

The background features a light gray gradient. On the left side, there is a complex network graph with dark gray nodes and edges. Scattered across the right side are several thin, light gray triangles of various sizes and orientations. The title 'Graph attention Networks (GAT)' is centered in a large, bold, dark gray font.

# Graph attention Networks (GAT)

---

Antonio Longa<sup>1,2</sup>

MobS<sup>1</sup> Lab, Fondazione Bruno Kessler, Trento, Italy  
SML<sup>2</sup> Lab, University of Trento, Italy



Recap **01**

Introduction **02**

Graph attention layer  
(GAT) **03**

## TABLE OF CONTENTS



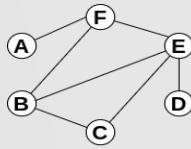
**04** Pros of GAT

**05** Message passing  
Implementation

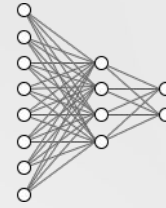
**06** Implement our  
GCNConv

**07** GAT implementation

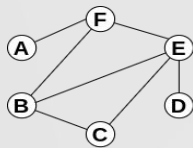
# 01 Recap



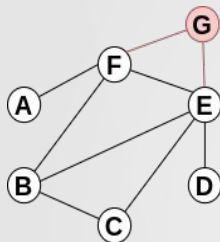
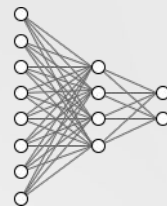
0	0	0	0	0	1
0	0	1	0	1	1
0	1	0	0	1	0
0	0	0	0	1	0
0	1	1	1	0	1
1	1	1	0	1	0



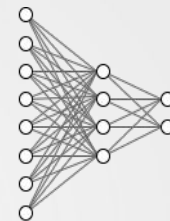
# 01 Recap



0	0	0	0	0	1
0	0	1	0	1	1
0	1	0	0	1	0
0	0	0	0	1	0
0	1	1	1	0	1
1	1	1	0	1	0



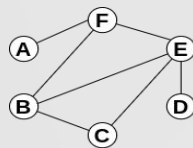
0	0	0	0	0	1	0
0	0	1	0	1	1	0
0	1	0	0	1	0	0
0	0	0	0	1	0	0
0	1	1	1	0	1	1
1	1	1	0	1	0	1
0	0	0	0	1	1	0

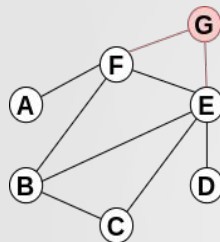
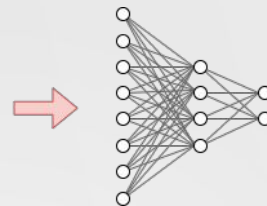


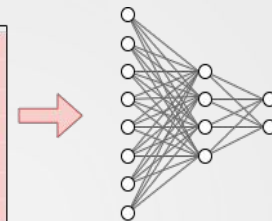
## PROBLEMS:

- Different sizes

# 01 Recap



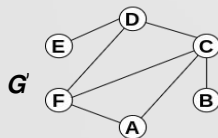
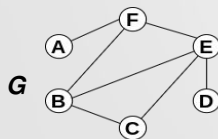
$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 & 1 & 0 \end{bmatrix}$$


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


## PROBLEMS:

- Different sizes

- NOT invariant to nodes ordering



$$G = G'$$

Adj( $G$ )

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

Adj( $G'$ )

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

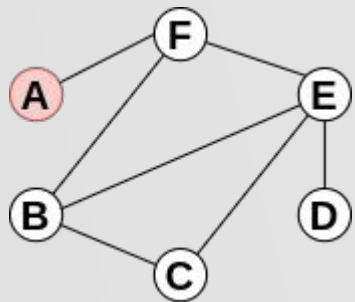
$$\text{Adj}(G) \neq \text{Adj}(G')$$

# 01 Recap

## COMPUTATION GRAPH

The neighbour of a node defines its computation graph

INPUT GRAPH

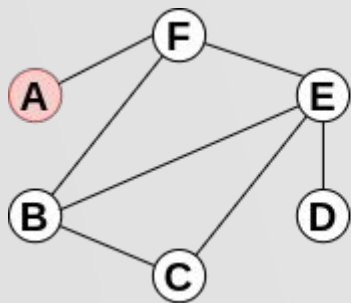


# 01 Recap

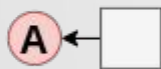
## COMPUTATION GRAPH

The neighbour of a node defines its computation graph

INPUT GRAPH



COMPUTATION GRAPH

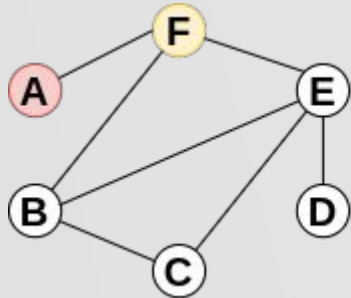


# 01 Recap

## COMPUTATION GRAPH

The neighbour of a node defines its computation graph

INPUT GRAPH



COMPUTATION GRAPH



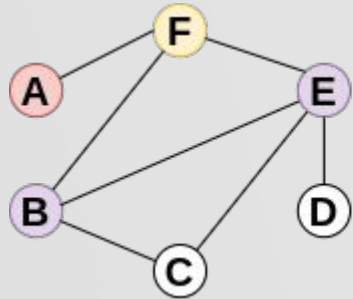


# 01 Recap

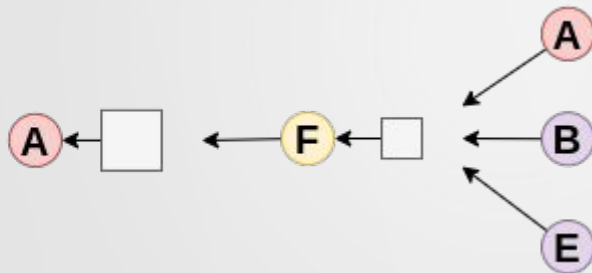
## COMPUTATION GRAPH

The neighbour of a node defines its computation graph

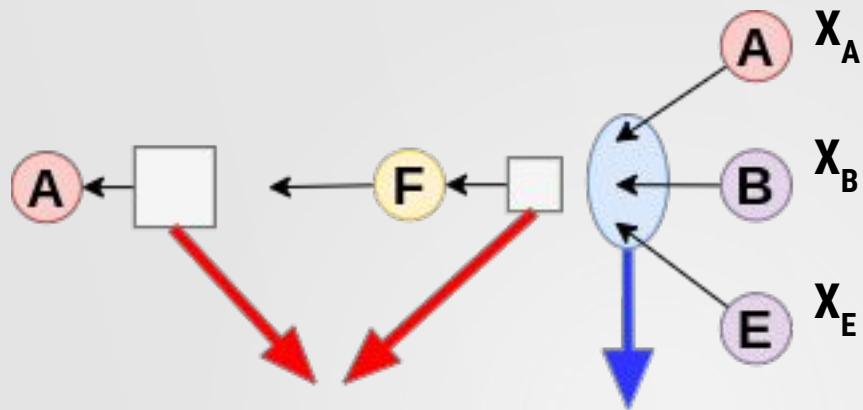
INPUT GRAPH



COMPUTATION GRAPH



# 01 Recap



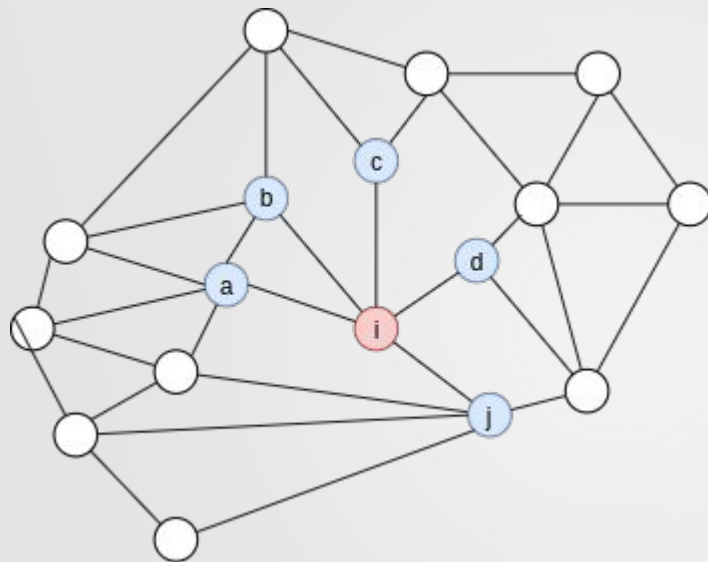
Neural Networks

Ordering invariant  
Aggregation

Sum  
Average

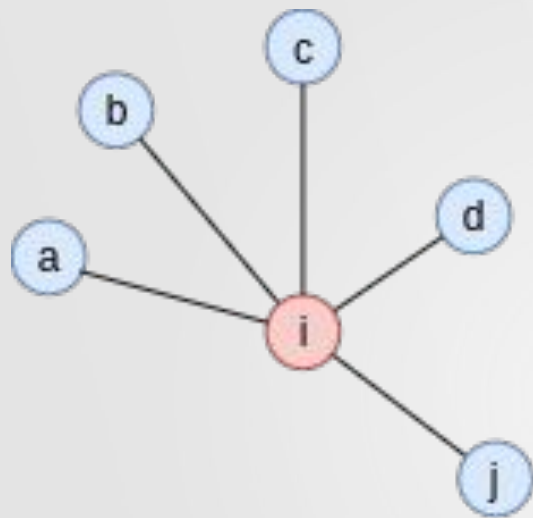
# 02 Introduction

---

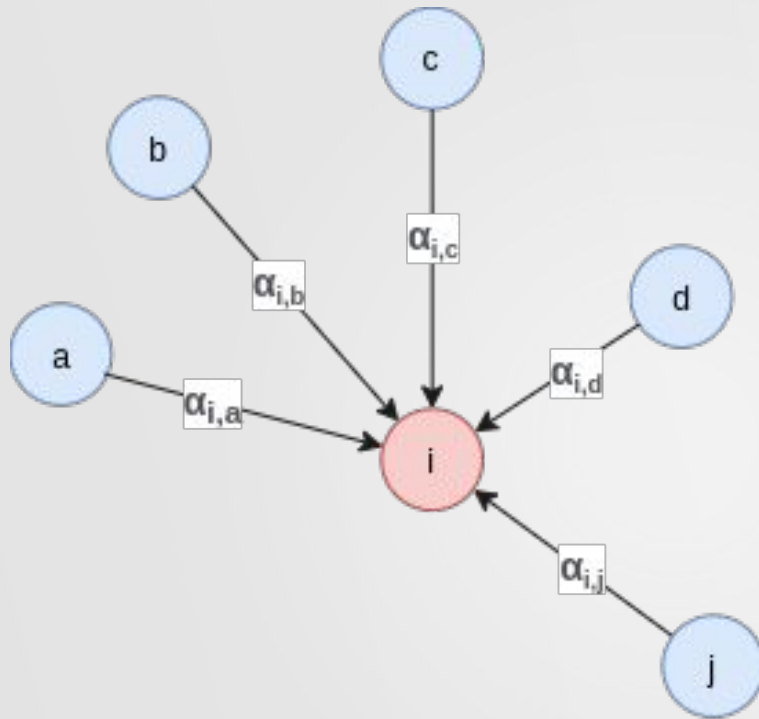


## 02 Introduction

---

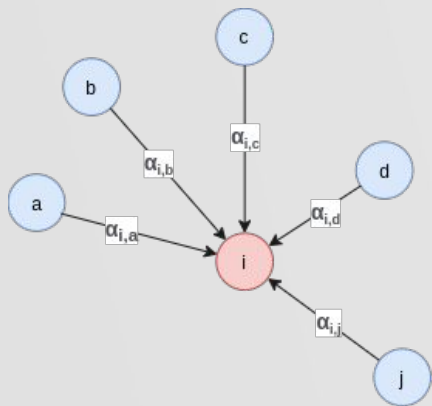


## 02 Introduction

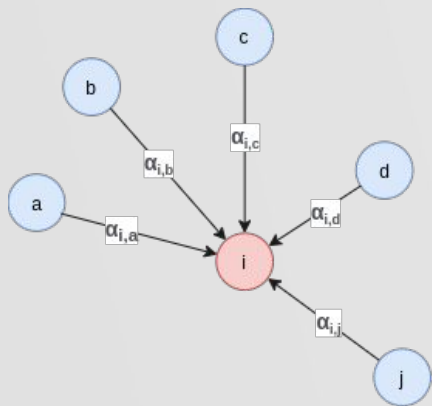


## 02 Introduction

How much features of node "c" are important to node "i"?



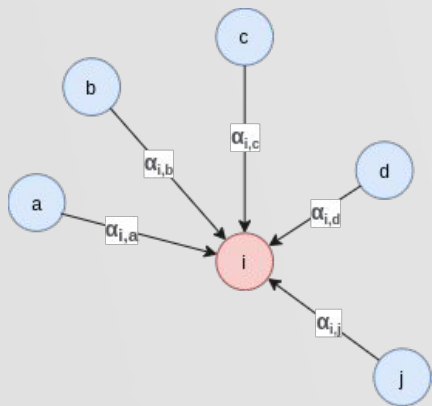
## 02 Introduction



How much features of node "c" are important to node "i"?

Can we learn such importance, in an automatic manner?

## 02 Introduction



How much features of node "c" are important to node "i"?

Can we learn such importance, in an automatic manner?

**YES, with GAT**



# 03 Graph Attention Networks GAT

Published as a conference paper at ICLR 2018

## GRAPH ATTENTION NETWORKS

**Petar Veličković\***

Department of Computer Science and Technology  
University of Cambridge  
petar.velickovic@cst.cam.ac.uk

**Guillem Cucurull\***

Centre de Visió per Computador, UAB  
gcucurull@gmail.com

**Arantxa Casanova\***

Centre de Visió per Computador, UAB  
ar.casanova.8@gmail.com

**Adriana Romero**

Montréal Institute for Learning Algorithms  
adriana.romero.soriano@umontreal.ca

**Pietro Liò**

Department of Computer Science and Technology  
University of Cambridge  
pietro.lio@cst.cam.ac.uk

**Yoshua Bengio**

Montréal Institute for Learning Algorithms  
yoshua.umontreal@gmail.com

**Petar  
Veličković**

Senior Research Scientist at DeepMind




## 03 Graph Attention layer

---

**INPUT:** a set of node features  $\mathbf{h} = \{\bar{h}_1, \bar{h}_2, \dots, \bar{h}_n\}$   $\bar{h}_i \in \mathbf{R}^F$

**OUTPUT:** a **new** set of node features  $\mathbf{h}' = \{\bar{h}'_1, \bar{h}'_2, \dots, \bar{h}'_n\}$   $\bar{h}'_i \in \mathbf{R}^{F'}$





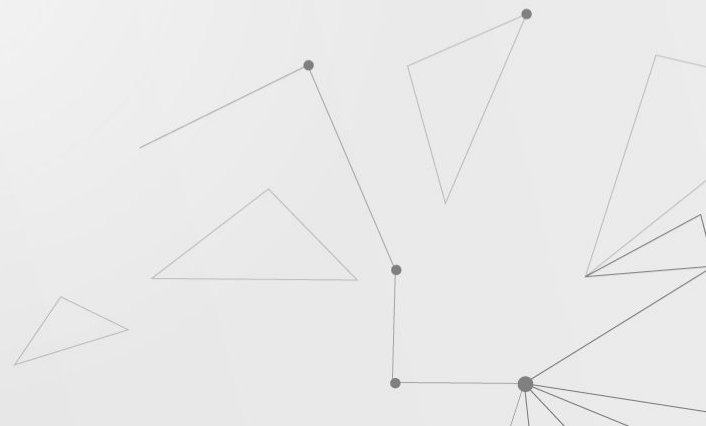
## 03 Graph Attention layer

---

1) apply a **parameterized linear transformation** to every node

$$\mathbf{W} \cdot \bar{h}_i$$

$$\mathbf{W} \in \mathbf{R}^{F' \times F}$$





## 03 Graph Attention layer

---

1) apply a **parameterized linear transformation** to every node

$$\mathbf{W} \cdot \bar{h}_i$$

$$\mathbf{W} \in \mathbf{R}^{F' \times F}$$

$$(F' \times F) \cdot F$$





## 03 Graph Attention layer

---

1) apply a **parameterized linear transformation** to every node

$$\mathbf{W} \cdot \bar{h}_i$$

$$\mathbf{W} \in \mathbf{R}^{F' \times F}$$

~~$$(F' \times F) \cdot F$$~~

$$F'$$





## 03 Graph Attention layer

---

2) Self attention

$$a : \mathbf{R}^{F'} \times \mathbf{R}^{F'} \rightarrow \mathbf{R}$$



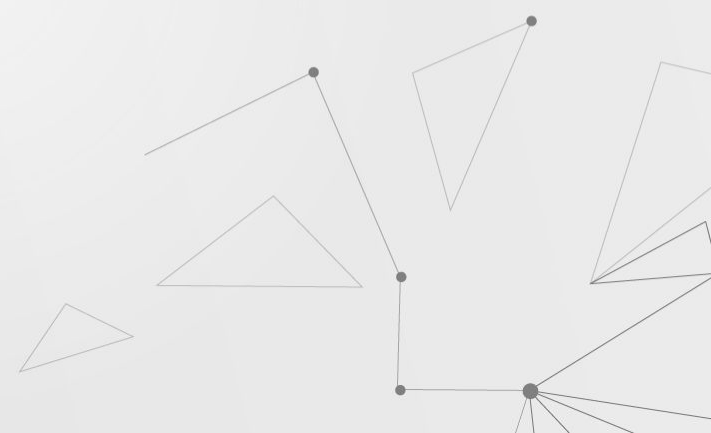


## 03 Graph Attention layer

---

2) Self attention

$$a : \mathbf{R}^{F'} \times \mathbf{R}^{F'} \rightarrow \mathbf{R}$$

$$e_{i,j} = a(\mathbf{W} \cdot \bar{h}_i, \mathbf{W} \cdot \bar{h}_j)$$


## 03 Graph Attention layer

2) Self attention

$$a : \mathbf{R}^{F'} \times \mathbf{R}^{F'} \rightarrow \mathbf{R}$$

$$e_{i,j} = a(\mathbf{W} \cdot \bar{h}_i, \mathbf{W} \cdot \bar{h}_j)$$

Specify the importance of node j's features to node i



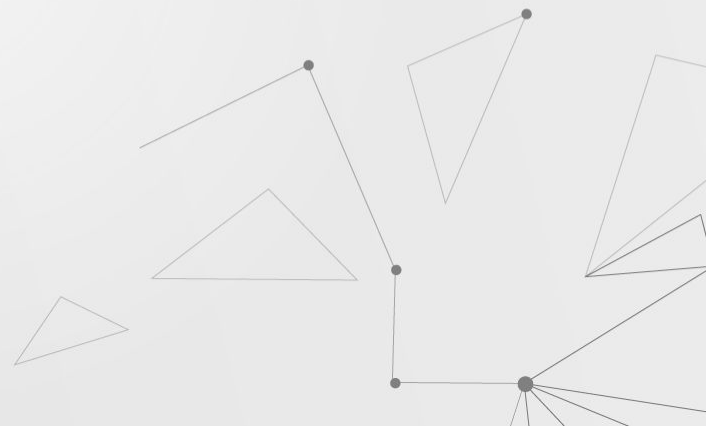


## 03 Graph Attention layer

---

### 3) Normalization

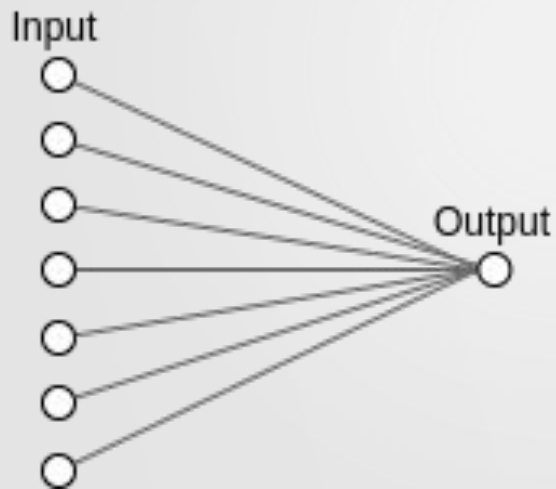
$$\alpha_{i,j} = \textit{softmax}_j(e_{i,j}) = \frac{\exp(e_{i,j})}{\sum_{k \in N(i)} \exp(e_{i,k})}$$



## 03 Graph Attention layer

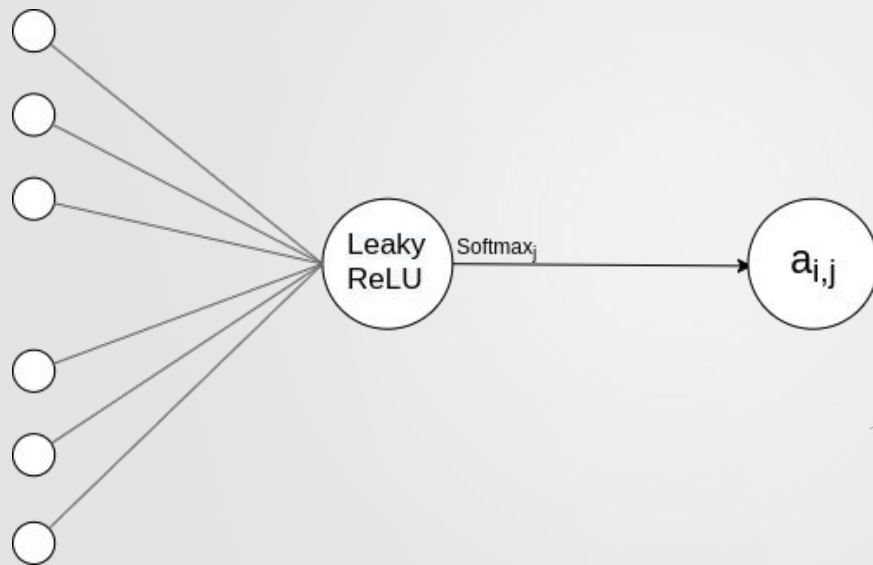
### 4) Attention mechanism $a$

Is a single-layer feed forward neural network



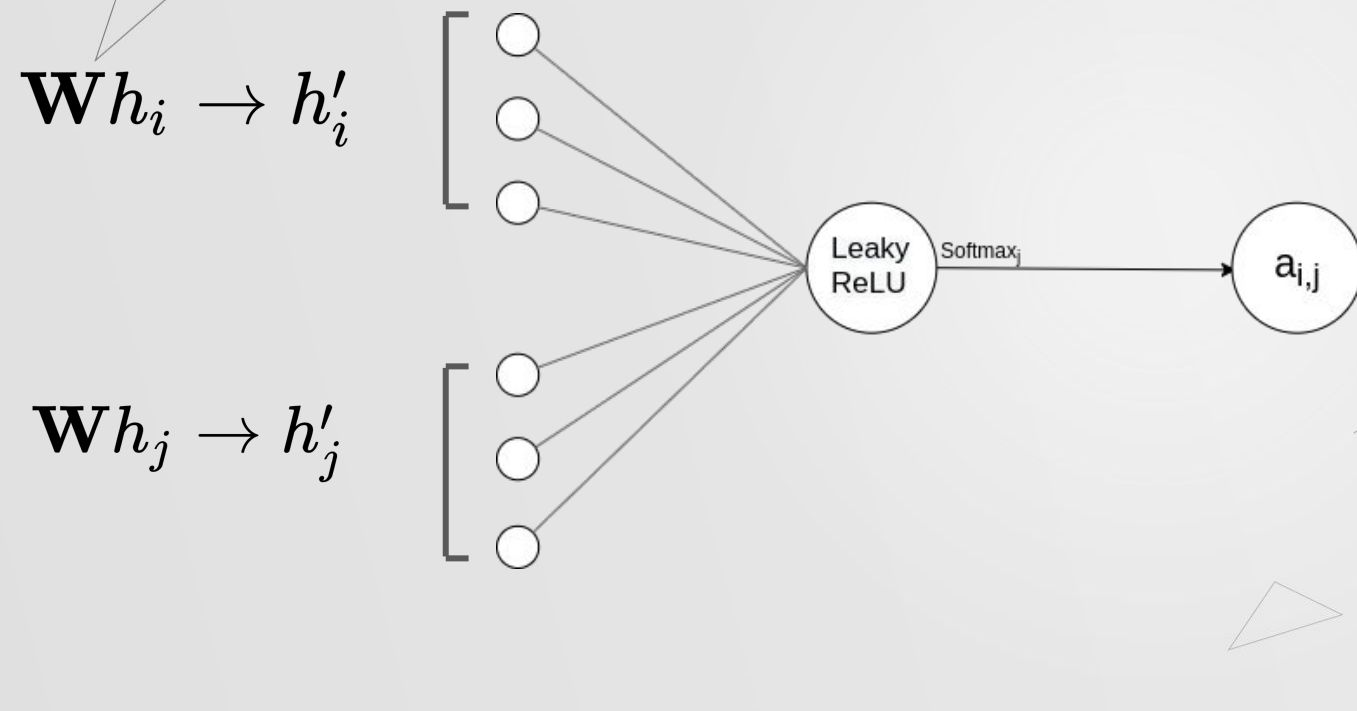
## 03 Graph Attention layer

### 4) Attention mechanism $a$



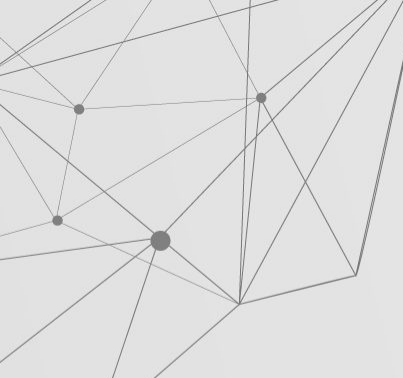
## 03 Graph Attention layer

### 4) Attention mechanism $a$

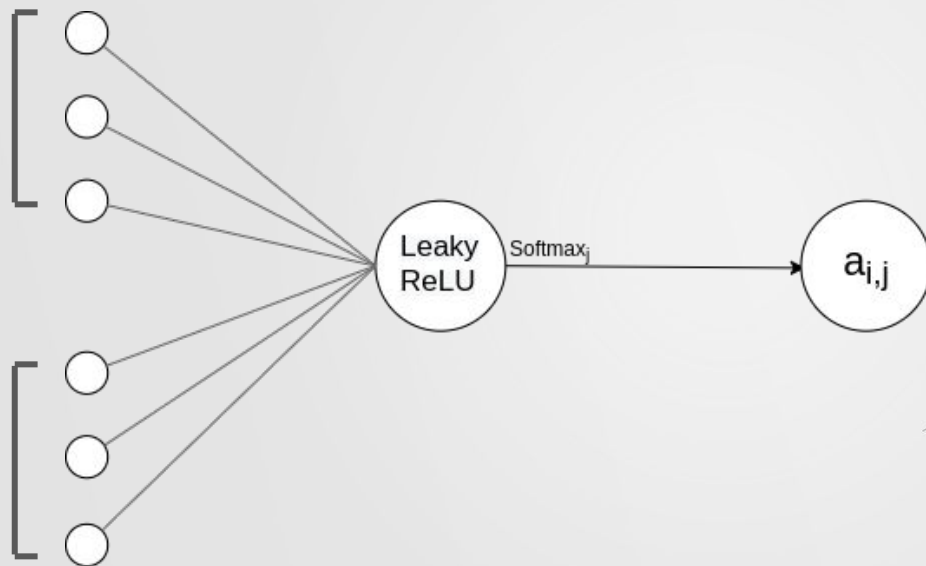


## 03 Graph Attention layer

### 4) Attention mechanism $a$

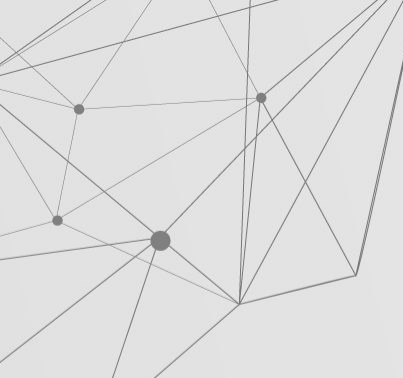

$$\mathbf{W}h_i \rightarrow h'_i$$
$$(F' \times F)F \rightarrow F'$$

$$\mathbf{W}h_j \rightarrow h'_j$$
$$(F' \times F)F \rightarrow F'$$

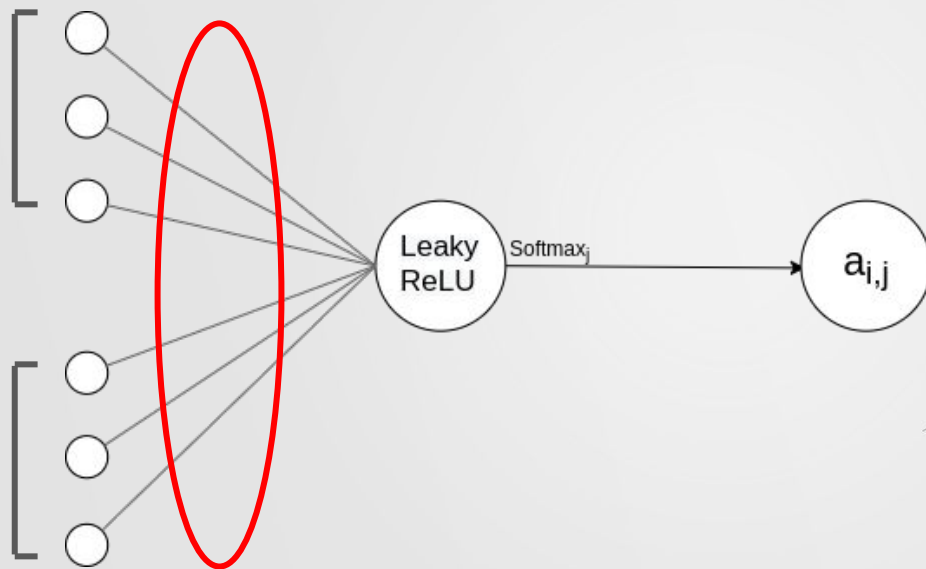


## 03 Graph Attention layer

### 4) Attention mechanism $a$


$$\mathbf{W}h_i \rightarrow h'_i$$
$$(F' \times F)F \rightarrow F'$$

$$\mathbf{W}h_j \rightarrow h'_j$$
$$(F' \times F)F \rightarrow F'$$

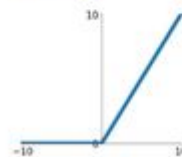


$$\bar{a} \in \mathbf{R}^{2F'}$$


# 03 Graph Attention layer

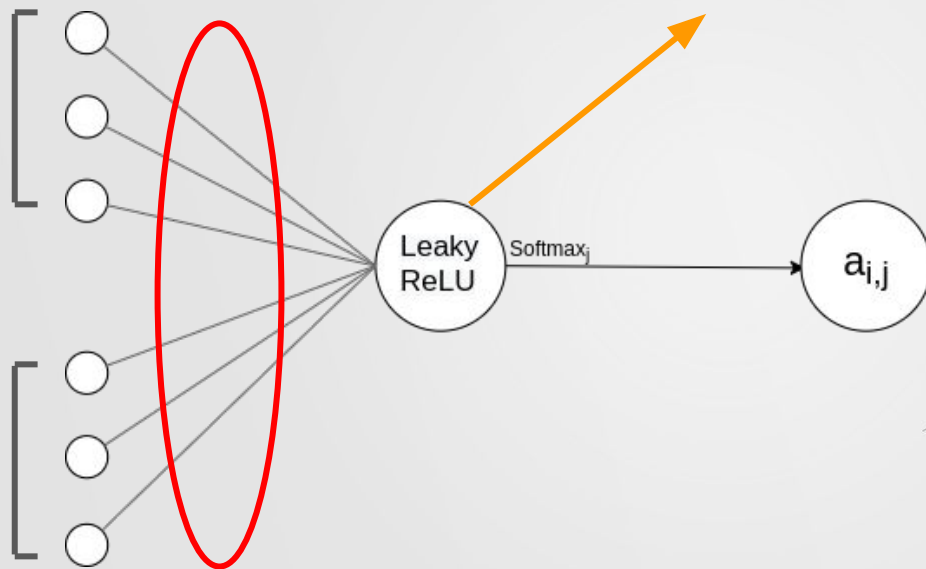
## 4) Attention mechanism $a$

ReLU  
 $\max(0, x)$



$$\mathbf{W}h_i \rightarrow h'_i$$
$$(F' \times F)F \rightarrow F'$$

$$\mathbf{W}h_j \rightarrow h'_j$$
$$(F' \times F)F \rightarrow F'$$



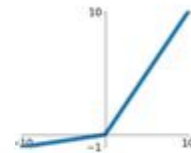
$$\bar{a} \in \mathbf{R}^{2F'}$$

## 03 Graph Attention layer

### 4) Attention mechanism $a$

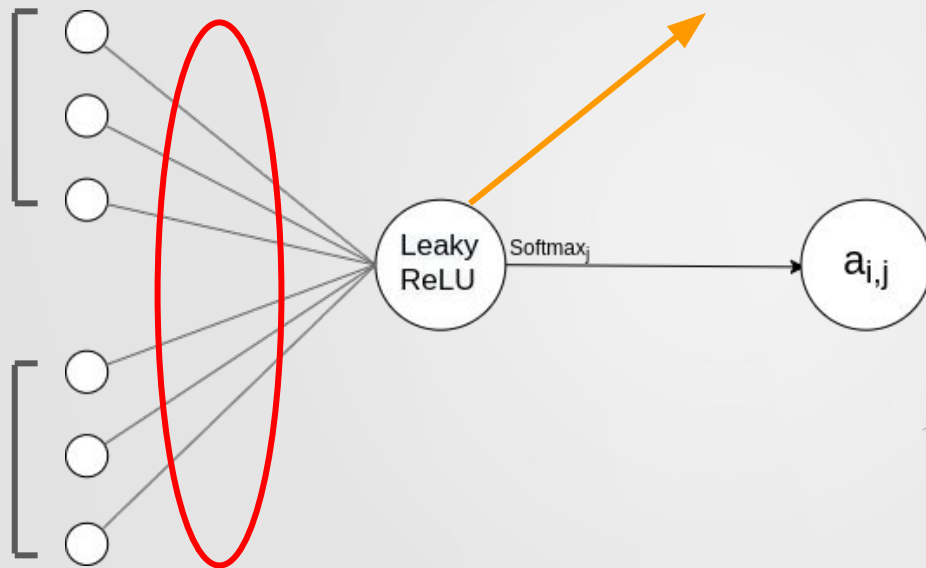
Leaky ReLU

$$\max(0.2x, x)$$



$$\mathbf{W}h_i \rightarrow h'_i$$
$$(F' \times F)F \rightarrow F'$$

$$\mathbf{W}h_j \rightarrow h'_j$$
$$(F' \times F)F \rightarrow F'$$



$$\bar{a} \in \mathbf{R}^{2F'}$$

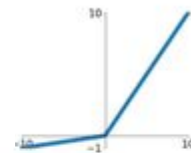


# 03 Graph Attention layer

## 4) Attention mechanism $a$

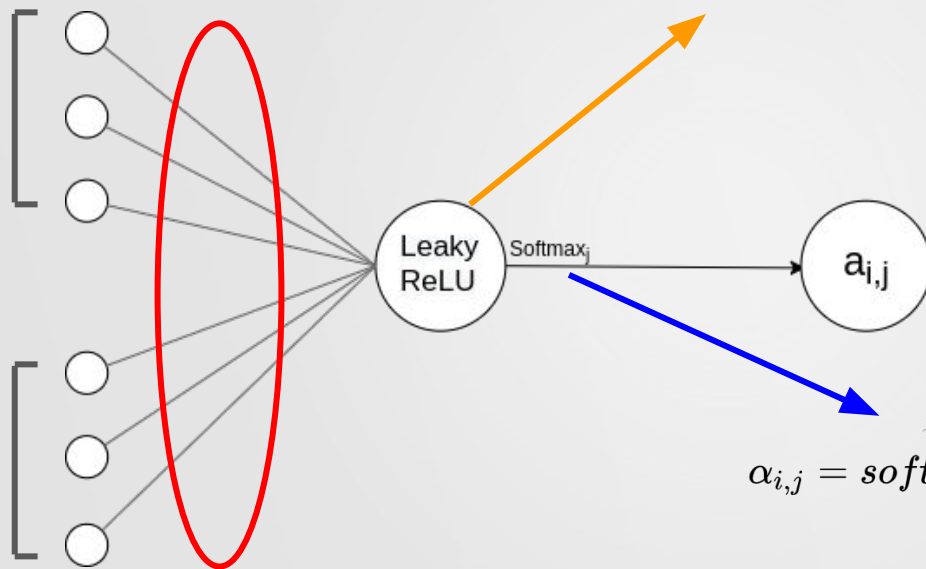
Leaky ReLU

$$\max(0.2x, x)$$



$$\mathbf{W}h_i \rightarrow h'_i$$
$$(F' \times F)F \rightarrow F'$$

$$\mathbf{W}h_j \rightarrow h'_j$$
$$(F' \times F)F \rightarrow F'$$



$$\bar{a} \in \mathbf{R}^{2F'}$$

$$\alpha_{i,j} = \text{softmax}_j(e_{i,j}) = \frac{\exp(e_{i,j})}{\sum_{k \in N(i)} \exp(e_{i,k})}$$

## 03 Graph Attention layer

### 4) Attention mechanism $a$

$$\alpha_{i,j} = \frac{\exp(\text{LeakyReLU}(\bar{a}^T [\mathbf{W}h_i || \mathbf{W}h_j]))}{\sum_{k \in N(i)} \exp(\text{LeakyReLU}(\bar{a}^T [\mathbf{W}h_i || \mathbf{W}h_k]))}$$

$(F' \times F)F$     $(F' \times F)F$

$\bar{a}^T \rightarrow \text{transpose}(a)$

$|| \rightarrow \text{concatenation}$

## 03 Graph Attention layer

### 4) Attention mechanism $a$

$$\alpha_{i,j} = \frac{\exp(\text{LeakyReLU}(\bar{a}^T [\overset{\substack{F' \\ (F' \times F)F}}{\mathbf{W}h_i}} || \overset{\substack{F' \\ (F' \times F)F}}{\mathbf{W}h_j}]))}{\sum_{k \in N(i)} \exp(\text{LeakyReLU}(\bar{a}^T [\mathbf{W}h_i || \mathbf{W}h_k]))}$$

$\bar{a}^T \rightarrow \text{transpose}(a)$

$|| \rightarrow \text{concatenation}$

### 03 Graph Attention layer

#### 4) Attention mechanism $a$

$$\alpha_{i,j} = \frac{\exp(\text{LeakyReLU}(\bar{a}^T [\mathbf{W}h_i || \mathbf{W}h_j]))}{\sum_{k \in N(i)} \exp(\text{LeakyReLU}(\bar{a}^T [\mathbf{W}h_i || \mathbf{W}h_k]))}$$

$(2F' \times 1)$   
 $\overbrace{\mathbf{W}h_i \parallel \mathbf{W}h_j}^{(F' \times F)F \quad (F' \times F)F}$

$\bar{a}^T \rightarrow \text{transpose}(a)$

$|| \rightarrow \text{concatenation}$

### 03 Graph Attention layer

#### 4) Attention mechanism $a$

$$\alpha_{i,j} = \frac{\exp(\text{LeakyReLU}(\bar{a}^T [\mathbf{W}h_i || \mathbf{W}h_j]))}{\sum_{k \in N(i)} \exp(\text{LeakyReLU}(\bar{a}^T [\mathbf{W}h_i || \mathbf{W}h_k]))}$$

Diagram illustrating the dimensions of the tensors in the attention mechanism:

- $\bar{a}$  has dimensions  $(1 \times 2F')$  (indicated by a blue bracket).
- $\mathbf{W}h_i$  and  $\mathbf{W}h_j$  have dimensions  $(F' \times F)F$  (indicated by red brackets).
- The concatenated vector  $[\mathbf{W}h_i || \mathbf{W}h_j]$  has dimensions  $(2F' \times 1)$  (indicated by a red bracket).

$\bar{a}^T \rightarrow \text{transpose}(a)$

$|| \rightarrow \text{concatenation}$

## 03 Graph Attention layer

### 4) Attention mechanism $a$

$$\alpha_{i,j} = \frac{\exp(\text{LeakyReLU}(\bar{a}^T [\mathbf{W}h_i || \mathbf{W}h_j]))}{\sum_{k \in N(i)} \exp(\text{LeakyReLU}(\bar{a}^T [\mathbf{W}h_i || \mathbf{W}h_k]))}$$

Diagram illustrating the dimensions of the tensors in the attention mechanism:

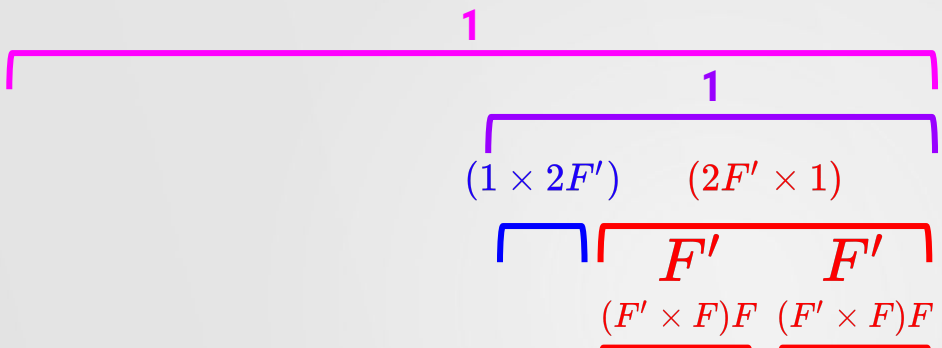
- $\bar{a}$  (purple bracket):  $(1 \times 2F')$
- $\mathbf{W}h_i$  (blue bracket):  $F'$
- $\mathbf{W}h_j$  (red bracket):  $F'$
- $\mathbf{W}h_i$  (red bracket):  $(F' \times F)F$
- $\mathbf{W}h_j$  (red bracket):  $(F' \times F)F$

$\bar{a}^T \rightarrow \text{transpose}(a)$

$|| \rightarrow \text{concatenation}$

# 03 Graph Attention layer

## 4) Attention mechanism $a$



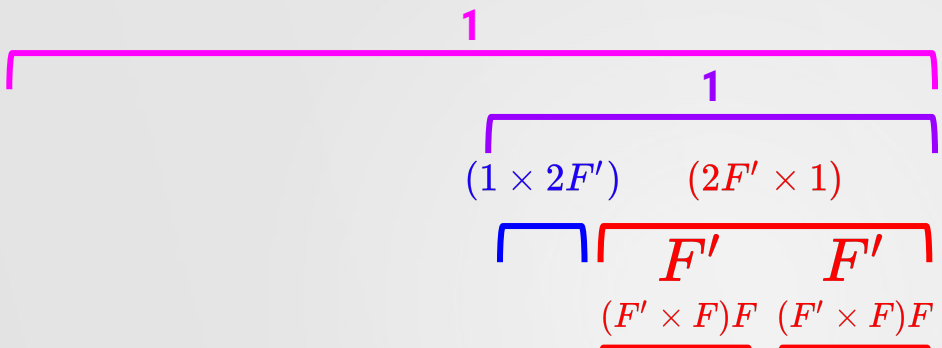
$$\alpha_{i,j} = \frac{\exp(\text{LeakyReLU}(\bar{a}^T [\mathbf{W}h_i || \mathbf{W}h_j]))}{\sum_{k \in N(i)} \exp(\text{LeakyReLU}(\bar{a}^T [\mathbf{W}h_i || \mathbf{W}h_k]))}$$

$\bar{a}^T \rightarrow \text{transpose}(a)$

$|| \rightarrow \text{concatenation}$

# 03 Graph Attention layer

## 4) Attention mechanism $a$



$$\alpha_{i,j} = \frac{\exp(\text{LeakyReLU}(\bar{a}^T [\mathbf{W}h_i || \mathbf{W}h_j]))}{\sum_{k \in N(i)} \exp(\text{LeakyReLU}(\bar{a}^T [\mathbf{W}h_i || \mathbf{W}h_k]))} = \frac{\mathbf{R}}{\mathbf{R}} = \mathbf{R}$$

$\bar{a}^T \rightarrow \text{transpose}(a)$

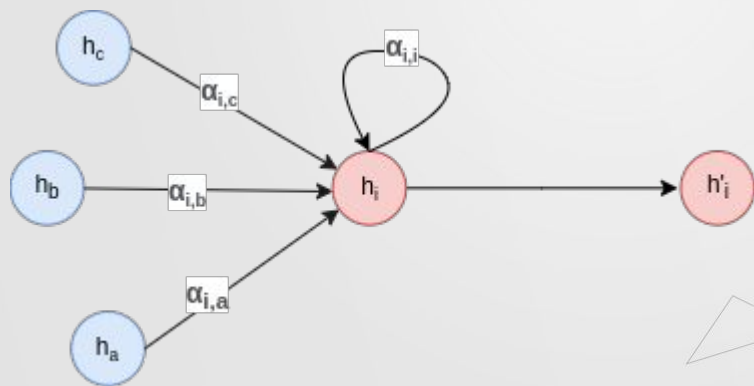
$|| \rightarrow \text{concatenation}$



## 03 Graph Attention layer

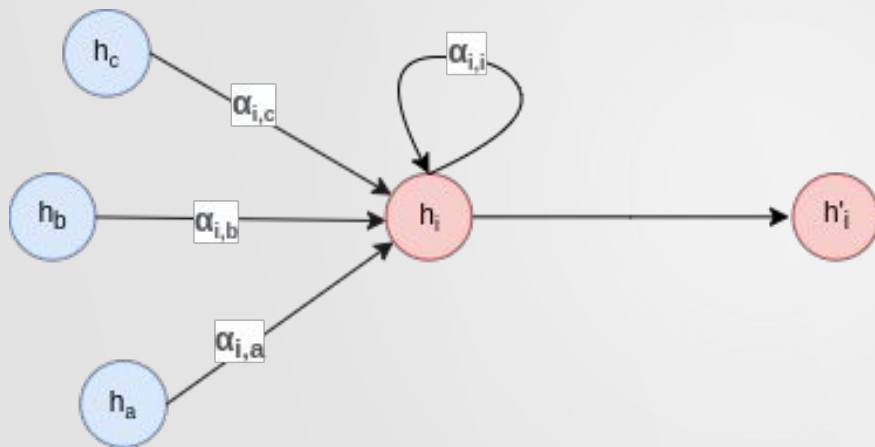
5) Use it :)

$$h'_i = \sigma(\sum_{j \in N(i)} \alpha_{i,j} \mathbf{W} h_j)$$



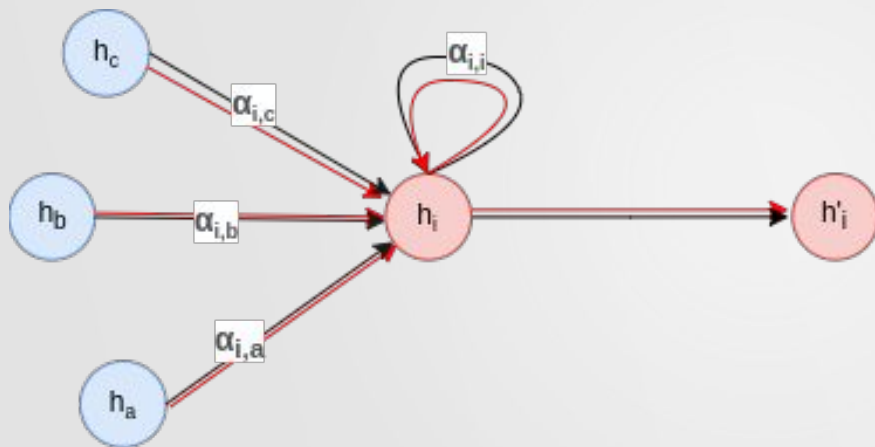
# 03 Graph Attention layer

## 6) Multi-head attention



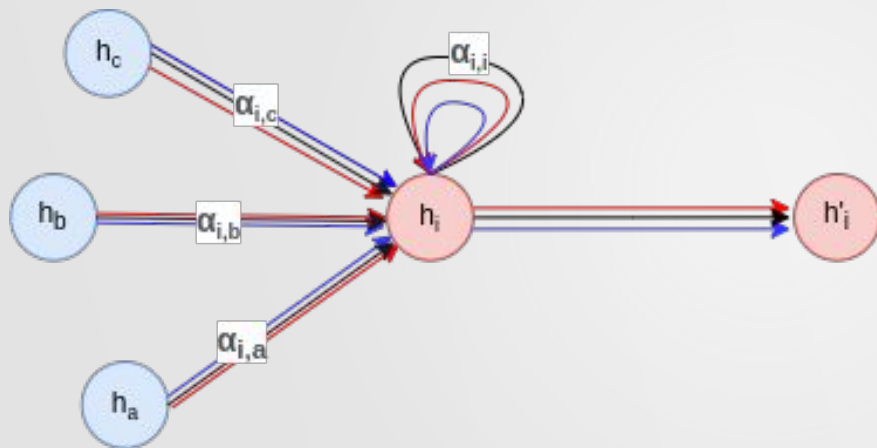
# 03 Graph Attention layer

## 6) Multi-head attention



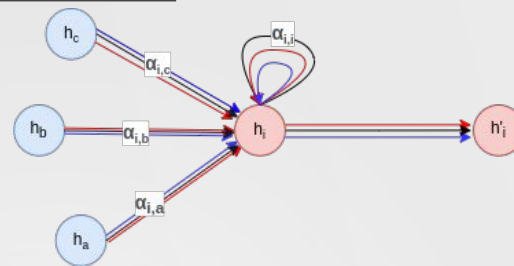
# 03 Graph Attention layer

## 6) Multi-head attention



# 03 Graph Attention layer

## 6) Multi-head attention



### Concatenation

$$h'_i = ||_{k=1}^K \sigma(\sum_{j \in N(i)} \alpha_{i,j}^k \mathbf{W}^k h_j)$$

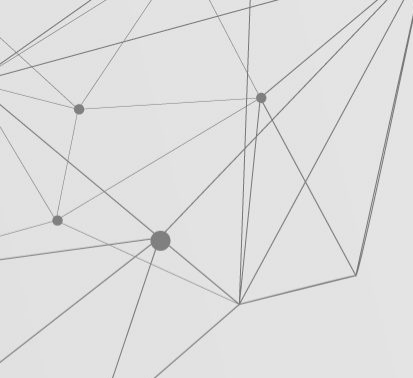
### Average

$$h'_i = \sigma(\frac{1}{K} \sum_{k=1}^K \sum_{j \in N(i)} \alpha_{i,j}^k \mathbf{W}^k h_j)$$

- On the final (prediction) layer of the network

## 04 Pros of GAT

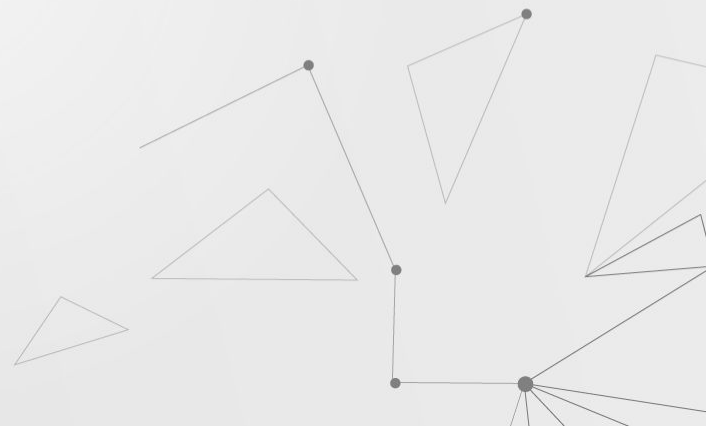
---



Computationally **efficient**

Self-attention layers can be **parallelized across edges**

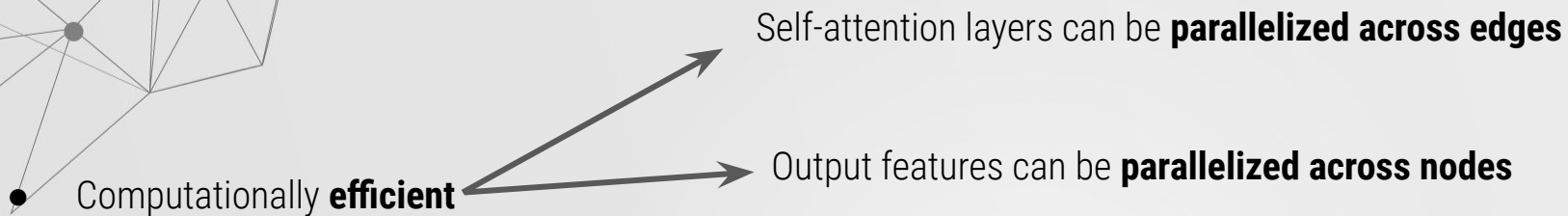
Output features can be **parallelized across nodes**



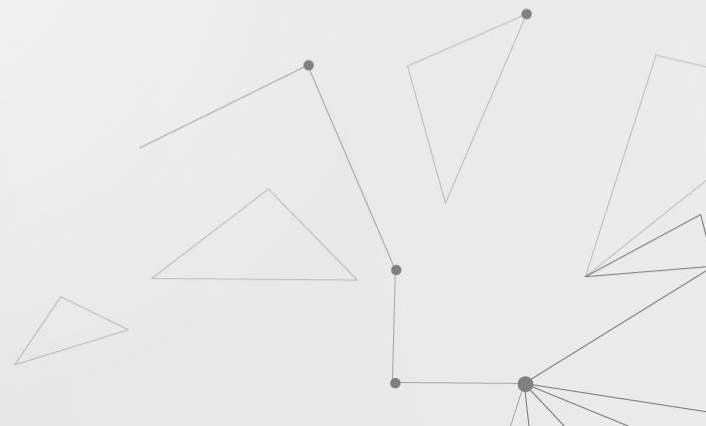


## 04 Pros of GAT

---



- Allows to assign **different importances to nodes** of a same neighborhood



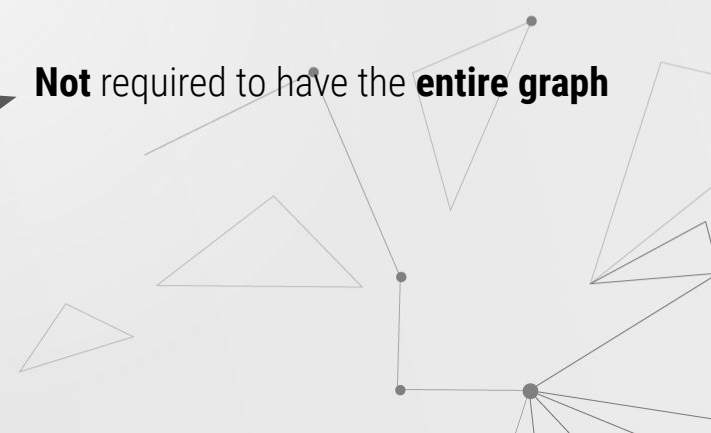


## 04 Pros of GAT

---

- Computationally **efficient**
  - Self-attention layers can be **parallelized across edges**
  - Output features can be **parallelized across nodes**

- Allows to assign **different importances to nodes** of a same neighborhood

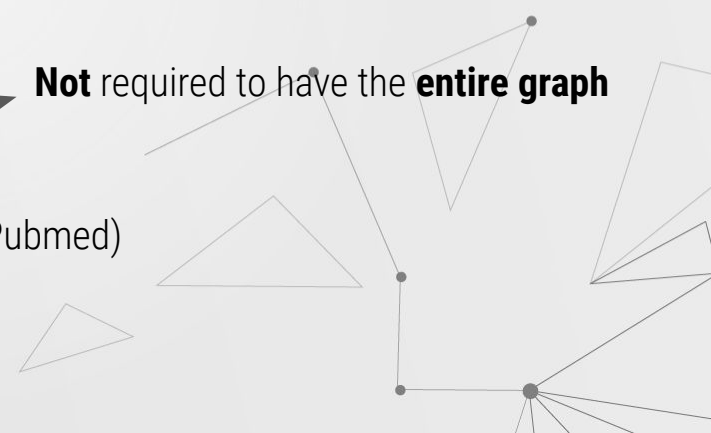
- It is applied in a **shared manner** to all edges in the graph
    - Not** required to have the **entire graph**
- 






## 04 Pros of GAT

---

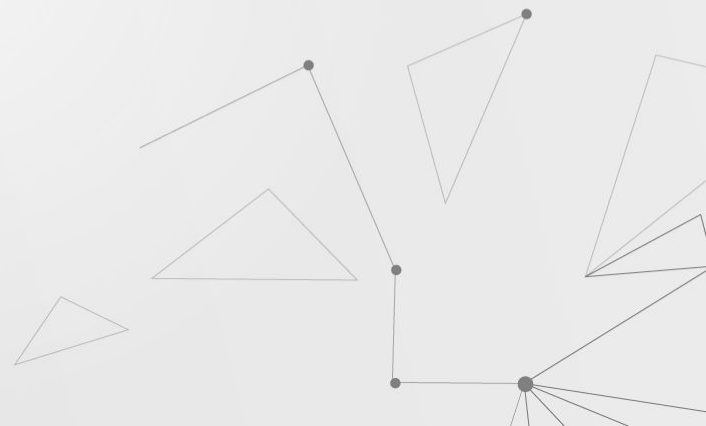
- Computationally **efficient**
    - Self-attention layers can be **parallelized across edges**
    - Output features can be **parallelized across nodes**
  - Allows to assign **different importances to nodes** of a same neighborhood
  - It is applied in a **shared manner** to all edges in the graph
    - **Not** required to have the **entire graph**
  - Works in **both**:
    - **Transductive** learning (Cora, Citeseer, Pubmed)
    - **Inductive** learning (PPI)
- 

## 05 Message passing implementation

$$\mathbf{x}_i^{(k)} = \gamma^{(k)} \left( \mathbf{x}_i^{(k-1)}, \bigoplus_{j \in \mathcal{N}(i)} \phi^{(k)} \left( \mathbf{x}_i^{(k-1)}, \mathbf{x}_j^{(k-1)}, \mathbf{e}_{j,i} \right) \right),$$




Features representations of  
node  $i$  at the  $k$ -th layer

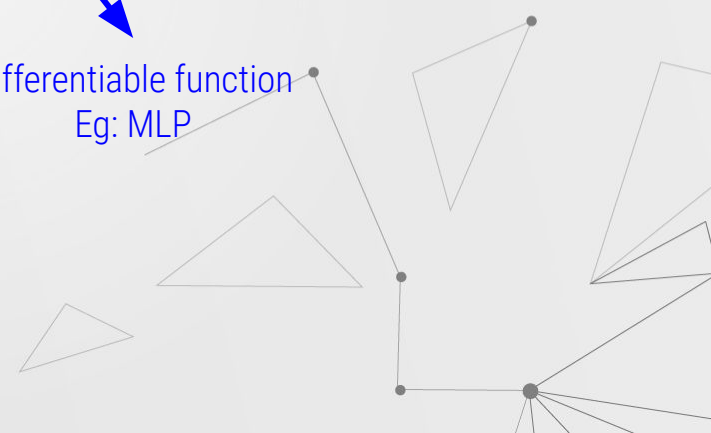


## 05 Message passing implementation

$$\mathbf{x}_i^{(k)} = \gamma^{(k)} \left( \mathbf{x}_i^{(k-1)}, \bigoplus_{j \in \mathcal{N}(i)} \phi^{(k)} \left( \mathbf{x}_i^{(k-1)}, \mathbf{x}_j^{(k-1)}, \mathbf{e}_{j,i} \right) \right),$$




Features representations of  
node  $i$  at the  $k$ -th layer



Differentiable function  
Eg: MLP

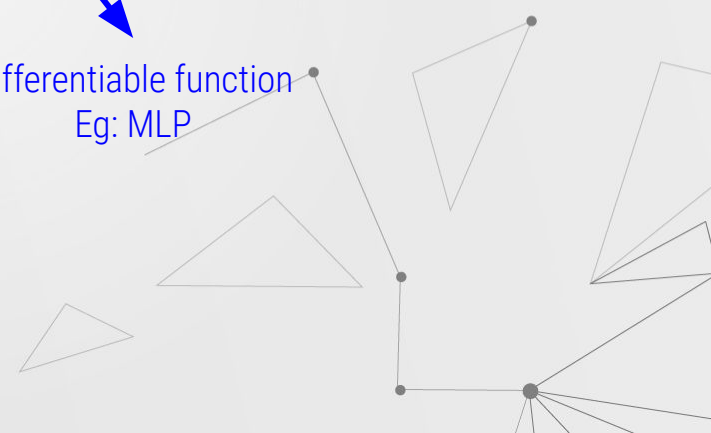
# 05 Message passing implementation


$$\mathbf{x}_i^{(k)} = \gamma^{(k)} \left( \mathbf{x}_i^{(k-1)}, \bigoplus_{j \in \mathcal{N}(i)} \phi^{(k)} \left( \mathbf{x}_i^{(k-1)}, \mathbf{x}_j^{(k-1)}, \mathbf{e}_{j,i} \right) \right),$$


Features representations of  
node i at the k-th layer

- Feature rep of node i at the (k-1)-th layer
- Feature rep of node j at the (k-1)-th layer
- **[optionally]** features of edge (i,j)

Differentiable function  
Eg: MLP



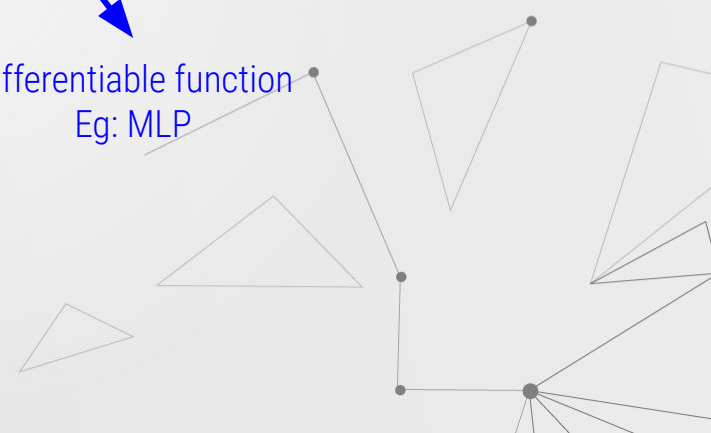
# 05 Message passing implementation


$$\mathbf{x}_i^{(k)} = \gamma^{(k)} \left( \mathbf{x}_i^{(k-1)}, \square_{j \in \mathcal{N}(i)} \phi^{(k)} \left( \mathbf{x}_i^{(k-1)}, \mathbf{x}_j^{(k-1)}, \mathbf{e}_{j,i} \right) \right),$$


Features representations of  
node i at the k-th layer

Differentiable, ordering  
invariant function.  
For every j in the  
neighbourhood of i.  
Eg: sum, average, etc...

Differentiable function  
Eg: MLP

- 
- Feature rep of node i at the (k-1)-th layer
  - Feature rep of node j at the (k-1)-th layer
  - **[optionally]** features of edge (i,j)

# 05 Message passing implementation

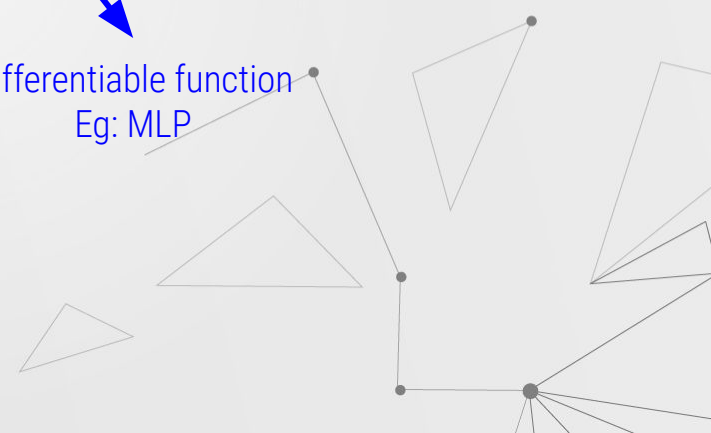

$$\mathbf{x}_i^{(k)} = \gamma^{(k)} \left( \mathbf{x}_i^{(k-1)}, \square_{j \in \mathcal{N}(i)} \phi^{(k)} \left( \mathbf{x}_i^{(k-1)}, \mathbf{x}_j^{(k-1)}, \mathbf{e}_{j,i} \right) \right),$$

Features representations of  
node  $i$  at the  $k$ -th layer

Differentiable function  
Eg: MLP

Differentiable, ordering  
invariant function.  
For every  $j$  in the  
neighbourhood of  $i$ .  
Eg: sum, average, etc...

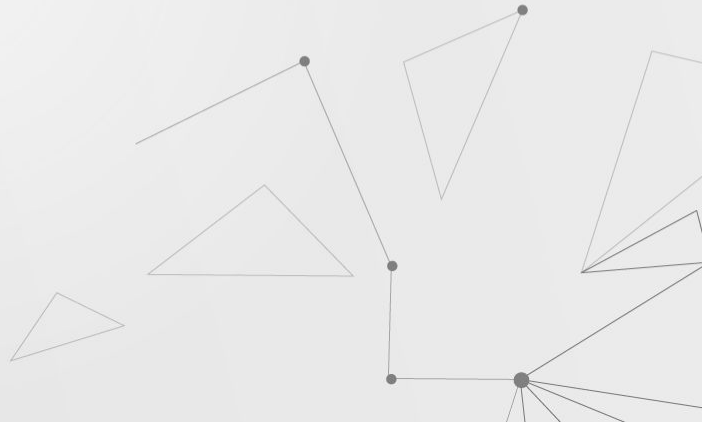
Differentiable function  
Eg: MLP

- 
- Feature rep of node  $i$  at the  $(k-1)$ -th layer
  - Feature rep of node  $j$  at the  $(k-1)$ -th layer
  - **[optionally]** features of edge  $(i,j)$



## 05 Message passing implementation

PyTorch Geometric provides the MessagePassing base class.



## 05 Message passing implementation

PyTorch Geometric provides the MessagePassing base class.

$$\mathbf{x}_i^{(k)} = \gamma^{(k)} \left( \mathbf{x}_i^{(k-1)}, \bigoplus_{j \in \mathcal{N}(i)} \phi^{(k)} \left( \mathbf{x}_i^{(k-1)}, \mathbf{x}_j^{(k-1)}, \mathbf{e}_{j,i} \right) \right),$$

Differentiable function  
Eg: MLP

message()



# 05 Message passing implementation

PyTorch Geometric provides the MessagePassing base class.

$$\mathbf{x}_i^{(k)} = \gamma^{(k)}\left(\mathbf{x}_i^{(k-1)}, \bigoplus_{j \in \mathcal{N}(i)} \phi^{(k)}\left(\mathbf{x}_i^{(k-1)}, \mathbf{x}_j^{(k-1)}, \mathbf{e}_{j,i}\right)\right),$$

Differentiable function

Eg: MLP



update()

Differentiable functions

Eg: MLP



message()

# 05 Message passing implementation

PyTorch Geometric provides the MessagePassing base class.

$$\mathbf{x}_i^{(k)} = \gamma^{(k)}\left(\mathbf{x}_i^{(k-1)}, \square_{j \in \mathcal{N}(i)} \phi^{(k)}\left(\mathbf{x}_i^{(k-1)}, \mathbf{x}_j^{(k-1)}, \mathbf{e}_{j,i}\right)\right),$$

Differentiable function

Eg: MLP

update()

Aggregation

Sum, avg, concat

Differentiable function

Eg: MLP

message()



## 05 Message passing implementation

PyTorch Geometric provides the MessagePassing base class.


### PARAMETERS

```
CLASS MessagePassing ( aggr: Optional[str] = 'add', flow: str = 'source_to_target', node_dim: int =  
- 2 ) [source]
```

Base class for creating message passing layers of the form

$$\mathbf{x}'_i = \gamma_{\Theta} \left( \mathbf{x}_i, \square_{j \in \mathcal{N}(i)} \phi_{\Theta} (\mathbf{x}_i, \mathbf{x}_j, \mathbf{e}_{j,i}) \right),$$

where  $\square$  denotes a differentiable, permutation invariant function, e.g., sum, mean or max, and  $\gamma_{\Theta}$  and  $\phi_{\Theta}$  denote differentiable functions such as MLPs. See [here](#) for the accompanying tutorial.



# 05 Message passing implementation

PyTorch Geometric provides the MessagePassing base class.

## PARAMETERS

**aggr** (*string, optional*) – The aggregation scheme to use ( "add" , "mean" , "max" OR None ).  
(default: "add" )

**CLASS MessagePassing** ( aggr: Optional[str] = 'add', flow: str = 'source\_to\_target', node\_dim: int = - 2 ) [\[source\]](#)

Base class for creating message passing layers of the form

$$\mathbf{x}'_i = \gamma_{\Theta} \left( \mathbf{x}_i, \square_{j \in \mathcal{N}(i)} \phi_{\Theta} \left( \mathbf{x}_i, \mathbf{x}_j, \mathbf{e}_{j,i} \right) \right),$$

where  $\square$  denotes a differentiable, permutation invariant function, e.g., sum, mean or max, and  $\gamma_{\Theta}$  and  $\phi_{\Theta}$  denote differentiable functions such as MLPs. See [here](#) for the accompanying tutorial.

# 05 Message passing implementation

PyTorch Geometric provides the MessagePassing base class.

## PARAMETERS

**aggr** (*string, optional*) – The aggregation scheme to use ( "add" , "mean" , "max" OR None ).  
(default: "add" )

**flow** (*string, optional*) – The flow direction of message passing ( "source\_to\_target" OR "target\_to\_source" ). (default: "source\_to\_target" )

**CLASS MessagePassing** ( aggr: Optional[str] = 'add', flow: str = 'source\_to\_target', node\_dim: int = - 2 ) [\[source\]](#)

Base class for creating message passing layers of the form

$$\mathbf{x}'_i = \gamma_{\Theta} \left( \mathbf{x}_i, \square_{j \in \mathcal{N}(i)} \phi_{\Theta} \left( \mathbf{x}_i, \mathbf{x}_j, \mathbf{e}_{j,i} \right) \right),$$

where  $\square$  denotes a differentiable, permutation invariant function, e.g., sum, mean or max, and  $\gamma_{\Theta}$  and  $\phi_{\Theta}$  denote differentiable functions such as MLPs. See [here](#) for the accompanying tutorial.



# 05 Message passing implementation

---

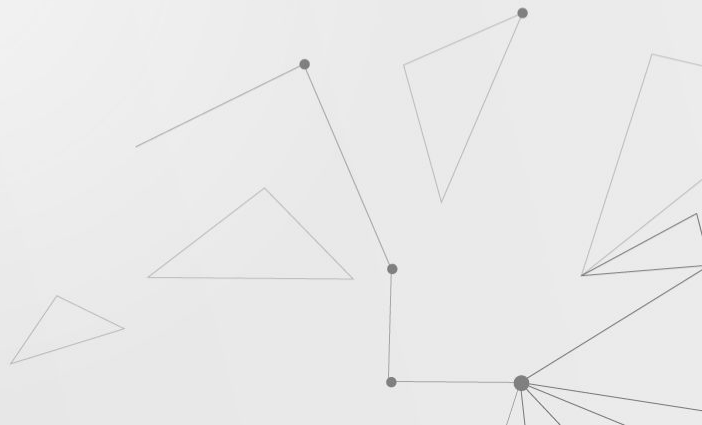
PyTorch Geometric provides the MessagePassing base class.

## METHODS

```
CLASS MessagePassing ( aggr: Optional[str] = 'add', flow: str = 'source_to_target', node_dim: int =  
- 2 ) [source]
```

Aggregates messages from neighbors  
(sum, mean, max)

```
aggregate ( inputs: torch.Tensor, index: torch.Tensor, ptr: Optional[torch.Tensor] = None, dim_size:  
Optional[int] = None ) → torch.Tensor [source]
```



# 05 Message passing implementation

PyTorch Geometric provides the MessagePassing base class.

## METHODS

```
CLASS MessagePassing ( aggr: Optional[str] = 'add', flow: str = 'source_to_target', node_dim: int =  
- 2 ) [source]
```

Aggregates messages from neighbors  
(sum, mean, max)

```
aggregate ( inputs: torch.Tensor, index: torch.Tensor, ptr: Optional[torch.Tensor] = None, dim_size:  
Optional[int] = None ) → torch.Tensor [source]
```

Constructs messages from node  $j$  to  
node  $i$  in analogy to  $\phi_{\theta}$

```
message ( x_j: torch.Tensor ) → torch.Tensor [source]
```

# 05 Message passing implementation

PyTorch Geometric provides the MessagePassing base class.

## METHODS

```
CLASS MessagePassing ( aggr: Optional[str] = 'add', flow: str = 'source_to_target', node_dim: int =  
- 2 ) \[source\]
```

Aggregates messages from neighbors  
(sum, mean, max)

```
aggregate ( inputs: torch.Tensor, index: torch.Tensor, ptr: Optional[torch.Tensor] = None, dim_size:  
Optional[int] = None ) → torch.Tensor \[source\]
```

Constructs messages from node  $j$  to  
node  $i$  in analogy to  $\phi_{\theta}$

```
message ( x_j: torch.Tensor ) → torch.Tensor \[source\]
```

Propagate messages

```
propagate ( edge_index: Union[torch.Tensor, torch_sparse.tensor.SparseTensor], size:  
Optional[Tuple[int, int]] = None, **kwargs ) \[source\]
```



# 05 Message passing implementation

PyTorch Geometric provides the MessagePassing base class.

## METHODS

```
CLASS MessagePassing ( aggr: Optional[str] = 'add', flow: str = 'source_to_target', node_dim: int =  
- 2 ) \[source\]
```

Aggregates messages from neighbors  
(sum, mean, max)

```
aggregate ( inputs: torch.Tensor, index: torch.Tensor, ptr: Optional[torch.Tensor] = None, dim_size:  
Optional[int] = None ) → torch.Tensor \[source\]
```

Constructs messages from node  $j$  to  
node  $i$  in analogy to  $\phi\theta$

```
message ( x_j: torch.Tensor ) → torch.Tensor \[source\]
```

Propagate messages

```
propagate ( edge_index: Union[torch.Tensor, torch_sparse.tensor.SparseTensor], size:  
Optional[Tuple[int, int]] = None, **kwargs ) \[source\]
```

Updates node embeddings in  
analogy to  $\gamma\theta$

```
update ( inputs: torch.Tensor ) → torch.Tensor \[source\]
```



# 05 Message passing implementation

---

## HOW TO USE IT?

### Layer Name

```
class GCNConv(MessagePassing):
    def __init__(self, in_channels, out_channels):
        super(GCNConv, self).__init__(aggr='add')

    def forward(self, x, edge_index):
        return self.propagate(edge_index, x=x, norm=norm)

    def message(self, ...):

        return ...
```



# 05 Message passing implementation

## HOW TO USE IT?

GCNConv inherits from MessagePassing

Layer Name

```
class GCNConv(MessagePassing):
    def __init__(self, in_channels, out_channels):
        super(GCNConv, self).__init__(aggr='add')

    def forward(self, x, edge_index):
        return self.propagate(edge_index, x=x, norm=norm)

    def message(self, ...):
        return ...
```

# 05 Message passing implementation

## HOW TO USE IT?

GCNConv inherits from MessagePassing

Layer Name

```
class GCNConv(MessagePassing):  
    def __init__(self, in_channels, out_channels):  
        super(GCNConv, self).__init__(aggr='add')  
  
    def forward(self, x, edge_index):  
        return self.propagate(edge_index, x=x, norm=norm)  
  
    def message(self, ...):  
        return ...
```

Initialize the class, call “super” specifying your aggregations (add,max,mean)

# 05 Message passing implementation

## HOW TO USE IT?

GCNConv inherits from MessagePassing

Layer Name

```
class GCNConv(MessagePassing):  
    def __init__(self, in_channels, out_channels):  
        super(GCNConv, self).__init__(aggr='add')
```

Initialize the class, call "super" specifying your aggregations (add,max,mean)

```
    def forward(self, x, edge_index):  
        return self.propagate(edge_index, x=x, norm=norm)
```

Forward and propagate

```
    def message(self, ...):  
        return ...
```

# 05 Message passing implementation

## HOW TO USE IT?

GCNConv inherits from MessagePassing

Layer Name

```
class GCNConv(MessagePassing):  
    def __init__(self, in_channels, out_channels):  
        super(GCNConv, self).__init__(aggr='add')
```

Initialize the class, call "super" specifying your aggregations (add,max,mean)

```
    def forward(self, x, edge_index):  
        return self.propagate(edge_index, x=x, norm=norm)
```

Forward and propagate

```
    def message(self, ...):  
        return ...
```

Compute the message

## 06 Implement our GCNConv

Simple example

$$\mathbf{x}_i^{(k)} = \sum_{j \in \mathcal{N}(i) \cup \{i\}} \frac{1}{\sqrt{\deg(i)} \cdot \sqrt{\deg(j)}} \cdot \left( \Theta \cdot \mathbf{x}_j^{(k-1)} \right)$$

## 06 Implement our GCNConv

Simple example

$$\mathbf{x}_i^{(k)} = \sum_{j \in \mathcal{N}(i) \cup \{i\}} \frac{1}{\sqrt{\deg(i)} \cdot \sqrt{\deg(j)}} \cdot \left( \Theta \cdot \mathbf{x}_j^{(k-1)} \right)$$

$$\mathbf{x}_i^{(k)} = \gamma^{(k)} \left( \mathbf{x}_i^{(k-1)}, \bigoplus_{j \in \mathcal{N}(i)} \phi^{(k)} \left( \mathbf{x}_i^{(k-1)}, \mathbf{x}_j^{(k-1)}, \mathbf{e}_{j,i} \right) \right),$$



## 06 Implement our GCNConv

Simple example

$$\mathbf{x}_i^{(k)} = \sum_{j \in \mathcal{N}(i) \cup \{i\}} \frac{1}{\sqrt{\deg(i)} \cdot \sqrt{\deg(j)}} \cdot (\Theta \cdot \mathbf{x}_j^{(k-1)})$$

$$\mathbf{x}_i^{(k)} = \gamma^{(k)} \left( \mathbf{x}_i^{(k-1)}, \square_{j \in \mathcal{N}(i)} \phi^{(k)} \left( \mathbf{x}_i^{(k-1)}, \mathbf{x}_j^{(k-1)}, \mathbf{e}_{j,i} \right) \right),$$



## 06 Implement our GCNConv

---

### Simple example

$$\mathbf{x}_i^{(k)} = \sum_{j \in \mathcal{N}(i) \cup \{i\}} \frac{1}{\sqrt{\deg(i)} \cdot \sqrt{\deg(j)}} \cdot \left( \Theta \cdot \mathbf{x}_j^{(k-1)} \right)$$

In steps:

1. Add self loops
2. A linear transformation to node feature matrix
3. Compute normalization coefficients
4. Normalize node features
5. Sum up neighboring node features



## 06 Implement our GCNConv

### Simple example

$$\mathbf{x}_i^{(k)} = \sum_{j \in \mathcal{N}(i) \cup \{i\}} \frac{1}{\sqrt{\deg(i)} \cdot \sqrt{\deg(j)}} \cdot \left( \Theta \cdot \mathbf{x}_j^{(k-1)} \right)$$

In steps:

1. Add self loops
2. A linear transformation to node feature matrix
3. Compute normalization coefficients
4. Normalize node features
5. Sum up neighboring node features

Forward method

Message method  
int



# 06 Implement our GCNConv

## GCNConv inherits from MessagePassing

```
class GCNConv(MessagePassing):
    def __init__(self, in_channels, out_channels):
        super(GCNConv, self).__init__(aggr='add') # "Add" aggregation (Step 5).
        self.lin = torch.nn.Linear(in_channels, out_channels)

    def forward(self, x, edge_index):
        # x has shape [N, in_channels]
        # edge_index has shape [2, E]

        # Step 1: Add self-loops to the adjacency matrix.
        edge_index, _ = add_self_loops(edge_index, num_nodes=x.size(0))

        # Step 2: Linearly transform node feature matrix.
        x = self.lin(x)

        # Step 3: Compute normalization.
        row, col = edge_index
        deg = degree(col, x.size(0), dtype=x.dtype)
        deg_inv_sqrt = deg.pow(-0.5)
        norm = deg_inv_sqrt[row] * deg_inv_sqrt[col]

        # Step 4-5: Start propagating messages.
        return self.propagate(edge_index, x=x, norm=norm)

    def message(self, x_j, norm):
        # x_j has shape [E, out_channels]

        # Step 4: Normalize node features.
        return norm.view(-1, 1) * x_j
```

# 06 Implement our GCNConv

GCNConv inherits from MessagePassing

```
class GCNConv(MessagePassing):
    def __init__(self, in_channels, out_channels):
        super(GCNConv, self).__init__(aggr='add') # "Add" aggregation (Step 5).
        self.lin = torch.nn.Linear(in_channels, out_channels)

    def forward(self, x, edge_index):
        # x has shape [N, in_channels]
        # edge_index has shape [2, E]

        # Step 1: Add self-loops to the adjacency matrix.
        edge_index, _ = add_self_loops(edge_index, num_nodes=x.size(0))

        # Step 2: Linearly transform node feature matrix.
        x = self.lin(x)

        # Step 3: Compute normalization.
        row, col = edge_index
        deg = degree(col, x.size(0), dtype=x.dtype)
        deg_inv_sqrt = deg.pow(-0.5)
        norm = deg_inv_sqrt[row] * deg_inv_sqrt[col]

        # Step 4-5: Start propagating messages.
        return self.propagate(edge_index, x=x, norm=norm)

    def message(self, x_j, norm):
        # x_j has shape [E, out_channels]

        # Step 4: Normalize node features.
        return norm.view(-1, 1) * x_j
```

1) Add self loops



# 06 Implement our GCNConv

GCNConv inherits from MessagePassing

```
class GCNConv(MessagePassing):
    def __init__(self, in_channels, out_channels):
        super(GCNConv, self).__init__(aggr='add') # "Add" aggregation (Step 5).
        self.lin = torch.nn.Linear(in_channels, out_channels)

    def forward(self, x, edge_index):
        # x has shape [N, in_channels]
        # edge_index has shape [2, E]

        # Step 1: Add self-loops to the adjacency matrix.
        edge_index, _ = add_self_loops(edge_index, num_nodes=x.size(0))

        # Step 2: Linearly transform node feature matrix.
        x = self.lin(x)

        # Step 3: Compute normalization.
        row, col = edge_index
        deg = degree(col, x.size(0), dtype=x.dtype)
        deg_inv_sqrt = deg.pow(-0.5)
        norm = deg_inv_sqrt[row] * deg_inv_sqrt[col]

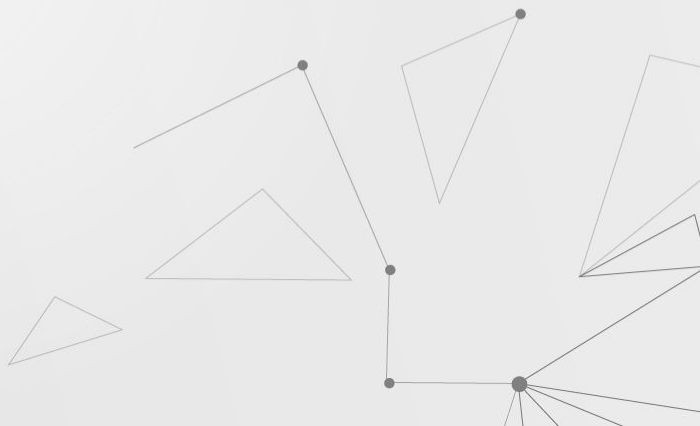
        # Step 4-5: Start propagating messages.
        return self.propagate(edge_index, x=x, norm=norm)

    def message(self, x_j, norm):
        # x_j has shape [E, out_channels]

        # Step 4: Normalize node features.
        return norm.view(-1, 1) * x_j
```

1) Add self loops

2) A linear transformation to node feature matrix





# 06 Implement our GCNConv

GCNConv inherits from MessagePassing

```
class GCNConv(MessagePassing):  
    def __init__(self, in_channels, out_channels):  
        super(GCNConv, self).__init__(aggr='add') # "Add" aggregation (Step 5).  
        self.lin = torch.nn.Linear(in_channels, out_channels)
```

```
    def forward(self, x, edge_index):  
        # x has shape [N, in_channels]  
        # edge_index has shape [2, E]  
  
        # Step 1: Add self-loops to the adjacency matrix.  
        edge_index, _ = add_self_loops(edge_index, num_nodes=x.size(0))
```

1) Add self loops

```
        # Step 2: Linearly transform node feature matrix.  
        x = self.lin(x)
```

2) A linear transformation to node feature matrix

```
        # Step 3: Compute normalization.  
        row, col = edge_index  
        deg = degree(col, x.size(0), dtype=x.dtype)  
        deg_inv_sqrt = deg.pow(-0.5)  
        norm = deg_inv_sqrt[row] * deg_inv_sqrt[col]
```

3) Compute normalization coefficients

```
        # Step 4-5: Start propagating messages.  
        return self.propagate(edge_index, x=x, norm=norm)
```

```
    def message(self, x_j, norm):  
        # x_j has shape [E, out_channels]  
  
        # Step 4: Normalize node features.  
        return norm.view(-1, 1) * x_j
```

# 06 Implement our GCNConv

GCNConv inherits from MessagePassing

```
class GCNConv(MessagePassing):  
    def __init__(self, in_channels, out_channels):  
        super(GCNConv, self).__init__(aggr='add') # "Add" aggregation (Step 5).  
        self.lin = torch.nn.Linear(in_channels, out_channels)
```

```
    def forward(self, x, edge_index):  
        # x has shape [N, in_channels]  
        # edge_index has shape [2, E]  
  
        # Step 1: Add self-loops to the adjacency matrix.  
        edge_index, _ = add_self_loops(edge_index, num_nodes=x.size(0))
```

1) Add self loops

```
        # Step 2: Linearly transform node feature matrix.  
        x = self.lin(x)
```

2) A linear transformation to node feature matrix

```
        # Step 3: Compute normalization.  
        row, col = edge_index  
        deg = degree(col, x.size(0), dtype=x.dtype)  
        deg_inv_sqrt = deg.pow(-0.5)  
        norm = deg_inv_sqrt[row] * deg_inv_sqrt[col]
```

3) Compute normalization coefficients

```
        # Step 4-5: Start propagating messages.  
        return self.propagate(edge_index, x=x, norm=norm)
```

```
    def message(self, x_j, norm):  
        # x_j has shape [E, out_channels]
```

```
        # Step 4: Normalize node features.  
        return norm.view(-1, 1) * x_j
```

4) Normalize node features



# 06 Implement our GCNConv

GCNConv inherits from MessagePassing

```
class GCNConv(MessagePassing):
    def __init__(self, in_channels, out_channels):
        super(GCNConv, self).__init__(aggr='add') # "Add" aggregation (Step 5).
        self.lin = torch.nn.Linear(in_channels, out_channels)
```

5) Sum up neighboring node features

```
def forward(self, x, edge_index):
    # x has shape [N, in_channels]
    # edge_index has shape [2, E]

    # Step 1: Add self-loops to the adjacency matrix.
    edge_index, _ = add_self_loops(edge_index, num_nodes=x.size(0))
```

1) Add self loops

```
    # Step 2: Linearly transform node feature matrix.
    x = self.lin(x)
```

2) A linear transformation to node feature matrix

```
    # Step 3: Compute normalization.
    row, col = edge_index
    deg = degree(col, x.size(0), dtype=x.dtype)
    deg_inv_sqrt = deg.pow(-0.5)
    norm = deg_inv_sqrt[row] * deg_inv_sqrt[col]
```

3) Compute normalization coefficients

```
    # Step 4-5: Start propagating messages.
    return self.propagate(edge_index, x=x, norm=norm)
```

```
def message(self, x_j, norm):
    # x_j has shape [E, out_channels]

    # Step 4: Normalize node features.
    return norm.view(-1, 1) * x_j
```

4) Normalize node features





# 06 GAT implementation

---

Jupyter-Notebook

