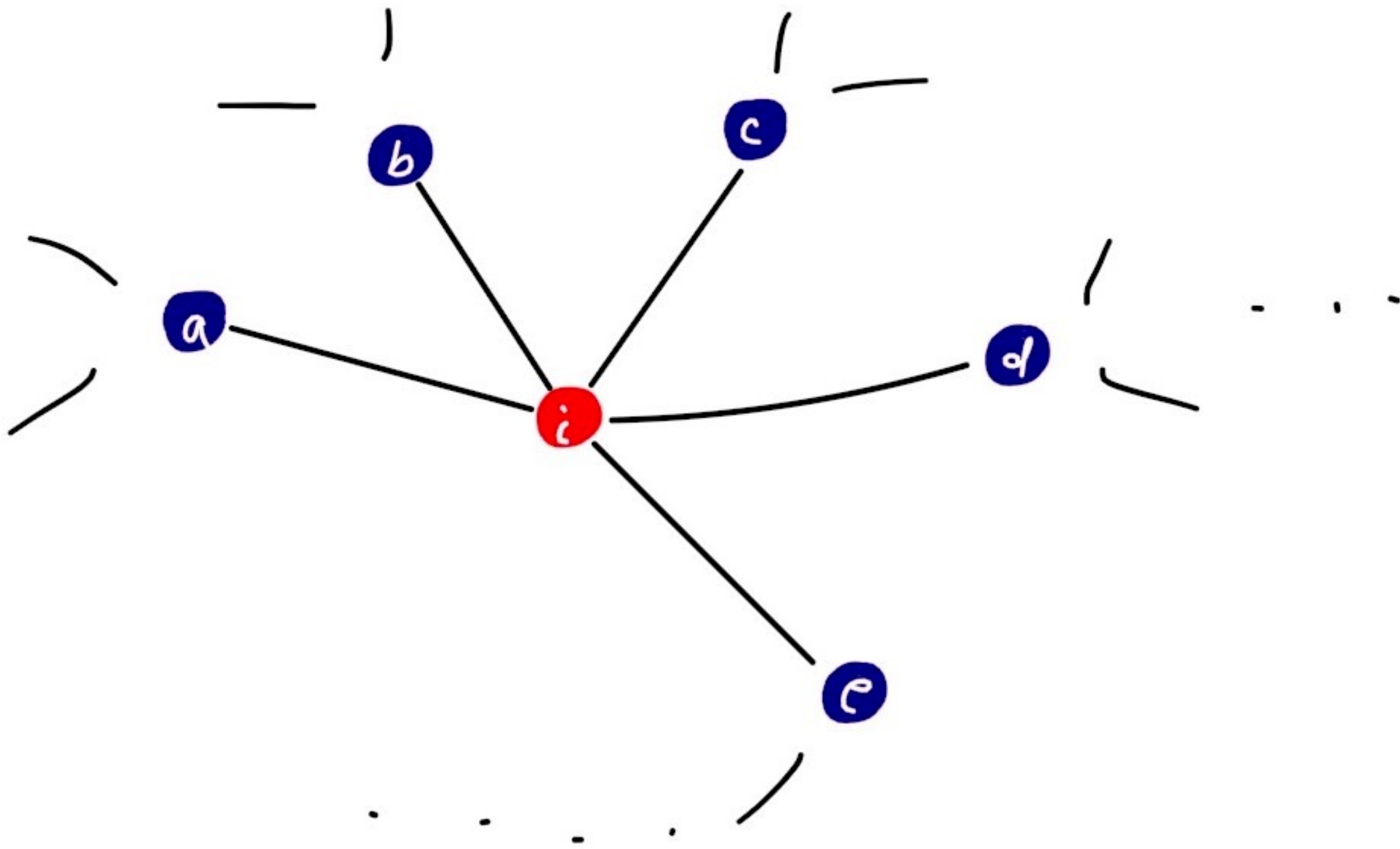
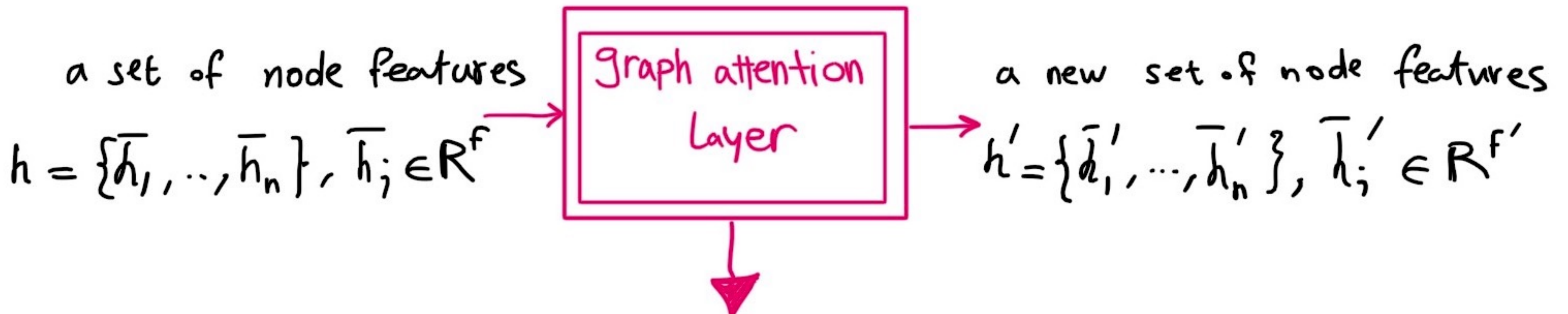


# Graph Attention Networks (GAT) in PyG ICLR'18



★ How much the features of  $a-e$  are important to node  $i$ ? → we'd like to learn this



① apply a parameterized linear transformation to each node  
 $W \cdot \bar{h}_i$ ,  $W \in \mathbb{R}^{F' \times F}$

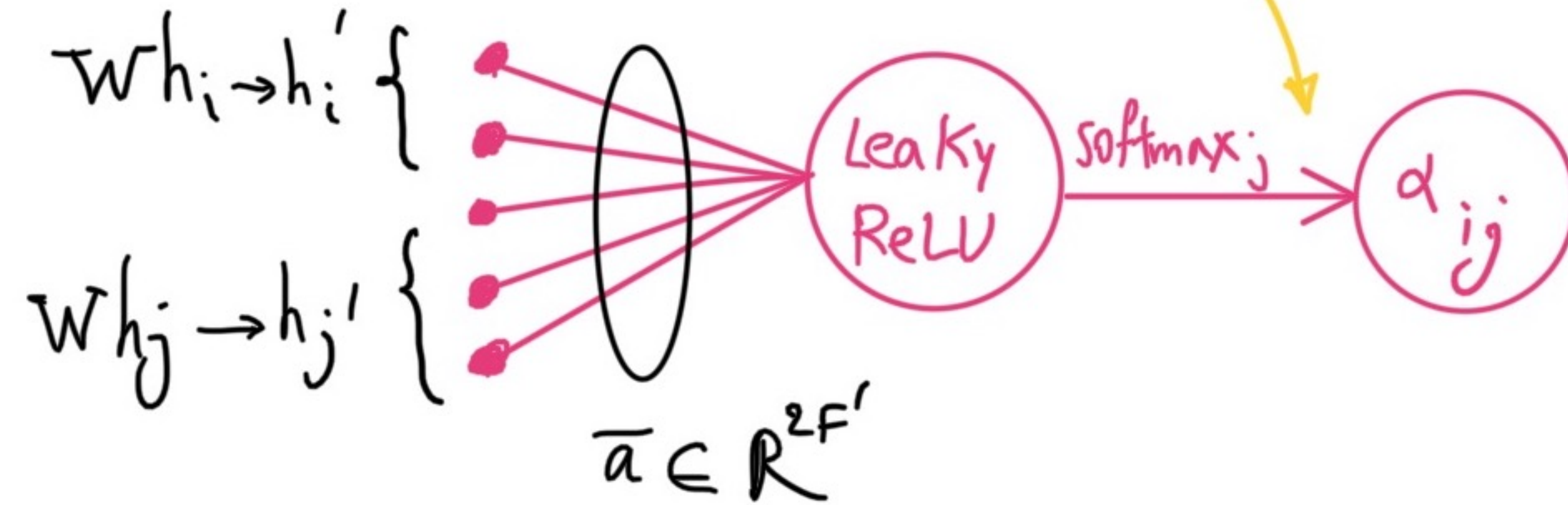
② do self attention  $a: \mathbb{R}^{F'} \times \mathbb{R}^{F'} \rightarrow \mathbb{R}$

importance of node  $j$ 's features to node  $i$ .  
 $e_{ij} = a(W \cdot \bar{h}_i, W \cdot \bar{h}_j)$



③ normalize.  $\alpha_{ij} = \text{softmax}_j(e_{ij}) = \frac{\exp(e_{ij})}{\underbrace{\sum_{k \in N(i)} \exp(e_{i,k})}_{\text{sum of neighbors of } i}}$

④ attention mechanism  $a$ : a single-layer FFNN.



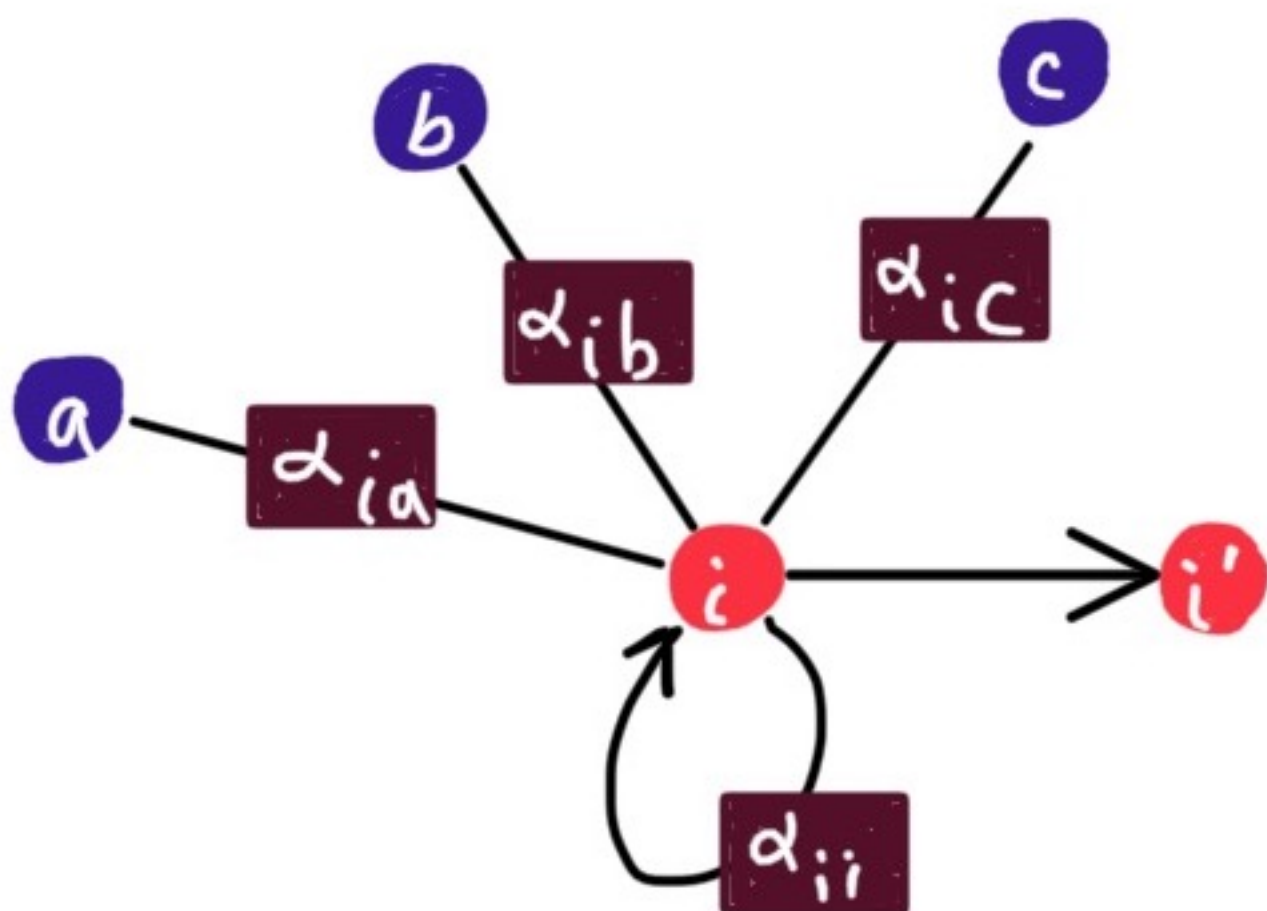
$\text{Leaky ReLU} = \max(0.2x, x)$

Complete formula:

$$\alpha_{ij} = \frac{\exp(\text{LeakyReLU}(\bar{a}^T [w_{h_i}; w_{h_j}]))}{\sum_{k \in N(i)} \exp(\text{LeakyReLU}(\bar{a}^T [w_{h_i}; w_{h_k}]))}$$

$\alpha_{ij} \in \mathbb{R}$

⑤ apply. 
$$h'_i = \sigma \left( \sum_{j \in N(i)} \alpha_{ij} W h_j \right)$$



⑥ multi-head attention

$$h'_i = \parallel_{k=1}^k \sigma \left( \sum_{j \in N(i)} \alpha_{ij}^k W^k h_j \right)$$

$\parallel$ : Concat

$$h'_i = \sigma \left( \frac{1}{k} \sum_{k=1}^k \sum_{j \in N(i)} \alpha_{ij}^k W h_j \right)$$

average on  
the last layer of  
the network



# Message Passing

$$x_i^{(k)} = \gamma^{(k)} \left( x_i^{(k-1)}, \bigoplus_{j \in \mathcal{N}(i)} \phi^{(k)} \left( x_i^{(k-1)}, x_j^{(k-1)}, e_{ji} \right) \right)$$

Diagram annotations:  $\gamma^{(k)}$  is circled in green with an arrow pointing to "MLP".  $\bigoplus_{j \in \mathcal{N}(i)}$  is circled in purple with an arrow pointing to "order-invariant function. For every neighbor  $j$  of node  $i$ . Sum, mean, ...".  $\phi^{(k)}$  is circled in red with an arrow pointing to "MLP".  $x_i^{(k-1)}$  is circled in red with an arrow pointing to "feature repr of node  $i$  at  $(k-1)$ th layer".  $e_{ji}$  is circled in red with an arrow pointing to "features of edge  $(i,j)$  if applicable".

feature repr of  
node  $i$  at  
the  $k$ -th layer

order-invariant  
function. For  
every neighbor  $j$   
of node  $i$ .  
Sum, mean, ...

feature repr  
of node  $i$   
at  $(k-1)$ th  
layer

features of  
edge  $(i,j)$   
if applicable



\* PyTorch Geometric provides MessagePassing class.

aggregate(): agg msgs from neighbors

message(): Construct msg from node  $j$  to  $i$  ( $\phi$ )

propagate(): propagate messages

update(): update node embeddings ( $\chi$ )

## Simple MessagePassing Usage

```
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 ...
```

Compute the message



# GCN

$$x_i^{(k)} = \sum_{j \in \mathcal{N}(i) \cup \{i\}} \frac{1}{\sqrt{\deg(i)} \cdot \sqrt{\deg(j)}} \cdot \left( \underbrace{\theta}_{\phi} \cdot x_j^{(k-1)} \right)$$

## steps to implement

- 1 add self-loops
- 2 a linear transformation to node feature matrix
- 3 Compute normalization coefficient
- 4 normalize node features
- 5 sum up neighboring node features

```
class GCNConv(MessagePassing):
```

```
    def __init__(self, in_channels, out_channels):
```

```
        super(GCNConv, self).__init__(aggr='add')
```

```
        self.lin = torch.nn.Linear(in_channels, out_channels)
```

$N \times \text{in\_channels}$   $\leftarrow$   $2 \times E$

```
    def forward(self, x, edge_idx):
```

- 1 `edge_idx, _ = add_self_loops(edge_idx, num_nodes = x.size(0))`
- 2 `x = self.lin(x)`



3 { row, Col = edge\_idx  
deg = degree(Col, x.size(0), dtype = x.dtype)  
deg\_inv\_sq = deg.pow(-0.5)  
norm = deg\_inv\_sq[row] \* deg\_inv\_sq[Col]

4 5 return self.propagate(edge\_idx, x=x, norm=norm)

→ Ex out - channels  
def message(self, x\_j, norm):  
4 return norm.view(-1, 1) \* x\_j

# GAN

## Simplified GAN layer:

```
class GATLayer(nn.Module):  
    def __init__(self, in_feats, out_feats, dropout, alpha,  
true for all but last layer ← Concat = True):  
        super(GATLayer, self).__init__() ↓  
slope of  
LeakyReLU  
        # get all input args ...  
  
        self.W = nn.Parameter(torch.zeros((in_feats, out_feats)))  
        nn.init.xavier_uniform_(self.W.data, gain=1.414)  
  
        self.a = nn.Parameter(torch.zeros((2 * out_feats, 1)))  
        nn.init.xavier_uniform_(self.a.data, gain=1.414)  
  
        self.leakyrelu = nn.LeakyReLU(self.alpha)
```

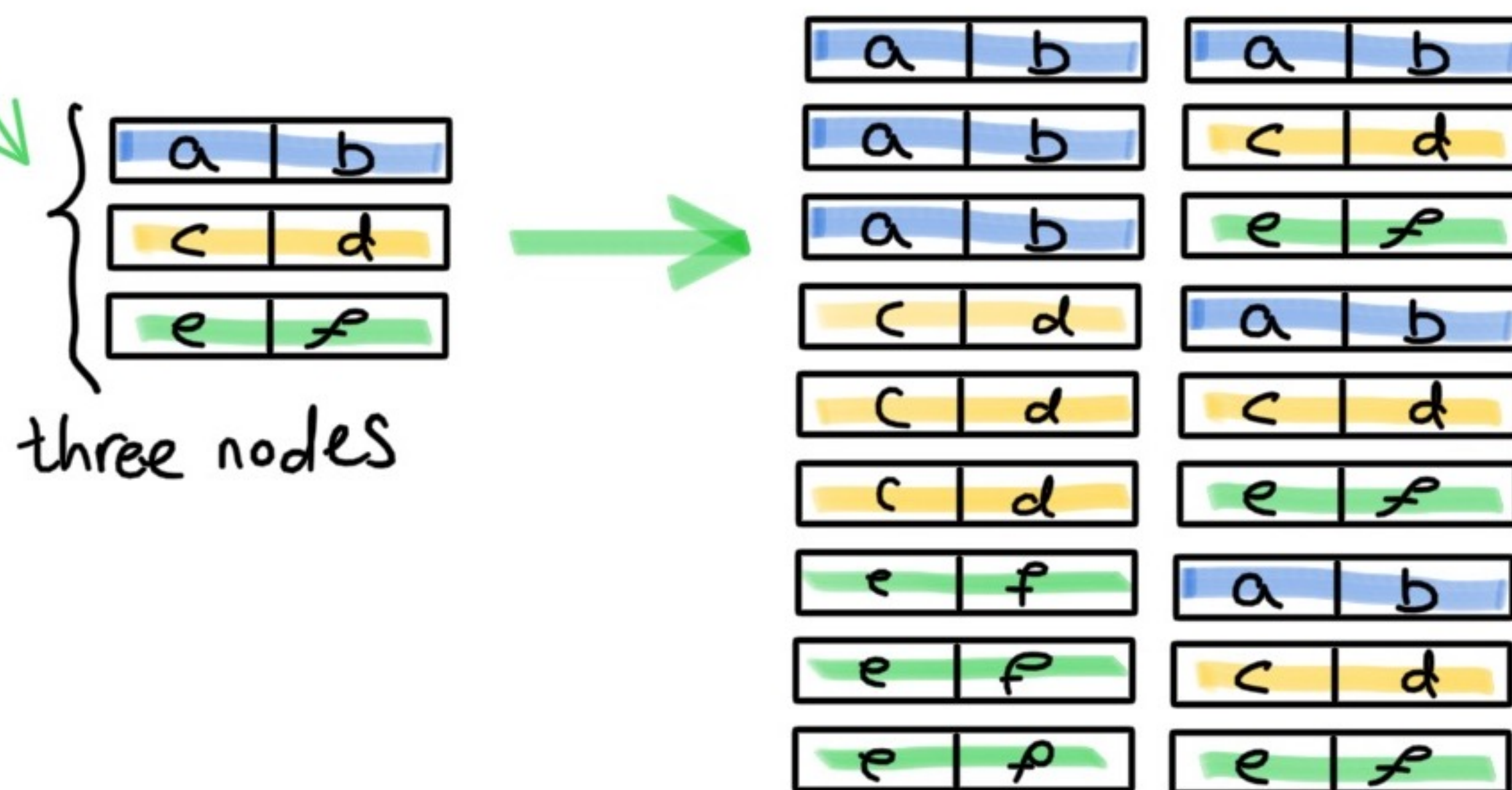


```
def forward(self, input, adj):
```

```
h = torch.mm(input, self.W) → linear transformation
```

```
a_inp = torch.cat([h.repeat(1, N).view(N * N, -1),  
                    h.repeat(N, 1)], dim = 1)  
                .view(N, -1, 2 * self.out_feats)
```

attn  
mechanism



```
e = self.leakyrelu(torch.matmul(a_inp, self.a).squeeze(2))
```

masked  
attn

```
zero_vec = -9e15 * torch.ones_like(e)  
attn = torch.where(adj > 0, e, zero_vec)
```

```
attn = F.softmax(attn, dim = 1)
```

```
attn = F.dropout(attn, self.dropout, training = self.training)
```

```
h_prime = torch.matmul(attn, h)
```

```
if self.Concat: return F.elu(h_prime)
```

```
else: return h_prime
```



Fortunately, we don't need to implement GAN layers.

```
from torch_geometric.data import Data
```

```
from torch_geometric.nn import GATConv
```

```
from torch_geometric.datasets import Planetoid
```

```
dataset = Planetoid(root='tmp/', name='Cora')
```

```
class GAT(torch.nn.Module):
```

```
    def __init__(self, GAT):
```

```
        self.hid = 8
```

```
        self.in-head = 8
```

```
        self.out-head = 1
```

```
        self.conv1 = GATConv(dataset.num_features, self.hid,  
                             heads=self.in-head, dropout=.6)
```

```
        self.conv2 = GATConv(self.hid * self.in-head,  
                             dataset.num_classes,  
                             Concat=False, heads=self.out-head,  
                             dropout=.6)
```



```
def forward(self, data):
```

```
    x, edge_idx = data.x, data.edge_index
```

```
    x = F.dropout(x, p=.6, training=self.training)
```

```
    x = self.conv1(x, edge_idx)
```

```
    x = F.relu(x)
```

```
    x = F.dropout(x, p=.6, training=self.training)
```

```
    x = self.conv2(x, edge_idx)
```

```
    return F.log_softmax(x, dim=1)
```

```
model = GAT().to(device)
```

```
data = dataset[0].to(device)
```

```
opt = torch.optim.Adam(model.parameters(), lr=.005,  
                        weight_decay=5e-4)
```



for epoch in range (epochs):

model.train()

opt.zero\_grad()

out = model(data)

loss = F.nll\_loss(out[data.train\_mask],  
data.y[data.train\_mask])

loss.backward()

opt.step()