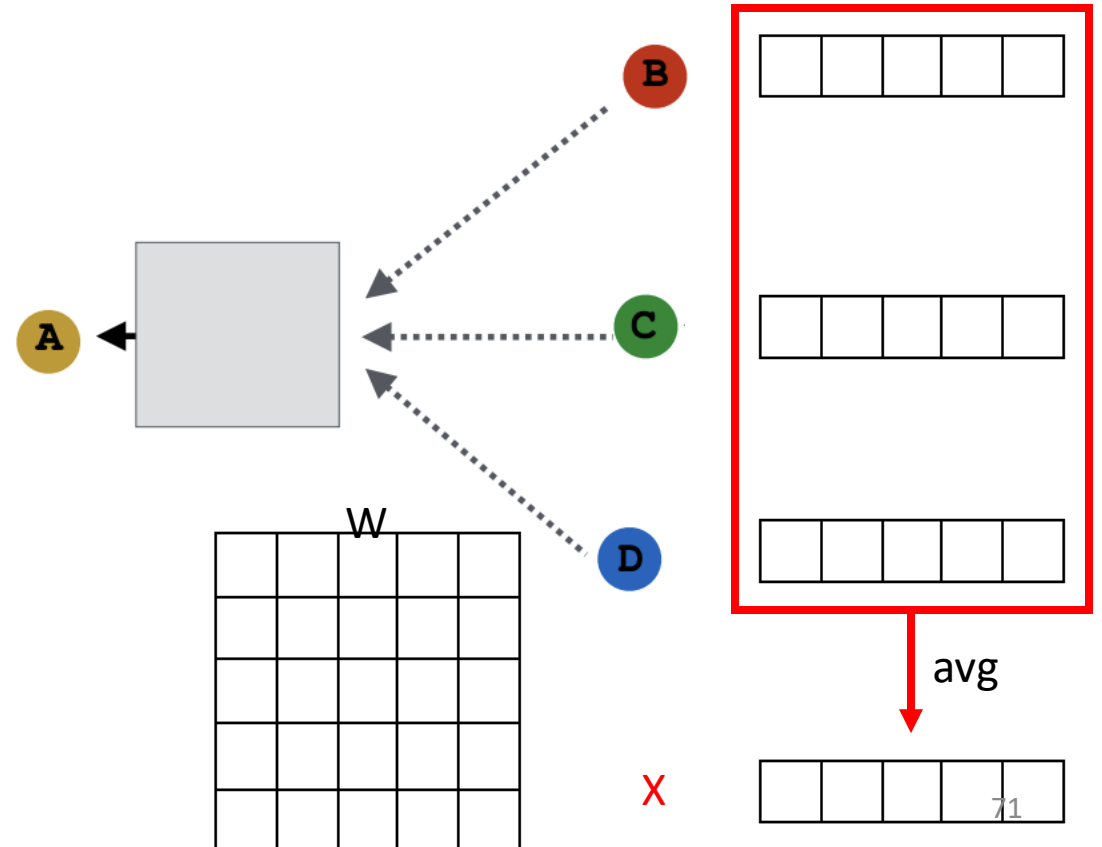
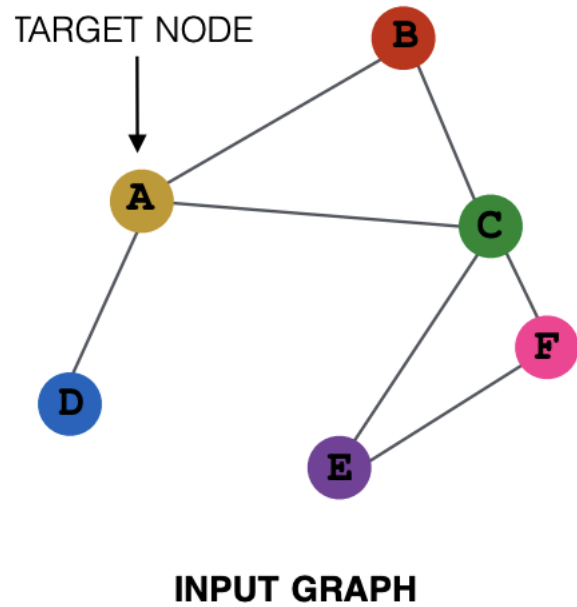
# Graph networks: aggregate neighbors



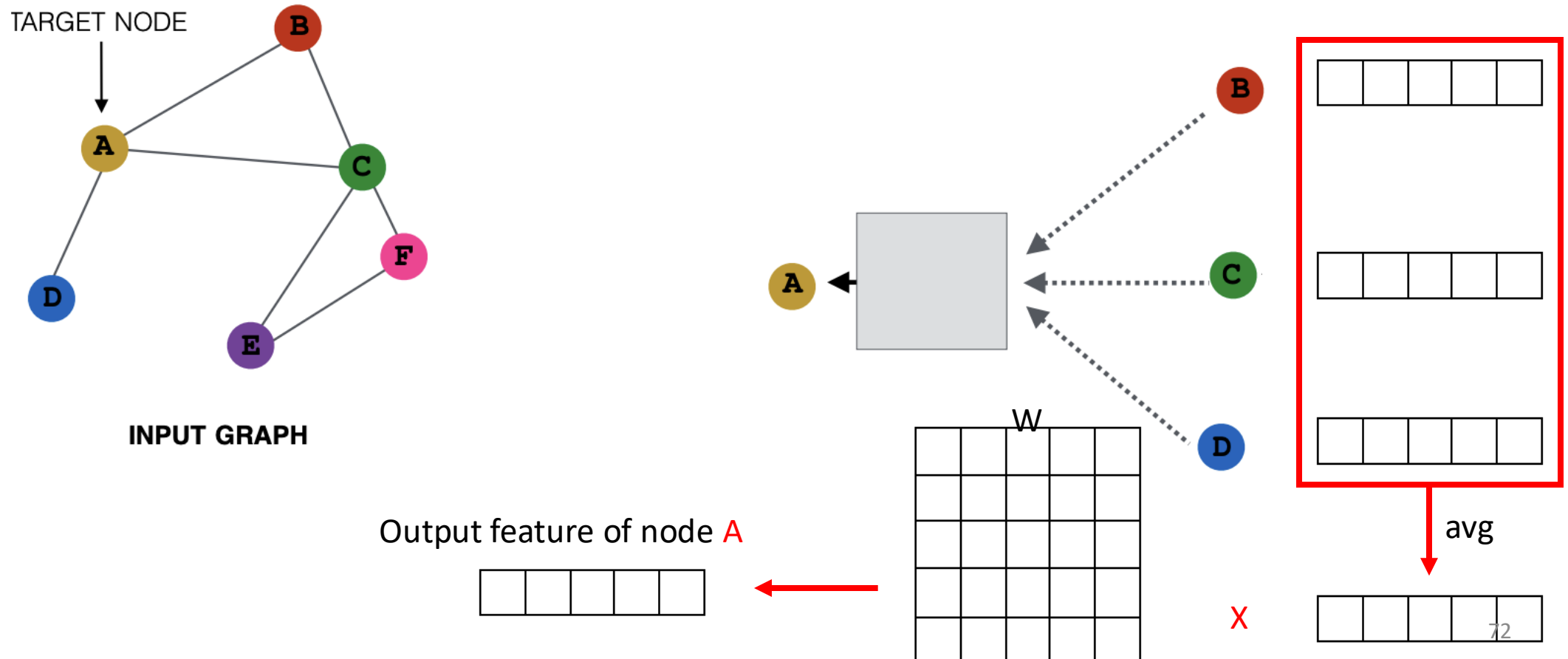
# Graph networks: aggregate neighbors



# Graph networks: aggregate neighbors

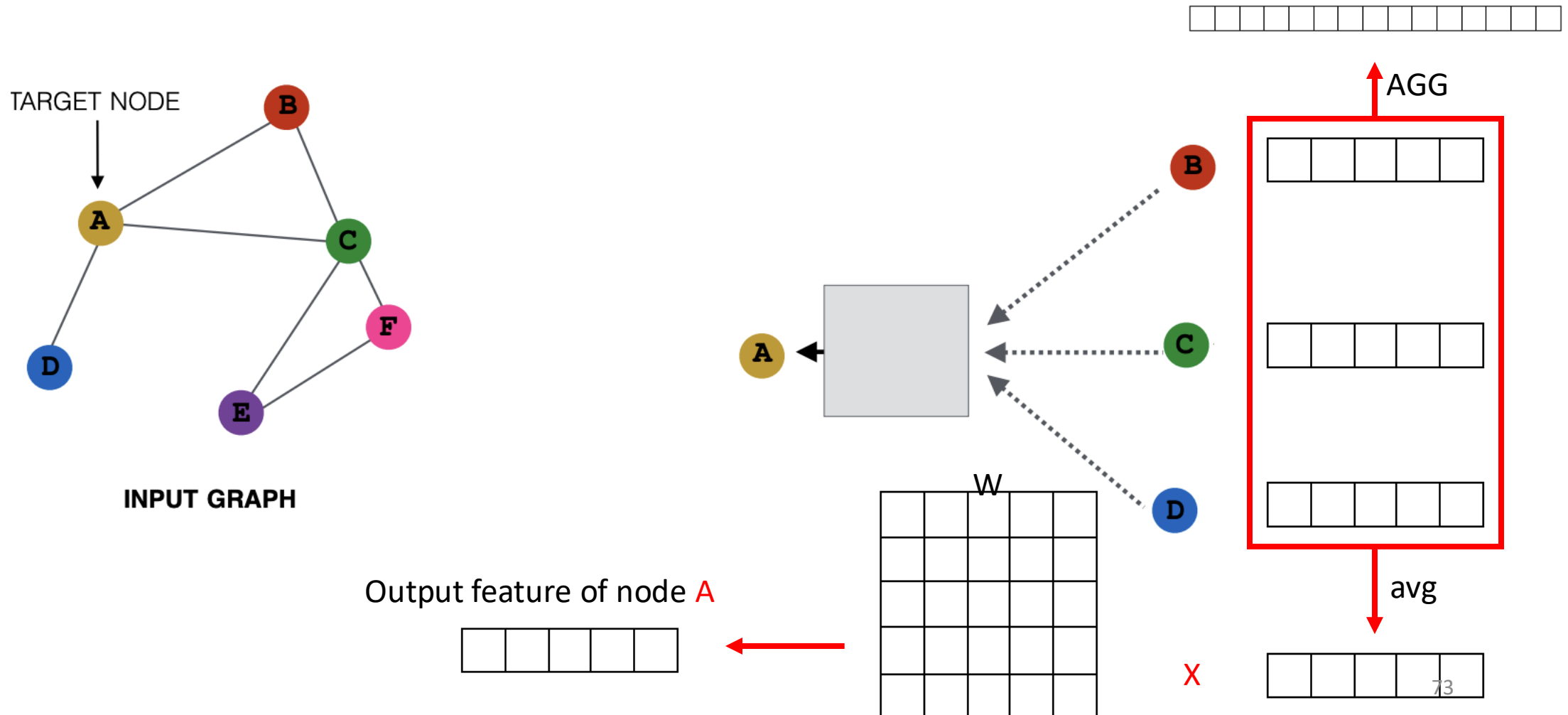


# Graph networks: aggregate neighbors

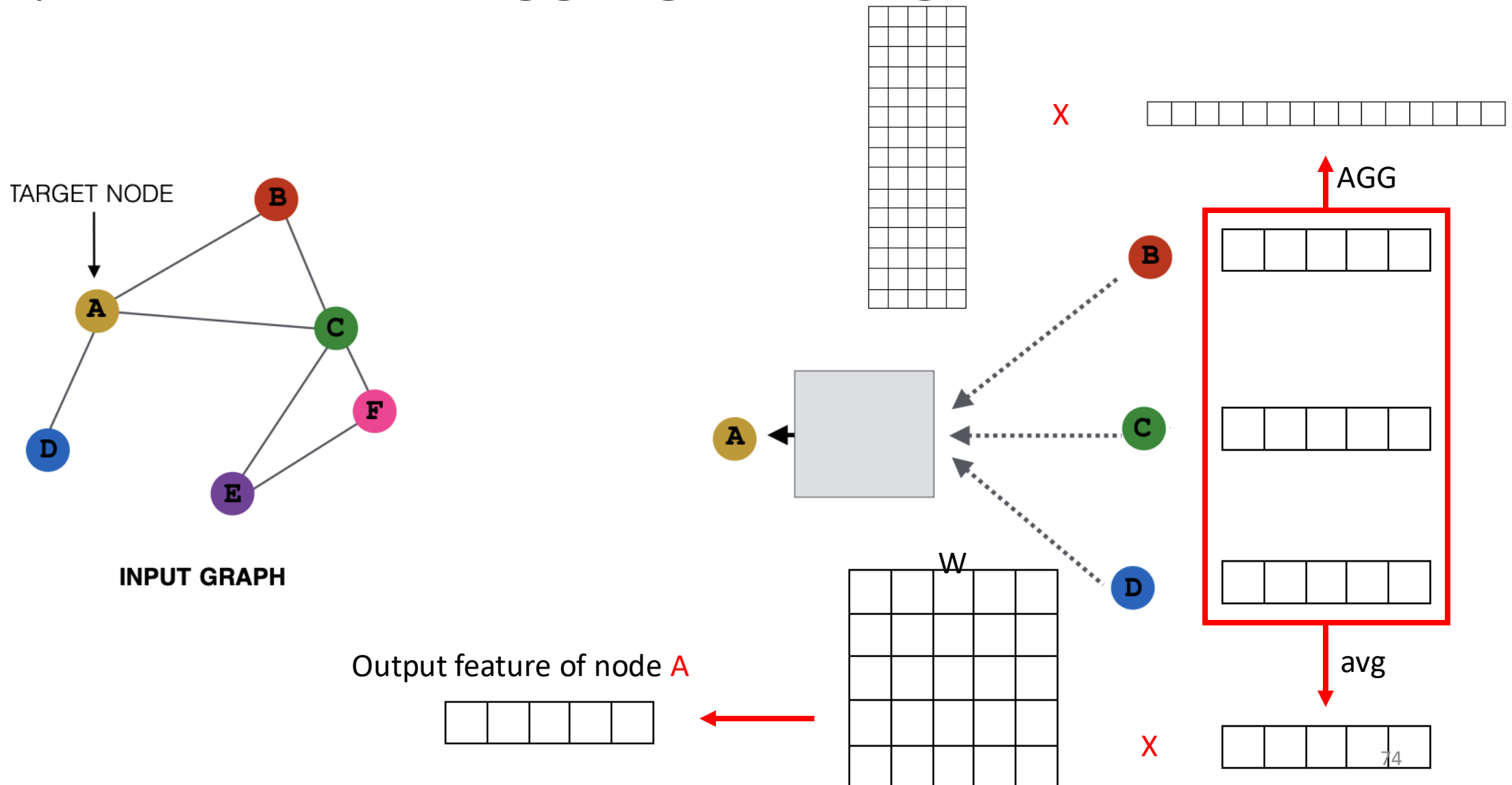




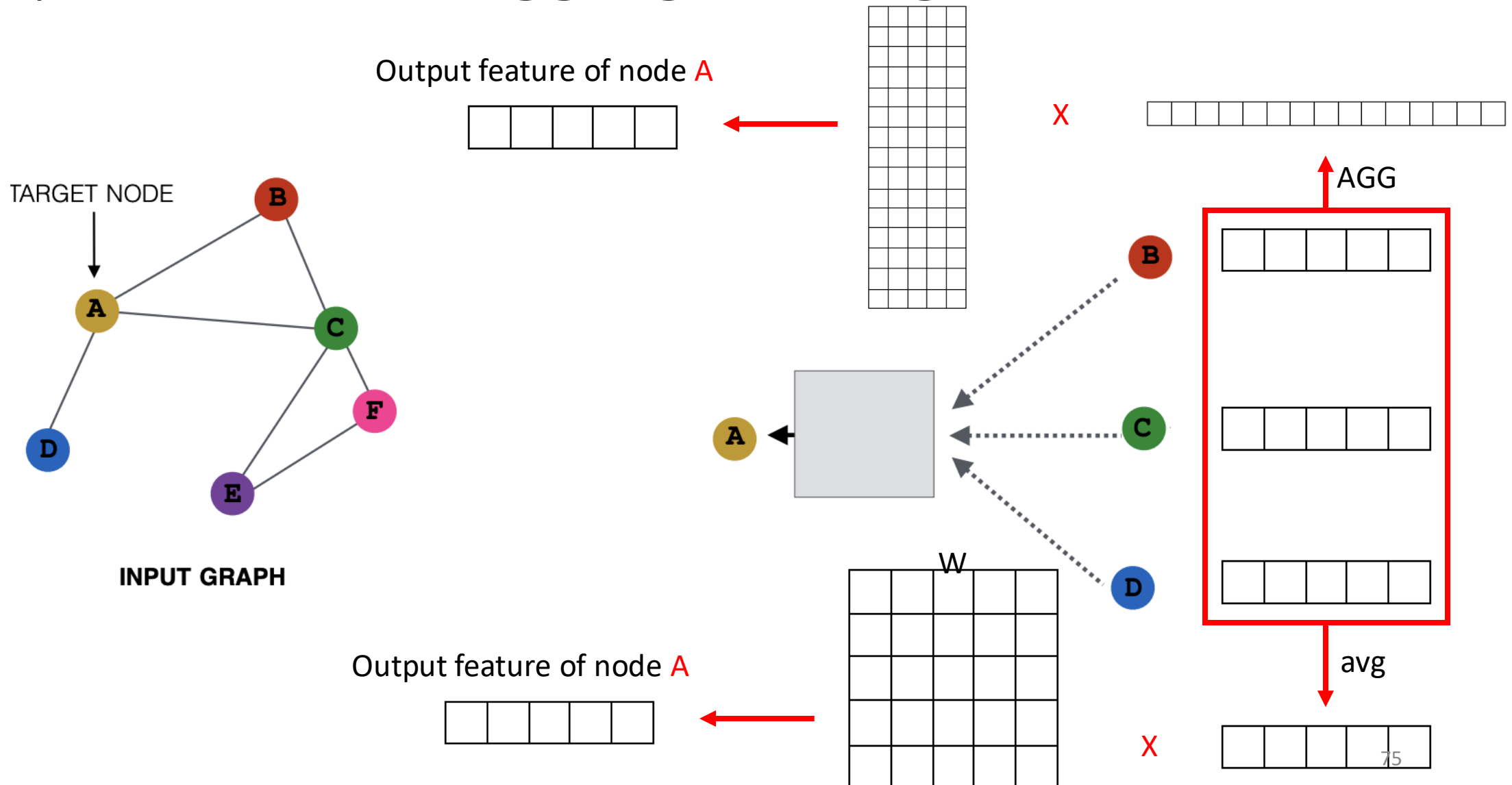
# Graph networks: aggregate neighbors



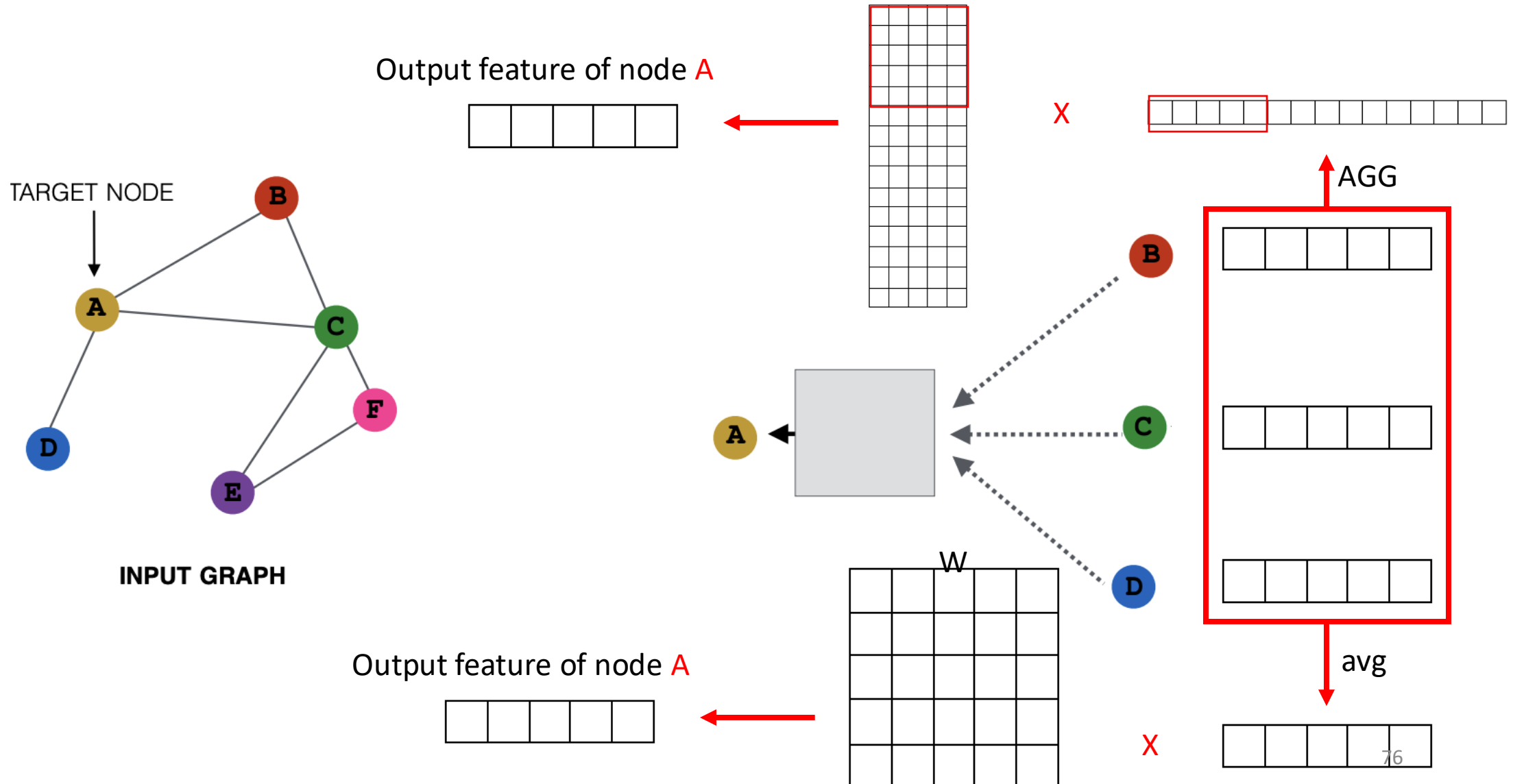
# Graph networks: aggregate neighbors



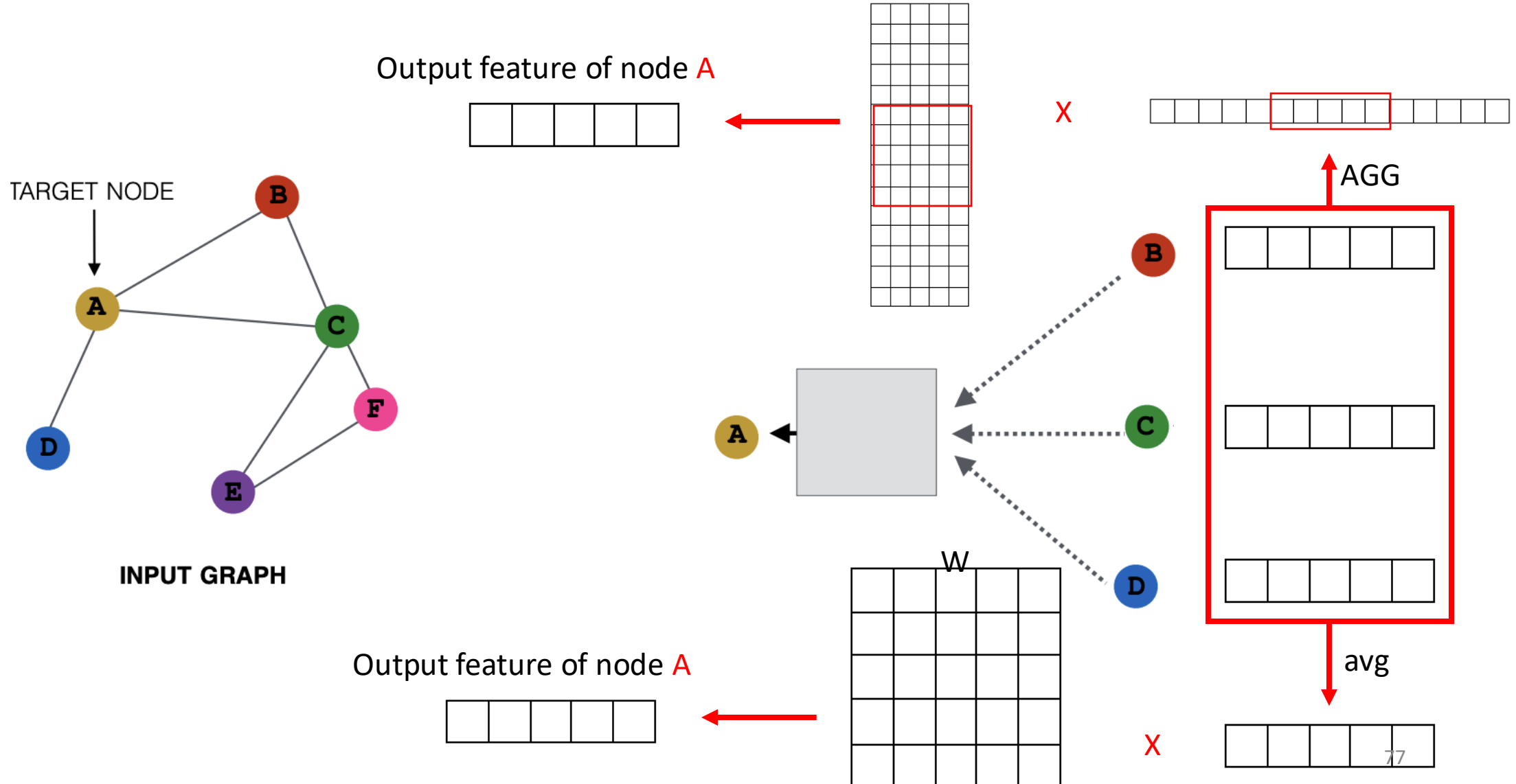
# Graph networks: aggregate neighbors



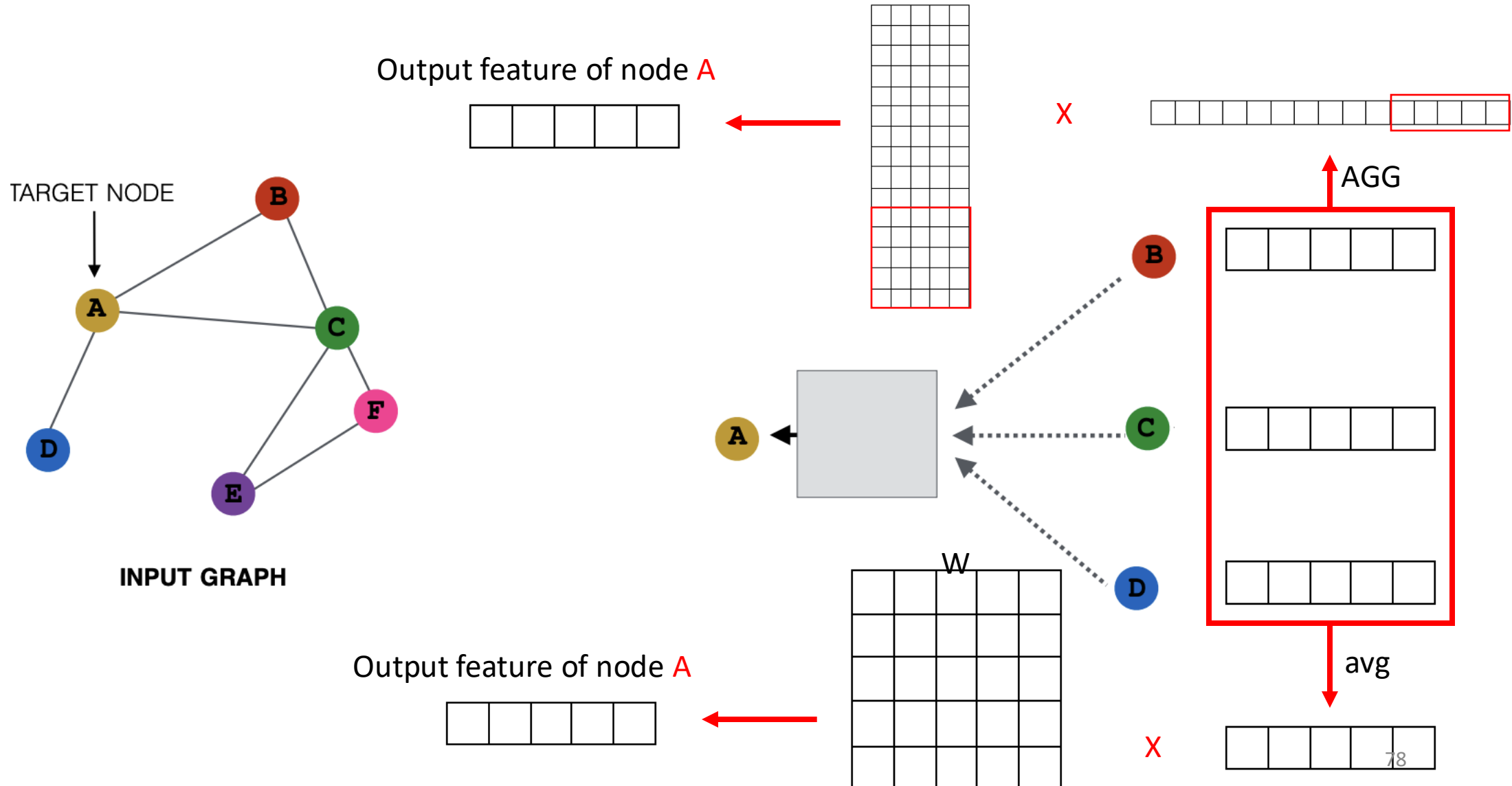
# Graph networks: aggregate neighbors



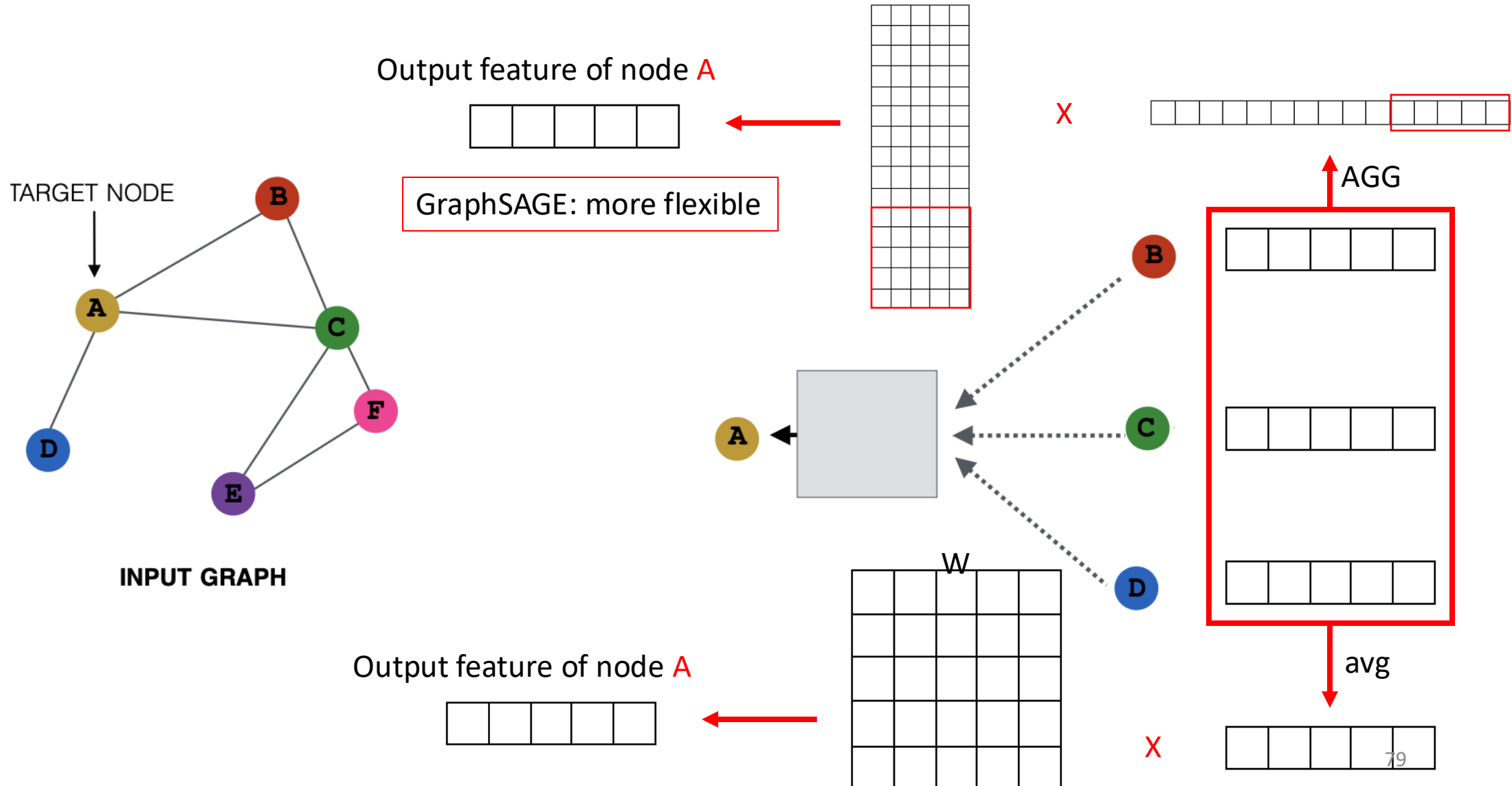
# Graph networks: aggregate neighbors



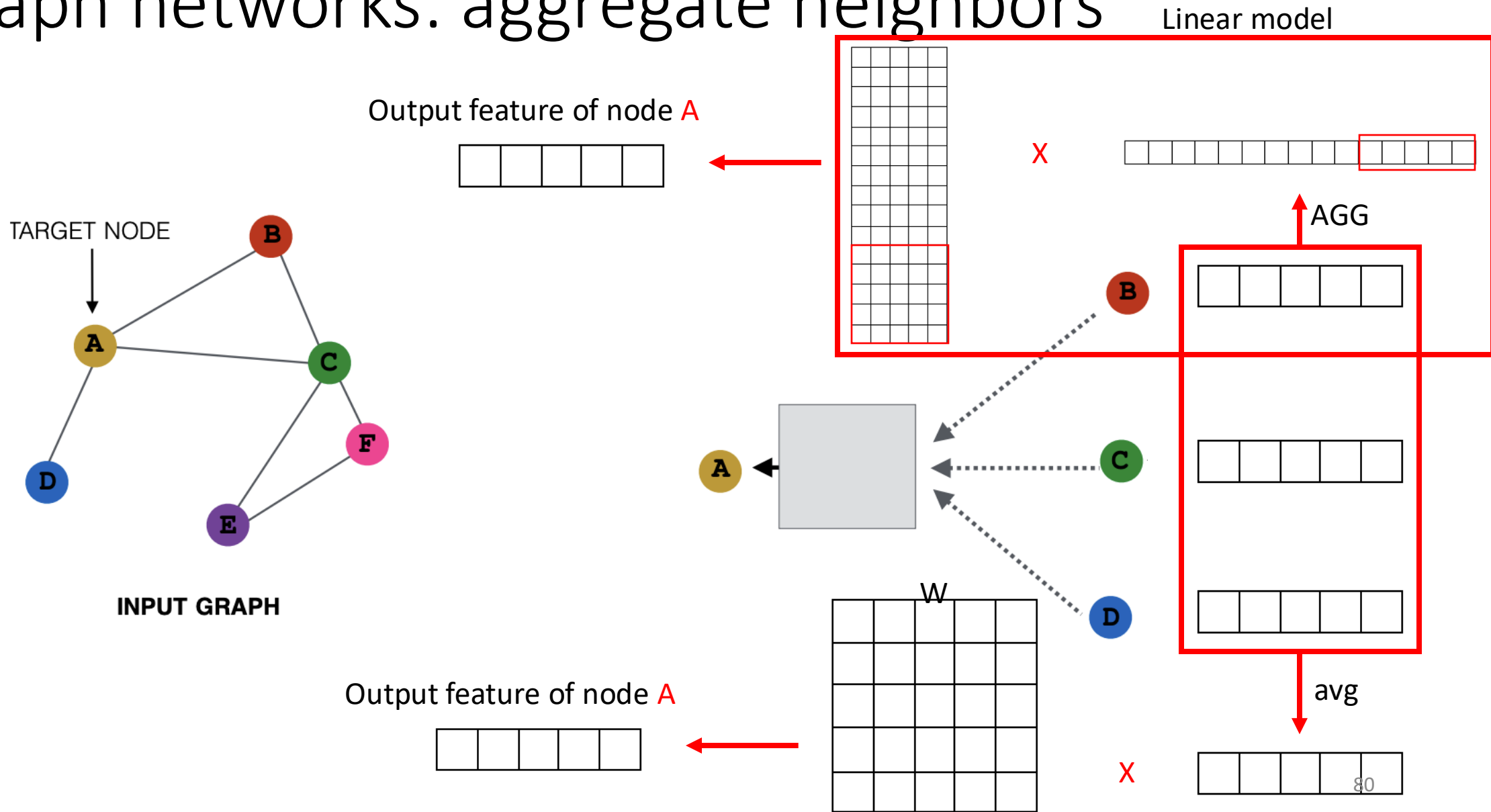
# Graph networks: aggregate neighbors



# Graph networks: aggregate neighbors

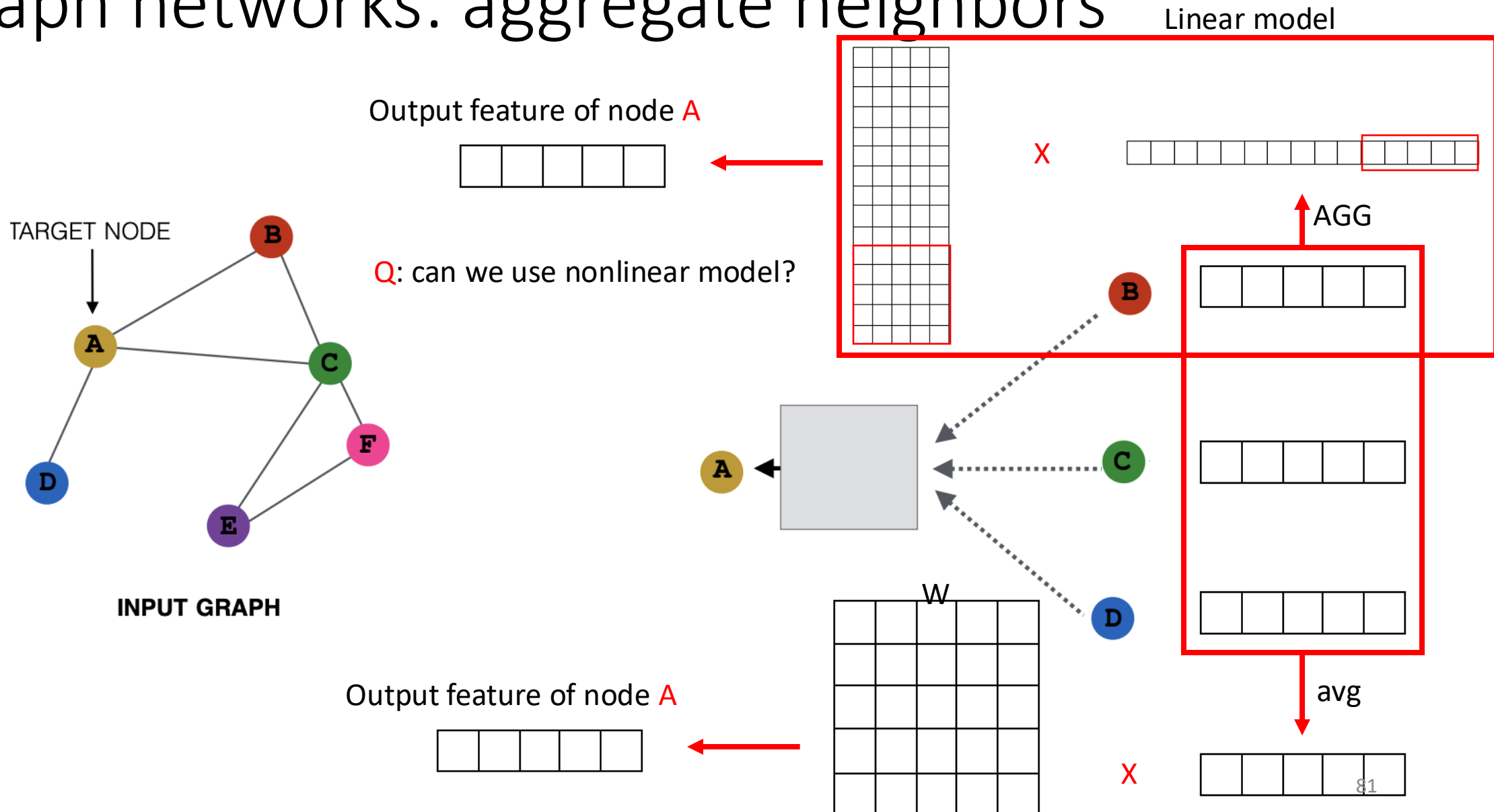


# Graph networks: aggregate neighbors

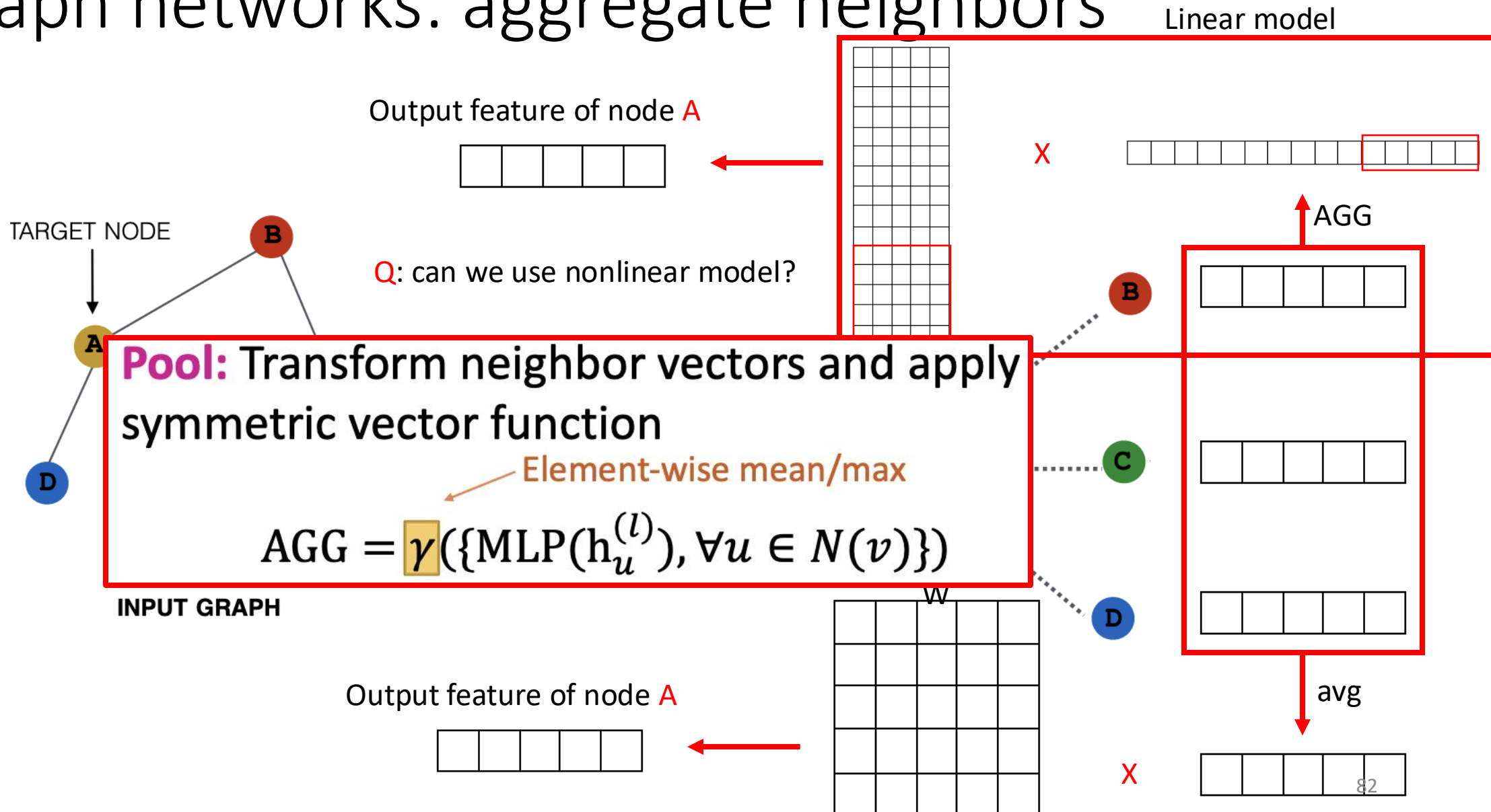




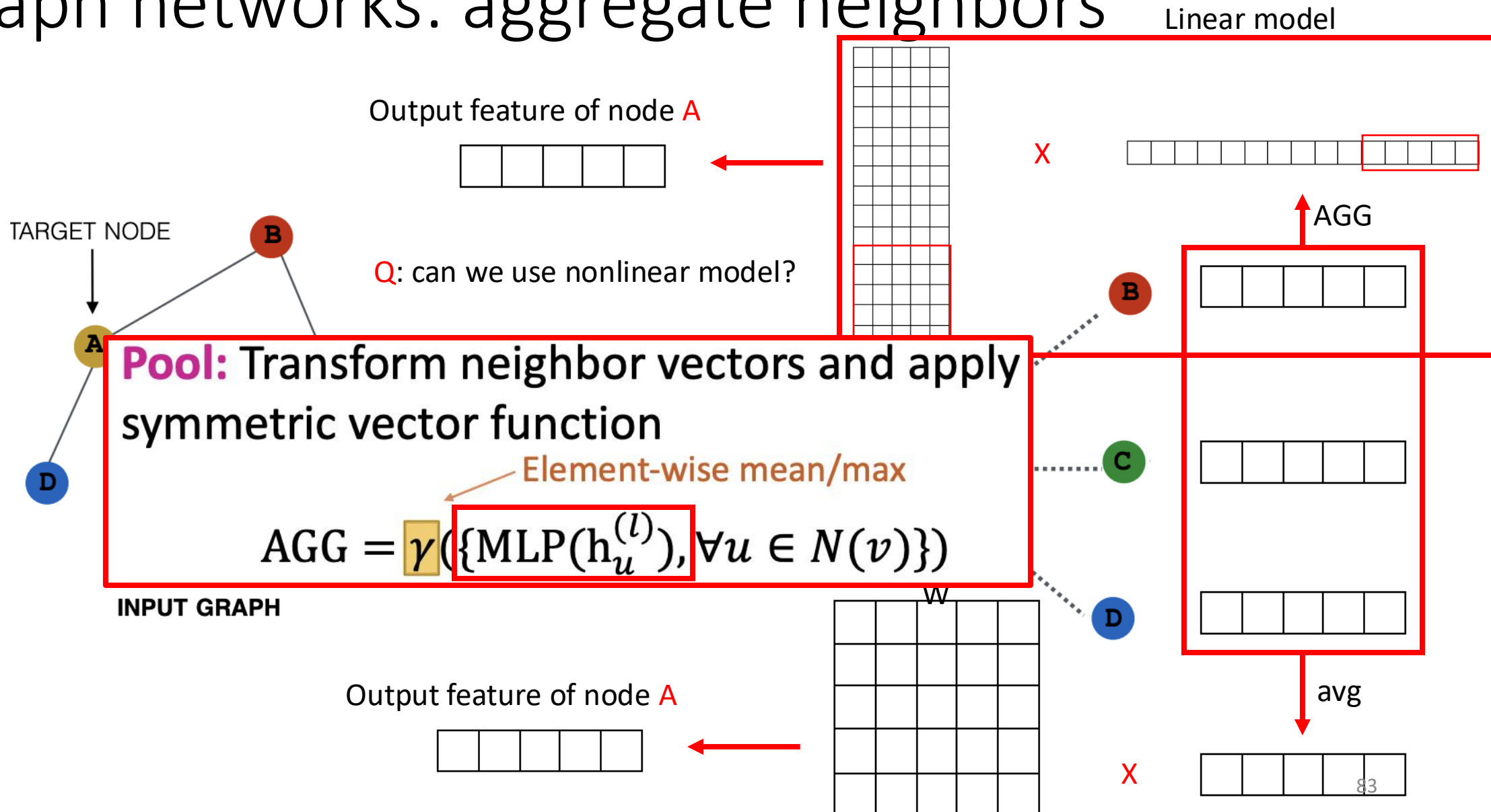
# Graph networks: aggregate neighbors



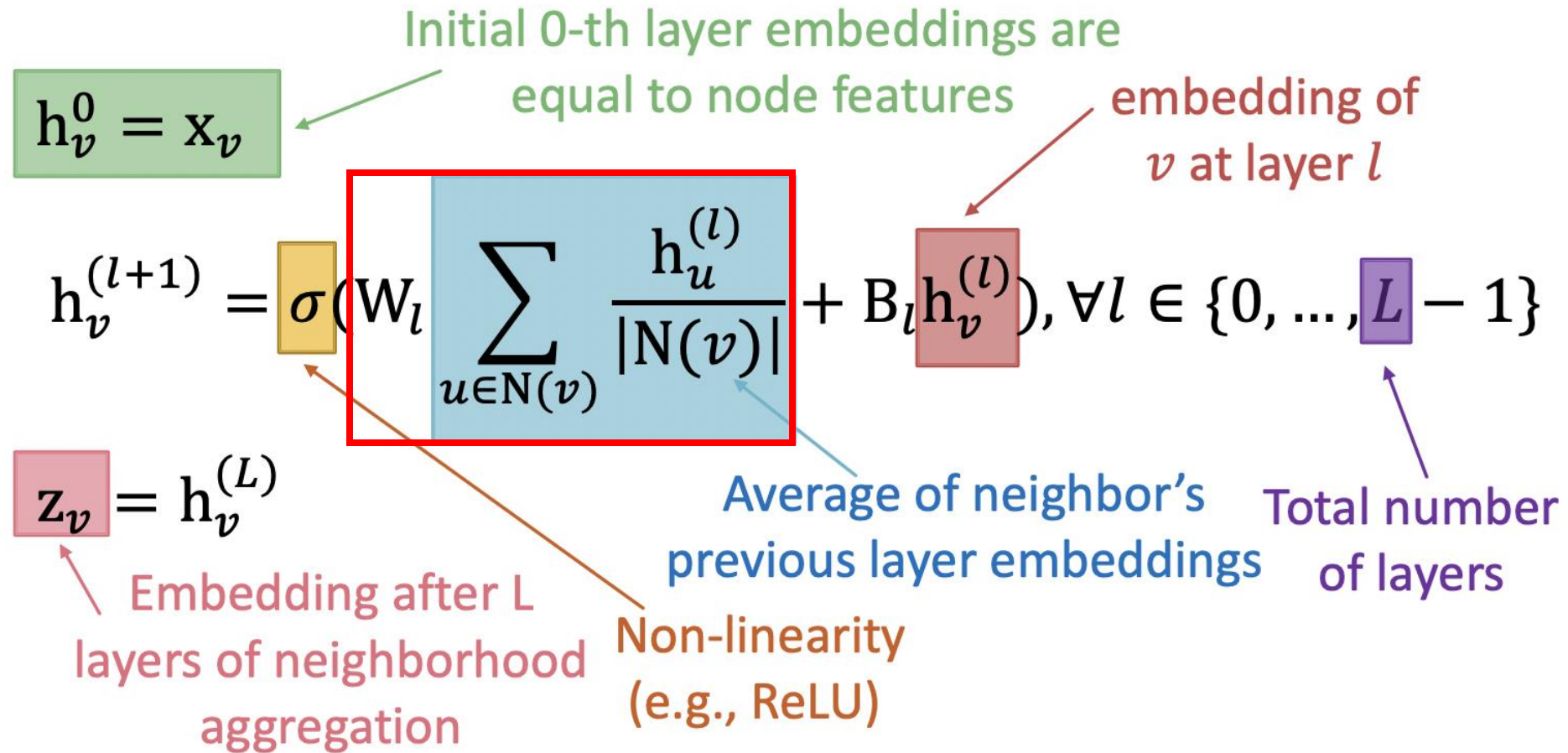
# Graph networks: aggregate neighbors



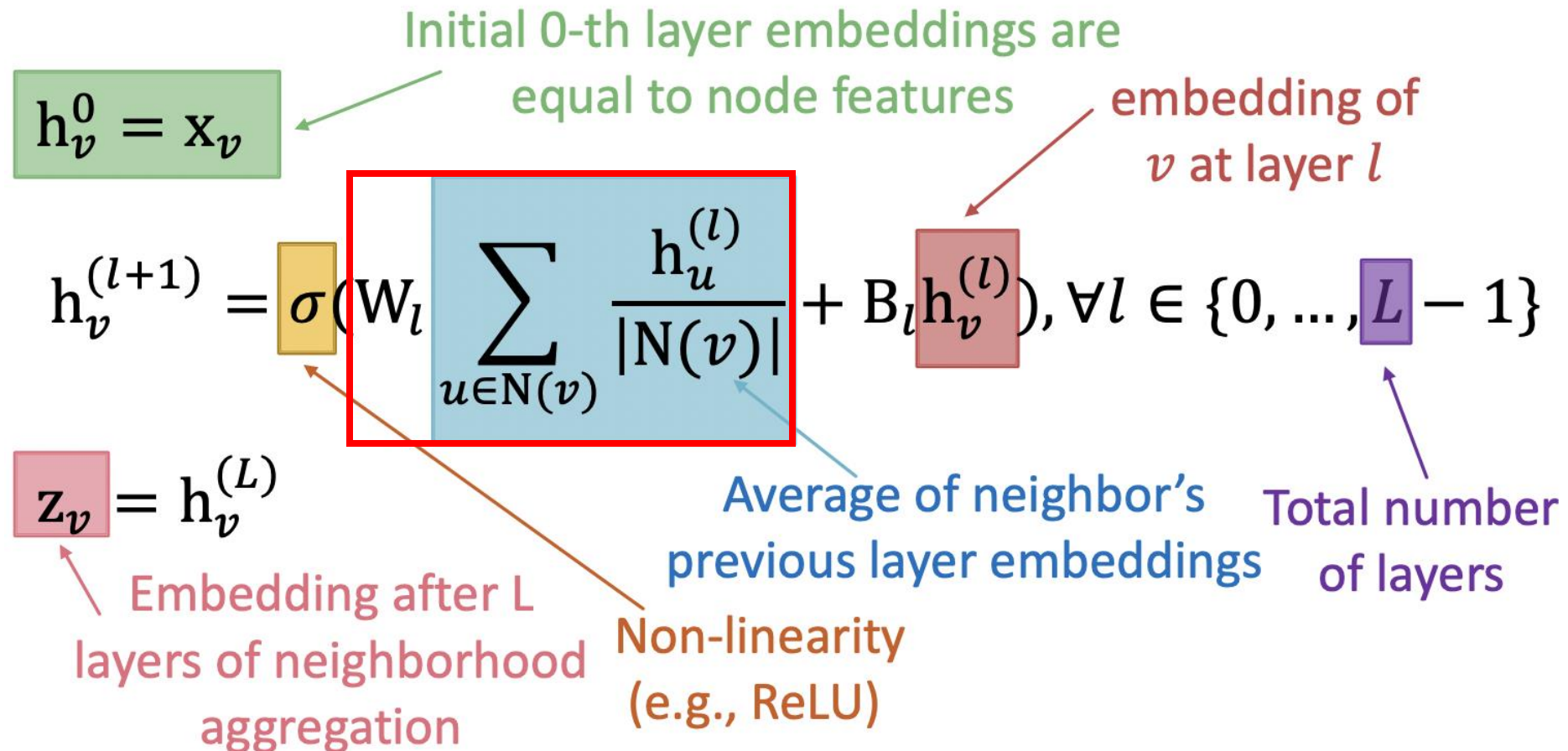
# Graph networks: aggregate neighbors



# Graph networks: aggregate neighbors

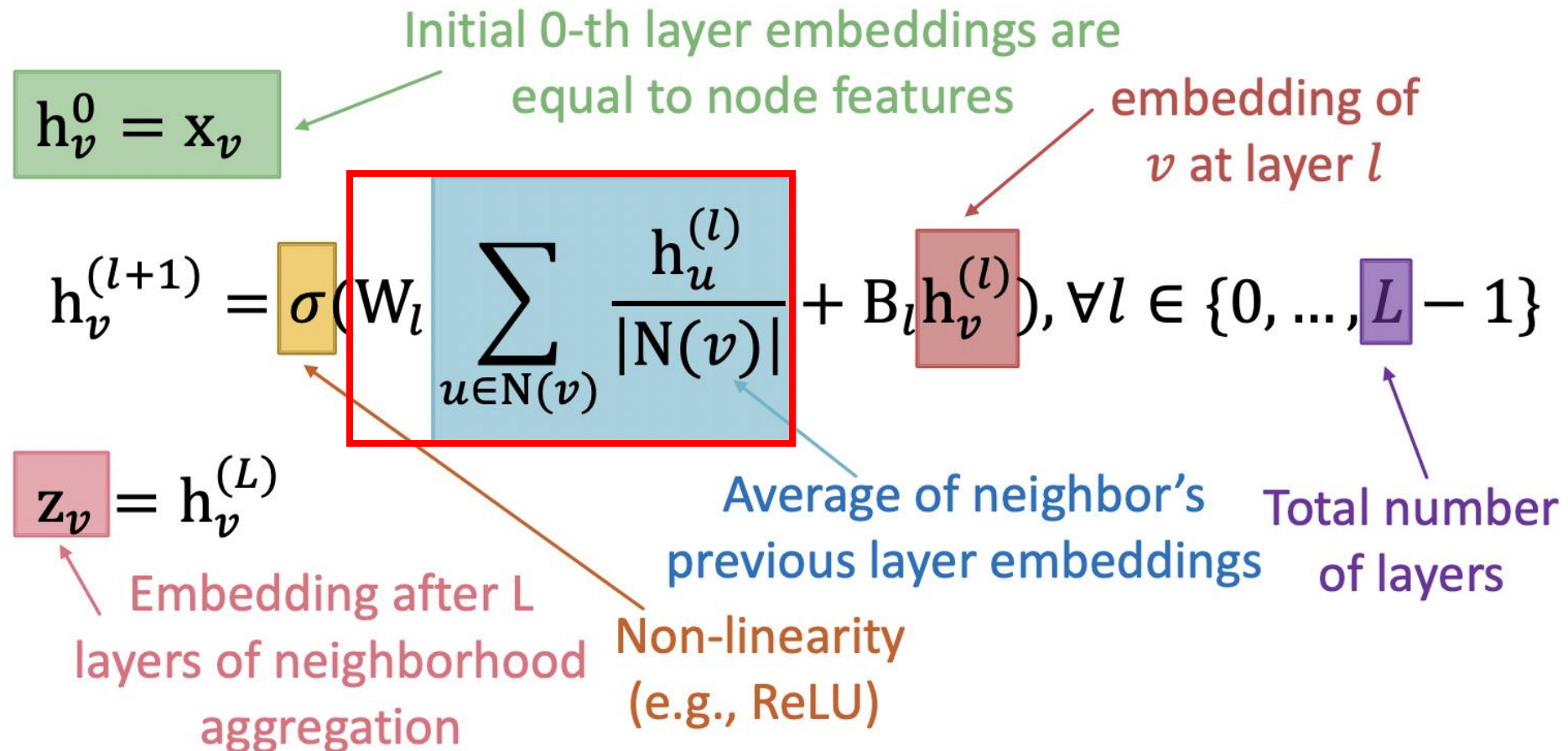


# Graph networks: aggregate neighbors



Q: can we use attention mechanism?

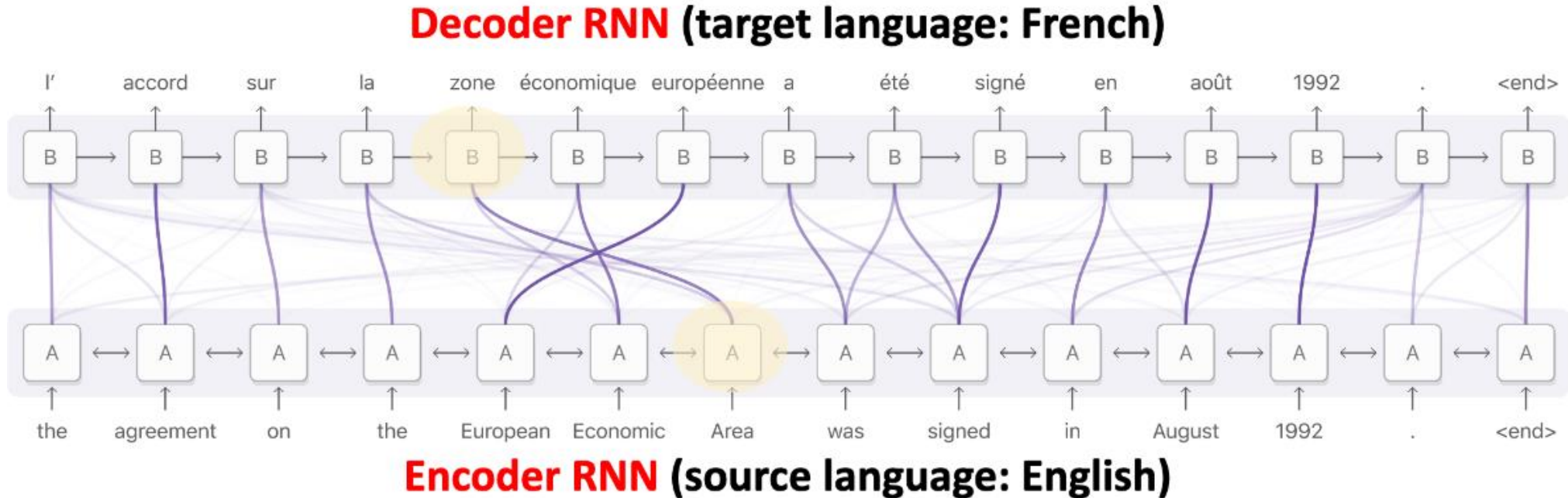
# Graph networks: aggregate neighbors



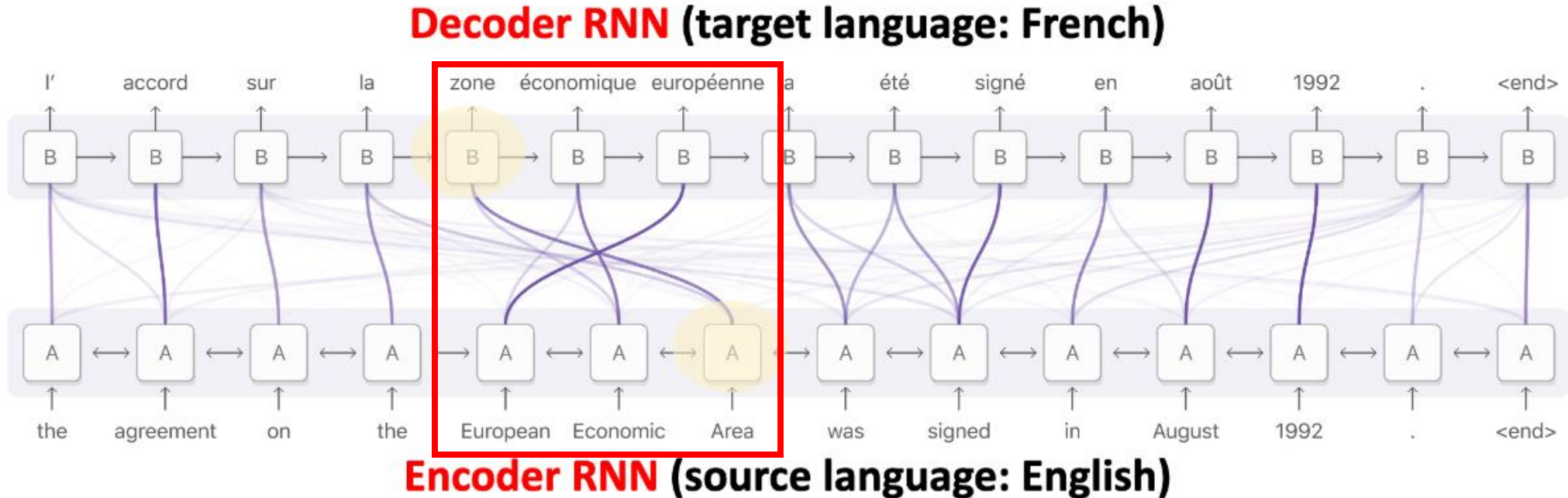
Q: can we use attention mechanism? ← Not all node's neighbors are equally important



# Input-output correlation

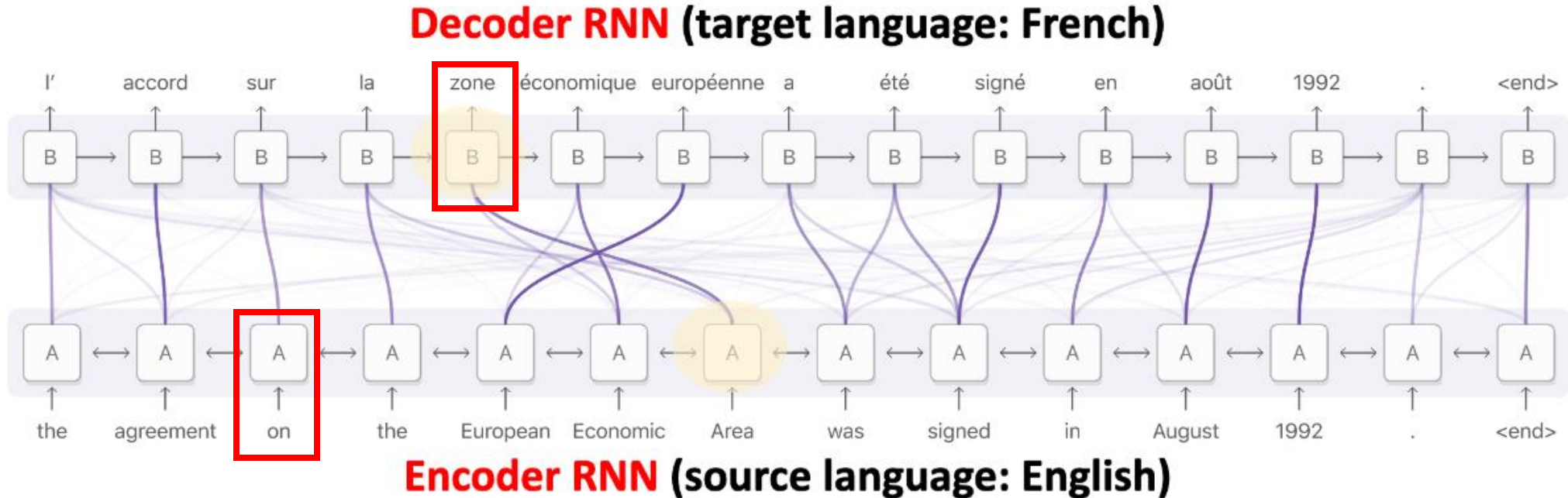


# Input-output correlation





# Input-output correlation



# Graph attention networks

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

Attention weights

# Graph attention networks

Q: what stand for importance of a node?

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

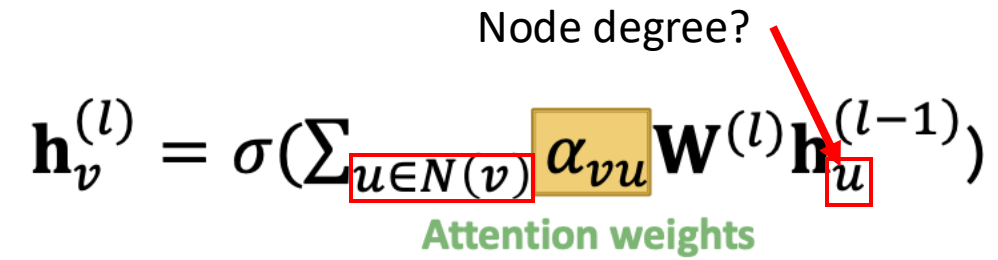
Attention weights

# Graph attention networks

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

Node degree?

Attention weights

The diagram shows the equation for the l-th layer of a graph attention network. The summation index u in N(v) is enclosed in a red box. The weight alpha\_vu is highlighted with a yellow box. The weight matrix W^(l) is in black. The input node index u in h\_u^(l-1) is also enclosed in a red box. A red arrow points from the text 'Node degree?' to the u in the denominator of the attention weights. The text 'Attention weights' is written in green below the alpha\_vu term.

# Graph attention networks

Learnable parameters?

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

Attention weights

# Graph attention networks

Learnable parameters?

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

Attention weights  
Between  $v$  and  $u$

# Graph attention networks

Learnable parameters?

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

Attention weights  
Between  $v$  and  $u$

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

# Graph attention networks

Learnable parameters?

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

Attention weights  
Between  $v$  and  $u$

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



# Graph attention networks

Learnable parameters?

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

Attention weights  
Between  $v$  and  $u$

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

Make use of last layer's output

# Graph attention networks

Learnable parameters?

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

Attention weights  
Between  $v$  and  $u$

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

Make use of last layer's output

# Graph attention networks

Learnable parameters?

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

Attention weights  
Between  $v$  and  $u$

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

Make use of last layer's output



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

# Graph attention networks

Learnable parameters?

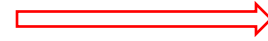
$$\mathbf{h}_v^{(l)} = \sigma\left(\sum_{u \in N(v)} \alpha_{vu} \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)}\right)$$

Attention weights  
Between  $v$  and  $u$

Softmax function

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

Make use of last layer's output



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

# Graph attention networks

Learnable parameters?

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

Attention weights  
Between  $v$  and  $u$

Softmax function  $\rightarrow$  normalized weights

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

Make use of last layer's output



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

# Graph attention networks

Learnable parameters?

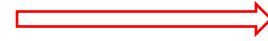
$$\mathbf{h}_v^{(l)} = \sigma\left(\sum_{u \in N(v)} \alpha_{vu} \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)}\right)$$

Attention weights  
Between  $v$  and  $u$

Softmax function  $\rightarrow$  normalized weights

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

Make use of last layer's output



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

# Graph attention networks

Learnable parameters?

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

Q: what is  $a$ 's structure?

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

Make use of last layer's output

Attention weights  
Between  $v$  and  $u$

Softmax function  $\rightarrow$  normalized weights

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

# Graph attention networks

Learnable parameters?

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

Q: what is  $a$ 's structure?

Attention weights  
Between  $v$  and  $u$

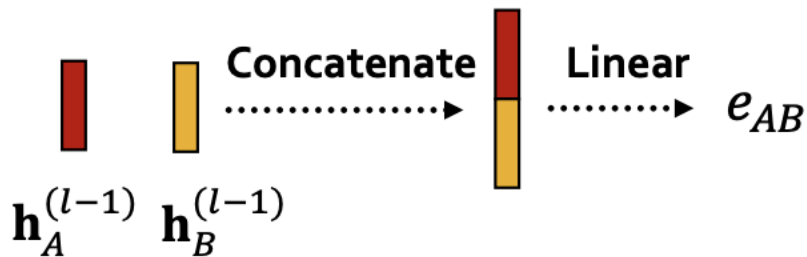
Softmax function  $\rightarrow$  normalized weights

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

Make use of last layer's output



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





# Graph attention networks

Learnable parameters?

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

Q: what is  $a$ 's structure?

Attention weights  
Between  $v$  and  $u$

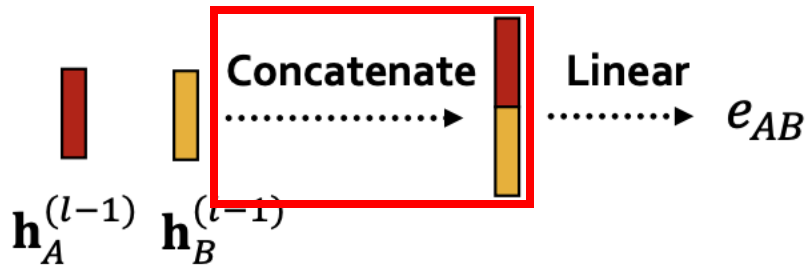
Softmax function  $\rightarrow$  normalized weights

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

Make use of last layer's output



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



# Graph attention networks

Learnable parameters?

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

Q: what is  $a$ 's structure?

Attention weights  
Between  $v$  and  $u$

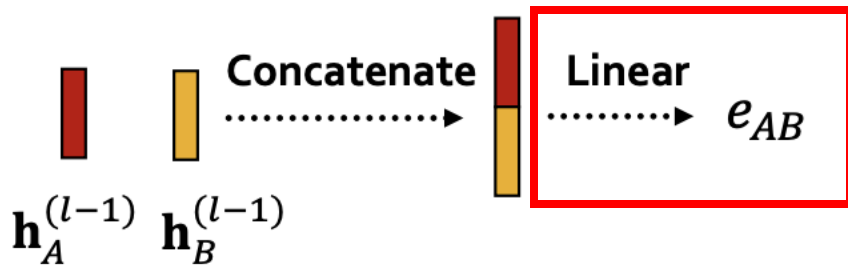
Softmax function  $\rightarrow$  normalized weights

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

Make use of last layer's output



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



# Graph attention networks

Learnable parameters?

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

Q: what is  $a$ 's structure?

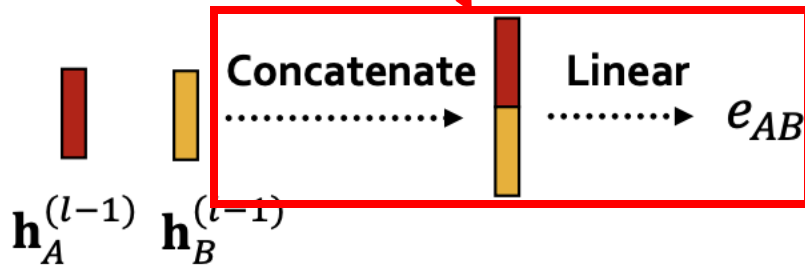
$$e_{vu} = a\left(\mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)}, \mathbf{W}^{(l)} \mathbf{h}_v^{(l-1)}\right)$$

Make use of last layer's output

Attention weights  
Between  $v$  and  $u$

Softmax function  $\rightarrow$  normalized weights

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



# Graph attention networks

Learnable parameters?

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

Attention weights  
Between  $v$  and  $u$

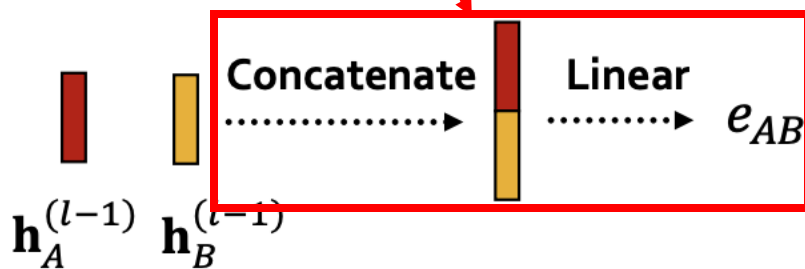
Softmax function  $\rightarrow$  normalized weights

Q: what is  $a$ 's structure?

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

Make use of last layer's output

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



$$\begin{aligned} e_{AB} &= a\left(\mathbf{W}^{(l)} \mathbf{h}_A^{(l-1)}, \mathbf{W}^{(l)} \mathbf{h}_B^{(l-1)}\right) \\ &= \text{Linear}\left(\text{Concat}\left(\mathbf{W}^{(l)} \mathbf{h}_A^{(l-1)}, \mathbf{W}^{(l)} \mathbf{h}_B^{(l-1)}\right)\right) \end{aligned}$$

# Graph attention networks

Learnable parameters?

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

Attention weights  
Between  $v$  and  $u$

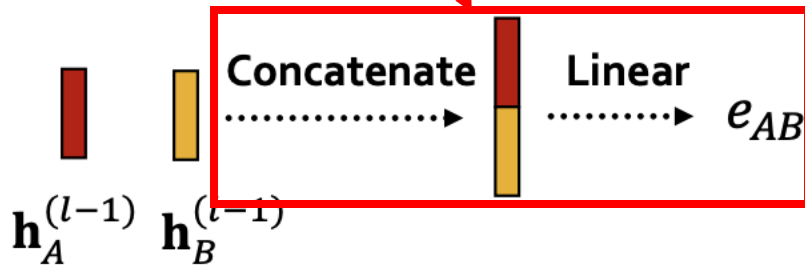
Softmax function  $\rightarrow$  normalized weights

Q: what is  $a$ 's structure?

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

Make use of last layer's output

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



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

$$= \text{Linear}\left(\text{Concat}\left(\mathbf{W}^{(l)} \mathbf{h}_A^{(l-1)}, \mathbf{W}^{(l)} \mathbf{h}_B^{(l-1)}\right)\right)$$

# Graph attention networks

Learnable parameters?

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

Attention weights  
Between  $v$  and  $u$

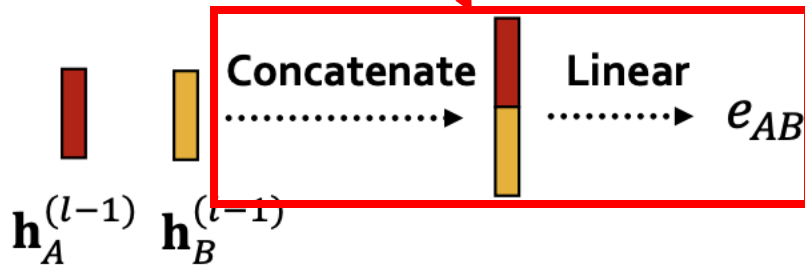
Softmax function  $\rightarrow$  normalized weights

Q: what is  $a$ 's structure?

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

Make use of last layer's output

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



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

$$= \text{Linear}\left(\text{Concat}\left(\mathbf{W}^{(l)} \mathbf{h}_A^{(l-1)}, \mathbf{W}^{(l)} \mathbf{h}_B^{(l-1)}\right)\right)$$

# Graph attention networks

Learnable parameters?

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

Attention weights  
Between  $v$  and  $u$

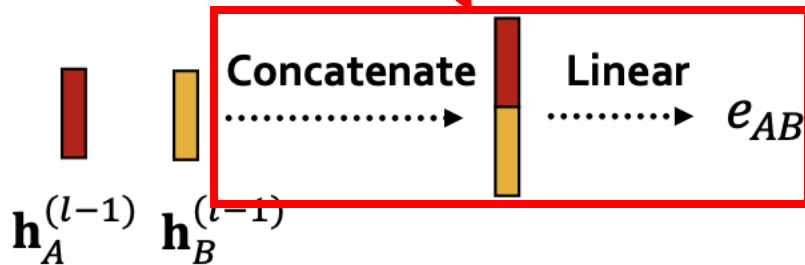
Softmax function  $\rightarrow$  normalized weights

Q: what is  $a$ 's structure?

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

Make use of last layer's output

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



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

$$= \text{Linear}\left(\text{Concat}\left(\mathbf{W}^{(l)} \mathbf{h}_A^{(l-1)}, \mathbf{W}^{(l)} \mathbf{h}_B^{(l-1)}\right)\right)$$