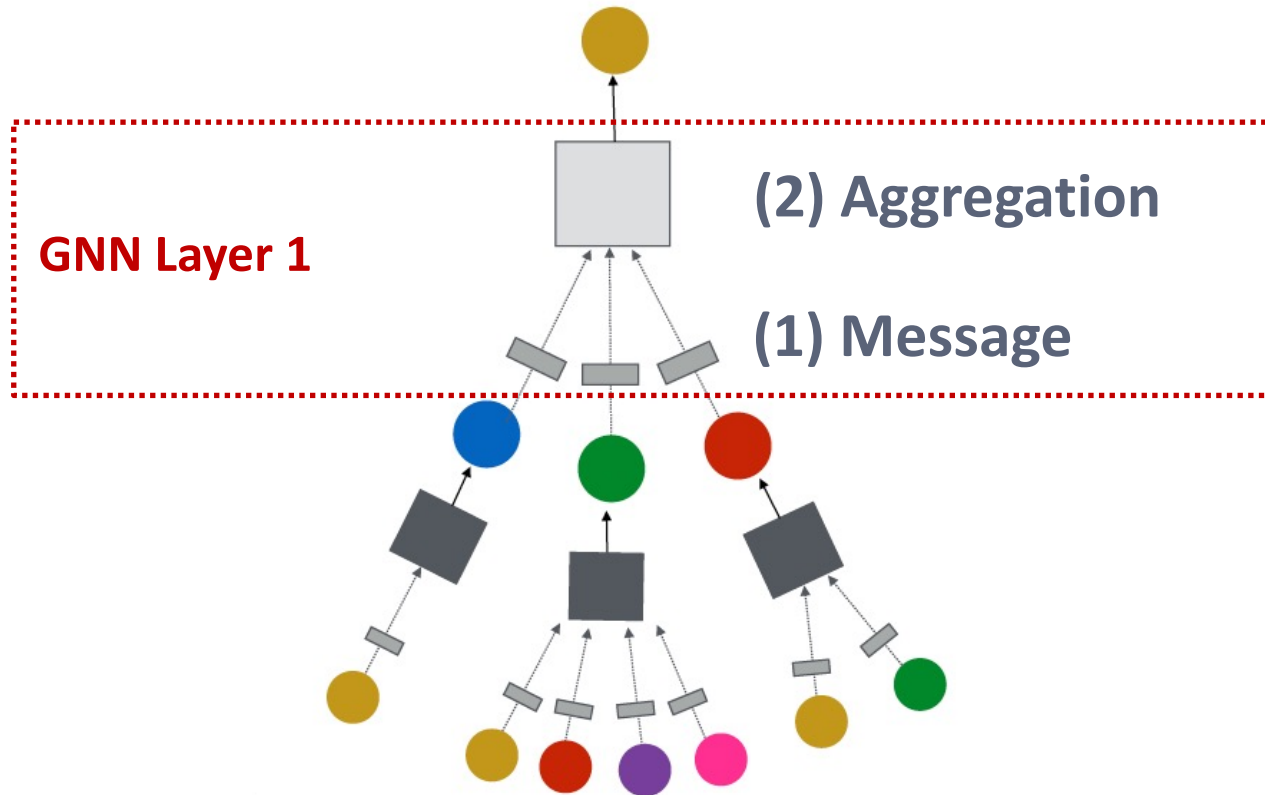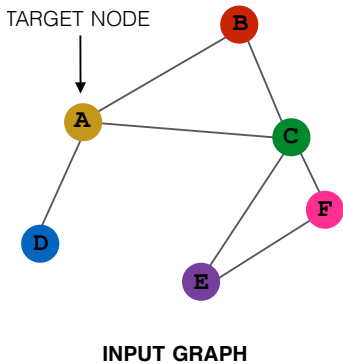# A Single Layer of a GNN

# A GNN Layer

**GNN Layer = Message + Aggregation**

- **Different instantiations under this perspective**
- **GCN, GraphSAGE, GAT, …**

TARGET NODE

INPUT GRAPH

**GNN Layer 1**

**(2) Aggregation**

**(1) Message**
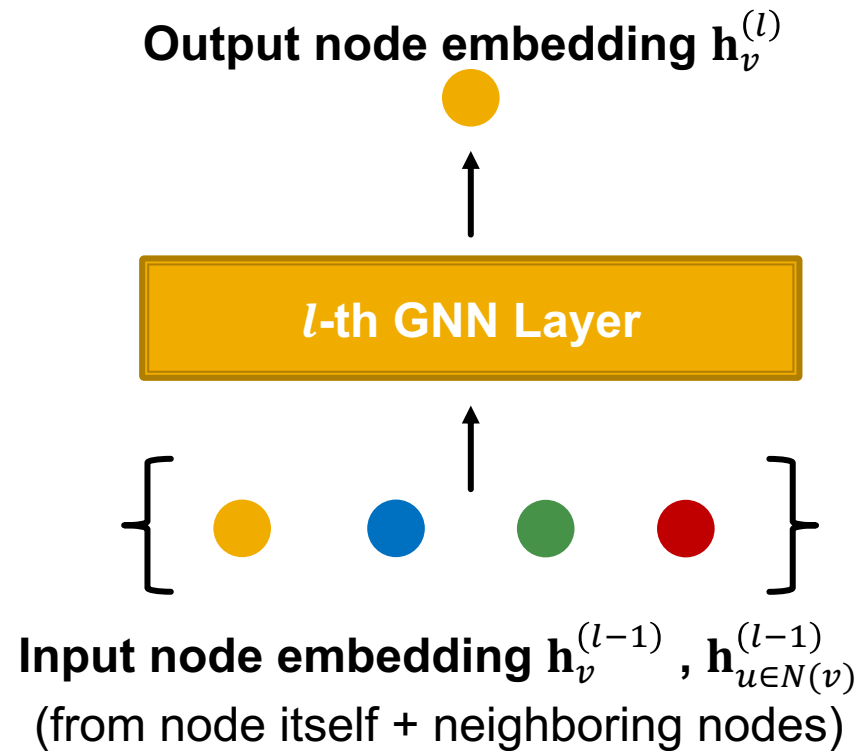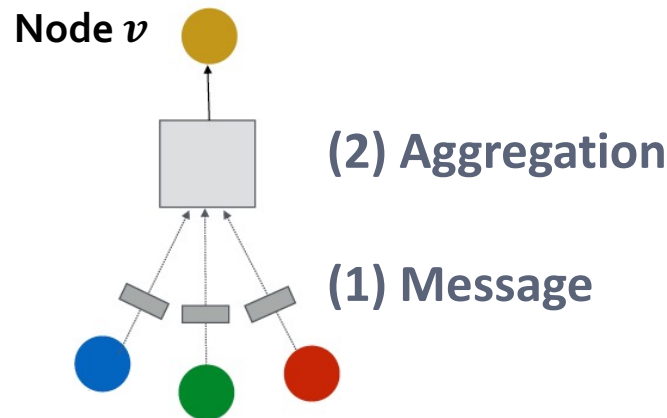
# A Single GNN Layer

- ## Idea of a GNN Layer:

  - Compress a set of vectors into a single vector

  - **Two-step process:**
    - **(1) Message**
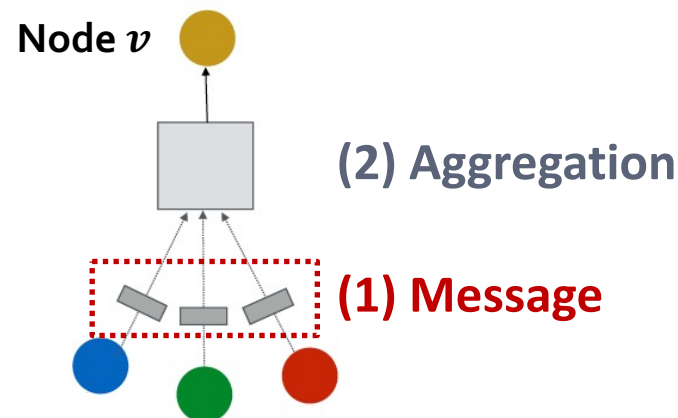    - **(2) Aggregation**

Node $v$

(2) Aggregation

(1) Message

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

$l$-th GNN Layer

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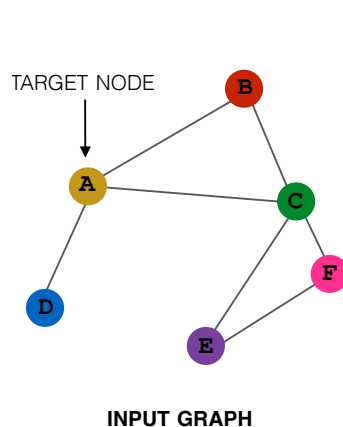(from node itself + neighboring nodes)

# Message Computation

- ## (1) Message computation

  - ### Message function: $\mathbf{m}_u^{(l)} = \mathrm{MSG}^{(l)}\left(\mathbf{h}_u^{(l-1)}\right)$

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

    - **Example:** A Linear layer $\mathbf{m}_u^{(l)} = \mathbf{W}^{(l)}\mathbf{h}_u^{(l-1)}$

      - Multiply node features with weight matrix $\mathbf{W}^{(l)}$



TARGET NODE

**INPUT GRAPH**

Node $v$

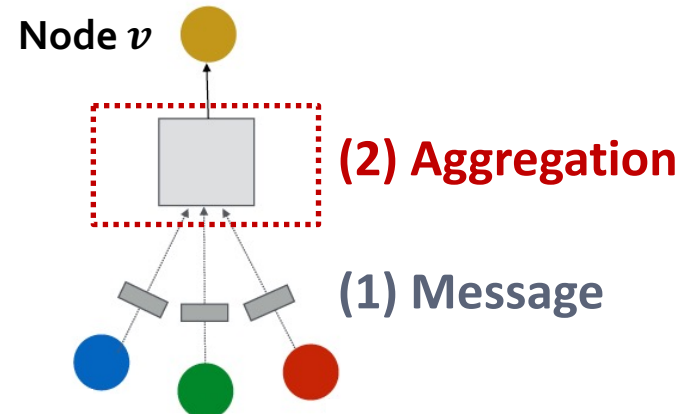(2) Aggregation

(1) Message

# Message Aggregation

- ## (2) Aggregation

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

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

  - **Example:** $\text{Sum}(\cdot), \text{Mean}(\cdot)$ or $\text{Max}(\cdot)$ aggregator

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



TARGET NODE

**INPUT GRAPH**

Node $v$

(2) Aggregation

(1) Message

# Message Aggregation: Issue

- **Issue:** Information from node $v$ itself **could get lost**

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

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

  - **(1) Message: compute message from node $v$ itself**

    - Usually, a **different message computation** will be performed

      🔵🟢🔴  $\mathbf{m}_u^{(l)} = \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)}$       🟠  $\mathbf{m}_v^{(l)} = \mathbf{B}^{(l)} \mathbf{h}_v^{(l-1)}$

  - **(2) Aggregation:** After aggregating from neighbors, we can **aggregate the message from node $v$ itself**

    - Via **concatenation** or **summation**

      **Then aggregate from node itself**

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

      **First aggregate from neighbors**

# A Single GNN Layer

- **Putting things together:**

  - **(1) Message**: each node computes a message
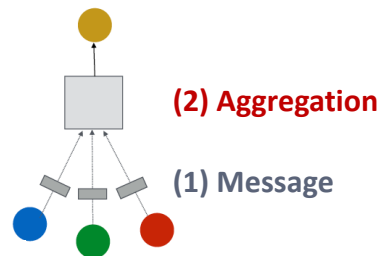
    $$\mathbf{m}_u^{(l)} = \mathrm{MSG}^{(l)}\left(\mathbf{h}_u^{(l-1)}\right), u \in \{N(v) \cup v\}$$

  - **(2) Aggregation**: aggregate messages from neighbors

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

  - **Nonlinearity (activation):** Adds expressiveness

    - Often written as $\sigma(\cdot)$: $\mathrm{ReLU}(\cdot)$, $\mathrm{Sigmoid}(\cdot)$, …

    - Can be added to **message or aggregation**



**(2) Aggregation**

**(1) Message**

# Classical GNN Layers: GCN (1)

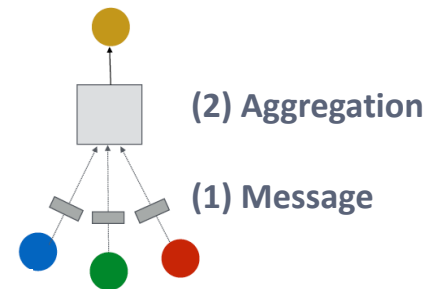- **(1) Graph Convolutional Networks (GCN)**

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

- **How to write this as Message + Aggregation?**

**Message**

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

**Aggregation**

(2) Aggregation

(1) Message
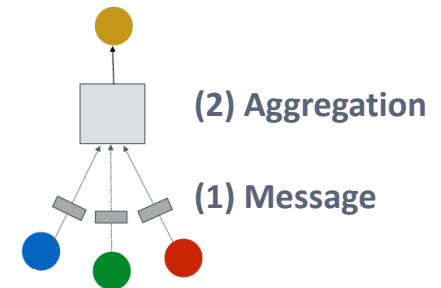
■ **(1) Graph Convolutional Networks (GCN)**

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

**(2) Aggregation**

**(1) Message**

■ **Message:**

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

**Normalized by node degree**
(In the GCN paper they use a slightly different normalization)

■ **Aggregation:**

■ **Sum** over messages from neighbors, then apply activation

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

In GCN graph is assumed to have self-edges that are included in the summation.

# Classical GNN Layers: GraphSAGE

- **(2) GraphSAGE**

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

- **How to write this as Message + Aggregation?**

  - **Message** is computed within the $\text{AGG}(\cdot)$

  - **Two-stage aggregation**

    - **Stage 1:** Aggregate from node neighbors
    $$\mathbf{h}_{N(v)}^{(l)} \leftarrow \text{AGG}\left(\left\{\mathbf{h}_u^{(l-1)}, \forall u \in N(v)\right\}\right)$$

    - **Stage 2:** Further aggregate over the node itself
    $$\mathbf{h}_v^{(l)} \leftarrow \sigma\left(\mathbf{W}^{(l)} \cdot \text{CONCAT}(\mathbf{h}_v^{(l-1)}, \mathbf{h}_{N(v)}^{(l)})\right)$$

# GraphSAGE Neighbor Aggregation

- **Mean:** Take a weighted average of neighbors

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

**Aggregation**     **Message computation**

- **Pool:** Transform neighbor vectors and apply symmetric vector function $\text{Mean}(\cdot)$ or $\text{Max}(\cdot)$

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

**Aggregation**     **Message computation**

- **LSTM:** Apply LSTM to reshuffled of neighbors

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

**Aggregation**

# GraphSAGE: L2 Normalization

- ### $\ell_2$ **Normalization:**

  - **Optional:** Apply $\ell_2$ normalization to $\mathbf{h}_v^{(l)}$ at every layer

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

  - Without $\ell_2$ normalization, the embedding vectors have different scales ($\ell_2$-norm) for vectors

  - In some cases (not always), normalization of embedding results in performance improvement

  - After $\ell_2$ normalization, all vectors will have the same $\ell_2$-norm

- **(3) Graph Attention Networks**

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

**Attention weights**

- **In GCN / GraphSAGE**

  - $\alpha_{vu} = \frac{1}{|N(v)|}$ is the **weighting factor (importance)** of node $u$'s message to node $v$

  - $\Longrightarrow \alpha_{vu}$ is defined **explicitly** based on the structural properties of the graph (node degree)

  - $\Longrightarrow$ All neighbors $u \in N(v)$ are equally important to node $v$

- **(3) Graph Attention Networks**

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

**Attention weights**

## Not all node's neighbors are equally important

- **Attention** is inspired by cognitive attention.
- The **attention** $\alpha_{vu}$ focuses on the important parts of the input data and fades out the rest.
  - **Idea:** the NN should devote more computing power on that small but important part of the data.
  - Which part of the data is more important depends on the context and is learned through training.

# Graph Attention Networks

**Can we do better than simple neighborhood aggregation?**

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

- **Goal:** Specify **arbitrary importance** to different neighbors of each node in the graph
- **Idea:** Compute embedding $\boldsymbol{h}_v^{(l)}$ of each node in the graph following an **attention strategy**:
  - Nodes attend over their neighborhoods' message
  - Implicitly specifying different weights to different nodes in a neighborhood
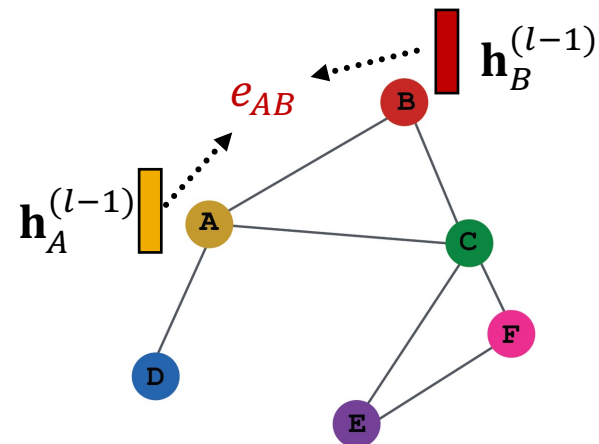
# Attention Mechanism (1)

- Let $\alpha_{vu}$ be computed as a byproduct of an **attention mechanism $a$:**

  - (1) Let $a$ compute **attention coefficients $e_{vu}$** across pairs of nodes $u$, $v$ based on their messages:

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

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

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

# Attention Mechanism (2)

- **Normalize** $e_{vu}$ into the **final attention weight** $\boldsymbol{\alpha_{vu}}$
  - Use the **softmax** function, so that $\sum_{u \in N(v)} \alpha_{vu} = 1$:

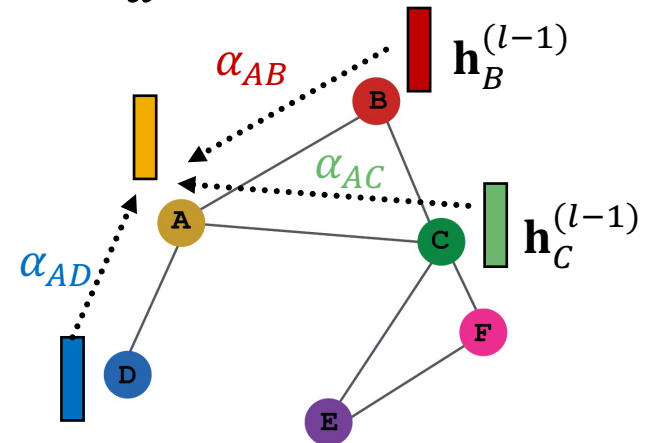$$\alpha_{vu} = \frac{\exp(e_{vu})}{\sum_{k \in N(v)} \exp(e_{vk})}$$

- **Weighted sum** based on the **final attention weight** $\boldsymbol{\alpha_{vu}}$

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

**Weighted sum using** $\alpha_{AB}, \alpha_{AC}, \alpha_{AD}$:
$\mathbf{h}_A^{(l)} = \sigma(\alpha_{AB} \mathbf{W}^{(l)} \mathbf{h}_B^{(l-1)} + \alpha_{AC} \mathbf{W}^{(l)} \mathbf{h}_C^{(l-1)} + \alpha_{AD} \mathbf{W}^{(l)} \mathbf{h}_D^{(l-1)})$
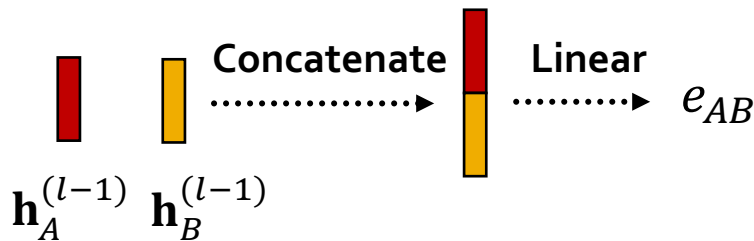
# Attention Mechanism (3)

- **What is the form of attention mechanism $a$?**

  - The approach is agnostic to the choice of $a$

    - E.g., use a simple single-layer neural network

      - $a$ have trainable parameters (weights in the Linear layer)



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

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

  - Parameters of $a$ are trained jointly:

    - Learn the parameters together with weight matrices (i.e., other parameter of the neural net $\mathbf{W}^{(l)}$) in an end-to-end fashion

# Attention Mechanism (4)

- **Multi-head attention:** Stabilizes the learning process of attention mechanism

  - **Create multiple attention scores** (each replica with a different set of parameters):

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

  $$\mathbf{h}_v^{(l)}[2] = \sigma(\sum_{u \in N(v)} \alpha_{vu}^2 \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)})$$

  $$\mathbf{h}_v^{(l)}[3] = \sigma(\sum_{u \in N(v)} \alpha_{vu}^3 \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)})$$

  - **Outputs are aggregated:**

    - By concatenation or summation

    - $\mathbf{h}_v^{(l)} = \text{AGG}(\mathbf{h}_v^{(l)}[1], \mathbf{h}_v^{(l)}[2], \mathbf{h}_v^{(l)}[3])$

# Benefits of Attention Mechanism

- **Key benefit:** Allows for (implicitly) specifying **different importance values $(\alpha_{vu})$ to different neighbors**

- **Computationally efficient**:
  - Computation of attentional coefficients can be parallelized across all edges of the graph
  - Aggregation may be parallelized across all nodes
- **Storage efficient**:
  - Sparse matrix operations do not require more than $O(V + E)$ entries to be stored
  - **Fixed** number of parameters, irrespective of graph size
- **Localized**:
  - Only **attends over local network neighborhoods**
- **Inductive capability**:
  - It is a shared *edge-wise* mechanism
  - It does not depend on the global graph structure