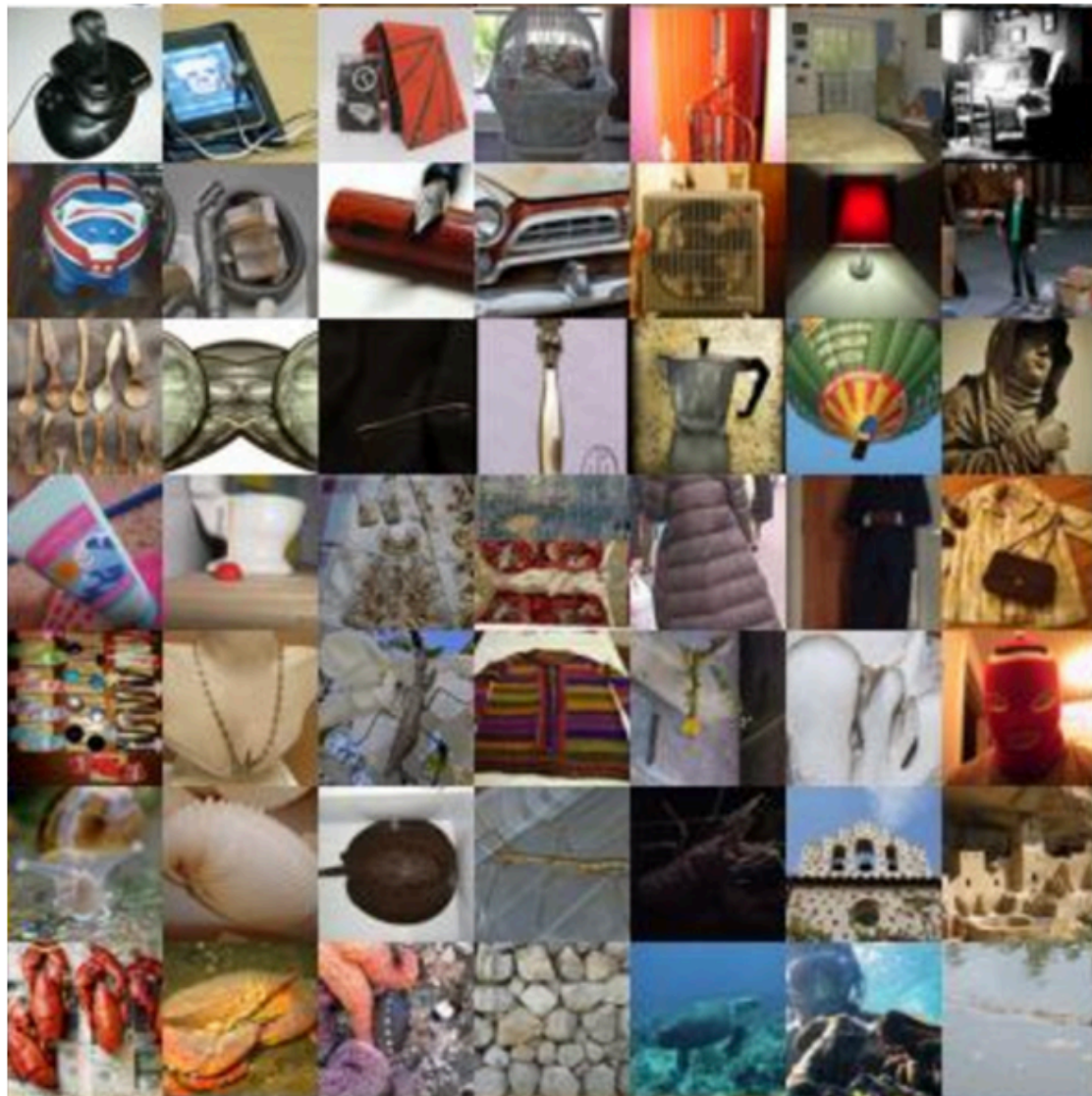# Topics in AI (CPSC 532S):
# Multimodal Learning with Vision, Language and Sound
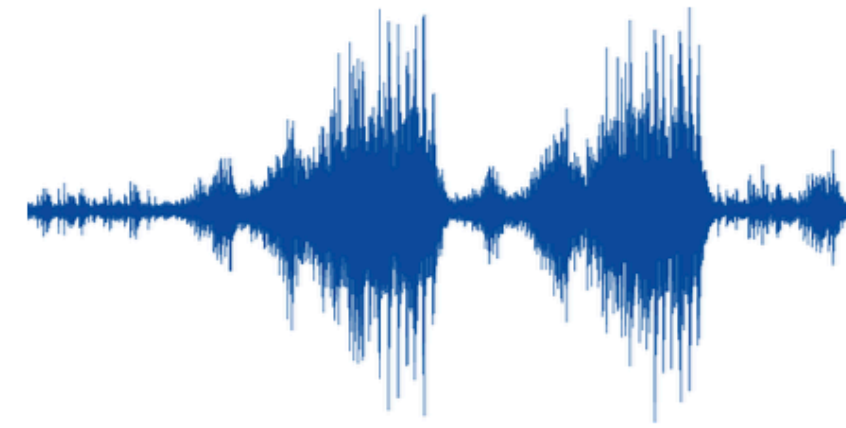
**Lecture 18: Graph Neural Networks**

# Traditional **Neural Networks**
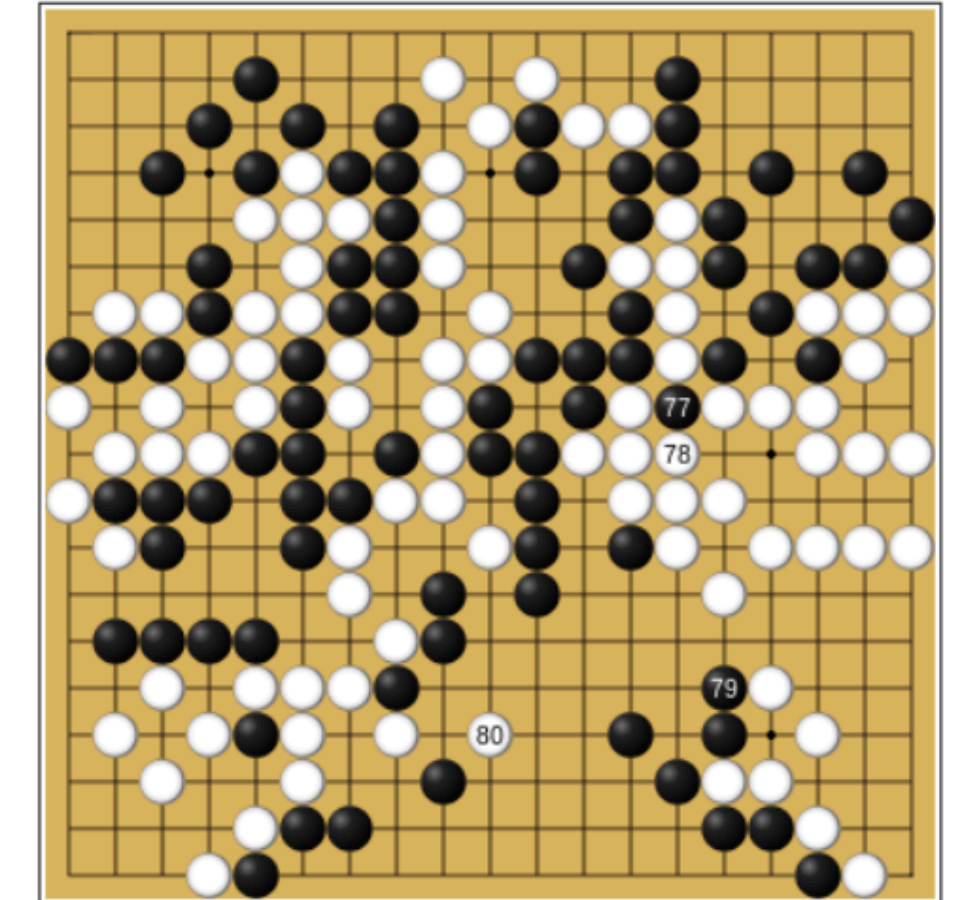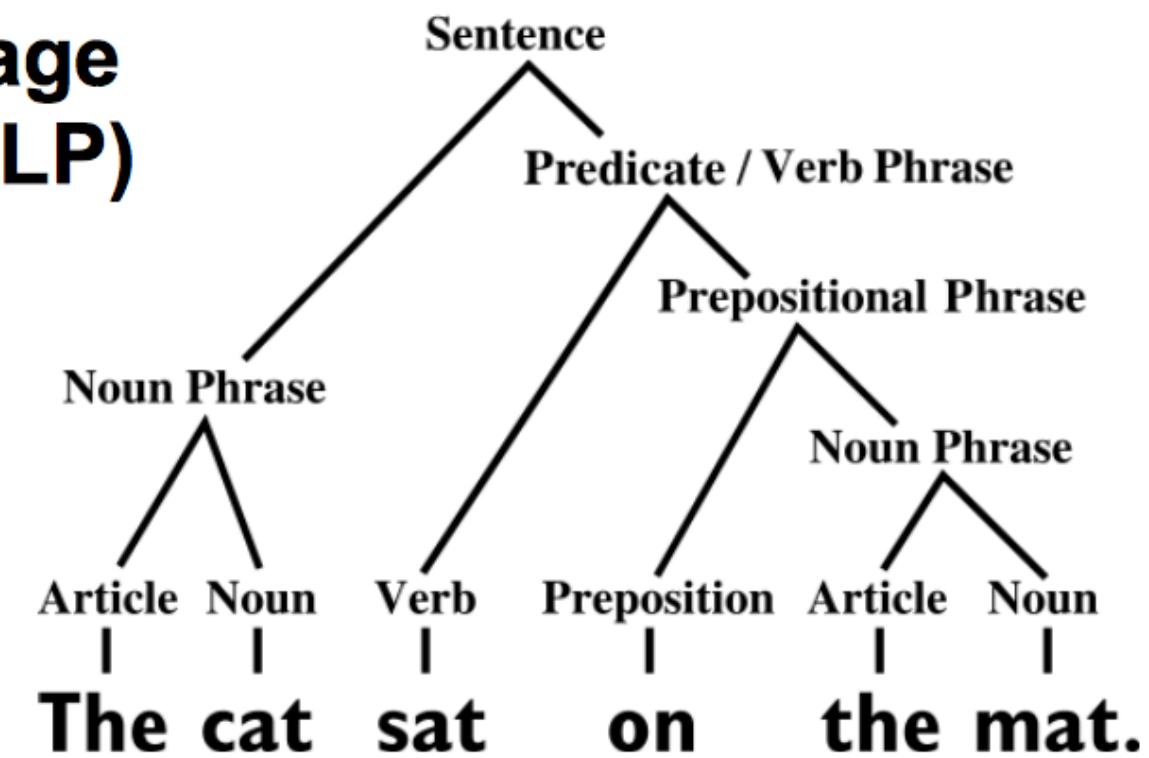
IM🔲GENET

**Speech data**

**Grid games**

**Natural language processing (NLP)**

Sentence

Predicate / Verb Phrase

Prepositional Phrase

Noun Phrase

Noun Phrase

Article  Noun  Verb  Preposition  Article  Noun

**The cat  sat  on  the mat.**

…

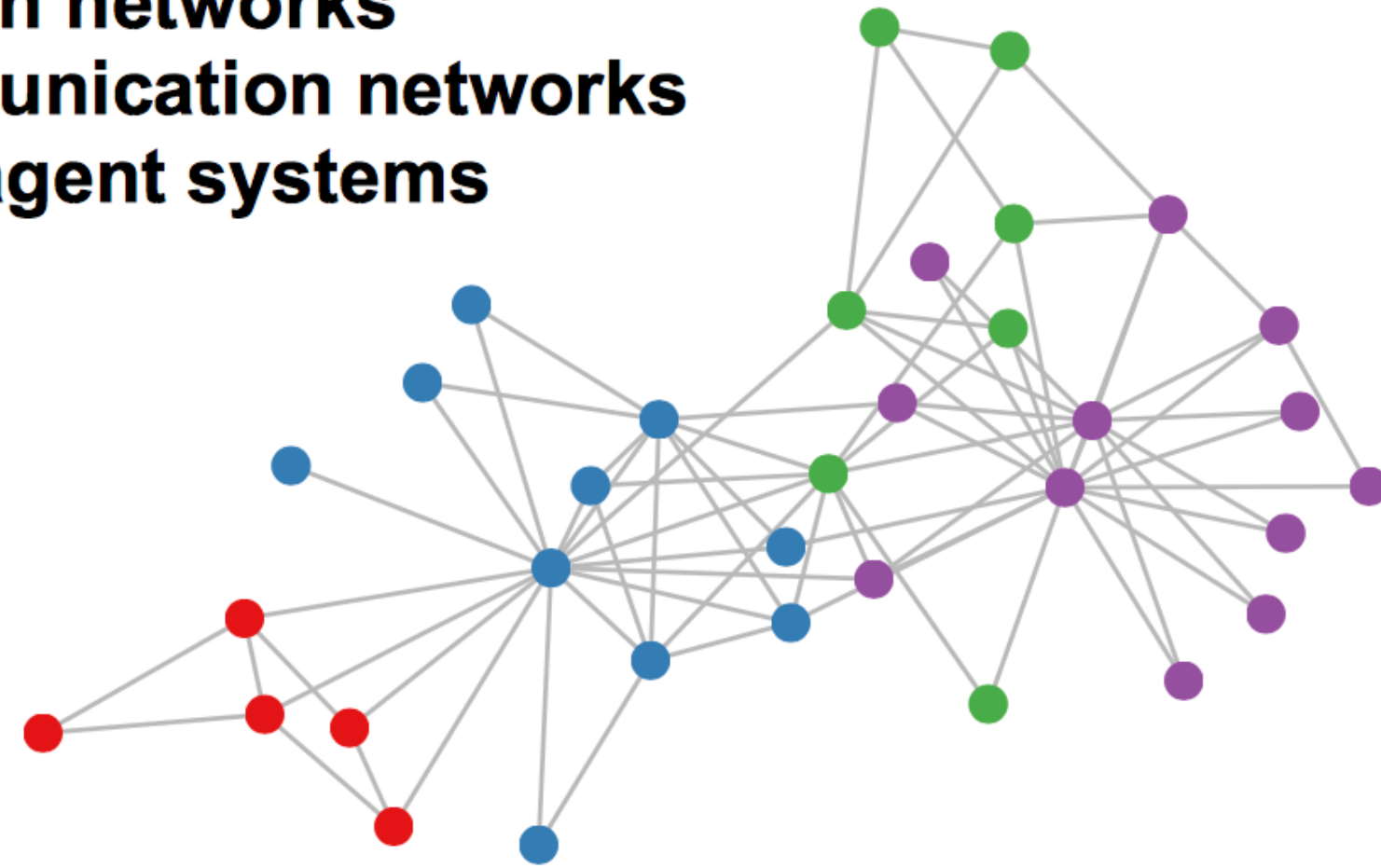**Deep neural nets that exploit:**

- translation equivariance (weight sharing)
- hierarchical compositionality

# **Graph**-structured Data

A lot of real-world data does not "live" on grids

**Social networks**
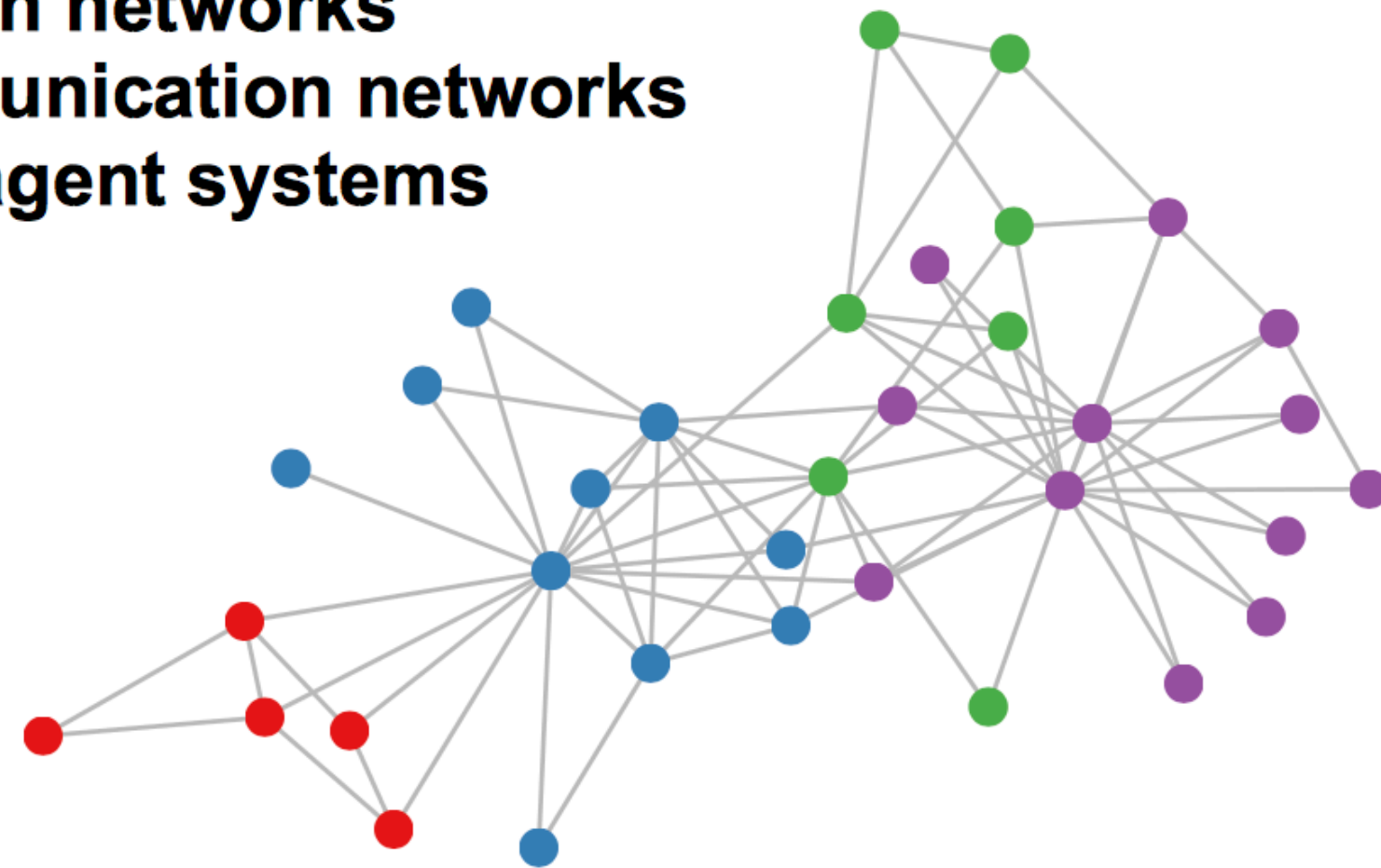**Citation networks**
**Communication networks**
**Multi-agent systems**

# **Graph**-structured Data

A lot of real-world data does not "live" on grids

**Social networks**
**Citation networks**
**Communication networks**
**Multi-agent systems**

**Protein interaction networks**

# **Graph**-structured Data

A lot of real-world data does not "live" on grids

**Social networks**
**Citation networks**
**Communication networks**
**Multi-agent systems**



Knowledge graphs

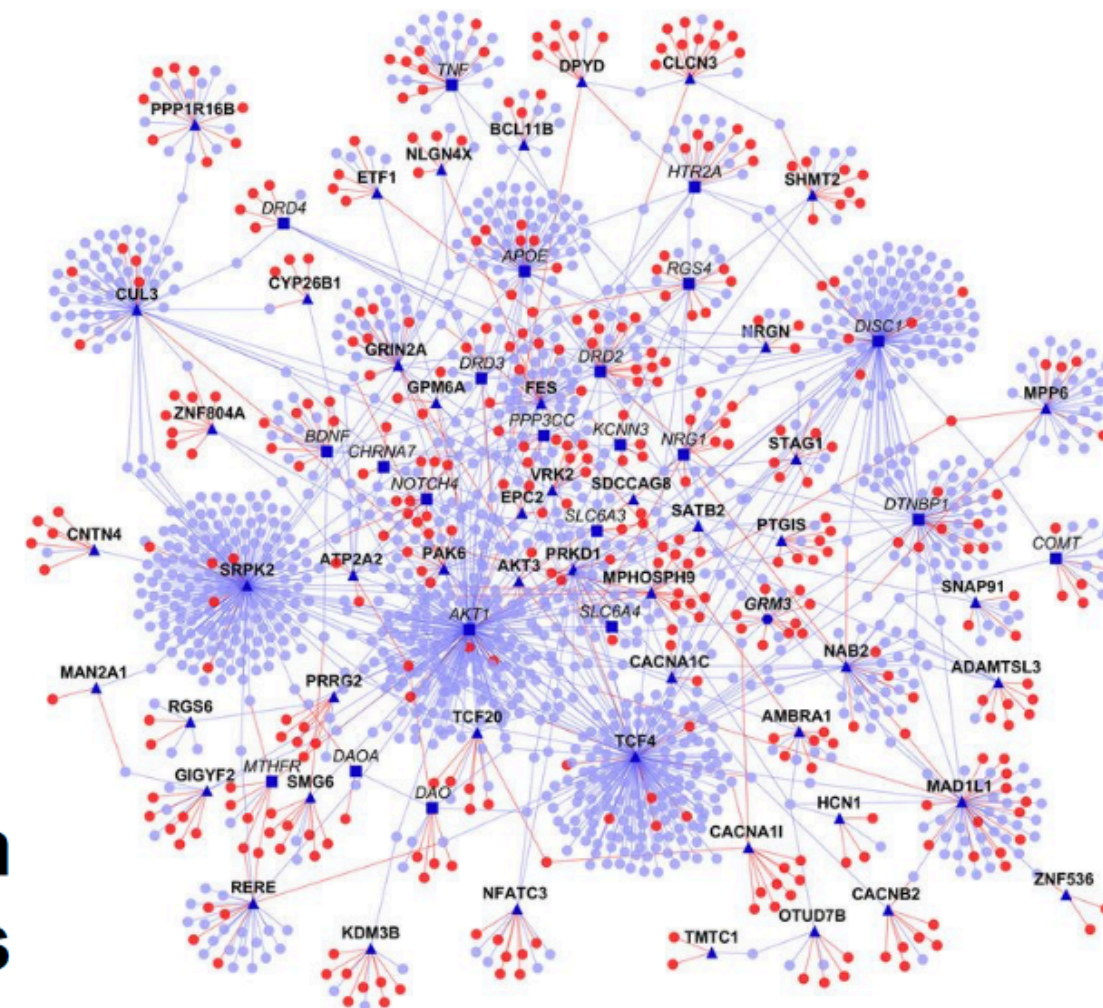Protein interaction
networks

# **Graph**-structured Data

A lot of real-world data does not "live" on grids

**Social networks**
**Citation networks**
**Communication networks**
**Multi-agent systems**

**Knowledge graphs**

**Molecules**

**Protein interaction networks**

**Road maps**

# **Graph**-structured Data

A lot of real-world data does not "live" on grids

**Social networks**
**Citation networks**
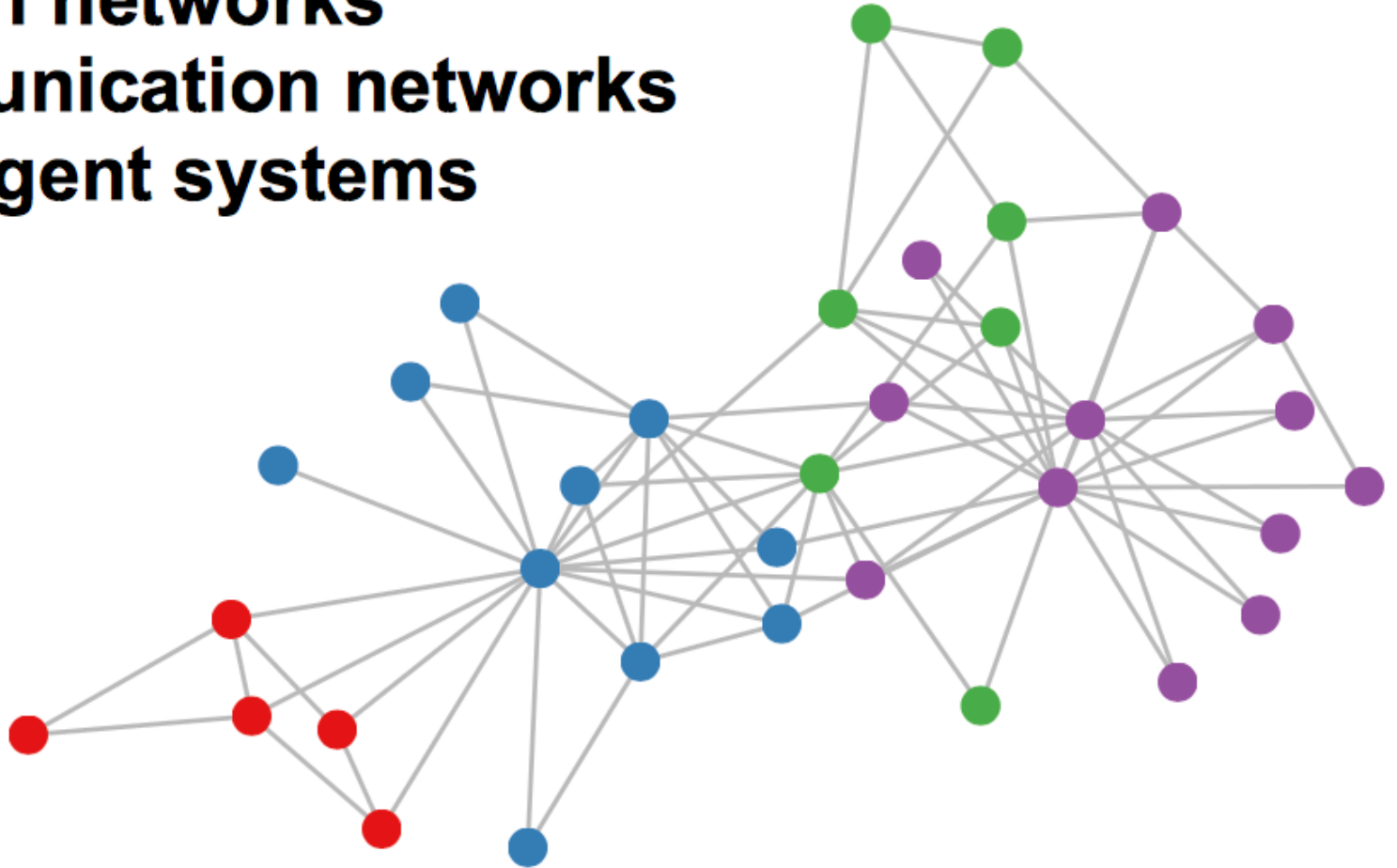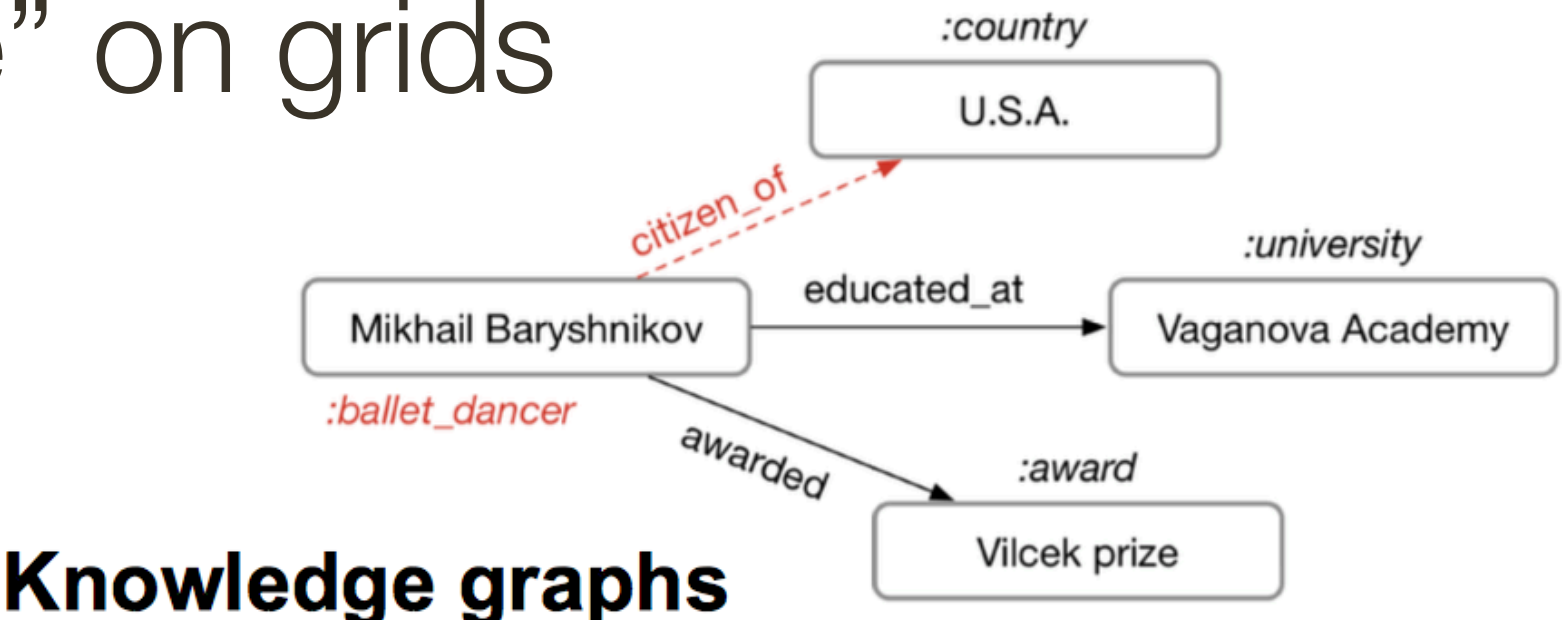**Communication networks**
**Multi-agent systems**

**Knowledge graphs**

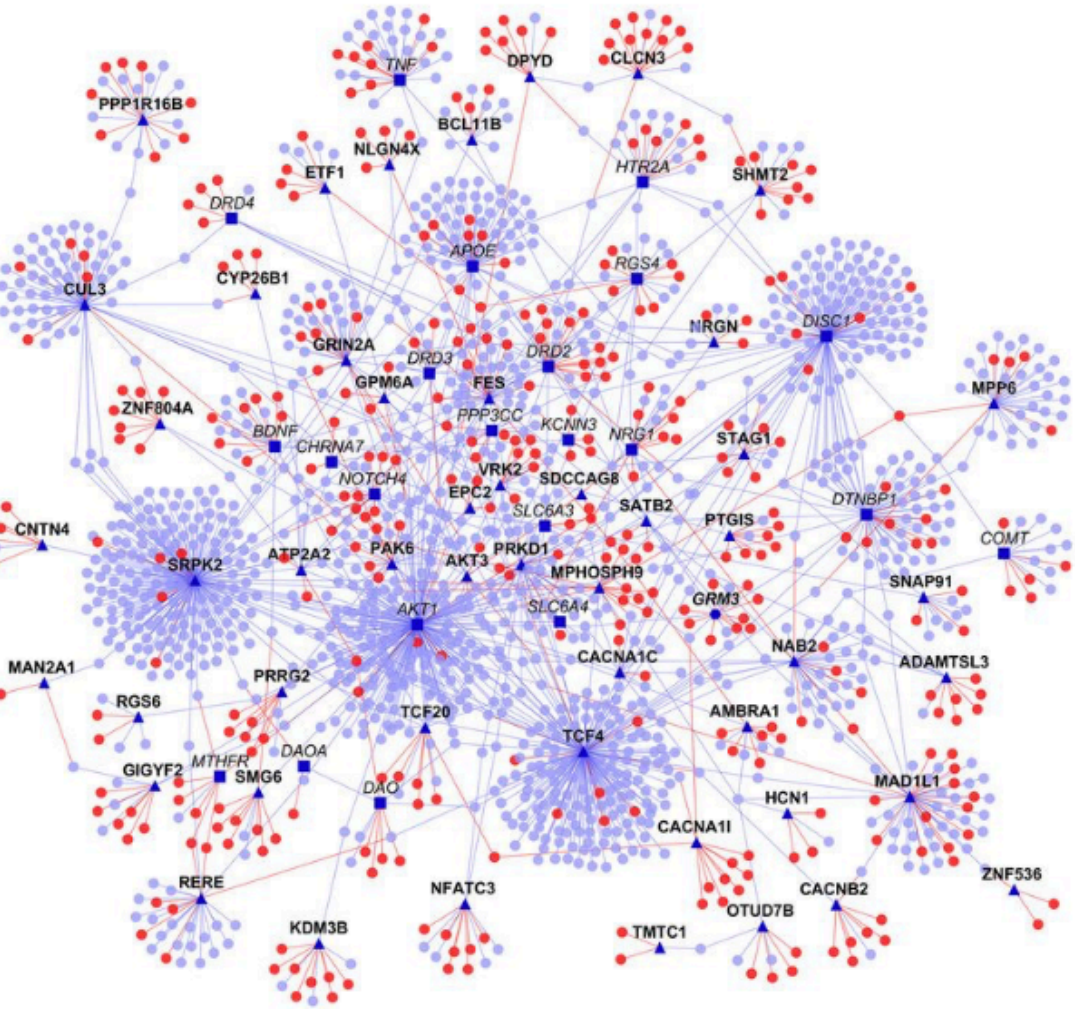**Molecules**

**Protein interaction networks**

**Road maps**

Standard **CNN** and **RNN** architectures don't work on this data

# Graph Neural Networks (GNNs)



**Main Idea**: Pass massages between pairs of nodes and agglomerate

**Alternative Interpretation**: Pass massages between nodes to refine node (and possibly edge) representations

# Graph Neural Networks (GNNs)

**Notation:** $\mathcal{G} = (\mathbf{A}, \mathbf{X})$

- Adjacency matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$
- Feature matrix $\mathbf{X} \in \mathbb{R}^{N \times F}$



**Main Idea**: Pass massages between pairs of nodes and agglomerate

**Alternative Interpretation**: Pass massages between nodes to refine node (and possibly edge) representations

* slide from Thomas Kipf, **University of Amsterdam**

# Recap: Convolutional Neural Networks (CNNs) on Grids

**Single CNN layer with 3x3 filter:**



(Animation by Vincent Dumoulin)

$\mathbf{h}_0$ $\mathbf{h}_1$ ...

$\mathbf{h}_i$

# Recap: Convolutional Neural Networks (CNNs) on Grids

**Single CNN layer with 3x3 filter:**



(Animation by Vincent Dumoulin)

$\mathbf{h}_0 \quad \mathbf{h}_1 \quad \ldots$

$\mathbf{h}_i$

$\mathbf{h}_i \in \mathbb{R}^F$ are (hidden layer) activations of a pixel/node

# Recap: Convolutional Neural Networks (CNNs) on Grids

**Single CNN layer with 3x3 filter:**



(Animation by Vincent Dumoulin)

$\mathbf{h}_0$     $\mathbf{h}_1$     ...

$\mathbf{h}_i$

**Update for a single pixel:**

- Transform messages individually $\mathbf{W}_i \mathbf{h}_i$

- Add everything up $\sum_i \mathbf{W}_i \mathbf{h}_i$

$\mathbf{h}_i \in \mathbb{R}^F$ are (hidden layer) activations of a pixel/node

# **Recap:** Convolutional Neural Networks (CNNs) on Grids

**Single CNN layer with 3x3 filter:**



(Animation by Vincent Dumoulin)

$\mathbf{h}_0 \quad \mathbf{h}_1 \quad \ldots$

$\mathbf{h}_i$

**Update for a single pixel:**

- Transform messages individually $\mathbf{W}_i \mathbf{h}_i$
- Add everything up $\sum_i \mathbf{W}_i \mathbf{h}_i$

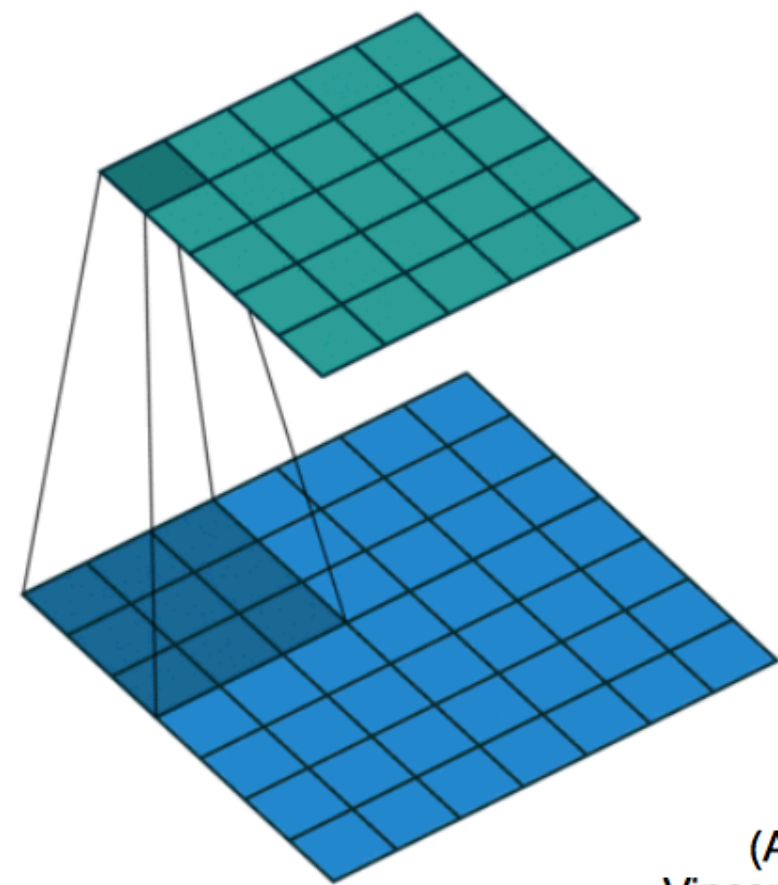$\mathbf{h}_i \in \mathbb{R}^F$ are (hidden layer) activations of a pixel/node

**Full update:**

$$\mathbf{h}_4^{(l+1)} = \sigma \left( \mathbf{W}_0^{(l)} \mathbf{h}_0^{(l)} + \mathbf{W}_1^{(l)} \mathbf{h}_1^{(l)} + \cdots + \mathbf{W}_8^{(l)} \mathbf{h}_8^{(l)} \right)$$

# Graph Convolutional Networks (GCNs)

Kipf & Welling (ICLR 2017), related previous works by Duvenaud et al. (NIPS 2015) and Li et al. (ICLR 2016)

Consider this
undirected graph:

# Graph Convolutional Networks (GCNs)

Kipf & Welling (ICLR 2017), related previous works by Duvenaud et al. (NIPS 2015) and Li et al. (ICLR 2016)



Consider this undirected graph:

Calculate update for node in red:

# Graph Convolutional Networks (GCNs)

Kipf & Welling (ICLR 2017), related previous works by Duvenaud et al. (NIPS 2015) and Li et al. (ICLR 2016)

Consider this undirected graph:

Calculate update for node in red:



**Update rule:**

$$\mathbf{h}_i^{(l+1)} = \sigma\left(\mathbf{h}_i^{(l)}\mathbf{W}_0^{(l)} + \sum_{j \in \mathcal{N}_i} \frac{1}{c_{ij}}\mathbf{h}_j^{(l)}\mathbf{W}_1^{(l)}\right)$$

**Scalability: subsample messages** [Hamilton et al., NIPS 2017]

$\mathcal{N}_i$ : neighbor indices       $c_{ij}$ : norm. constant (fixed/trainable)

# Graph Convolutional Networks (GCNs)

Kipf & Welling (ICLR 2017), related previous works by Duvenaud et al. (NIPS 2015) and Li et al. (ICLR 2016)

Consider this
undirected graph:

Calculate update
for node in red:

**Desirable properties:**

- Weight sharing over all locations
- Invariance to permutations
- Linear complexity O(E)
- Applicable both in transductive
  and inductive settings
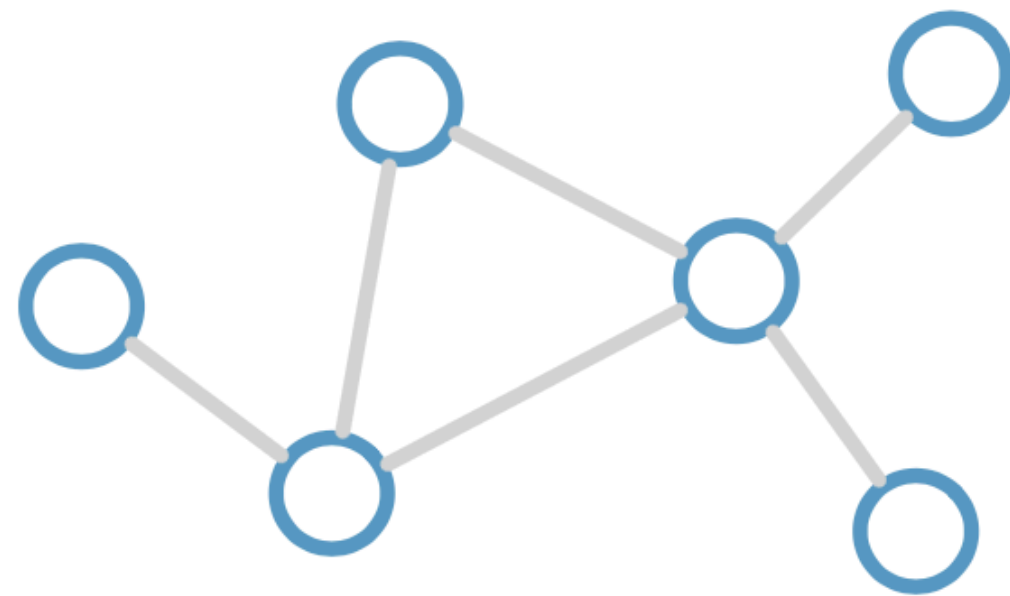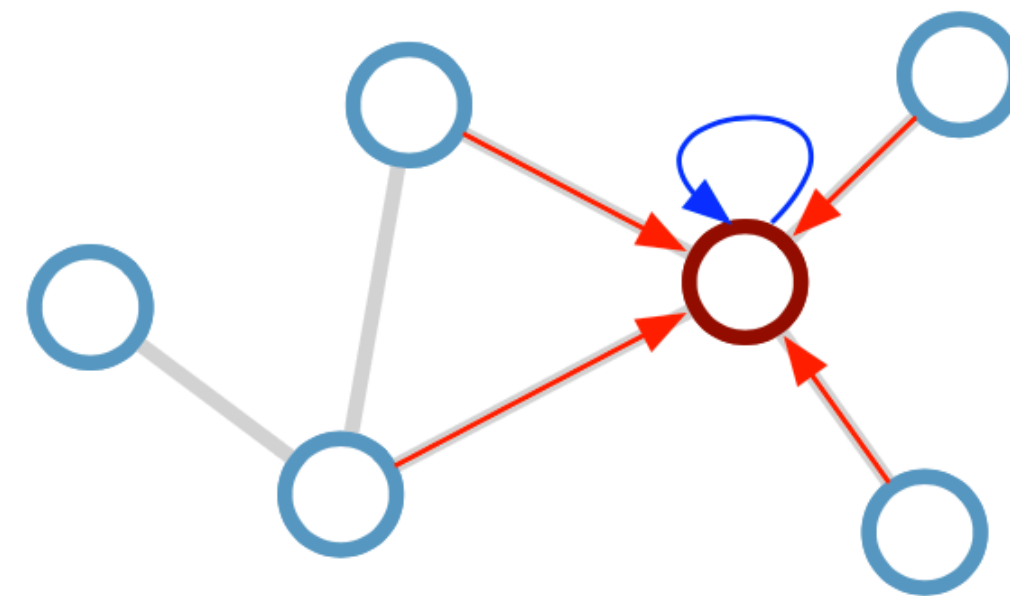


**Update
rule:**

$$\mathbf{h}_i^{(l+1)} = \sigma \left( \mathbf{h}_i^{(l)} \mathbf{W}_0^{(l)} + \sum_{j \in \mathcal{N}_i} \frac{1}{c_{ij}} \mathbf{h}_j^{(l)} \mathbf{W}_1^{(l)} \right)$$

**Scalability: subsample messages** [Hamilton et al., NIPS 2017]

$\mathcal{N}_i$ : neighbor indices    $c_{ij}$ : norm. constant
(fixed/trainable)

# GNNs with **Edge** Embeddings

Battaglia et al. (NIPS 2016), Gilmer et al. (ICML 2017), Kipf et al. (ICML 2018)



**Legend:** ▮ : Node embedding  ▮▮ : Edge embedding  → : MLP

Node-to-edge ($v \rightarrow e$)   Edge-to-node ($e \rightarrow v$)

**Formally:**

$$v \rightarrow e : \quad \mathbf{h}_{(i,j)}^l = f_e^l([\mathbf{h}_i^l, \mathbf{h}_j^l, \mathbf{x}_{(i,j)}])$$

$$e \rightarrow v : \quad \mathbf{h}_j^{l+1} = f_v^l([\textstyle\sum_{i \in \mathcal{N}_j} \mathbf{h}_{(i,j)}^l, \mathbf{x}_j])$$

# GNNs with **Edge** Embeddings

Battaglia et al. (NIPS 2016), Gilmer et al. (ICML 2017), Kipf et al. (ICML 2018)



Legend: ■ : Node embedding  ■■ : Edge embedding  → : MLP

Node-to-edge $(v \rightarrow e)$          Edge-to-node $(e \rightarrow v)$

**Pros:**

- Supports edge features
- More expressive than GCN
- As general as it gets (?)
- Supports sparse matrix ops

**Formally:**

$$v \rightarrow e : \quad \mathbf{h}^l_{(i,j)} = f^l_e([\mathbf{h}^l_i, \mathbf{h}^l_j, \mathbf{x}_{(i,j)}])$$

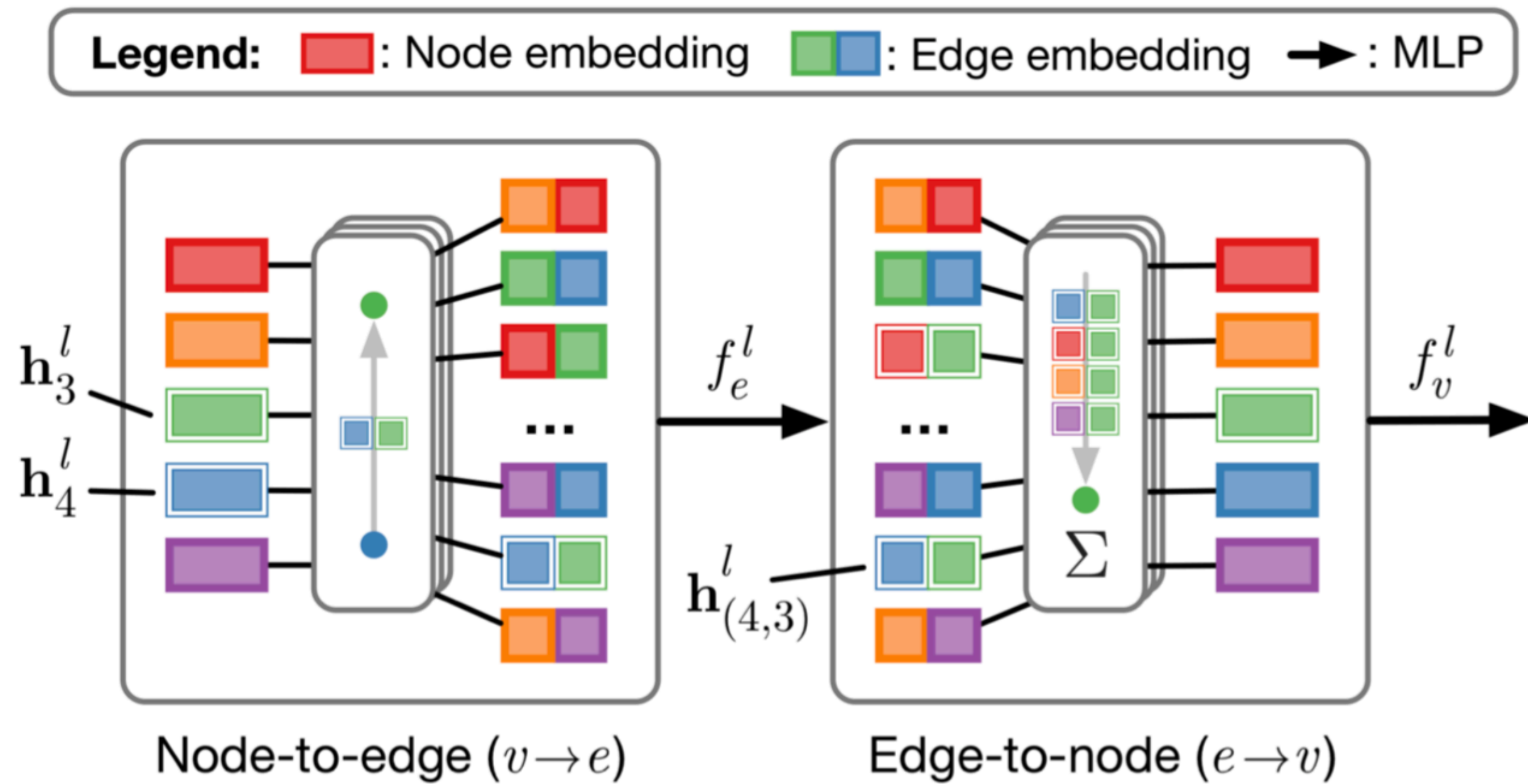$$e \rightarrow v : \quad \mathbf{h}^{l+1}_j = f^l_v([\textstyle\sum_{i \in \mathcal{N}_j} \mathbf{h}^l_{(i,j)}, \mathbf{x}_j])$$

# GNNs with **Edge** Embeddings

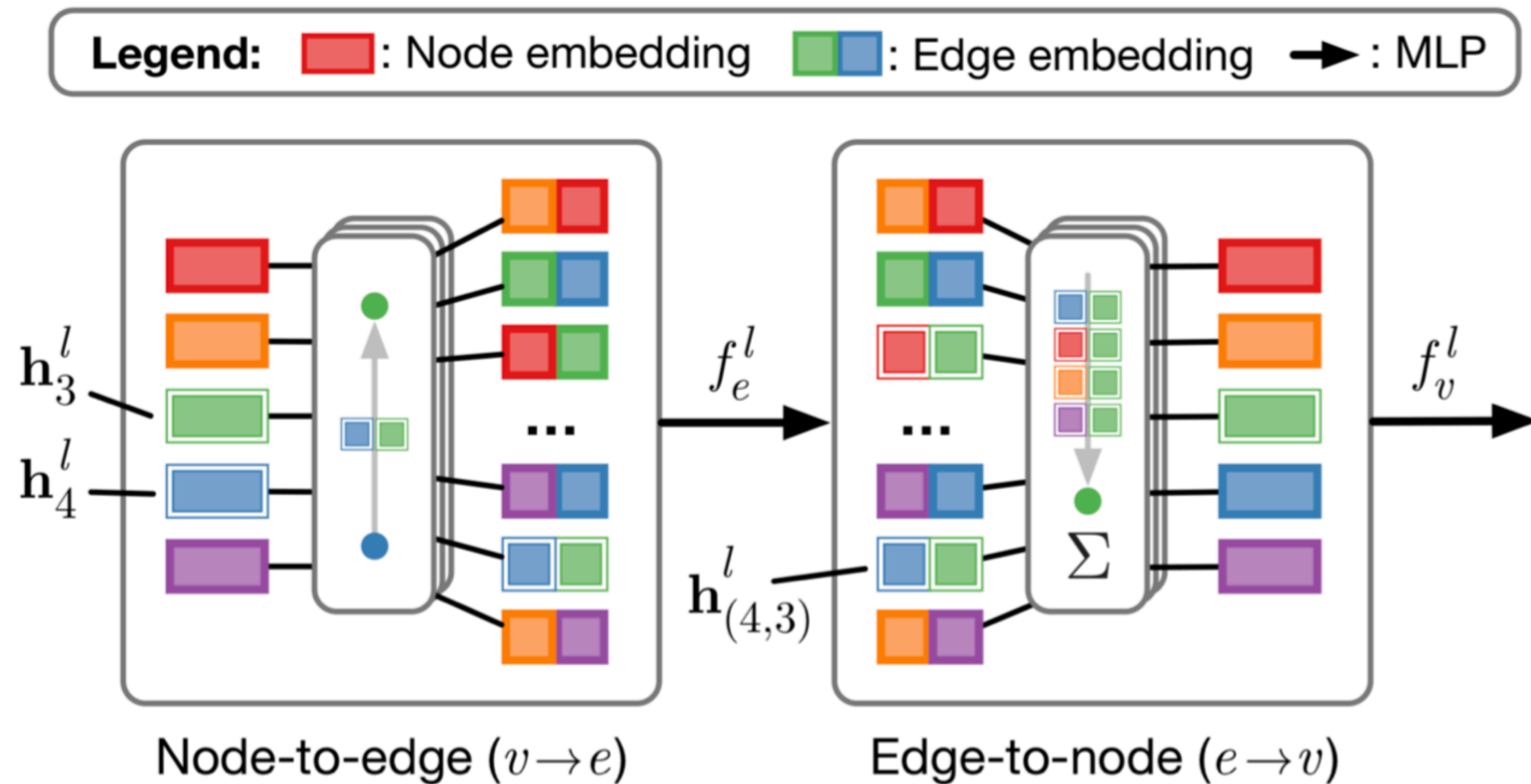Battaglia et al. (NIPS 2016), Gilmer et al. (ICML 2017), Kipf et al. (ICML 2018)



**Pros:**

- Supports edge features
- More expressive than GCN
- As general as it gets (?)
- Supports sparse matrix ops

**Cons:**

- Need to store intermediate edge-based activations
- Difficult to implement with subsampling
- ➡ In practice limited to small graphs

**Formally:**

$$v \to e: \quad \mathbf{h}^l_{(i,j)} = f^l_e([\mathbf{h}^l_i, \mathbf{h}^l_j, \mathbf{x}_{(i,j)}])$$

$$e \to v: \quad \mathbf{h}^{l+1}_j = f^l_v([\textstyle\sum_{i \in \mathcal{N}_j} \mathbf{h}^l_{(i,j)}, \mathbf{x}_j])$$

# Graph Neural Networks (GNNs) with **Attention**

Monti et al. (CVPR 2017), Hoshen (NIPS 2017), Veličković et al. (ICLR 2018)



[Figure from Veličković et al. (ICLR 2018)]

$$\vec{h}'_i = \sigma\left(\frac{1}{K}\sum_{k=1}^{K}\sum_{j \in \mathcal{N}_i} \alpha_{ij}^k \mathbf{W}^k \vec{h}_j\right)$$

# Graph Neural Networks (GNNs) with **Attention**

Monti et al. (CVPR 2017), Hoshen (NIPS 2017), Veličković et al. (ICLR 2018)
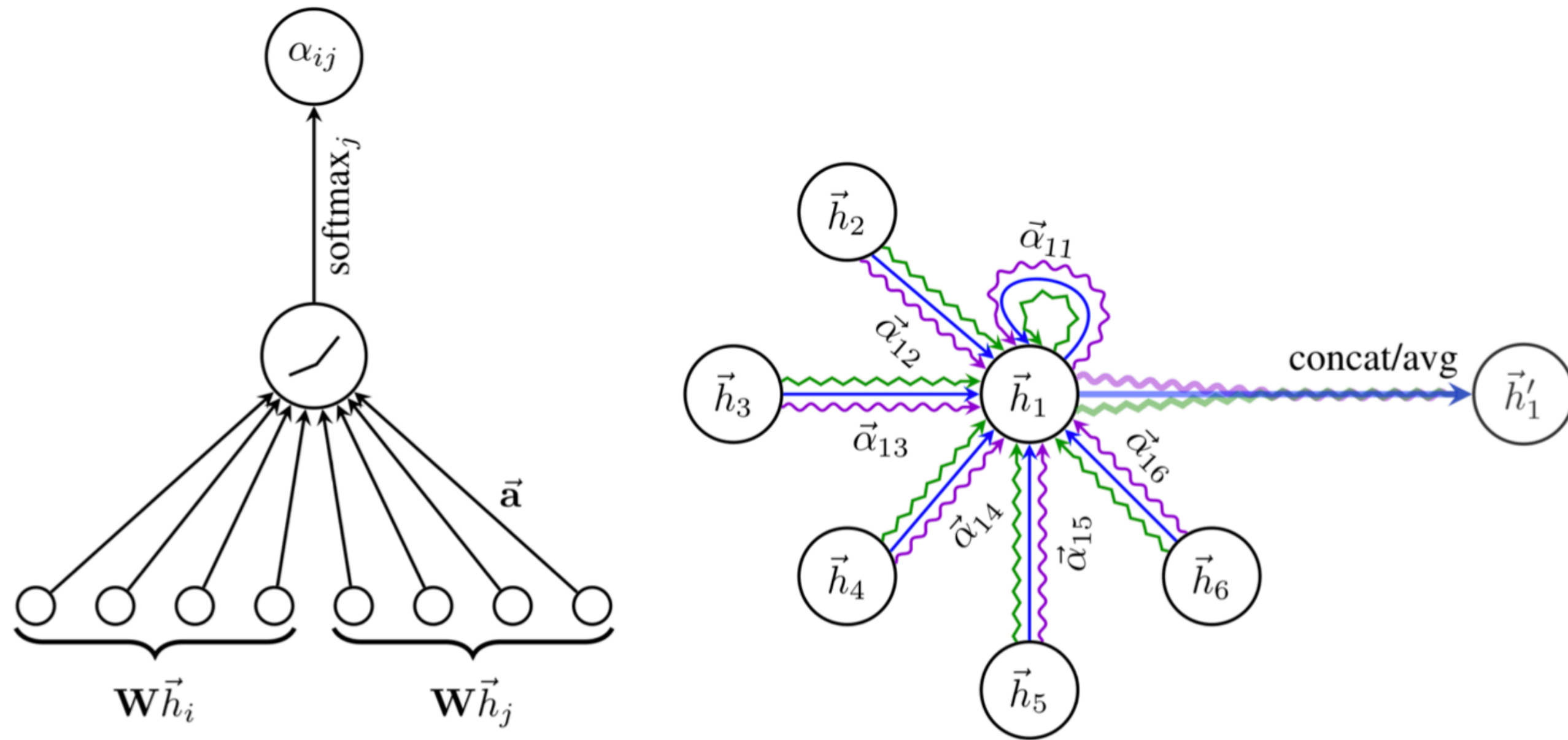


[Figure from Veličković et al. (ICLR 2018)]

$$\vec{h}'_i = \sigma \left( \frac{1}{K} \sum_{k=1}^{K} \sum_{j \in \mathcal{N}_i} \alpha_{ij}^k \mathbf{W}^k \vec{h}_j \right) \qquad \alpha_{ij} = \frac{\exp \left( \text{LeakyReLU} \left( \vec{\mathbf{a}}^T [\mathbf{W}\vec{h}_i \| \mathbf{W}\vec{h}_j] \right) \right)}{\sum_{k \in \mathcal{N}_i} \exp \left( \text{LeakyReLU} \left( \vec{\mathbf{a}}^T [\mathbf{W}\vec{h}_i \| \mathbf{W}\vec{h}_k] \right) \right)}$$

# Graph Neural Networks (GNNs) with **Attention**

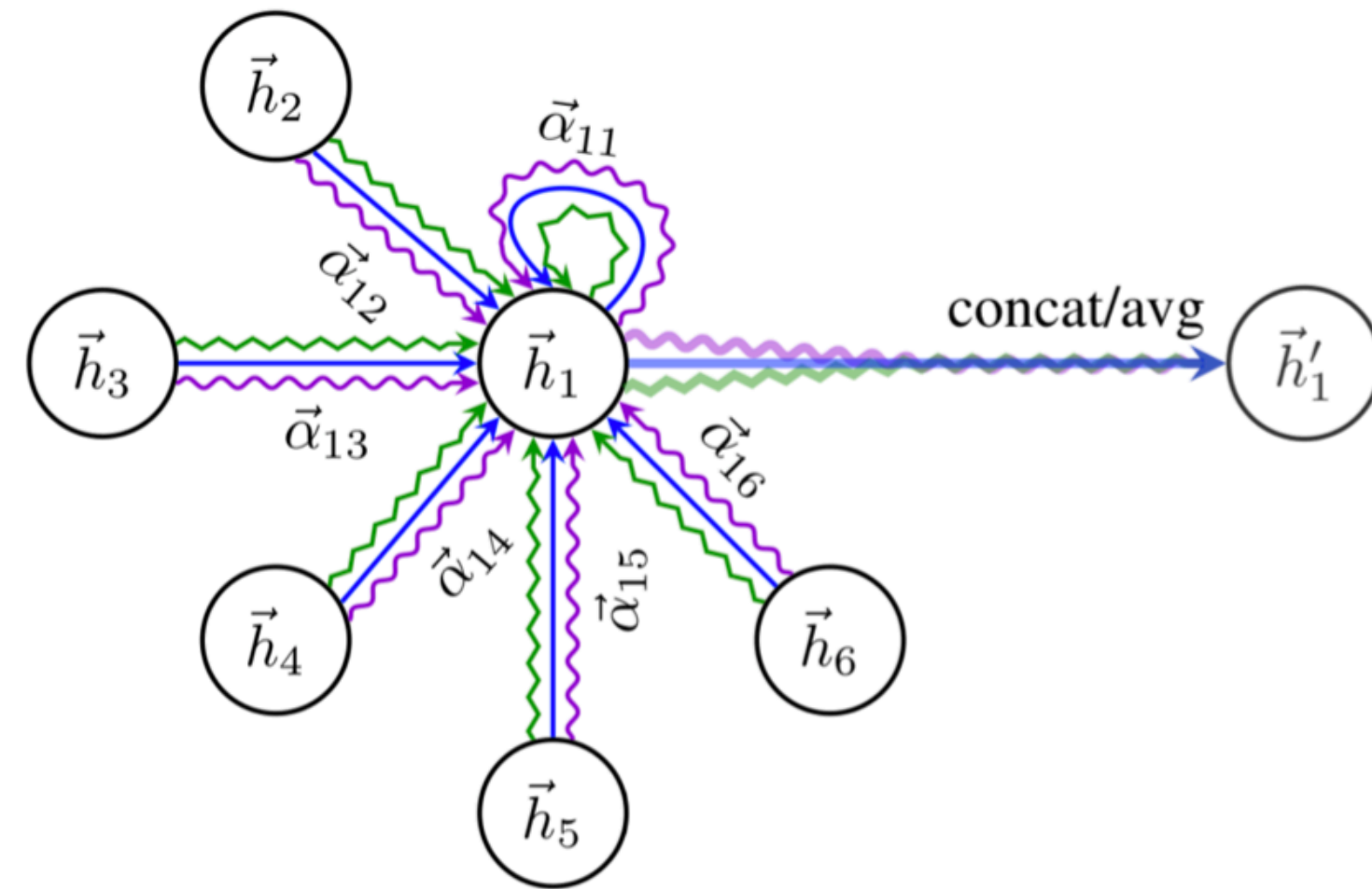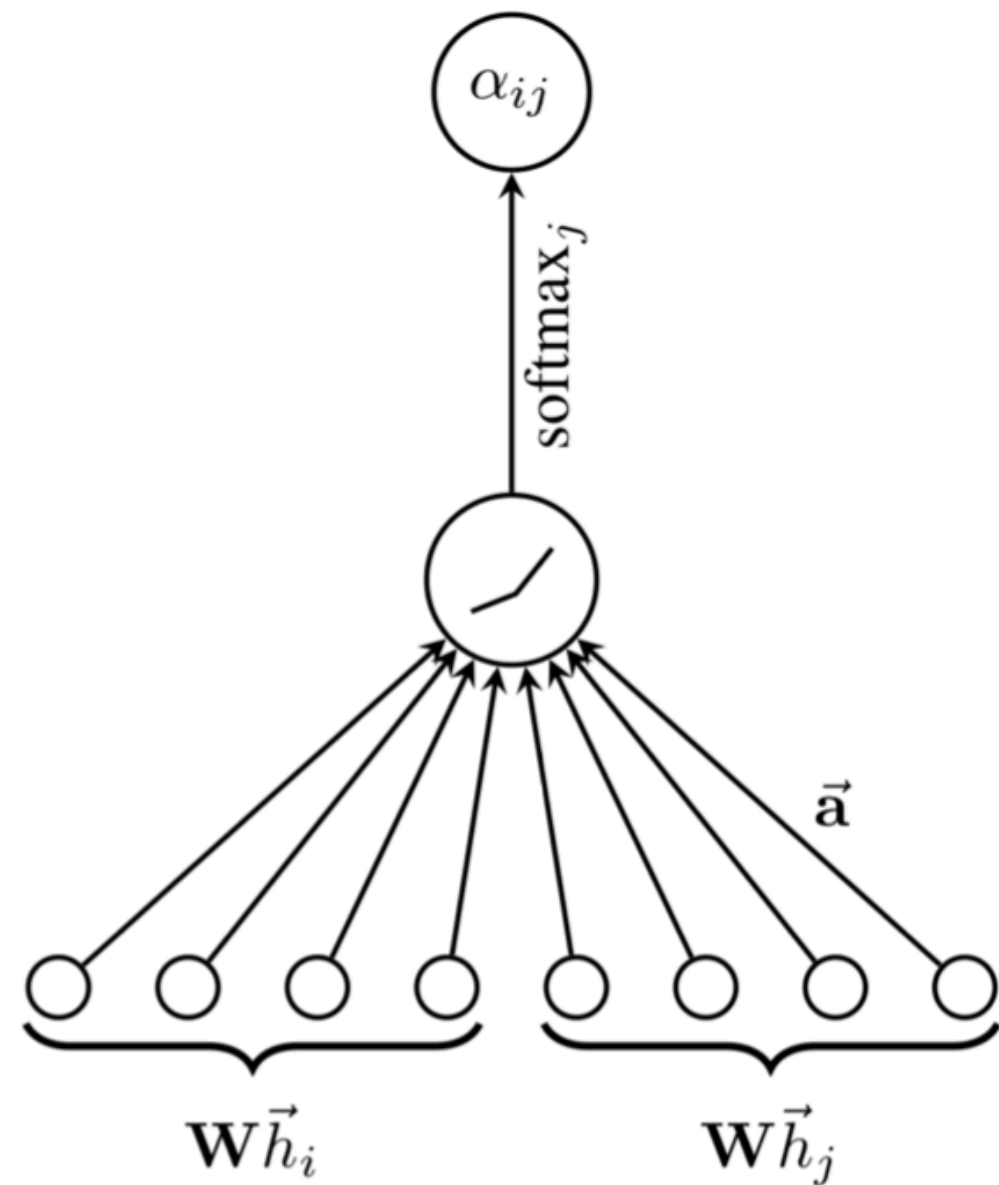Monti et al. (CVPR 2017), Hoshen (NIPS 2017), Veličković et al. (ICLR 2018)



[Figure from Veličković et al. (ICLR 2018)]

**Pros:**

- No need to store intermediate edge-based activation vectors (when using dot-product attn.)
- Slower than GCNs but faster than GNNs with edge embeddings

$$\vec{h}_i' = \sigma \left( \frac{1}{K} \sum_{k=1}^{K} \sum_{j \in \mathcal{N}_i} \alpha_{ij}^k \mathbf{W}^k \vec{h}_j \right) \qquad \alpha_{ij} = \frac{\exp\left( \text{LeakyReLU}\left( \vec{\mathbf{a}}^T [\mathbf{W}\vec{h}_i \| \mathbf{W}\vec{h}_j] \right) \right)}{\sum_{k \in \mathcal{N}_i} \exp\left( \text{LeakyReLU}\left( \vec{\mathbf{a}}^T [\mathbf{W}\vec{h}_i \| \mathbf{W}\vec{h}_k] \right) \right)}$$

# Graph Neural Networks (GNNs) with **Attention**

Monti et al. (CVPR 2017), Hoshen (NIPS 2017), Veličković et al. (ICLR 2018)



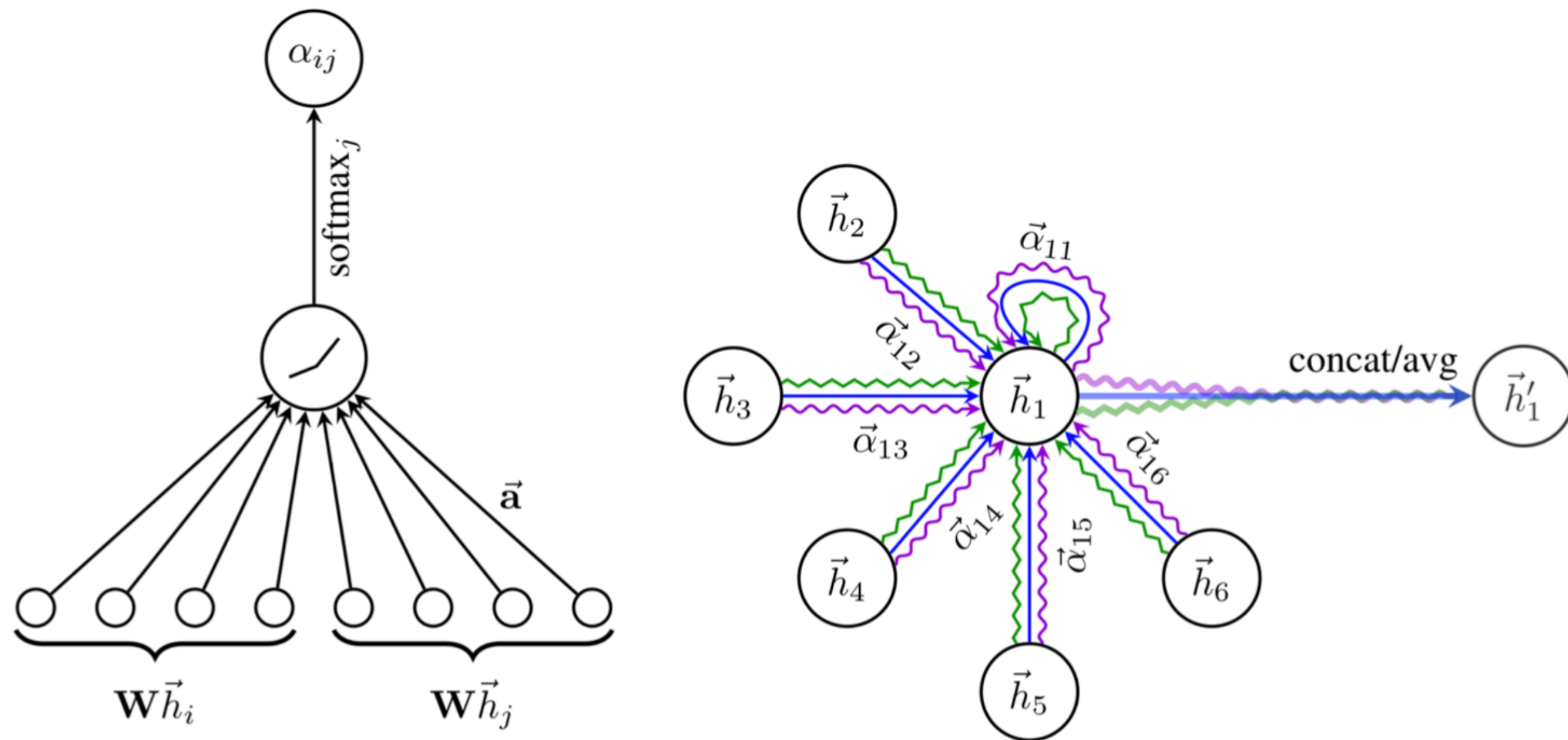[Figure from Veličković et al. (ICLR 2018)]

**Pros:**

- No need to store intermediate edge-based activation vectors (when using dot-product attn.)
- Slower than GCNs but faster than GNNs with edge embeddings

**Cons:**

- (Most likely) less expressive than GNNs with edge embeddings
- Can be more difficult to optimize

$$\vec{h}_i' = \sigma \left( \frac{1}{K} \sum_{k=1}^{K} \sum_{j \in \mathcal{N}_i} \alpha_{ij}^k \mathbf{W}^k \vec{h}_j \right) \qquad \alpha_{ij} = \frac{\exp\left( \text{LeakyReLU}\left( \vec{\mathbf{a}}^T [\mathbf{W}\vec{h}_i \| \mathbf{W}\vec{h}_j] \right) \right)}{\sum_{k \in \mathcal{N}_i} \exp\left( \text{LeakyReLU}\left( \vec{\mathbf{a}}^T [\mathbf{W}\vec{h}_i \| \mathbf{W}\vec{h}_k] \right) \right)}$$