

Graph Attention Networks

2020. 09. 11

Data Mining & Quality Analytics Lab.

강현규



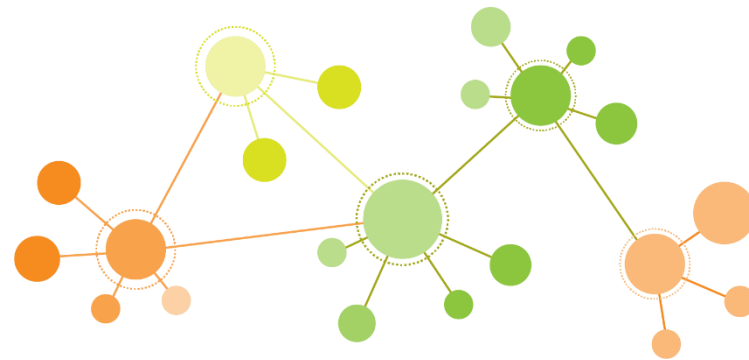
- **Introduction**
- **Attention**
- **Graph Neural Networks**
- **Graph Attention Networks**

01 | Introduction

■ 발표자 소개

- 강현규 (Hyeongyu Kang)
 - ✓ Data Mining & Quality Analytics Lab.
 - ✓ M.S. Student (2019.03 ~ Present)
 - ✓ 2021.02 졸업 예정
- Research Interest
 - ✓ Machine learning / Deep learning
 - ✓ Natural Language Processing / Graph Neural Networks

Text mining makes information extraction from huge volumes of data easier and structures the information as important facts, key terms or persons.



Hyeongyu Kang
(강현규)

M.S. Student
Industrial Management
Engineering,
Korea University

■ Seminar Topic

- Deep learning
 - ✓ 깊고 복잡한 신경망 구조 -> End-to-End learning
- Attention
 - ✓ 모델이 집중해서 학습해야 하는 곳 까지도 모델이 학습
 - ✓ Explainability + model performance

종료

Attention Network

Structure

Attention

Pass one-layer MLP with tanh activation

$$u_i = \tanh(W_1 h_i + b_1)$$

Hidden representation of sentence

Dot product

Similarity between sentence and sentence context vector

$$u_i^T u_j$$

Softmax

$$\alpha_i = \frac{\exp(u_i^T u_w)}{\sum_j \exp(u_j^T u_w)}$$

Normalized importance weight = Attention score of sentence

Hierarchical Attention Network for Document Classification

발표자: 강현규

📅 2019년 7월 12일

🕒 오후 1시 ~

📍 고려대학교 신공학관 218호

세미나 정보 보기 →

종료

Visual Attention

output

learned weights

$$y_{ij} = \sum_{a,b \in \mathcal{M}(i,j)} W_{i-a,j-b} x_{ab}$$

Convolution

matrix multiplication

learned transform

$$y_{ij} = \sum_{a,b \in \mathcal{M}(i,j)} \text{softmax}_{ab}(q_i^T k_{ab}) v_{ab}$$

$$q_i = W_q x_i \quad k_{ab} = W_k x_{ab} \quad v_{ab} = W_v x_{ab}$$

Self Attention

Visual Attention

발표자: 강현규

📅 2020년 2월 14일

🕒 오후 1시 ~

📍 고려대학교 신공학관 2층 컴퓨터실

세미나 정보 보기 →

D-1

softmax

Graph Attention Networks

발표자: 강현규

📅 2020년 9월 11일

🕒 오후 1시 ~

📍 온라인 비디오 시청(YouTube)

세미나 정보 보기 →



- Seminar Topic
 - ✓ Graph Neural Networks

Graph Neural Networks

Node-Edge로 구성된 그래프 데이터의

구조를 학습하는 딥러닝 모델

- Seminar Topic
 - Graph Attention Networks
 - ✓ Graph Neural Networks + Attention

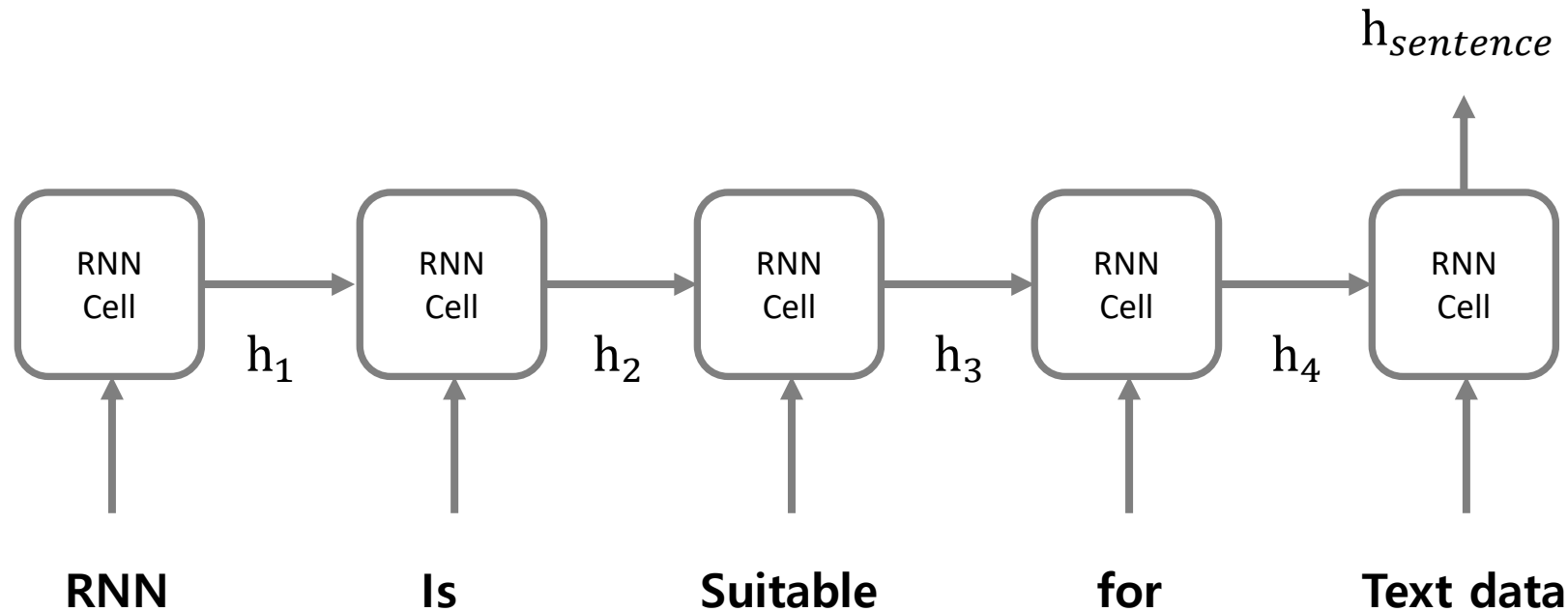
중요 Node에 가중치를 부여하는 어텐션 메커니즘을 사용하여

Graph Attention Networks

Node-Edge로 구성된 그래프 데이터의

구조를 학습하는 딥러닝 모델

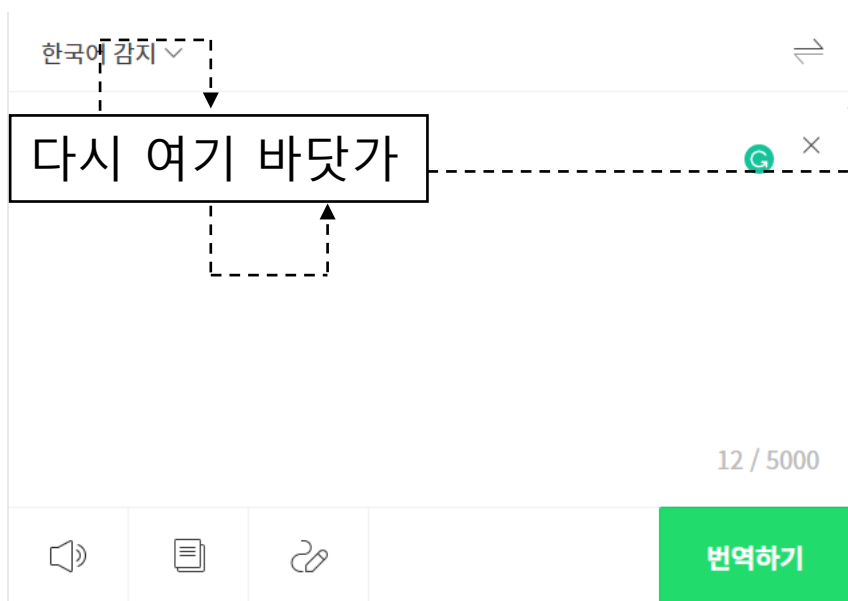
- Text data & Recurrent Neural Network (RNN)
 - RNN은 sequential data의 학습에 적합한 신경망 모델
 - Text data는 co-occurrence와 sequential pattern을 고려한 분석이 필요



02 | Attention

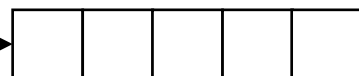
- Seq2seq
 - RNN encoder + RNN decoder
 - Machine translation

Encoder: RNN



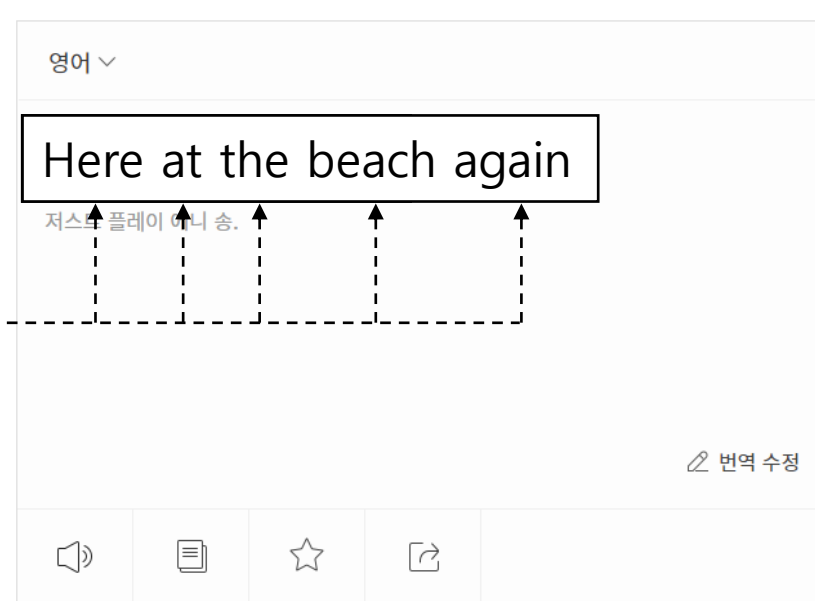
Input

Context vector



papago

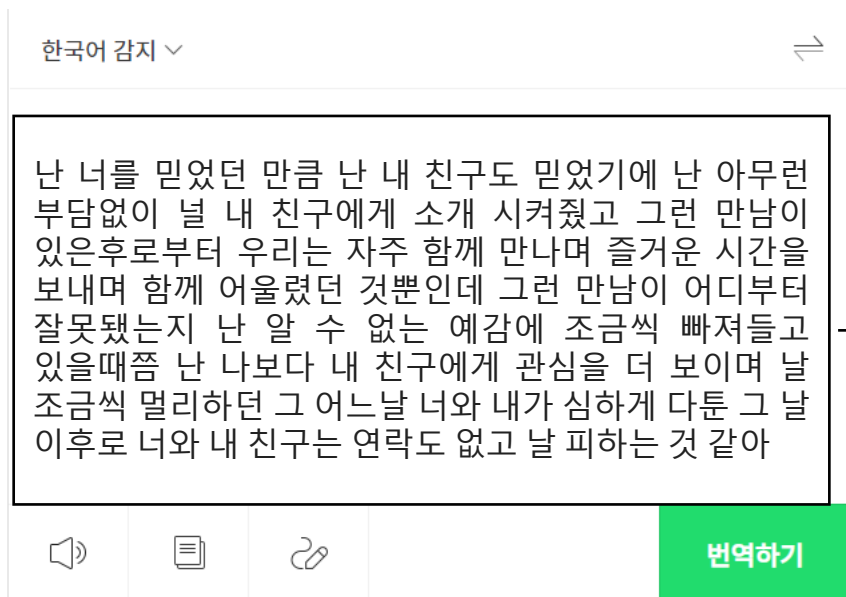
Decoder: RNN



Output

- Seq2seq
 - Long term dependency
 - Vanishing/Exploding gradient

Encoder: RNN



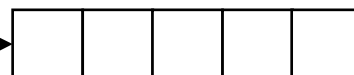
한국어 감지 ▾

난 너를 믿었던 만큼 난 내 친구도 믿었기에 난 아무런 부담없이 널 내 친구에게 소개 시켜줬고 그런 만남이 있은후로부터 우리는 자주 함께 만나며 즐거운 시간을 보내며 함께 어울렸던 것뿐인데 그런 만남이 어디부터 잘못됐는지 난 알 수 없는 예감에 조금씩 빠져들고 있을때쯤 난 나보다 내 친구에게 관심을 더 보이며 날 조금씩 멀리하던 그 어느날 너와 내가 심하게 다툰 그 날 이후로 너와 내 친구는 연락도 없고 날 피하는 것 같아

번역하기

Input

Context vector



papago

Decoder: RNN



영어 ▾

As much as I

fought
met
trusted

번역 수정

Output

- Seq2seq with Attention

- Relieve the encoder from the burden of having to encode all information into a fixed length vector

Encoder: RNN

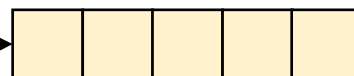
한국어 감지

난 너를 믿었던 만큼 난 내 친구도 믿었기에 난 아무런 부담없이 널 내 친구에게 소개 시켜줬고 그런 만남이 있은후로부터 우리는 자주 함께 만나며 즐거운 시간을 보내며 함께 어울렸던 것뿐인데 그런 만남이 어디부터 잘못됐는지 난 알 수 없는 예감에 조금씩 빠져들고 있을때쯤 난 나보다 내 친구에게 관심을 더 보이며 날 조금씩 멀리하던 그 어느날 너와 내가 심하게 다툰 그 날 이후로 너와 내 친구는 연락도 없고 날 피하는 것 같아

번역하기

Input

Context vector



papago

Decoder: RNN

영어

As much as I

fought
met
trusted

번역 수정

Output

02 | Attention

▪ Key, Query, Value

- Dictionary 자료 구조
- Query와 Key가 일치하면 Value를 return



Query: 유재석

Key	Value
이효리	천옥
엄정화	만옥
제시	은비
화사	실비
유재석	지미 유

부캐 Dictionary



Output: 지미 유

02 | Attention

▪ Key, Query, Value

- Dictionary 자료 구조
- Query와 Key가 일치하면 Value를 return

```
① def Similarity(Query, Key)  
    if Query = Key:  
        return 1  
    else:  
        return 0
```



Query: 유재석

Key	Sim	Value
이효리	0	천옥
엄정화	0	만옥
제시	0	은비
화사	0	실비
유재석	1	지미 유

부캐 Dictionary

02 | Attention

Key, Query, Value

- Dictionary 자료 구조
- Query와 Key가 일치하면 Value를 return

```
① def Similarity(Query, Key)  
    if Query = Key:  
        return 1  
    else:  
        return 0
```

```
② def SimXValue(Sim, Value)  
    output = Sim X int(value)  
    return output
```



Query: 유재석

Key	Sim	Value	SimXValue
이효리	0	천옥	0
엄정화	0	만옥	0
제시	0	은비	0
화사	0	실비	0
유재석	1	지미 유	지미 유

부캐 Dictionary

02 | Attention

Key, Query, Value

- Dictionary 자료 구조
- Query와 Key가 일치하면 Value를 return

```
① def Similarity(Query, Key)  
    if Query = Key:  
        return 1  
    else:  
        return 0
```

```
② def SimXValue(Sim, Value)  
    output = Sim X int(value)  
    return output
```



Query: 유재석

Key	Value
이효리	천옥
엄정화	만옥
제시	은비
화사	실비
유재석	지미 유

부캐 Dictionary

SimXValue
0
0
0
0
지미 유

Result
지미 유

```
③ def Result(outputs)  
    return sum(outputs)
```

▪ Key, Query, Value

- Dictionary 자료형의 결과 리턴 과정

- ① Similarity(key, value)
 - ✓ Key와 value의 유사도를 계산한다
- ② SimXValue(sim, value)
 - ✓ 유사도와 value를 곱한다
- ③ Result(outputs)
 - ✓ 유사도와 value를 곱한 값의 합을 리턴한다

$$result = \sum_i similarity(key, query) * value$$

```
① def Similarity(Query, Key)
    if Query = Key:
        return 1
    else:
        return 0
```



```
② def SimXValue(Sim, Value)
    output = Sim X int(value)
    return output
```



```
③ def Result(outputs)
    return sum(outputs)
```

▪ Key, Query, Value in Attention

- Attention: Query와 key의 유사도를 계산한 후 value의 가중합을 계산하는 과정
- Attention score: Value에 곱해지는 가중치
- Considerations
 - ✓ Key, Query, Value = Vectors (Matrix/Tensor)
 - ✓ Similarity function

$$output = \sum_i similarity(key, query) * value$$

Dictionary 자료 구조



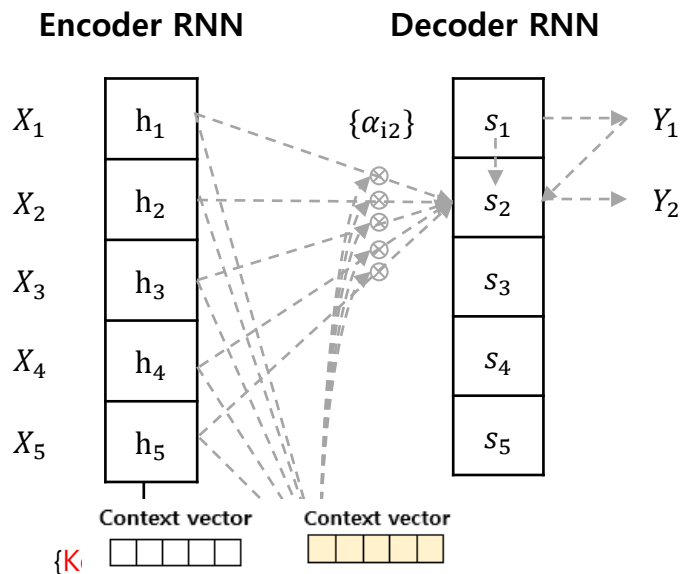
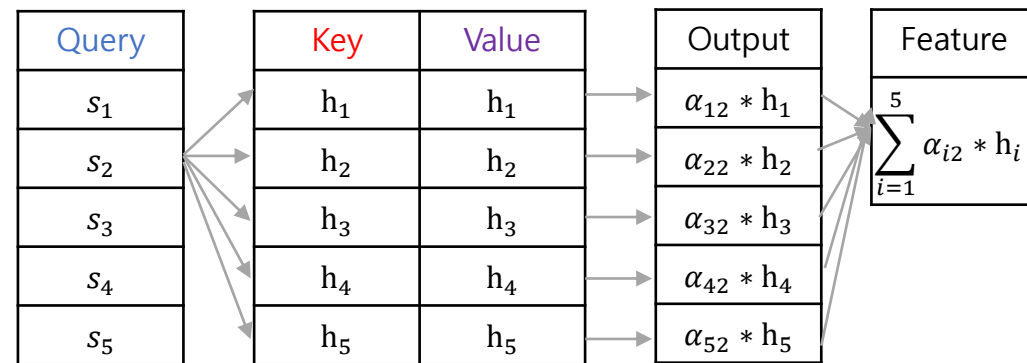
$$A(q, K, V) = \sum_i softmax(f(K, q)) V$$

Attention

02 | Attention

Attention in Seq2seq Machine Translation

- Attention: Query와 key의 유사도를 계산한 후 value의 가중합을 계산하는 과정
- Key, Value = Hidden states of encoder h_i
- Query = Hidden state of decoder s_2
- Feature = Context vector at time step 2



$$\{\alpha_{i2}\} = \text{softmax}(f(h_i, s_2))$$

$$\text{feature} = \sum_{i=1}^5 \alpha_{i2} * h_i$$



$$A(q, K, V) = \sum_i \text{softmax}(f(K, q)) V$$



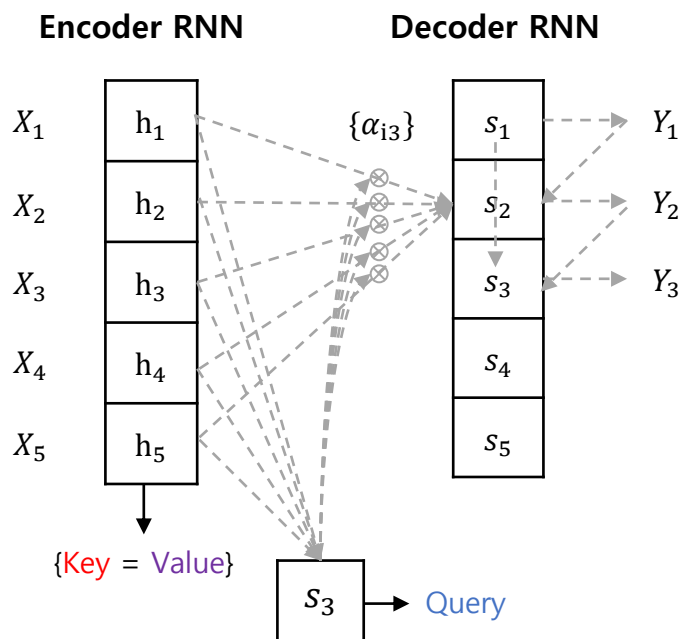
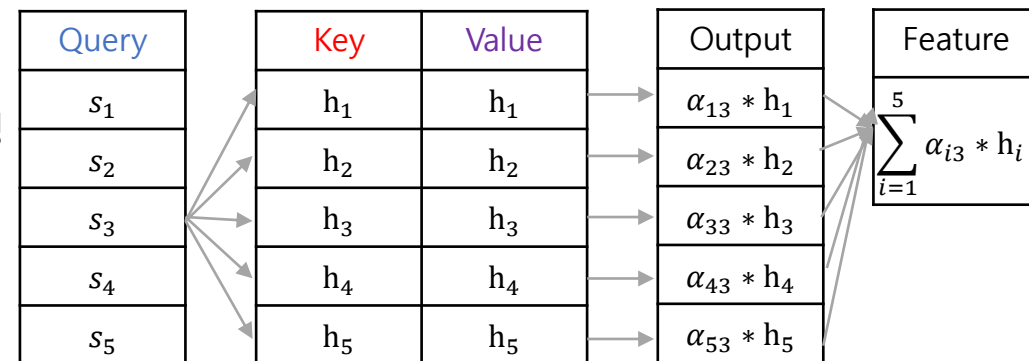
Machine Translation



02 | Attention

Attention in Seq2seq Machine Translation

- Attention: Query와 key의 유사도를 계산한 후 value의 가중합을 계산하는 과정
- Key, Value = Hidden states of encoder h_i
- Query = Hidden state of decoder s_3
- Feature = Context vector at time step 3



$$\{\alpha_{i3}\} = \text{softmax}(f(h_i, s_3))$$

$$\text{feature} = \sum_{i=1}^5 \alpha_{i3} * h_i$$



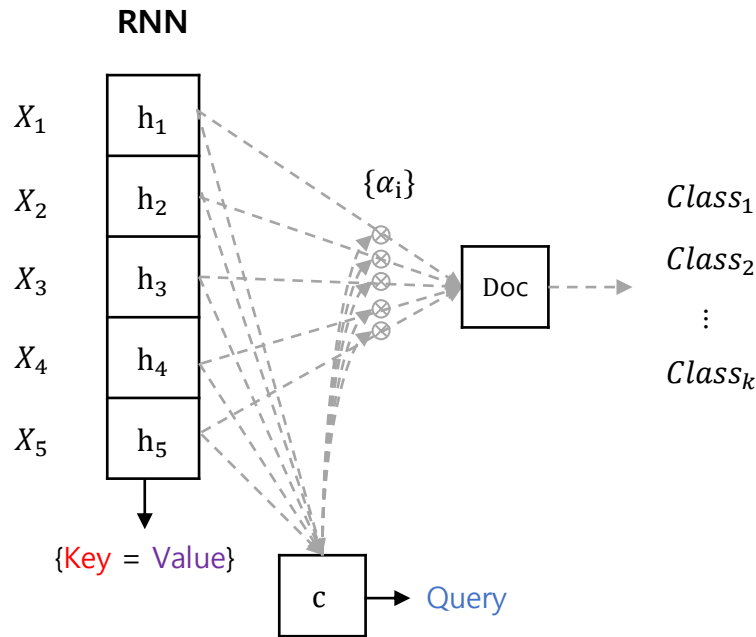
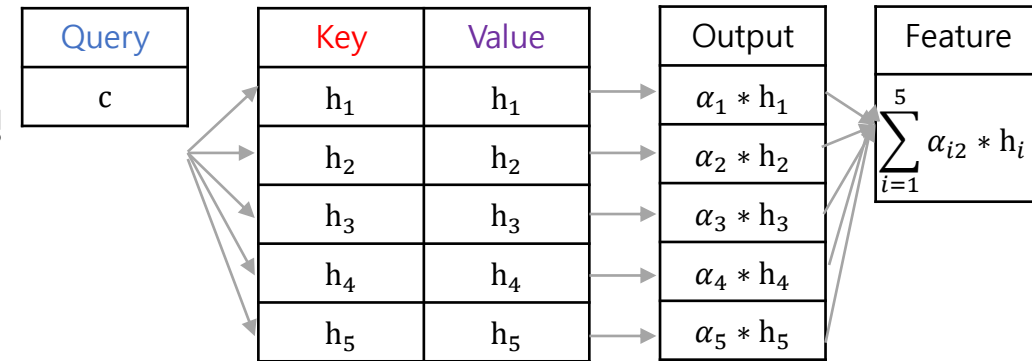
$$A(q, K, V) = \sum_i \text{softmax}(f(K, q)) V$$

Machine Translation



Attention in RNN-based Document Classification

- Attention: Query와 key의 유사도를 계산한 후 value의 가중합을 계산하는 과정
- Feature = Document vector
- Key, Value = Hidden states of RNN h_i
- Query = Learnable parameter vector c (context vector)



$$\{\alpha_i\} = \text{softmax}(f(h_i, c))$$

$$feature = \sum_{i=1}^5 \alpha_i * h_i$$



$$A(q, K, V) = \sum_i \text{softmax}(f(K, q)) V$$

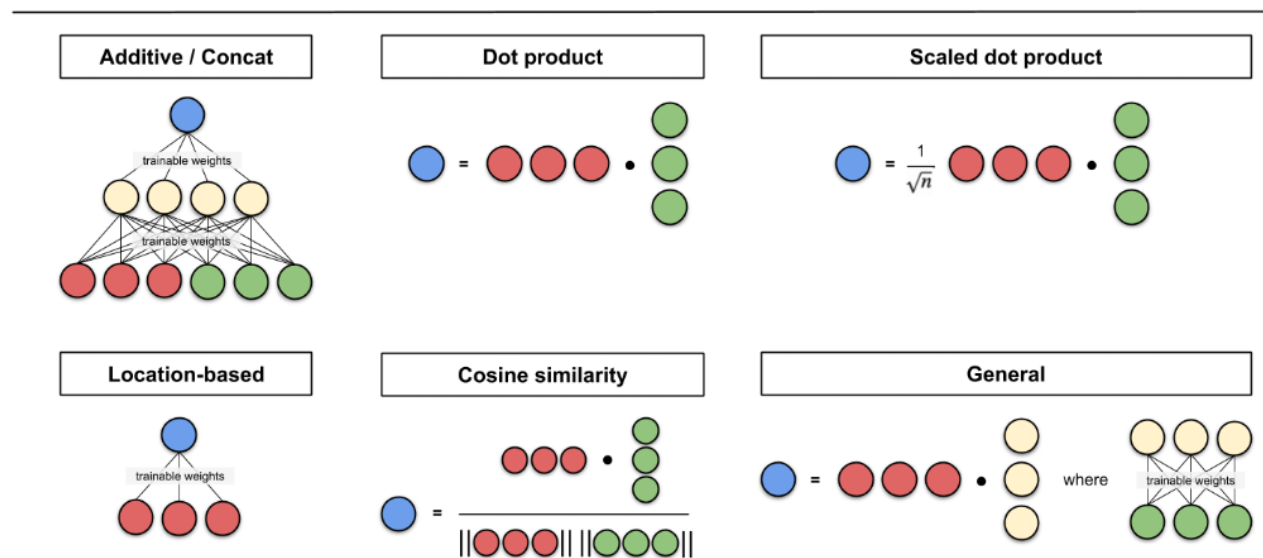
Document Classification

Similarity Function (Alignment model)

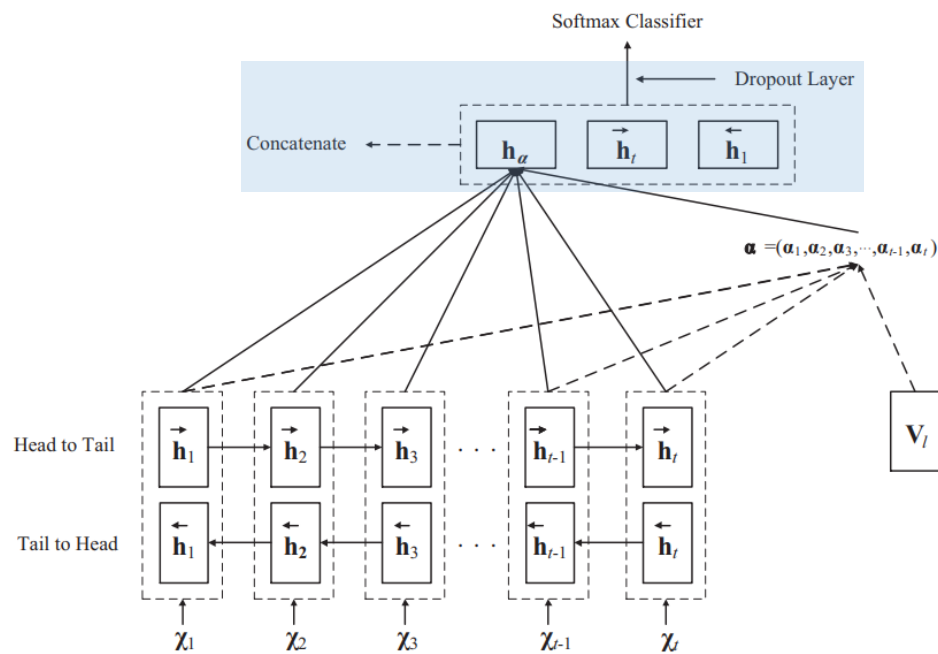
- Vector간 유사도를 계산하는 다양한 방법을 similarity function으로 사용 가능
- Bahdanau Attention (2014), Graph Attention Network (2018) -> Additive / Concat
- Luong (2015) -> 다양한 similarity function을 제시
- Transformer (2017) -> Scaled-dot product

$$A(q, K, V) = \sum_i softmax(f(K, q)) V$$

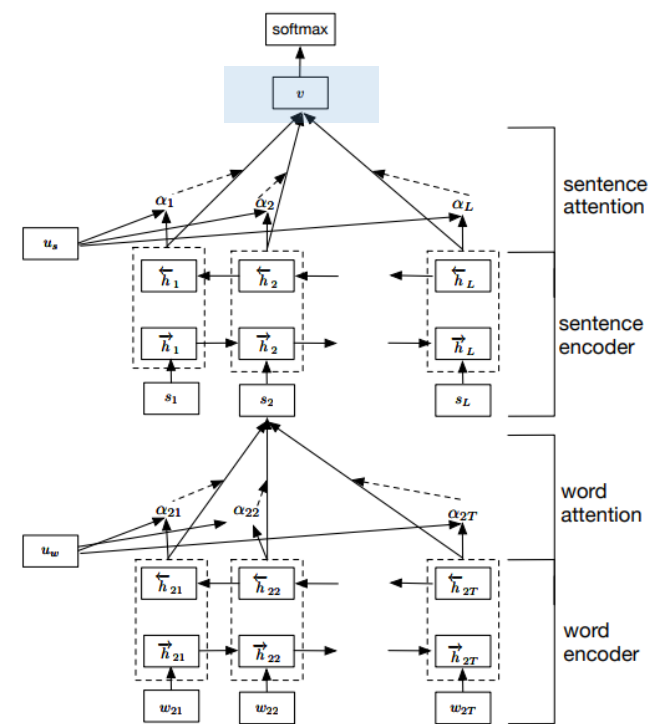
$$f(Key, Query) = score$$



- Feature Representation by RNN-based Network
 - Bi-RNN with Attention
 - Hierarchical Attention Network (2016)



Bidirectional RNN with Attention

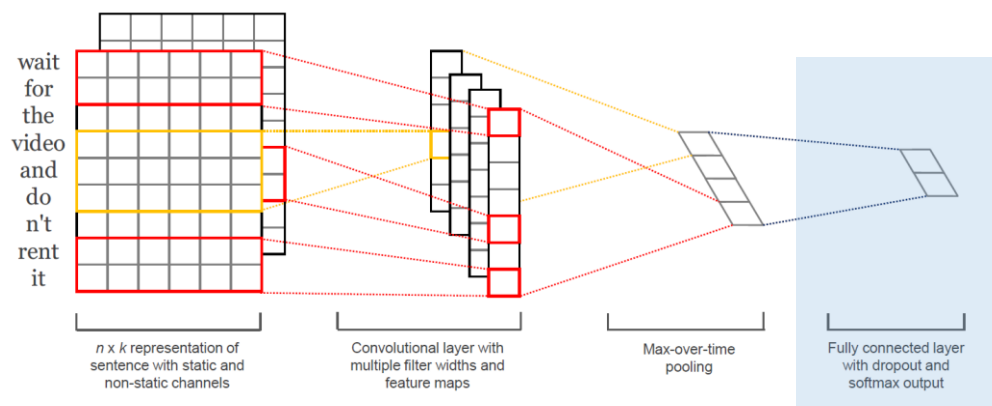


Hierarchical Attention Network

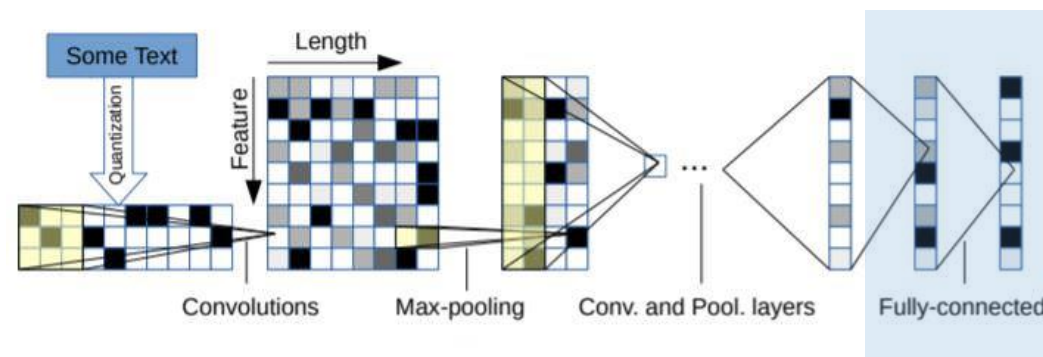
Liu, Tengfei, et al. "Recurrent networks with attention and convolutional networks for sentence representation and classification." *Applied Intelligence* 48.10 (2018): 3797-3806.

Yang, Zichao, et al. "Hierarchical attention networks for document classification." *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics* 2016.

- Feature Representation by CNN-based Network
 - TextCNN (2014)
 - Character-level CNN (2015)



TextCNN



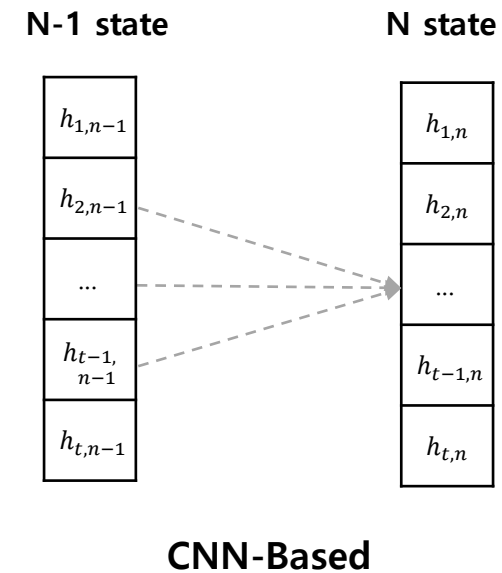
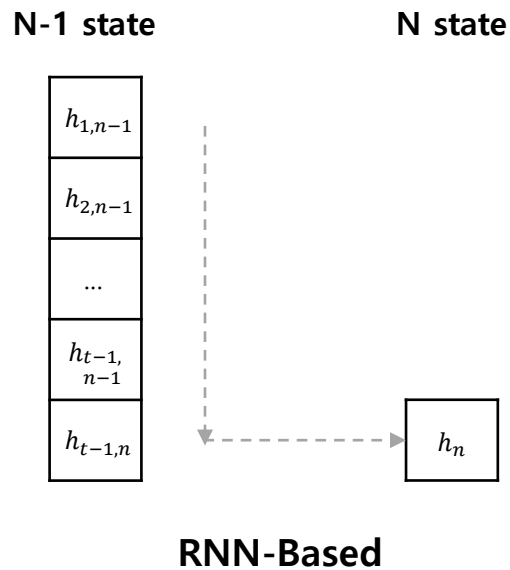
Character-level CNN

Kim, Yoon. "Convolutional neural networks for sentence classification." *arXiv preprint arXiv:1408.5882* (2014).

Zhang, Xiang, Junbo Zhao, and Yann LeCun. "Character-level convolutional networks for text classification." *Advances in neural information processing systems*. 2015.

Attention to Self-Attention

- RNN-based Network
 - ✓ Sequential data → Parallel computing X
 - ✓ Calculation time and complexity ↑
 - ✓ Vanishing gradient / Long term dependency
- CNN-based Network
 - ✓ Long path length between long-range dependencies



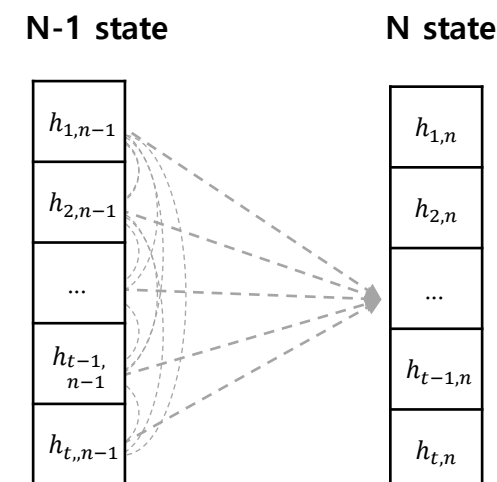
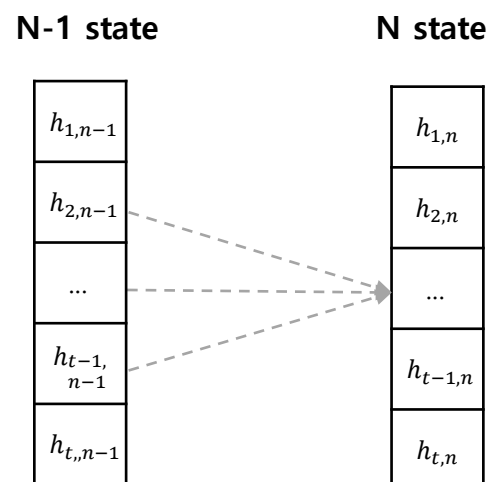
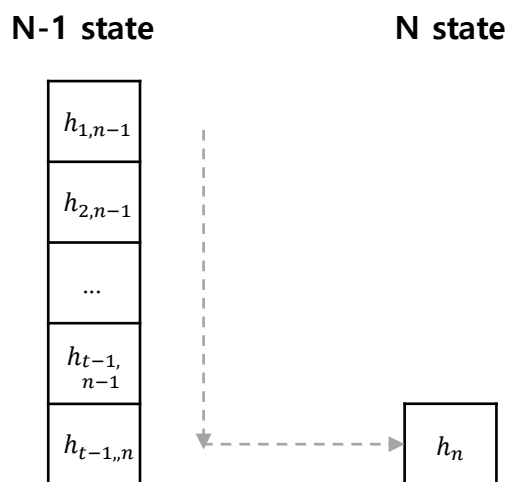
02 | Attention

Self-Attention

- RNN, CNN 구조를 사용하지 않고 attention만을 사용하여 feature representation
 - ✓ Key = Query = Value = Hidden state of word embedding vector
 - ✓ Scaled dot-product attention
 - ✓ Multi-head attention

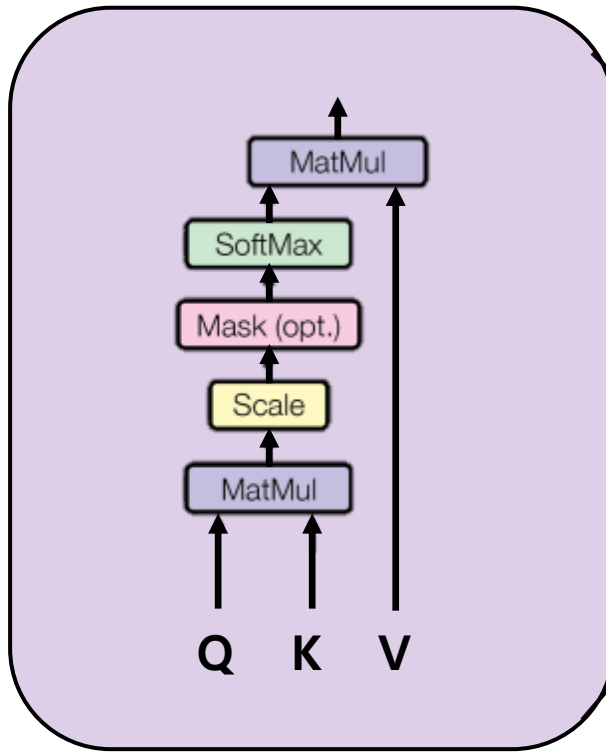
$$A(q, K, V) = \sum_i softmax(f(K, q)) V$$

Key = Query = Value

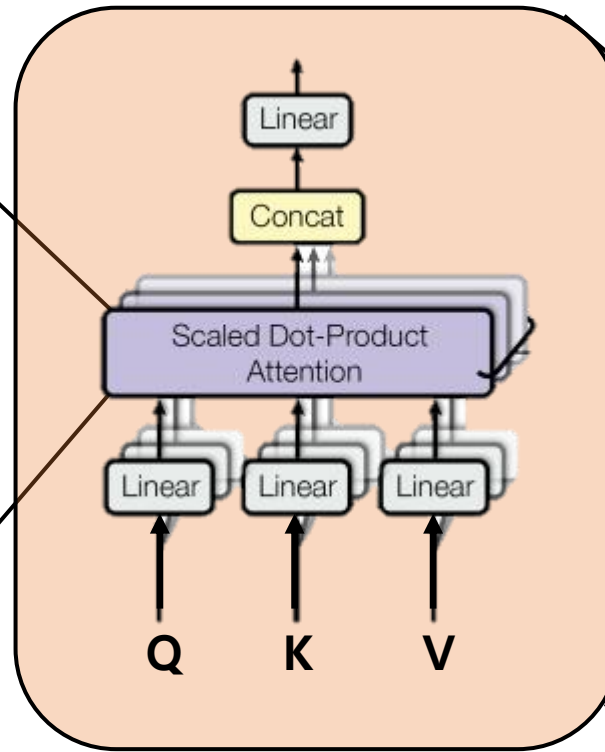


02 | Attention

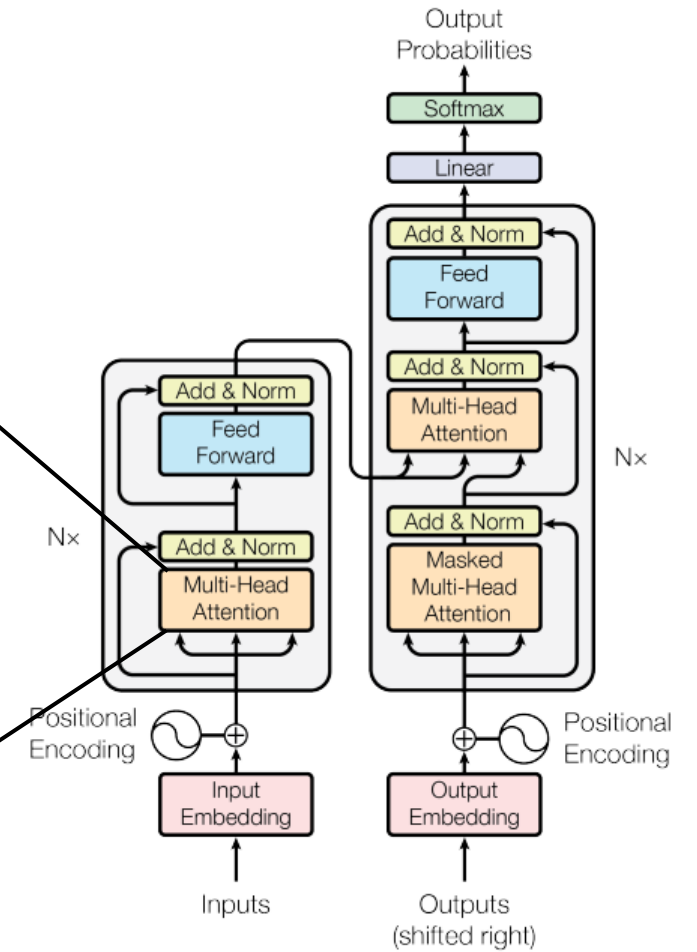
- Transformer



Scaled Dot-Product Attention



Multi-Head Attention

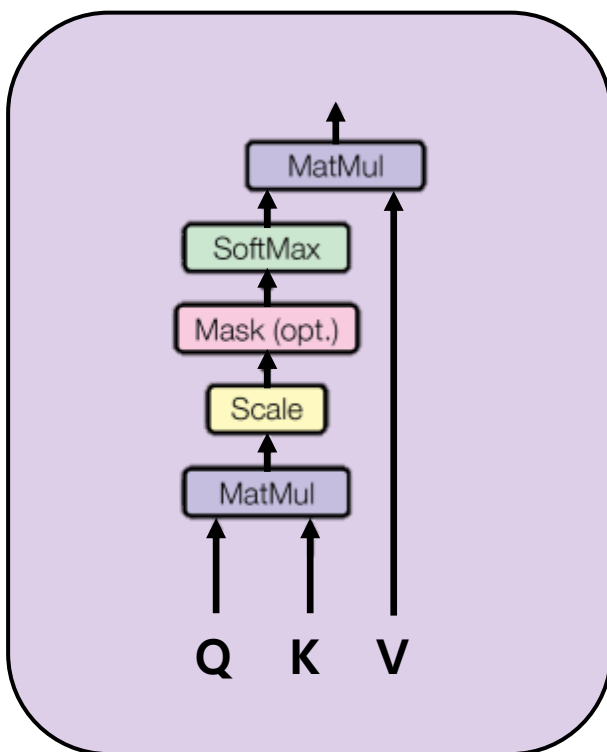


Transformer

Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems*. 2017..

Transformer

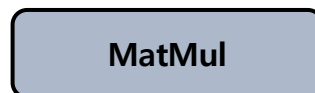
- Scale-dot product attention (Self-Attention)
 - ✓ Key = Query = Value = Hidden state of word embedding vector (X)
 - ✓ Similarity function = Dot-product



Scaled Dot-Product Attention

Generalized
Attention Form

①



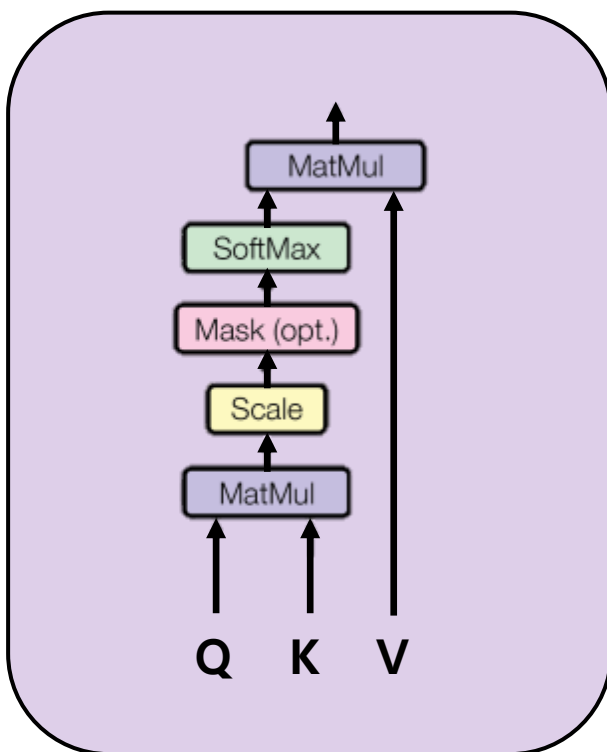
$$A(q, K, V) = \sum_i \text{softmax}(f(K, q)) V$$

$$f(K, Q) = QK^T \quad (K = XW^K, Q = XW^Q, V = XW^V)$$

While the two are similar in theoretical complexity, dot-product attention is much faster and more space-efficient in practice, since it can be implemented using highly optimized matrix multiplication code. (...)

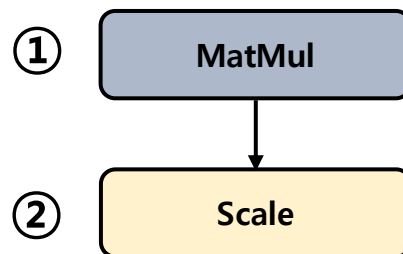
- Transformer

- Scale-dot product attention (Self-Attention)
 - ✓ Similarity function: Scaled-dot product



Scaled Dot-Product Attention

Generalized Attention Form



$$A(q, K, V) = \sum_i \text{softmax}(f(K, q)) V$$

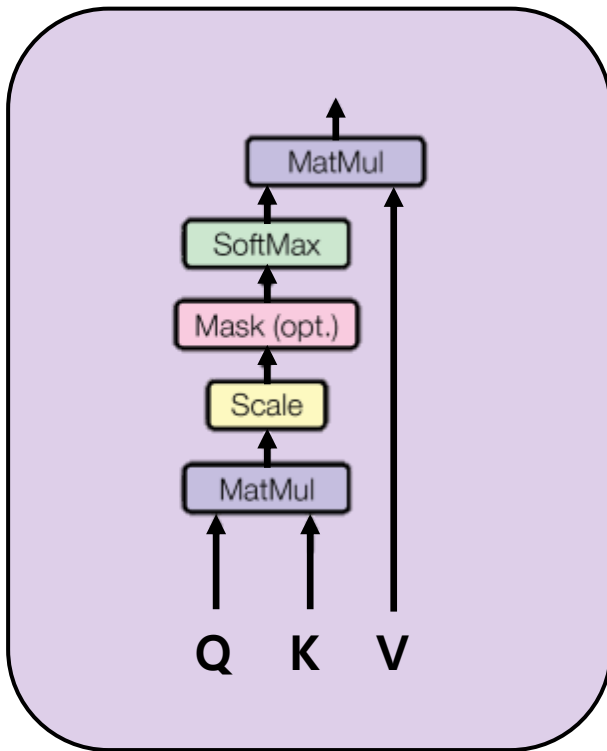
$$f(K, Q) = QK^T \quad (K = XW^K, Q = XW^Q, V = XW^V)$$

$$\frac{QK^T}{\sqrt{d_k}}$$

We suspect that for large values of d_k , the dot products grow large in magnitude, pushing the softmax function into regions where it has extremely small gradients. To counteract this effect, we scale the dot products (...)

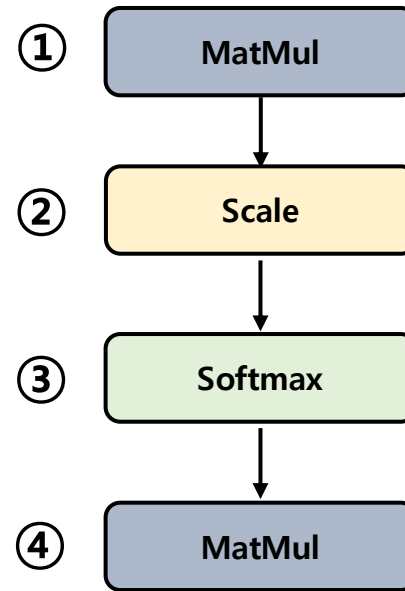
Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems*. 2017..

- Transformer
 - Scale-dot product attention (Self-Attention)
 - Weight-sum of value vectors



Scaled Dot-Product Attention

Generalized Attention Form



$$A(q, K, V) = \sum_i \text{softmax}(f(K, q)) V$$

$$f(K, Q) = QK^T \quad (K = KW^K, Q = QW^Q, V = QW^V)$$

$$\frac{QK^T}{\sqrt{d_k}}$$

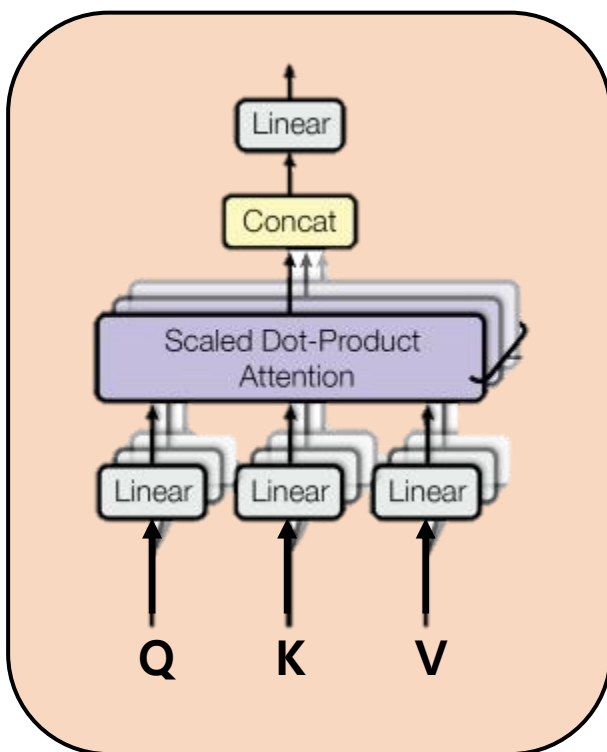
$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)$$

$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

- Transformer

- Multi-head Attention

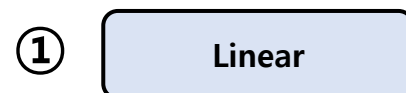
- ✓ Learning diverse input features



Multi-Head Attention

Self-Attention

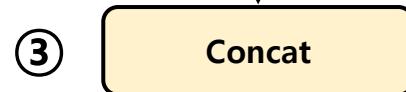
$$SA(q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



$$Q' = QW_i^Q \quad K' = KW_i^K \quad V' = VW_i^V \quad (i = 1 \dots h)$$



$$\text{head}_i = SA(Q', K', V')$$



$$[\text{head}_1, \text{head}_2, \dots, \text{head}_h]$$



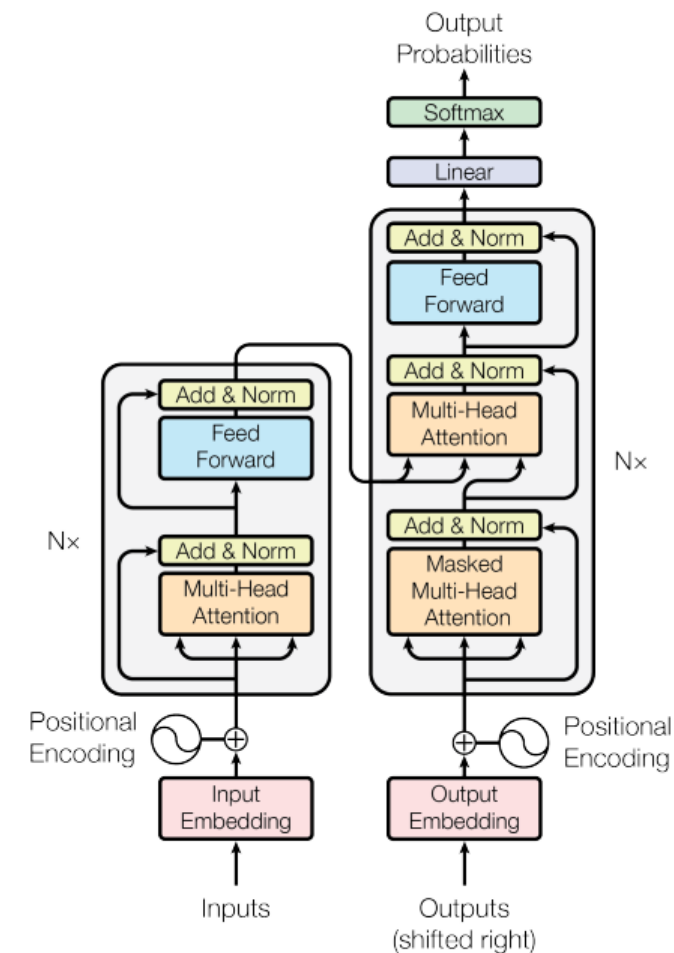
$$[\text{head}_1, \text{head}_2, \dots, \text{head}_h]W^O \\ = \text{MultiHead}(Q, K, V)$$

Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems*. 2017..

Transformer

- Contributions

- One is the **total computational complexity** per layer. (...)
- Another is the amount of computation that can **be parallelized**, (...)
- The third is the path length between long-range dependencies in the network. **Learning long-range dependencies** is a key challenge in many sequence transduction tasks. (...)
- As side benefit, self-attention could yield **more interpretable models**.



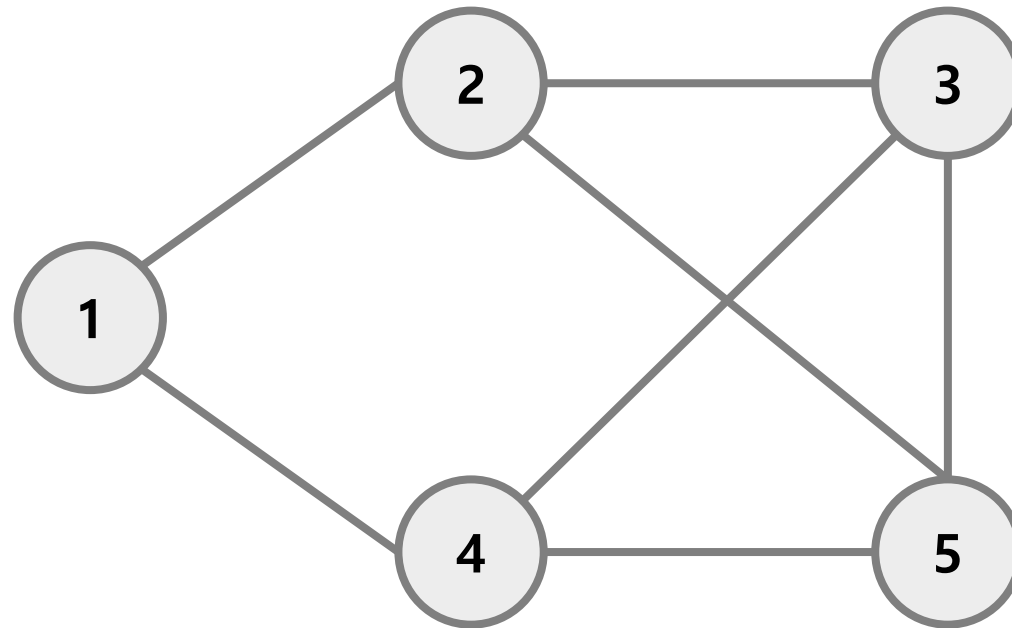
Transformer

Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems*. 2017..

03 | Graph Neural Networks

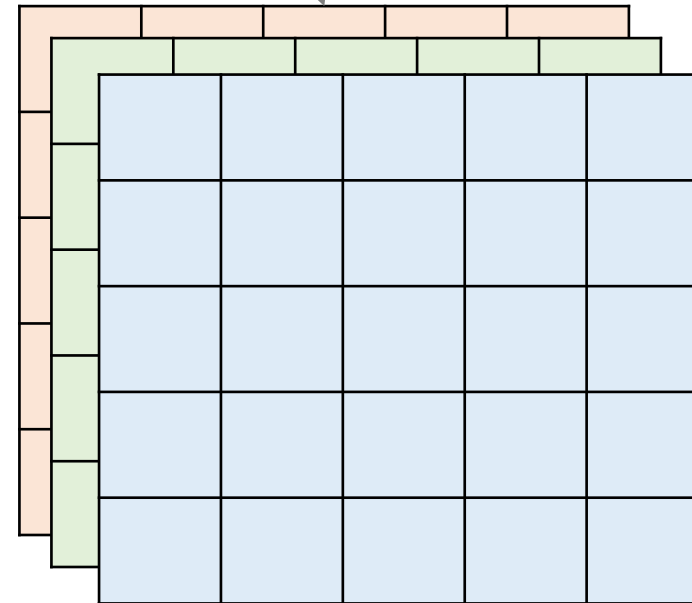
■ Graph Data Structure

- Node = Vertex
 - ✓ Represent elements of a system
- Edge
 - ✓ Relationship or interaction between nodes



03 | Graph Neural Networks

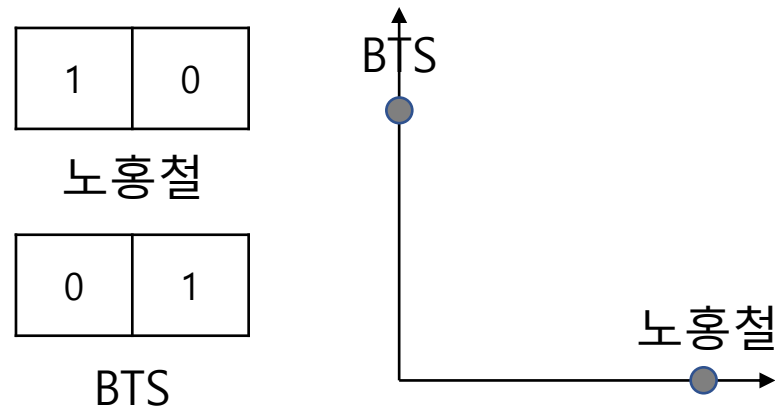
- Graph Data Structure
 - Image data: Euclidean space



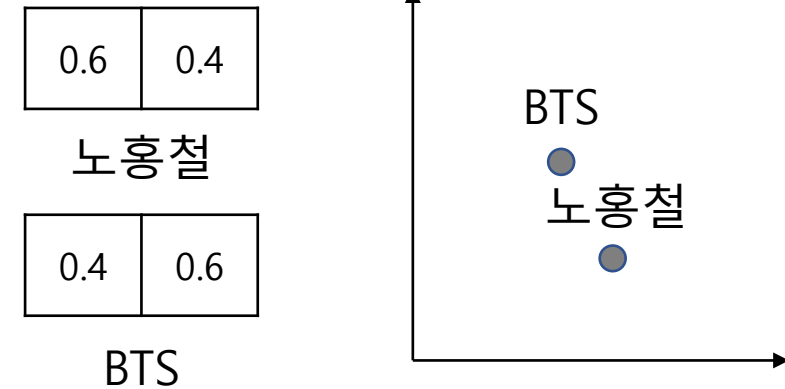
[3 X W X H]
dimension

03 | Graph Neural Networks

- Graph Data Structure
 - Text data: Euclidean space



One-hot encoding



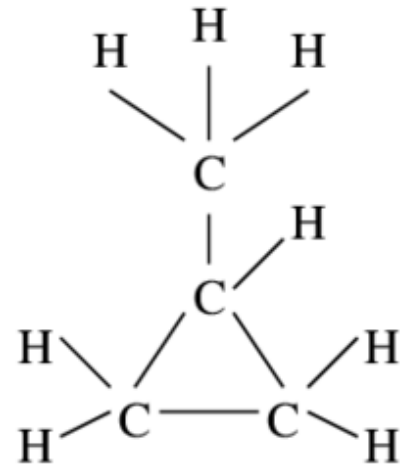
Distributed representation

03 | Graph Neural Networks

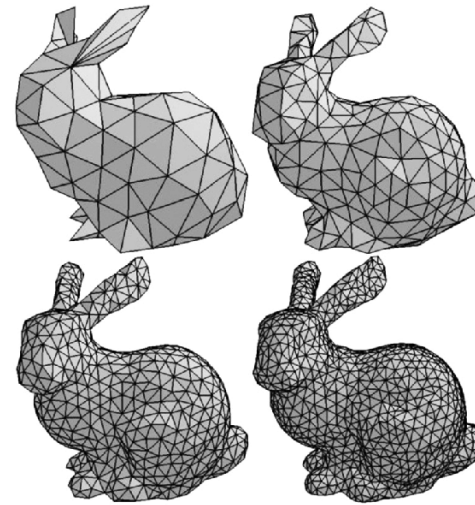
- Graph Data Structure
 - Graph data: Non-Euclidean Space



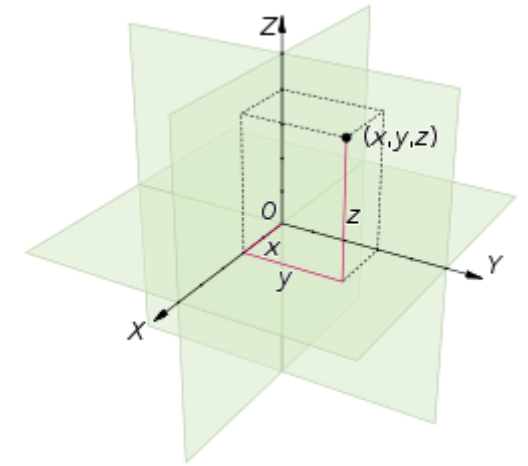
Social Network



Molecular Graph



3D Mesh

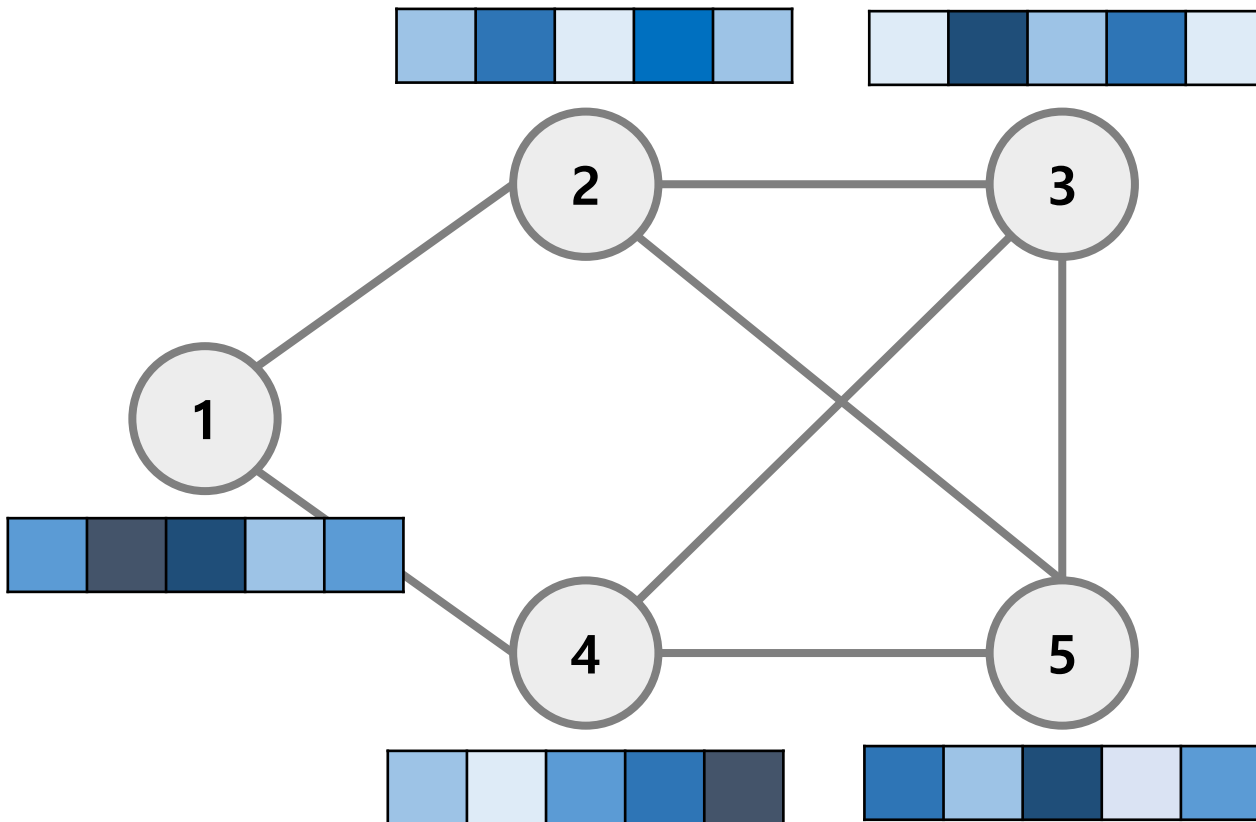


Euclidean Space

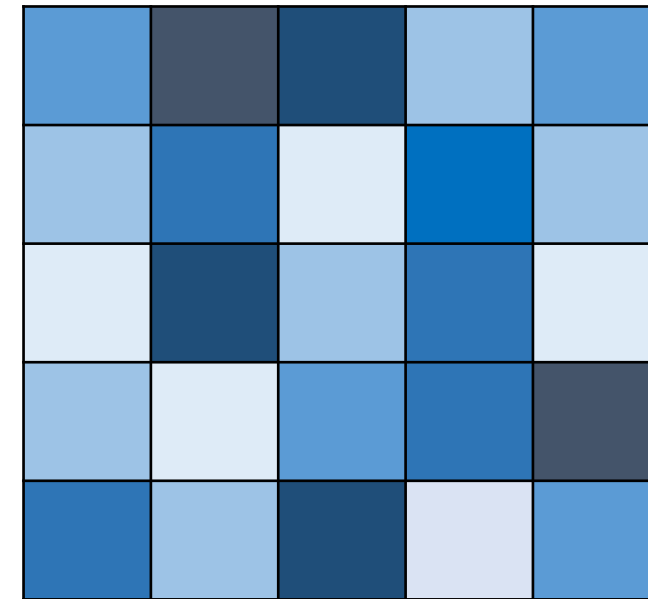
03 | Graph Neural Networks

Matrix Representation of Graph

- Node-Feature matrix
 - ✓ N by D dimension

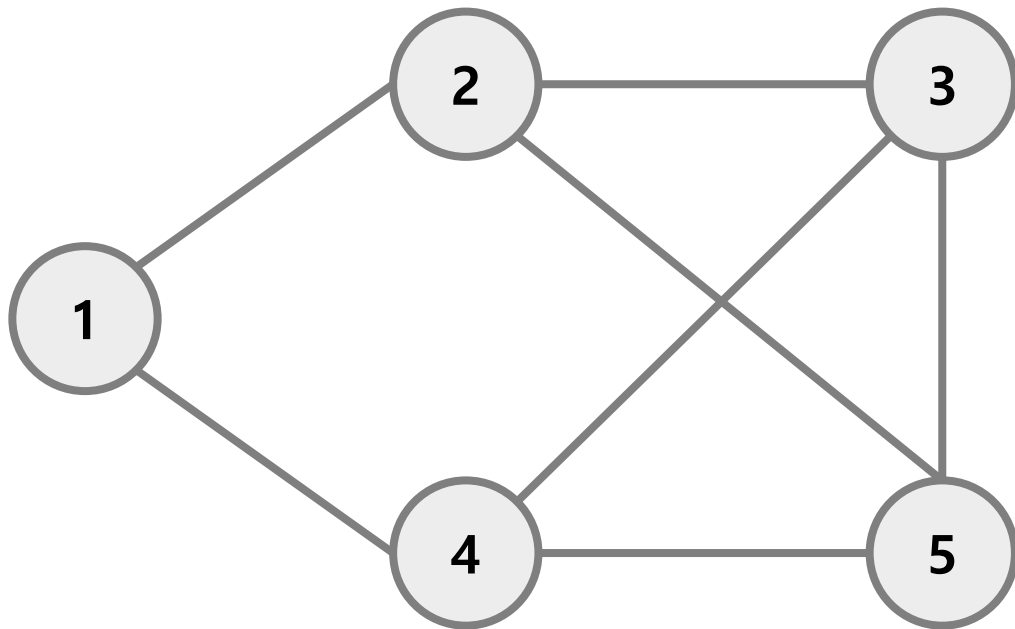


[Node-Feature Matrix]



03 | Graph Neural Networks

- Matrix Representation of Graph
 - Adjacency matrix: Undirected graph
 - ✓ N by N square matrix
 - ✓ Symmetric



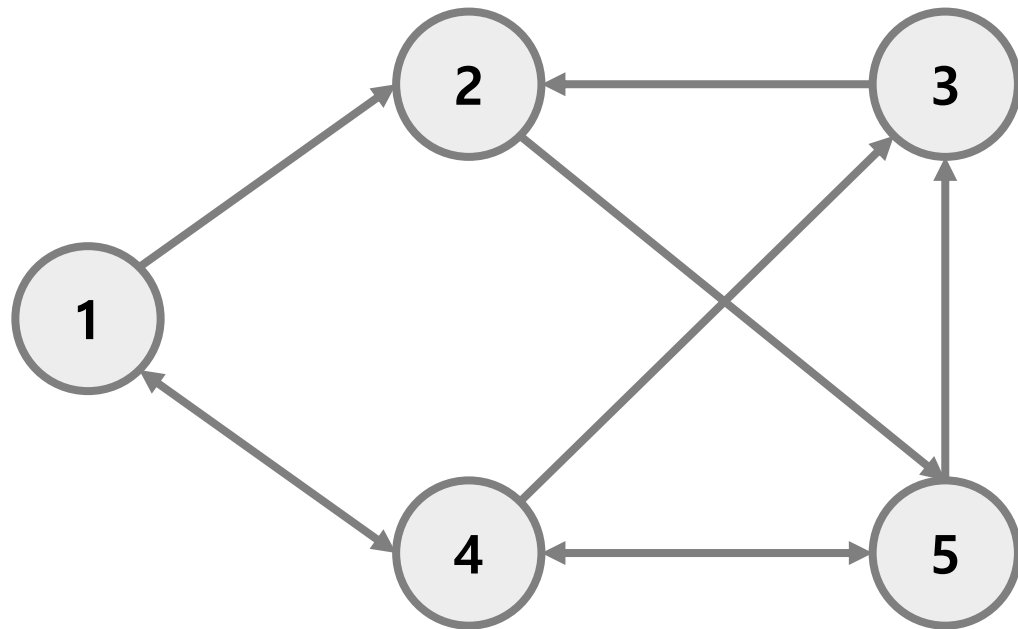
[Adjacency Matrix]

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

03 | Graph Neural Networks

Matrix Representation of Graph

- Adjacency matrix: Directed graph
 - ✓ N by N
 - ✓ Asymmetric

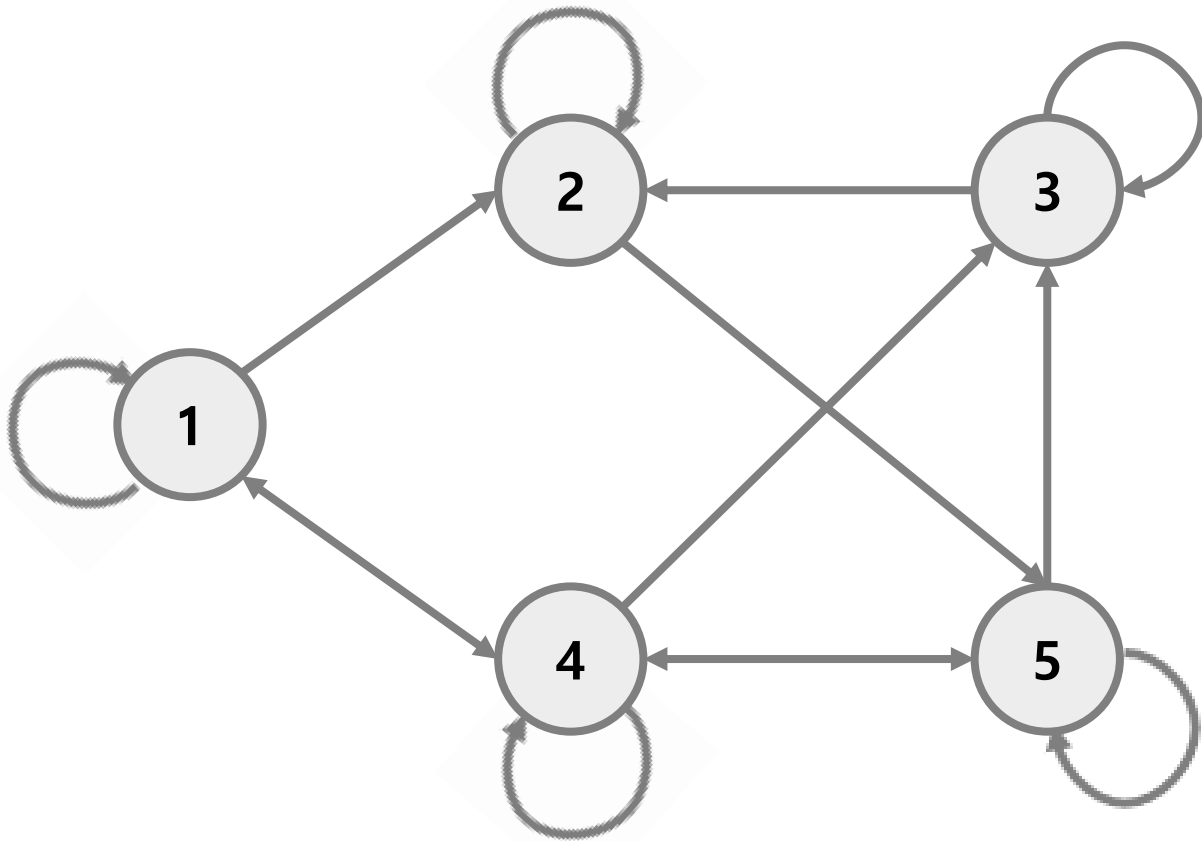


[Adjacency Matrix]

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

03 | Graph Neural Networks

- Matrix Representation of Graph
 - Adjacency matrix: Directed graph
 - ✓ Adjacency matrix + Identity matrix

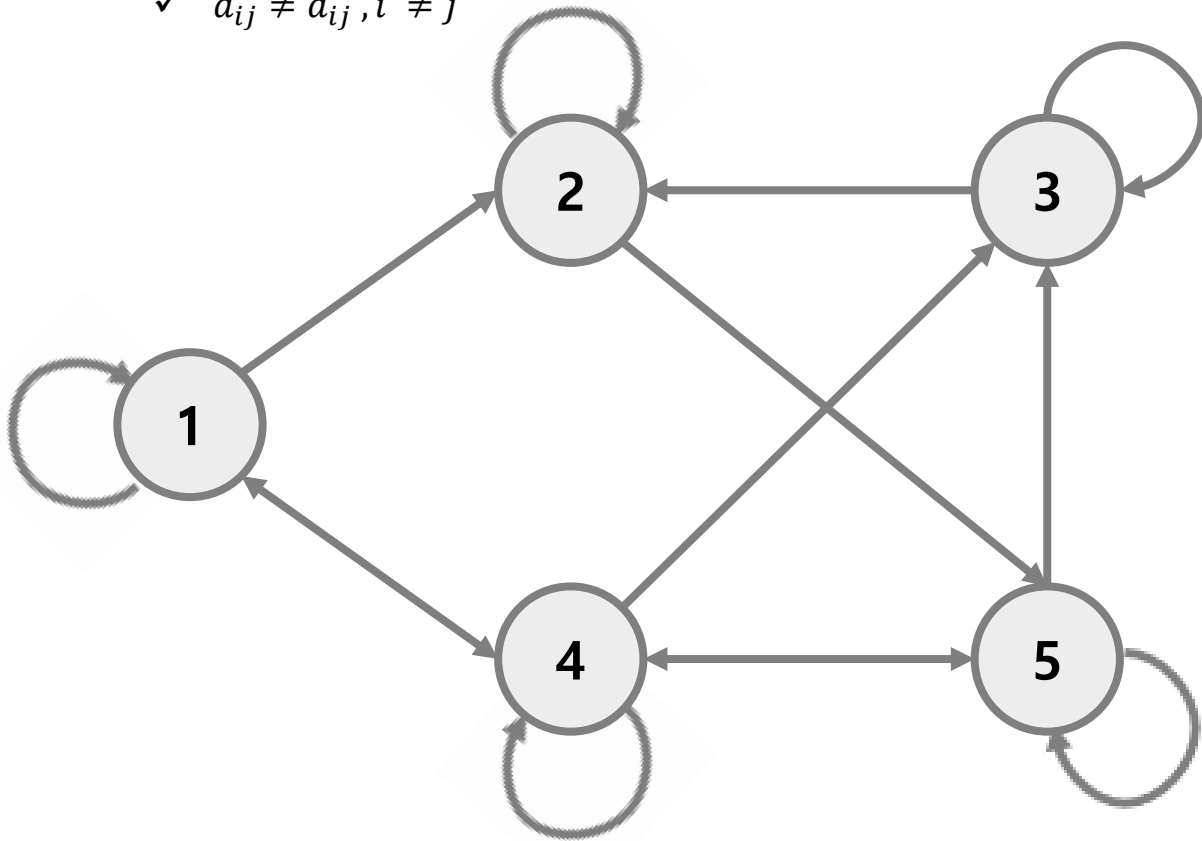


[Adjacency Matrix]

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

03 | Graph Neural Networks

- Matrix Representation of Graph
 - Adjacency matrix: Weighted directed graph
 - ✓ Edge information
 - ✓ $a_{ij} \neq a_{ji}, i \neq j$

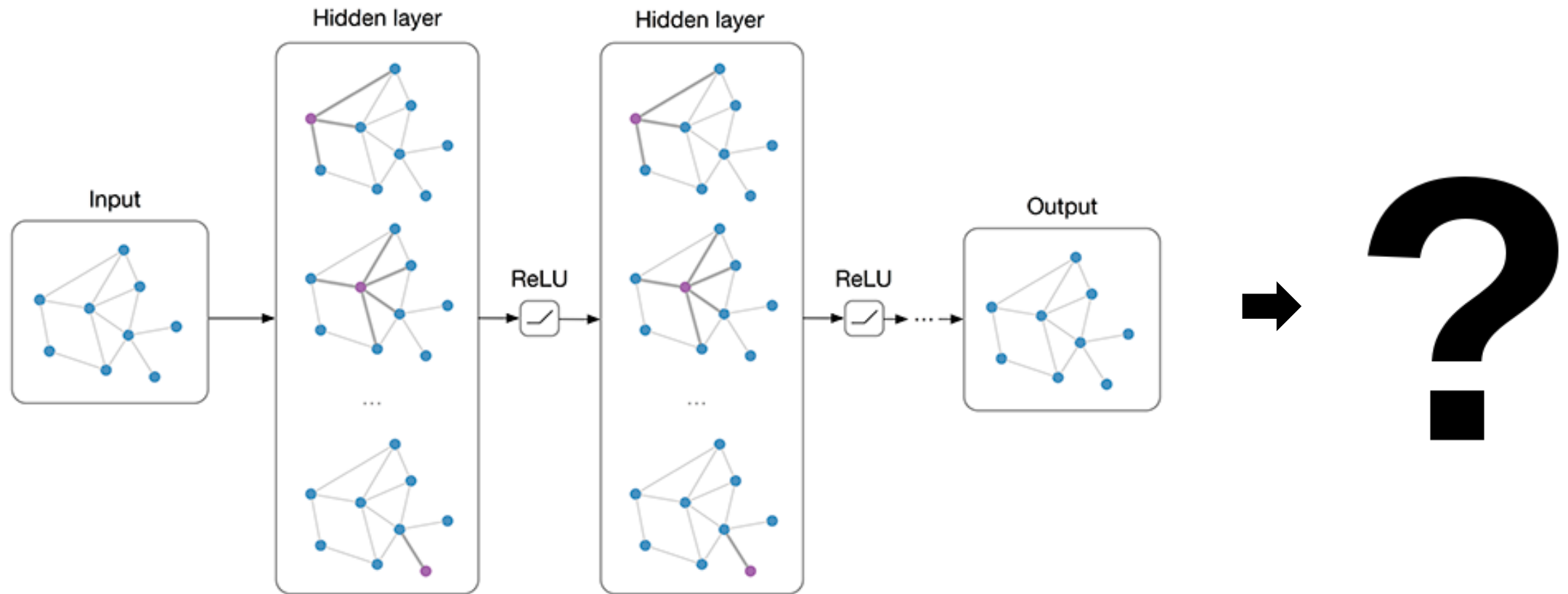


[Adjacency Matrix]

a_{11}	a_{12}	0	a_{14}	0
0	a_{22}	0	0	a_{25}
0	a_{32}	a_{33}	0	0
a_{41}	0	a_{43}	a_{44}	a_{45}
0	0	a_{53}	a_{54}	a_{55}

03 | Graph Neural Networks

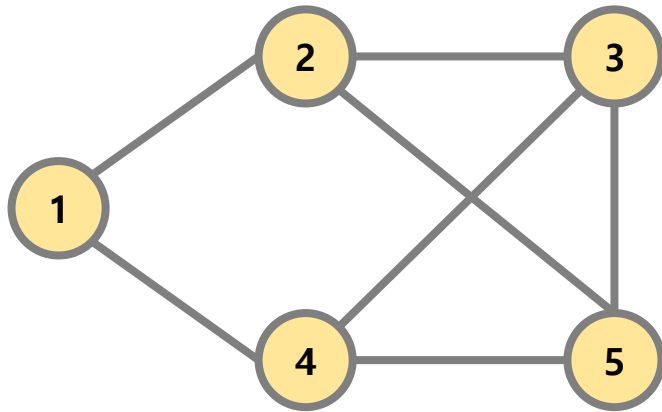
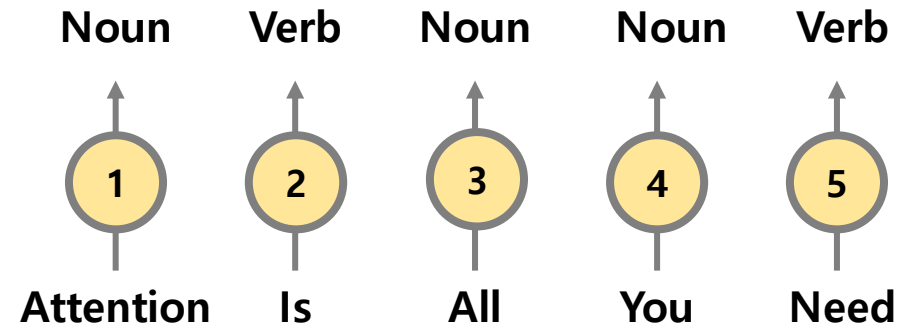
- Graph Neural Networks
 - Neural Networks for learning the structures of graphs



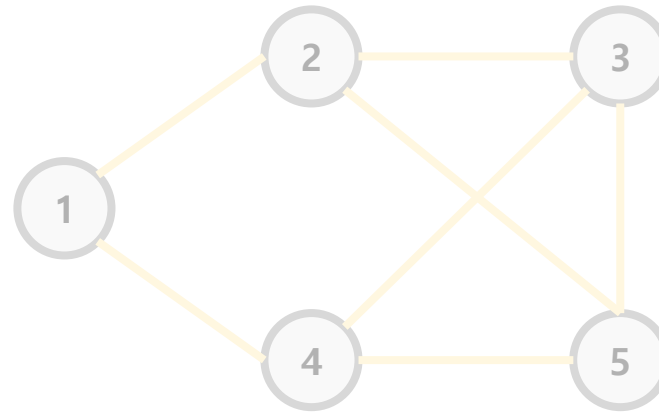
03 | Graph Neural Networks

▪ GNN Tasks

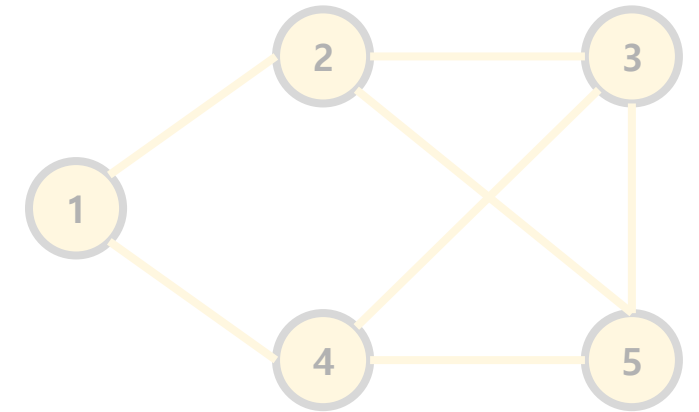
- Node Level
- Edge Level
- Graph Level



Node Level



Edge Level

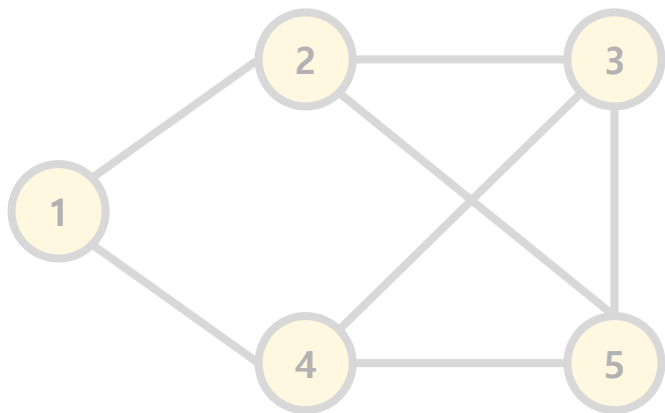
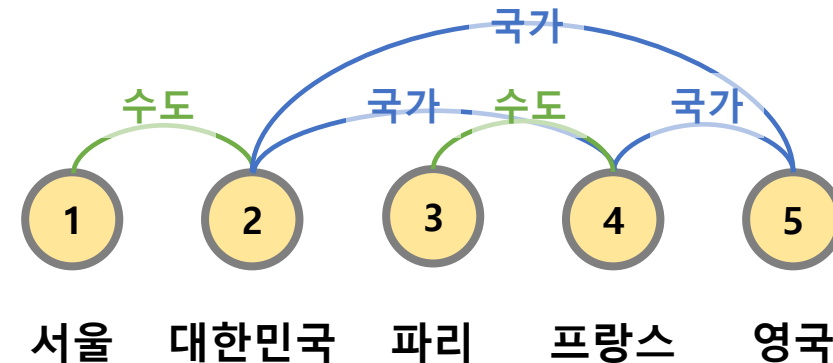
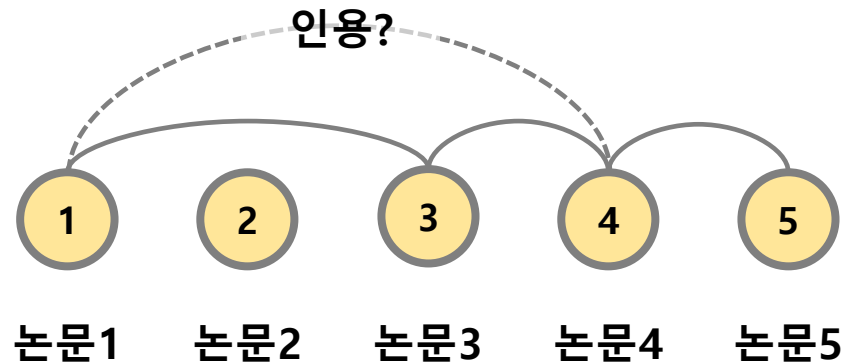


Graph Level

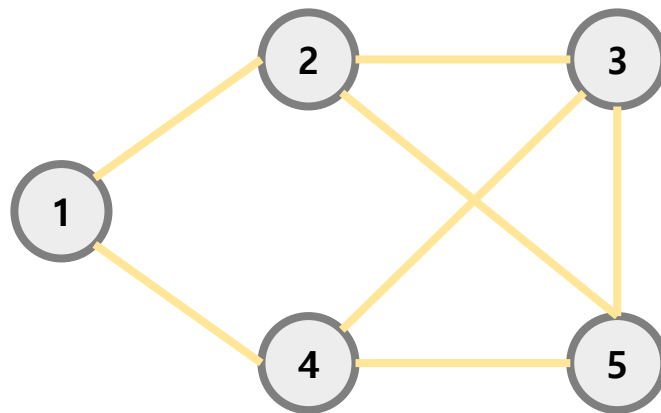
03 | Graph Neural Networks

▪ GNN Tasks

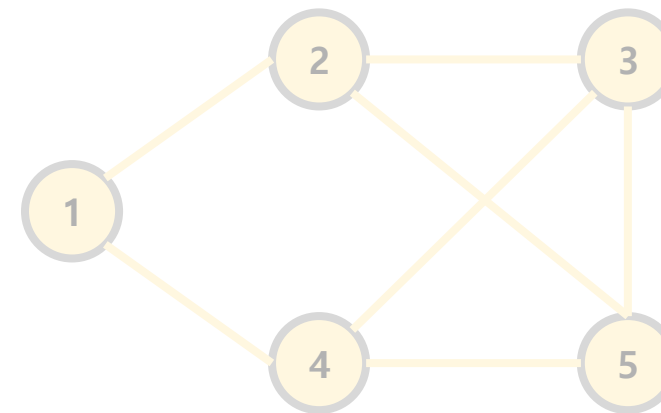
- Node Level
- Edge Level
- Graph Level



Node Level



Edge Level

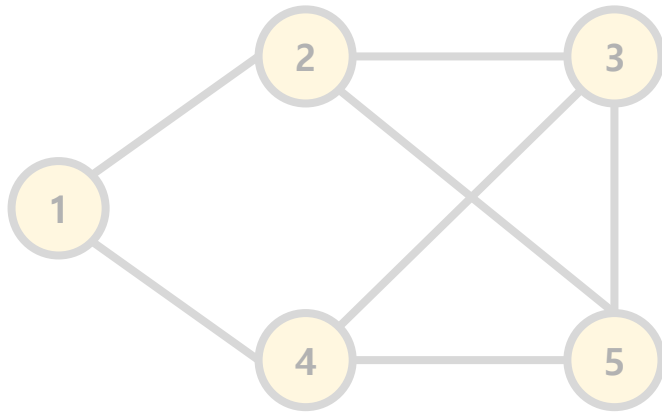
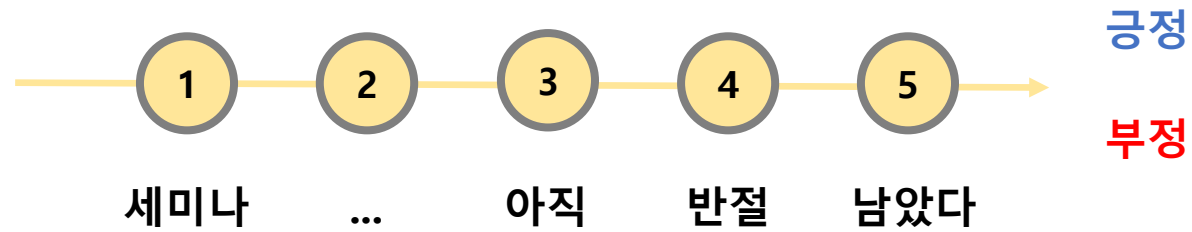


Graph Level

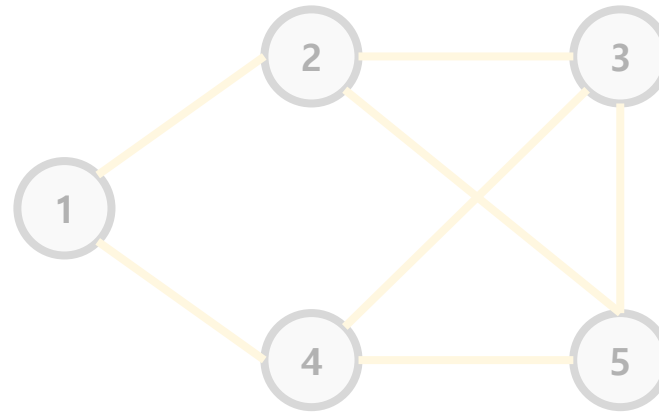
03 | Graph Neural Networks

▪ GNN Tasks

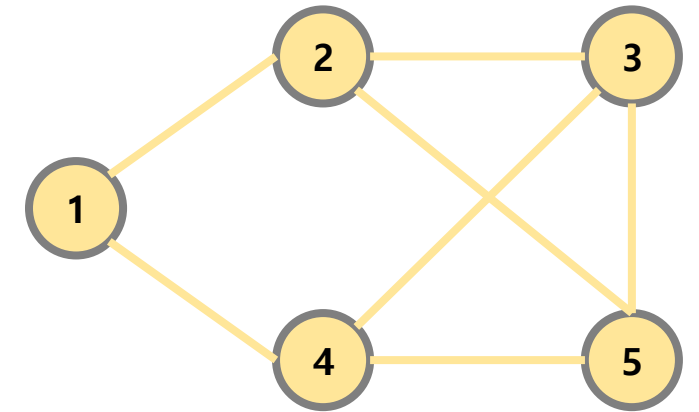
- Node Level
- Edge Level
- Graph Level



Node Level



Edge Level



Graph Level

03 | Graph Neural Networks

- Graph Task
 - Computer Vision
 - Natural Language Processing
 - Etc.

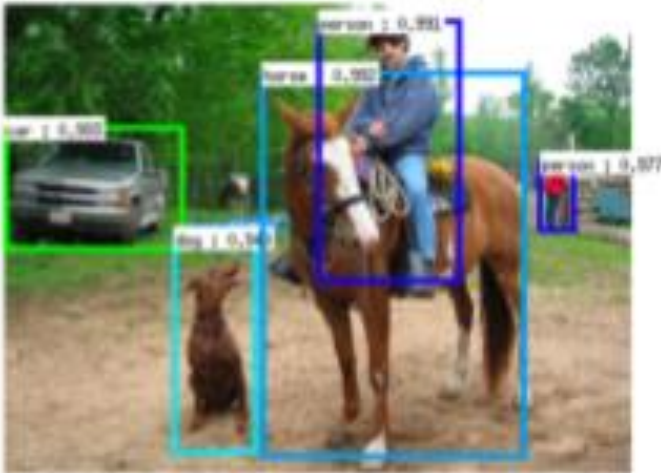
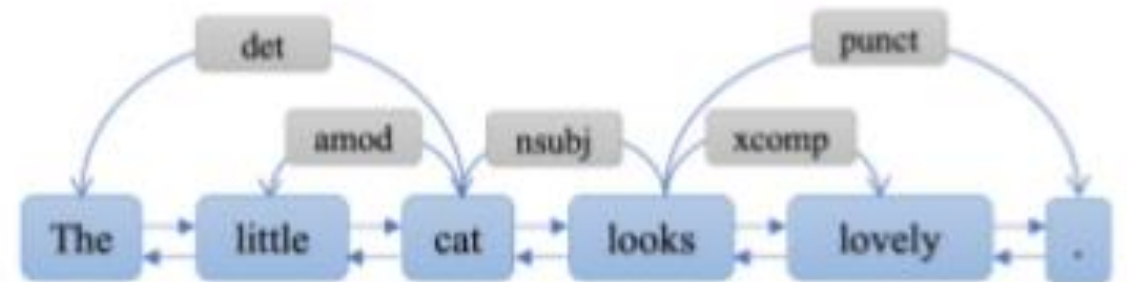


Image Graph

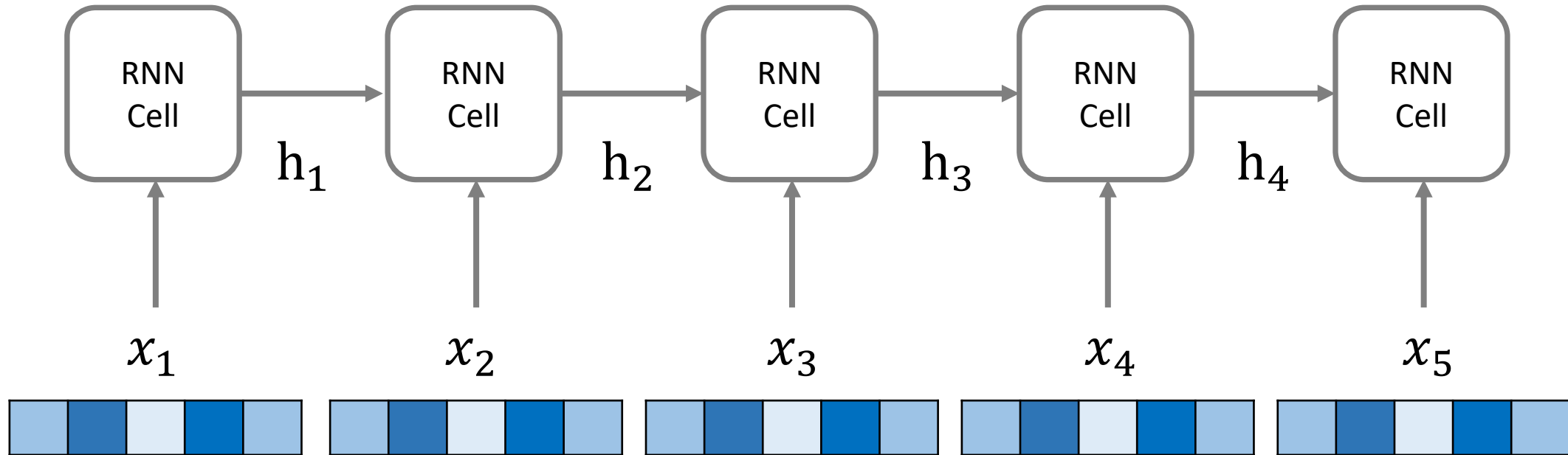


Text Graph

Zhou, Jie, et al. "Graph neural networks: A review of methods and applications." *arXiv preprint arXiv:1812.08434* (2018).

03 | Graph Neural Networks

- GNN Learning Process
 - RNN learning process



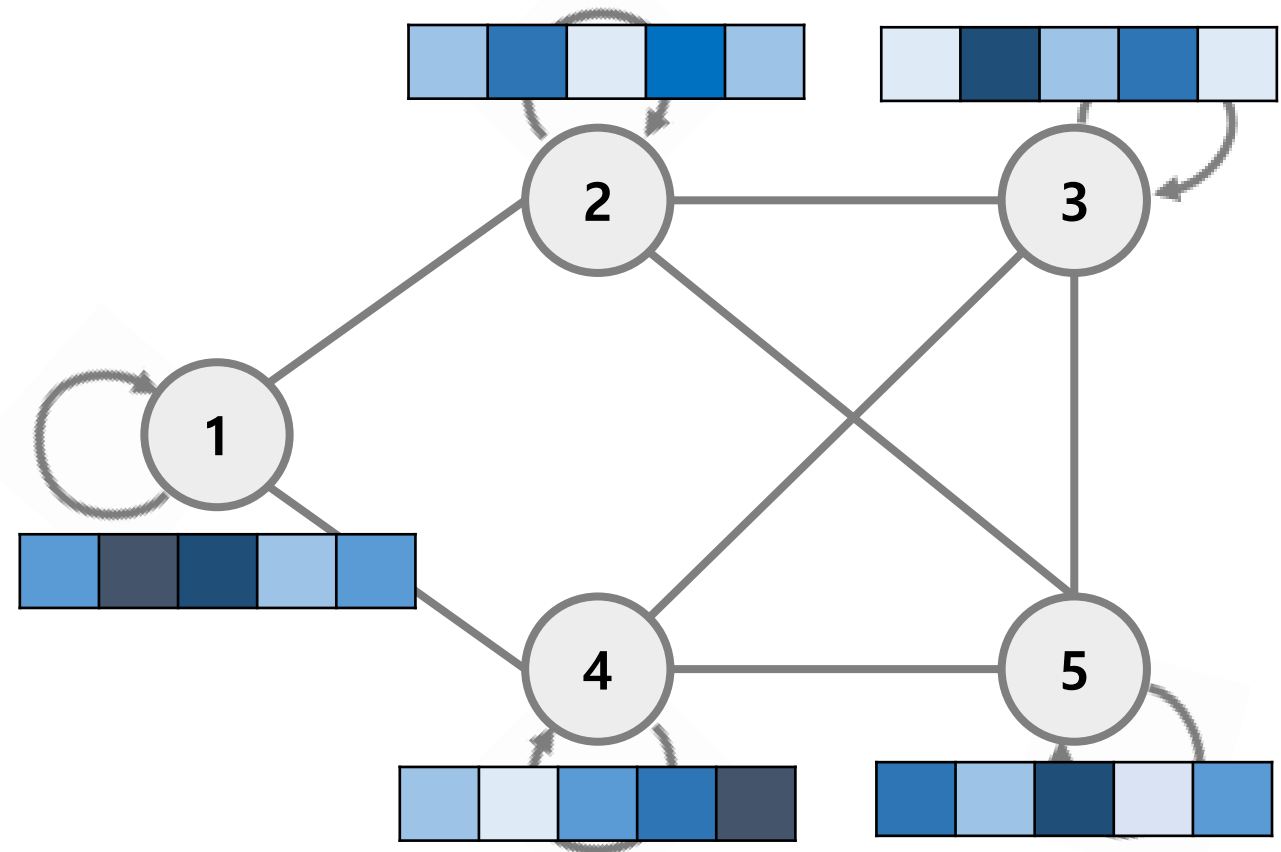
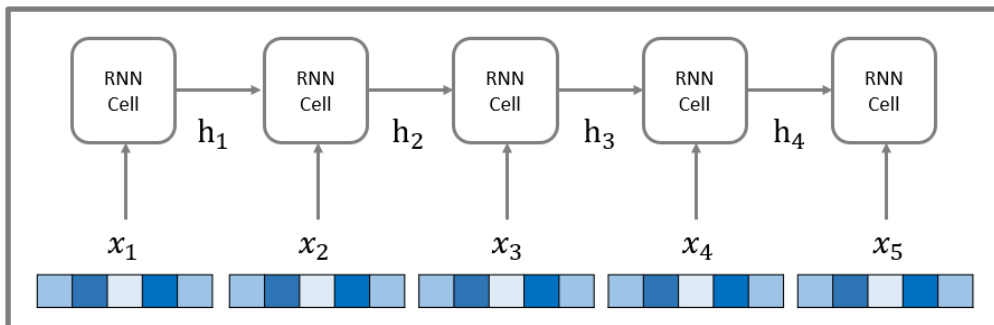
$$h_t = \text{combine}(h_{t-1}, x_t)$$

$$\text{combine} \in \{RNN, LSTM, GRU\}$$

03 | Graph Neural Networks

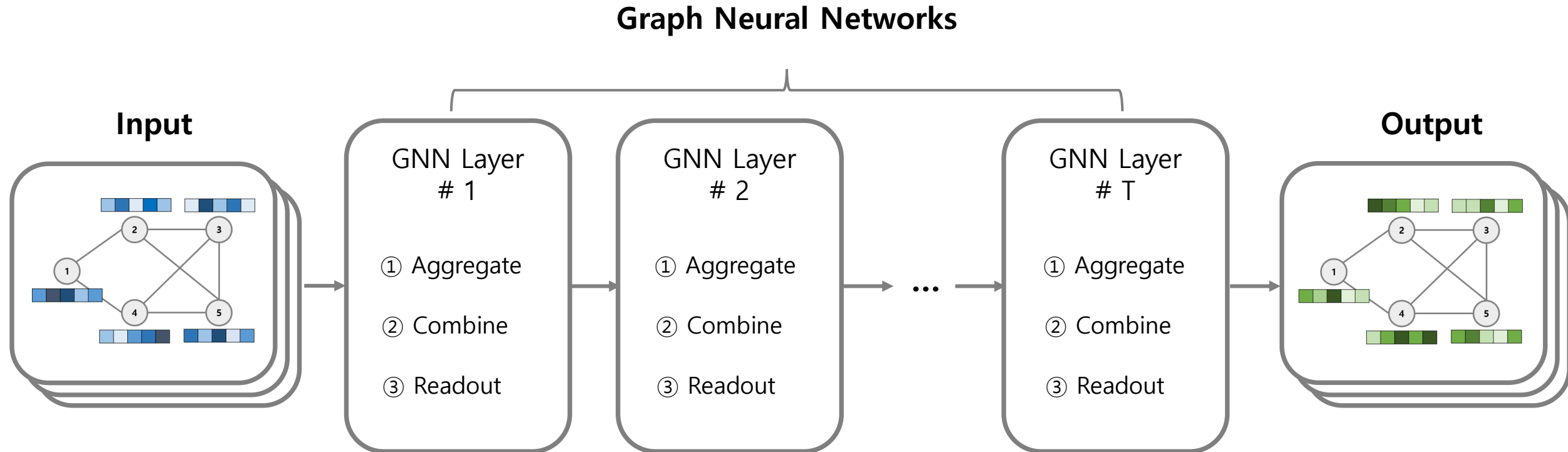
■ GNN Learning Process

- Graph: different with sequential data
 - ✓ No sequences(ordering)
 - ✓ Various graph structures
 - ✓ Multiple in-edge per node
- Considerations for encoding graphs
 - ✓ Information passes along edges
 - ✓ Information passes in parallel
 - ✓ Target nodes affected by multiple nodes



03 | Graph Neural Networks

- GNN Layer
 - Node feature update reflecting graph structures
 - ① Aggregate / Message passing
 - ② Combine / Update
 - ③ Readout



03 | Graph Neural Networks

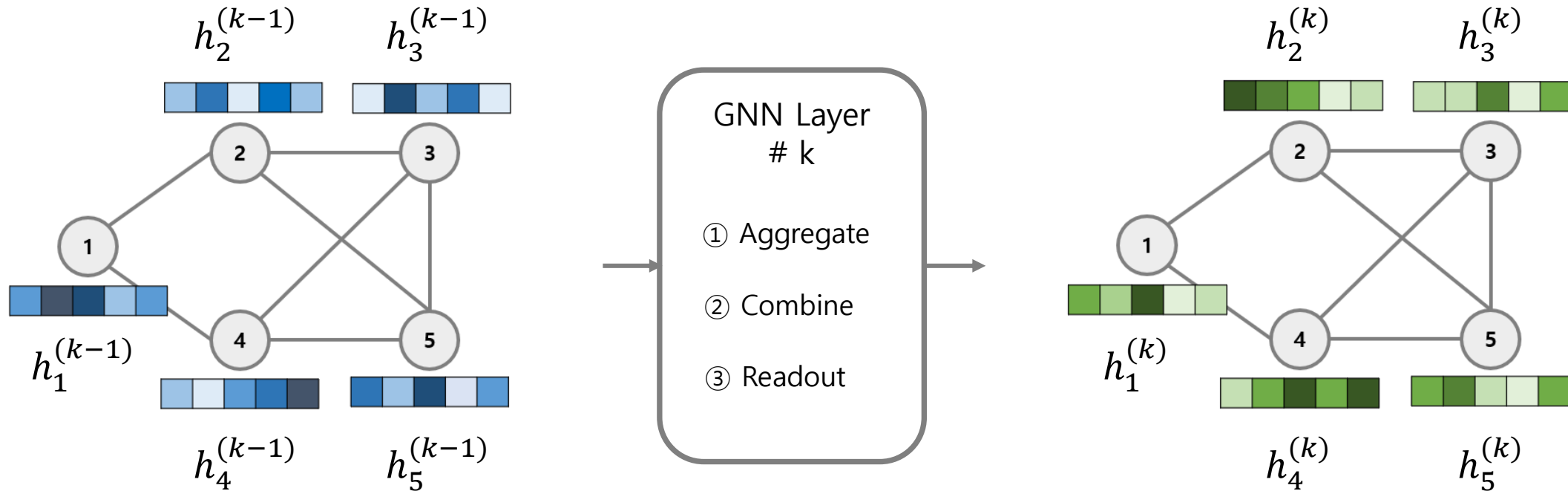
▪ GNN Notation

- $h_v^{(k)}$: hidden embedding node v at k th GNN layer
- $v = \text{target node}$
- $N(v) = \text{neighbor nodes of } v$
- $u = \text{neighbor node} \in N(v)$

$$\text{graph} = G(A, X)$$

$$X = \text{Node - Feature Matrix}$$

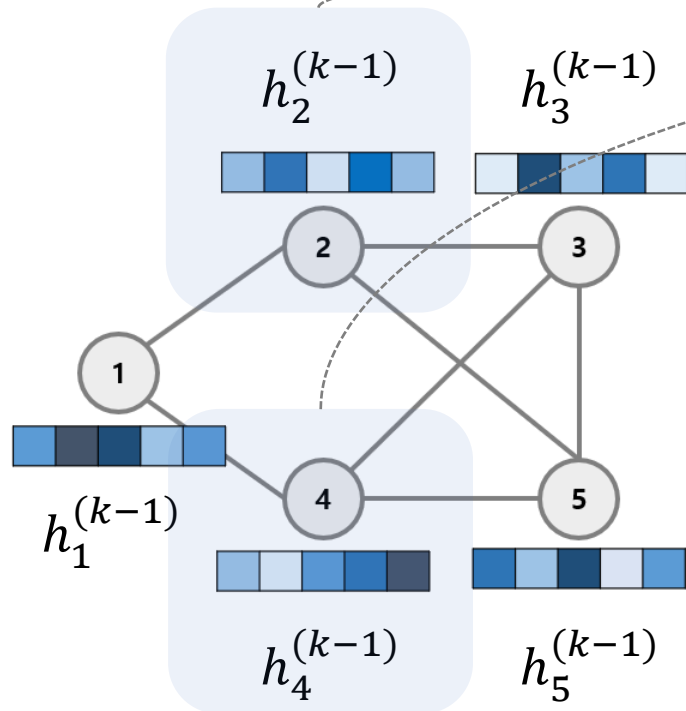
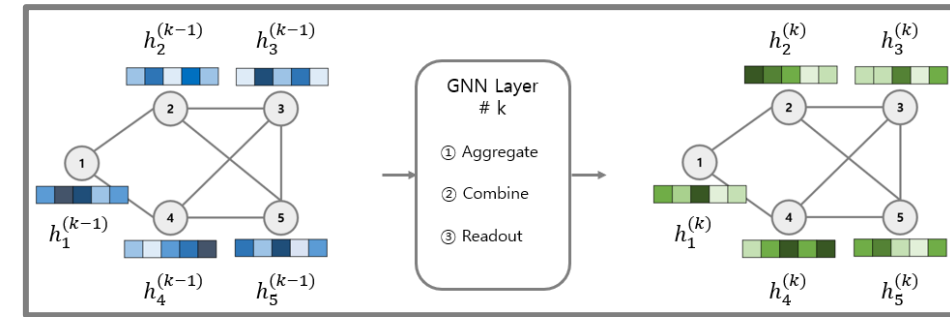
$$A = \text{Adjacency Matrix}$$



03 | Graph Neural Networks

■ GNN: ① Aggregate

- 타겟 노드의 이웃 노드들의 k-1 시점의 hidden state를 결합



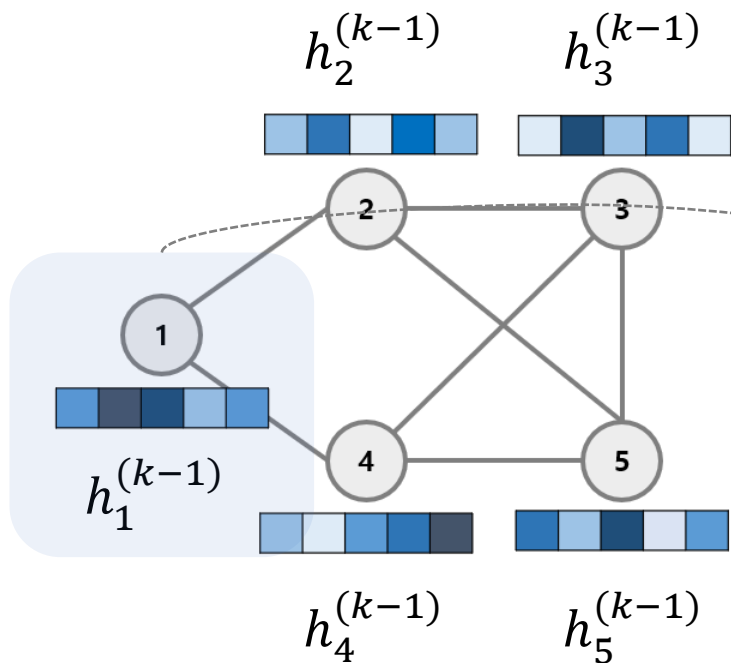
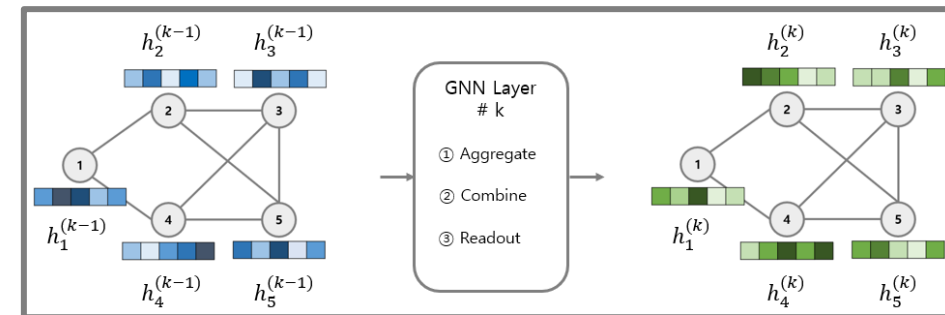
$$a_1^{(k-1)} = \text{aggregate}^k(\{ \text{[blue bar 1]}, \text{[blue bar 2]}, \text{[blue bar 3]}, \text{[blue bar 4]}, \text{[blue bar 5]} \})$$

$$a_v^{(k-1)} = \text{aggregate}^k(\{h_u^{(k-1)}, \forall u \in N(v)\})$$

03 | Graph Neural Networks

■ GNN: ② Combine

- k-1 시점 target node의 hidden state와 aggregated information을 사용하여 k 시점의 target node의 hidden state를 update



$$a_1^{(k-1)} = \text{aggregate}^k(\{ \text{[blue bar 1]}, \text{[blue bar 2]}, \text{[blue bar 3]}, \text{[blue bar 4]}, \text{[blue bar 5]} \})$$

$$h_1^{(k)} = \text{combine}^k(a_1^{(k-1)}, \text{[blue bar 1]})$$

$$= \text{[green bar 1]}$$

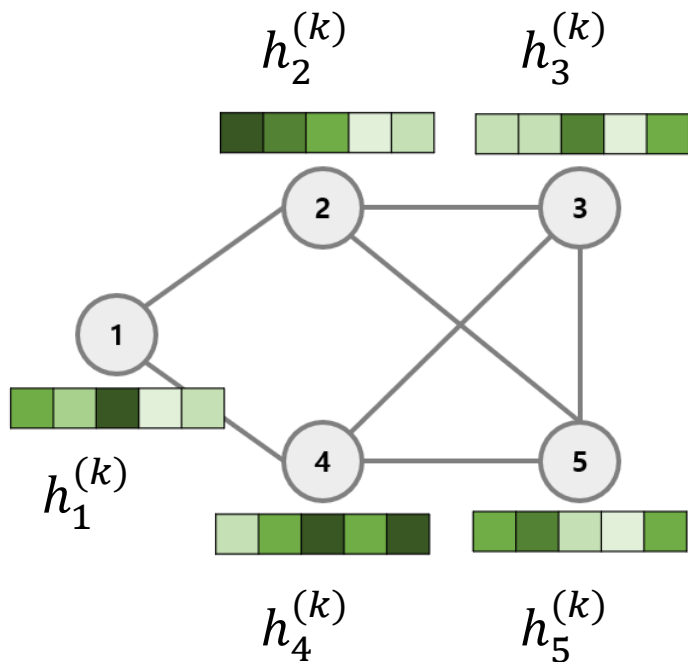
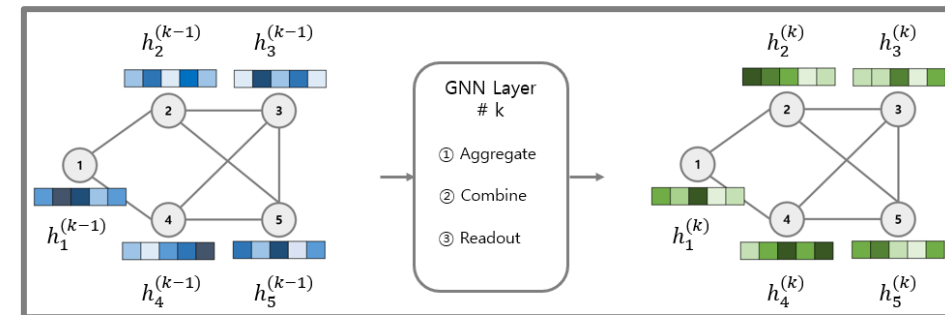
$$h_v^{(k)} = \text{combine}^k(a_u^{(k-1)}, h_v^{(k-1)})$$



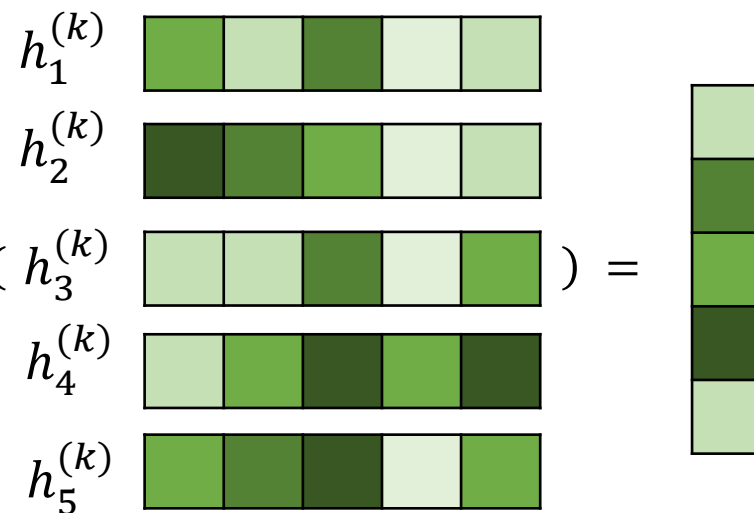
03 | Graph Neural Networks

■ GNN: ③ Readout

- K 시점의 모든 Node들의 hidden state를 결합하여 graph의 hidden state 생성
- Graph level classification



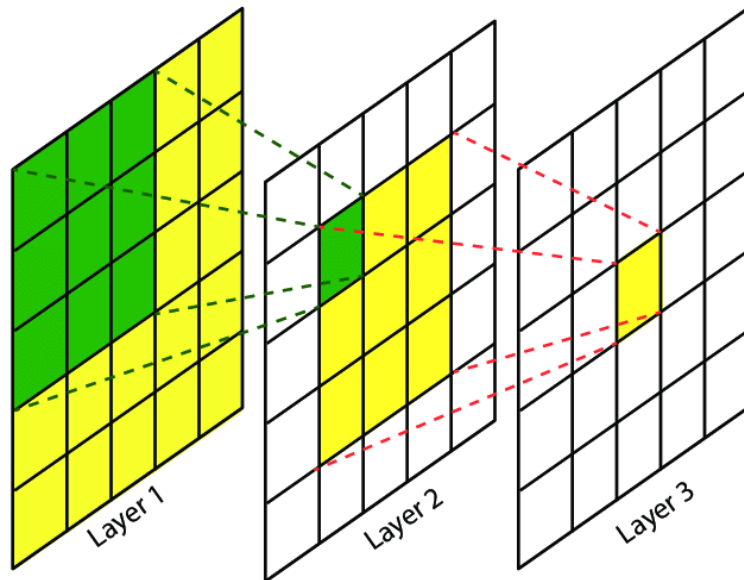
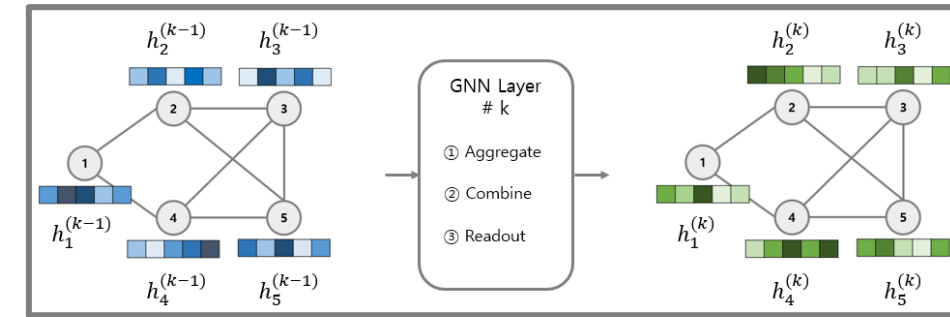
$$h_G^{(k)} = \text{readout}^k ($$



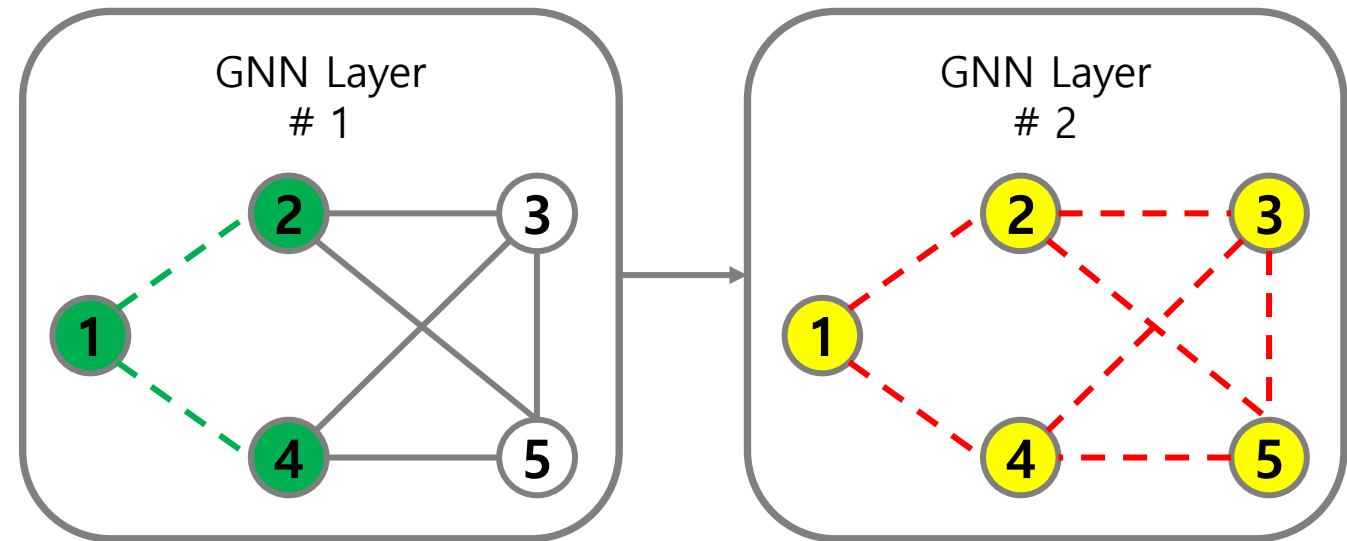
$$h_G^{(k)} = \text{readout}^k (h_v^{(k)}, \forall v \in G)$$

03 | Graph Neural Networks

- Stacking GNN Layer
 - CNN: Increase the receptive field
 - GNN: Increase hop of the graph



Convolutional Neural Networks



Graph Neural Networks

▪ GNN: Summary

- Aggregate

- ✓ $a_v^{(k-1)} = \text{aggregate}^k(\{h_u^{(k-1)}, \forall u \in N(v)\})$

- Combine

- ✓ $h_v^{(k)} = \text{combine}^k(a_u^{(k-1)}, h_v^{(k-1)})$

- Readout

- ✓ For graph level task

- ✓ $h_G^{(k)} = \text{readout}^k(h_v^{(k)}, \forall v \in G)$

- Stacking GNN Layer

- ✓ Increase hop of the graph



03 | Graph Neural Networks

■ GNN Variants

- Aggregate / Combine function의 정의에 따라 다양한 방식의 모델이 존재함
- Differentiable function

Name	Variant	Aggregator	Updater
Spectral Methods	ChebNet	$N_k = \mathbf{T}_k(\bar{\mathbf{L}})\mathbf{X}$	$\mathbf{H} = \sum_{k=0}^K N_k \Theta_k$
	1 st -order model	$N_0 = \mathbf{X}$ $N_1 = \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \mathbf{X}$	$\mathbf{H} = N_0 \Theta_0 + N_1 \Theta_1$
	Single parameter	$\mathbf{N} = (\mathbf{I}_N + \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}}) \mathbf{X}$	$\mathbf{H} = \mathbf{N} \Theta$
	GCN	$\mathbf{N} = \bar{\mathbf{D}}^{-\frac{1}{2}} \bar{\mathbf{A}} \bar{\mathbf{D}}^{-\frac{1}{2}} \mathbf{X}$	$\mathbf{H} = \mathbf{N} \Theta$
Non-spectral Methods	Convolutional networks in [33]	$\mathbf{h}_{\mathcal{N}_v}^t = \mathbf{h}_v^{t-1} + \sum_{k=1}^{\mathcal{N}_v} \mathbf{h}_k^{t-1}$	$\mathbf{h}_v^t = \sigma(\mathbf{h}_{\mathcal{N}_v}^t \mathbf{W}_L^{\mathcal{N}_v})$
	DCNN	Node classification: $\mathbf{N} = \mathbf{P} * \mathbf{X}$ Graph classification: $\mathbf{N} = \frac{1}{N} \mathbf{P} * \mathbf{X} / N$	$\mathbf{H} = f(\mathbf{W}^c \odot \mathbf{N})$
	GraphSAGE	$\mathbf{h}_{\mathcal{N}_v}^t = \text{AGGREGATE}_t(\{\mathbf{h}_u^{t-1}, \forall u \in \mathcal{N}_v\})$	$\mathbf{h}_v^t = \sigma(\mathbf{W}^t \cdot [\mathbf{h}_v^{t-1} \parallel \mathbf{h}_{\mathcal{N}_v}^t])$
Graph Attention Networks	GAT	$\alpha_{vk} = \frac{\exp(\text{LeakyReLU}(a^T [\mathbf{W} \mathbf{h}_v \parallel \mathbf{W} \mathbf{h}_k]))}{\sum_{j \in \mathcal{N}_v} \exp(\text{LeakyReLU}(a^T [\mathbf{W} \mathbf{h}_v \parallel \mathbf{W} \mathbf{h}_j]))}$ $\mathbf{h}_{\mathcal{N}_v}^t = \sigma(\sum_{k \in \mathcal{N}_v} \alpha_{vk} \mathbf{W} \mathbf{h}_k)$ Multi-head concatenation: $\mathbf{h}_{\mathcal{N}_v}^t = \parallel_{m=1}^M \sigma(\sum_{k \in \mathcal{N}_v} \alpha_{vk}^m \mathbf{W}^m \mathbf{h}_k)$ Multi-head average: $\mathbf{h}_{\mathcal{N}_v}^t = \sigma\left(\frac{1}{M} \sum_{m=1}^M \sum_{k \in \mathcal{N}_v} \alpha_{vk}^m \mathbf{W}^m \mathbf{h}_k\right)$	$\mathbf{h}_v^t = \mathbf{h}_{\mathcal{N}_v}^t$

Gated Graph Neural Networks	GGNN	$\mathbf{h}_{\mathcal{N}_v}^t = \sum_{k \in \mathcal{N}_v} \mathbf{h}_k^{t-1} + \mathbf{b}$	$\mathbf{z}_v^t = \sigma(\mathbf{W}^z \mathbf{h}_{\mathcal{N}_v}^t + \mathbf{U}^z \mathbf{h}_v^{t-1})$ $\mathbf{r}_v^t = \sigma(\mathbf{W}^r \mathbf{h}_{\mathcal{N}_v}^t + \mathbf{U}^r \mathbf{h}_v^{t-1})$ $\tilde{\mathbf{h}}_v^t = \tanh(\mathbf{W} \mathbf{h}_{\mathcal{N}_v}^t + \mathbf{U}(\mathbf{r}_v^t \odot \mathbf{h}_v^{t-1}))$ $\mathbf{h}_v^t = (1 - \mathbf{z}_v^t) \odot \mathbf{h}_v^{t-1} + \mathbf{z}_v^t \odot \tilde{\mathbf{h}}_v^t$
	Tree LSTM (Child sum)	$\mathbf{h}_{\mathcal{N}_v}^t = \sum_{k \in \mathcal{N}_v} \mathbf{h}_k^{t-1}$	$\mathbf{i}_v^t = \sigma(\mathbf{W}^i \mathbf{x}_v^t + \mathbf{U}^i \mathbf{h}_{\mathcal{N}_v}^t + \mathbf{b}^i)$ $\mathbf{f}_{vk}^t = \sigma(\mathbf{W}^f \mathbf{x}_v^t + \mathbf{U}^f \mathbf{h}_k^{t-1} + \mathbf{b}^f)$ $\mathbf{o}_v^t = \sigma(\mathbf{W}^o \mathbf{x}_v^t + \mathbf{U}^o \mathbf{h}_{\mathcal{N}_v}^t + \mathbf{b}^o)$ $\mathbf{u}_v^t = \tanh(\mathbf{W}^u \mathbf{x}_v^t + \mathbf{U}^u \mathbf{h}_{\mathcal{N}_v}^t + \mathbf{b}^u)$ $\mathbf{c}_v^t = \mathbf{i}_v^t \odot \mathbf{u}_v^t + \sum_{k \in \mathcal{N}_v} \mathbf{f}_{vk}^t \odot \mathbf{c}_k^{t-1}$ $\mathbf{h}_v^t = \mathbf{o}_v^t \odot \tanh(\mathbf{c}_v^t)$
Graph LSTM	Tree LSTM (N-ary)	$\mathbf{h}_{\mathcal{N}_v}^{ti} = \sum_{l=1}^K \mathbf{U}_l^i \mathbf{h}_{vl}^{t-1}$ $\mathbf{h}_{\mathcal{N}_v}^{tf} = \sum_{l=1}^K \mathbf{U}_l^f \mathbf{h}_{vl}^{t-1}$ $\mathbf{h}_{\mathcal{N}_v}^{to} = \sum_{l=1}^K \mathbf{U}_l^o \mathbf{h}_{vl}^{t-1}$ $\mathbf{h}_{\mathcal{N}_v}^{tu} = \sum_{l=1}^K \mathbf{U}_l^u \mathbf{h}_{vl}^{t-1}$	$\mathbf{i}_v^t = \sigma(\mathbf{W}^i \mathbf{x}_v^t + \mathbf{h}_{\mathcal{N}_v}^{ti} + \mathbf{b}^i)$ $\mathbf{f}_{vk}^t = \sigma(\mathbf{W}^f \mathbf{x}_v^t + \mathbf{h}_{\mathcal{N}_v}^{tf} + \mathbf{b}^f)$ $\mathbf{o}_v^t = \sigma(\mathbf{W}^o \mathbf{x}_v^t + \mathbf{h}_{\mathcal{N}_v}^{to} + \mathbf{b}^o)$ $\mathbf{u}_v^t = \tanh(\mathbf{W}^u \mathbf{x}_v^t + \mathbf{h}_{\mathcal{N}_v}^{tu} + \mathbf{b}^u)$ $\mathbf{c}_v^t = \mathbf{i}_v^t \odot \mathbf{u}_v^t + \sum_{l=1}^K \mathbf{f}_{vl}^t \odot \mathbf{c}_{vl}^{t-1}$ $\mathbf{h}_v^t = \mathbf{o}_v^t \odot \tanh(\mathbf{c}_v^t)$
	Graph LSTM in [34]	$\mathbf{h}_{\mathcal{N}_v}^{ti} = \sum_{k \in \mathcal{N}_v} \mathbf{U}_{m(v,k)}^i \mathbf{h}_k^{t-1}$ $\mathbf{h}_{\mathcal{N}_v}^{to} = \sum_{k \in \mathcal{N}_v} \mathbf{U}_{m(v,k)}^o \mathbf{h}_k^{t-1}$ $\mathbf{h}_{\mathcal{N}_v}^{tu} = \sum_{k \in \mathcal{N}_v} \mathbf{U}_{m(v,k)}^u \mathbf{h}_k^{t-1}$	$\mathbf{i}_v^t = \sigma(\mathbf{W}^i \mathbf{x}_v^t + \mathbf{h}_{\mathcal{N}_v}^{ti} + \mathbf{b}^i)$ $\mathbf{f}_{vk}^t = \sigma(\mathbf{W}^f \mathbf{x}_v^t + \mathbf{U}_{m(v,k)}^f \mathbf{h}_k^{t-1} + \mathbf{b}^f)$ $\mathbf{o}_v^t = \sigma(\mathbf{W}^o \mathbf{x}_v^t + \mathbf{h}_{\mathcal{N}_v}^{to} + \mathbf{b}^o)$ $\mathbf{u}_v^t = \tanh(\mathbf{W}^u \mathbf{x}_v^t + \mathbf{h}_{\mathcal{N}_v}^{tu} + \mathbf{b}^u)$ $\mathbf{c}_v^t = \mathbf{i}_v^t \odot \mathbf{u}_v^t + \sum_{k \in \mathcal{N}_v} \mathbf{f}_{vk}^t \odot \mathbf{c}_k^{t-1}$ $\mathbf{h}_v^t = \mathbf{o}_v^t \odot \tanh(\mathbf{c}_v^t)$

Zhou, Jie, et al. "Graph neural networks: A review of methods and applications." *arXiv preprint arXiv:1812.08434* (2018).

03 | Graph Neural Networks

▪ GNN

- Aggregate

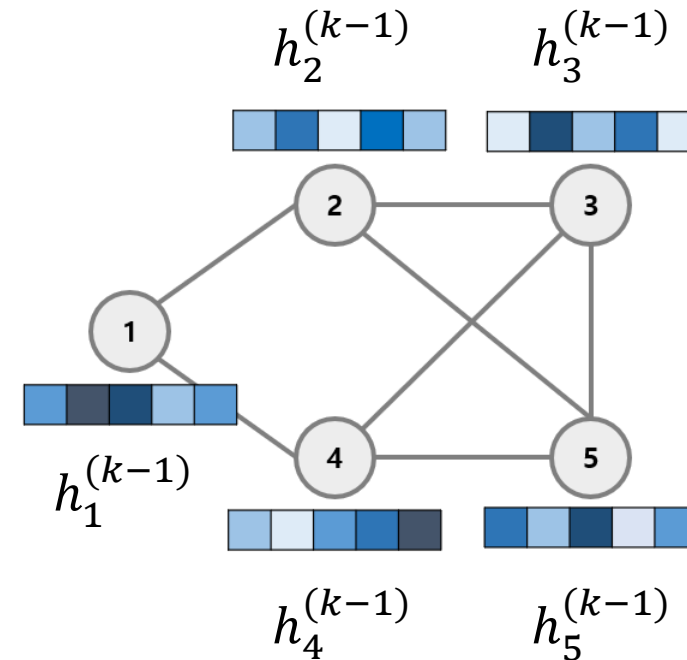
$$\checkmark a_v^{(k-1)} = \sum_{u \in N(v)} h_u^{(k-1)}$$

- Combine

$$\checkmark h_v^{(k)} = \text{Relu}(W_{\text{self}}h_v^{(k-1)} + W_{\text{neigh}}a_v^{(k-1)})$$

- Matrix form

$$\checkmark H^{(t)} = \text{Relu}(H^{(t-1)}W_{\text{self}} + AH^{(t-1)}W_{\text{neigh}}a_v^{(k-1)})$$



03 | Graph Neural Networks

▪ GNN (Self-loop)

- When aggregating, self information is needed

- $W_{\text{self}} = W_{\text{neigh}} = W$

- Aggregate

$$\checkmark a_v^{(k-1)} = \sum_{u \in N(v)} h_u^{(k-1)} + h_v^{(k-1)}$$

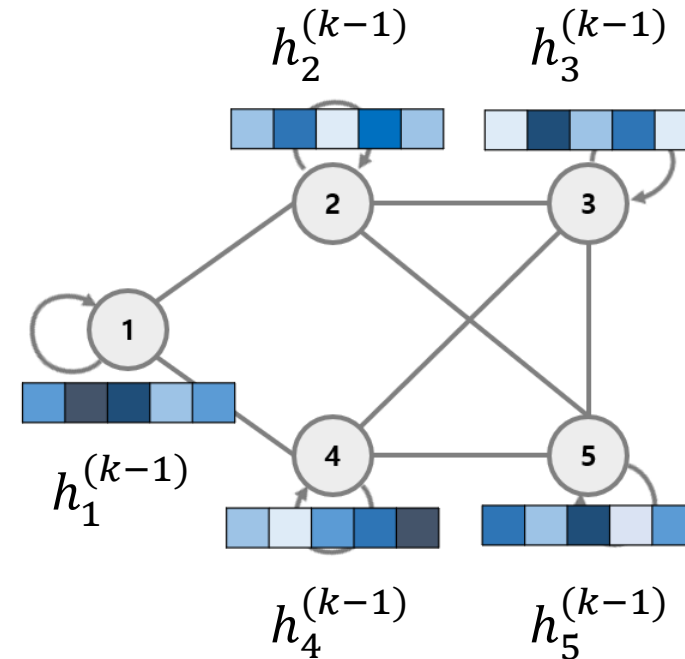
- Combine

$$\checkmark h_v^{(k)} = \text{Relu}(W a_v^{(k-1)})$$

$$\checkmark h_v^{(k)} = \text{Relu}(W \sum_{u \in N(v) \cup \{v\}} h_u^{(k-1)})$$

- Matrix form

$$\checkmark H^{(t)} = \text{Relu}((A + I)H^{(t-1)}W)$$



Graph Convolutional Networks (2016)

- Amsterdam Univ., CIFAR
- If graph size is too large
 - ✓ Unstable and sensitive to node degrees
 - ✓ Degree = # of neighbor nodes for each node
- Normalized aggregate function

$$a_v^{(k-1)} = \sum_{u \in N(v) \cup \{v\}} \frac{h_u^{(k-1)}}{\sqrt{|N(v)| |N(u)|}}$$

- Combine
 - ✓ $h_v^{(k)} = \text{Relu}(W a_v^{(k-1)})$
- Matrix form

$$H^{(t)} = \text{Relu}((D^{-\frac{1}{2}}(A + I) D^{-\frac{1}{2}} H^{(t-1)} W)$$

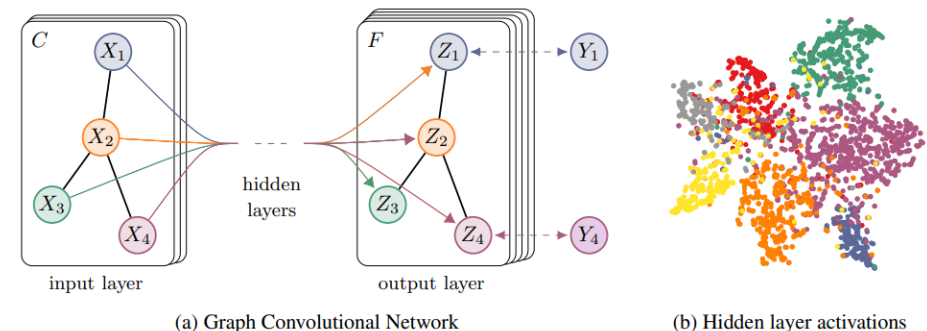
SEMI-SUPERVISED CLASSIFICATION WITH GRAPH CONVOLUTIONAL NETWORKS

Thomas N. Kipf
University of Amsterdam
T.N.Kipf@uva.nl

Max Welling
University of Amsterdam
Canadian Institute for Advanced Research (CIFAR)
M.Welling@uva.nl

ABSTRACT

We present a scalable approach for semi-supervised learning on graph-structured data that is based on an efficient variant of convolutional neural networks which operate directly on graphs. We motivate the choice of our convolutional architecture via a localized first-order approximation of spectral graph convolutions. Our model scales linearly in the number of graph edges and learns hidden layer representations that encode both local graph structure and features of nodes. In a number of experiments on citation networks and on a knowledge graph dataset we demonstrate that our approach outperforms related methods by a significant margin.



03 | Graph Neural Networks

▪ Gated Graph Neural Networks (2016)

- Stacking deep layers lead overfitting / vanishing gradient
- ICLR, Toronto Univ., Microsoft
- Aggregate

$$\checkmark a_v^{(k-1)} = \sum_{u \in N(v)} \sum h_u^{(k-1)}$$

- Combine

$$\checkmark h_v^{(k)} = GRU(h_v^{(k-1)}, a_v^{(k-1)})$$

- Matrix form

$$\checkmark H^{(t)} = GRU(((A + I)W, H^{(t-1)})$$

GATED GRAPH SEQUENCE NEURAL NETWORKS

Yujia Li* & Richard Zemel

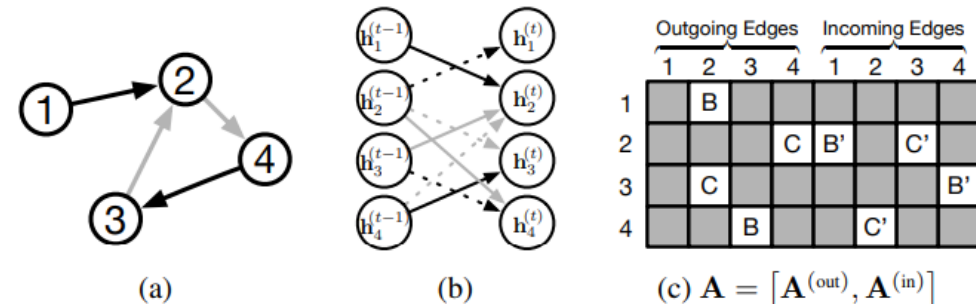
Department of Computer Science, University of Toronto
Toronto, Canada
{yujiali, zemel}@cs.toronto.edu

Marc Brockschmidt & Daniel Tarlow

Microsoft Research
Cambridge, UK
{mabrocks, dtarlow}@microsoft.com

ABSTRACT

Graph-structured data appears frequently in domains including chemistry, natural language semantics, social networks, and knowledge bases. In this work, we study feature learning techniques for graph-structured inputs. Our starting point is previous work on Graph Neural Networks (Scarselli et al., 2009), which we modify to use gated recurrent units and modern optimization techniques and then extend to output sequences. The result is a flexible and broadly useful class of neural network models that has found applications in a variety of domains including



Li, Yujia, et al. "Gated graph sequence neural networks." *arXiv preprint arXiv:1511.05493* (2015).

GraphSAGE (2017)

- NIPS, Stanford Univ.
- Consider node importance or ordering
- Aggregate

- Mean aggregate

$$\checkmark a_v^{(k-1)} = \sum_{u \in N(v)} \frac{h_u^{(k-1)}}{|N(u)|}$$

- LSTM aggregate (Random permutation of neighbors)

$$\checkmark a_v^{(k-1)} = LSTM(\sum\{W_{agg}h_u^{(k-1)}, \forall u \in N(v)\})$$

- Pooling aggregate

$$\checkmark a_v^{(k-1)} = Pool(\{W_{pool}h_u^{(k-1)}, \forall u \in N(v)\})$$

$$\checkmark Pool = \textit{element-wise mean or max}$$

- Combine

$$\checkmark h_v^{(k)} = Relu(W[h_v^{(k-1)}, a_v^{(k-1)}])$$

- Residual connection / Skip connection

Inductive Representation Learning on Large Graphs

William L. Hamilton*
wleif@stanford.edu

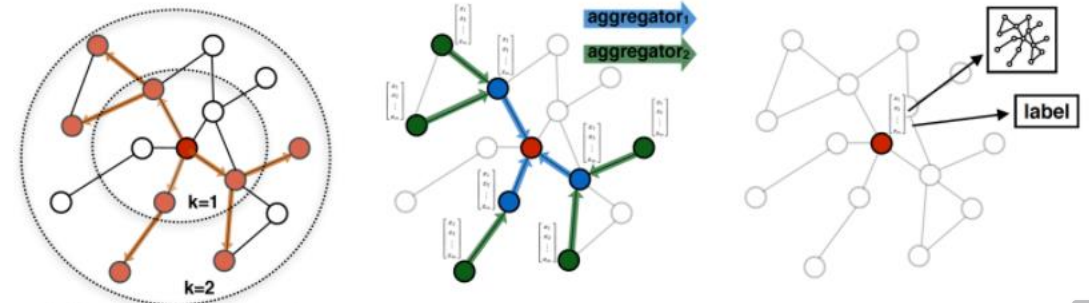
Rex Ying*
rexying@stanford.edu

Jure Leskovec
jure@cs.stanford.edu

Department of Computer Science
Stanford University
Stanford, CA, 94305

Abstract

Low-dimensional embeddings of nodes in large graphs have proved extremely useful in a variety of prediction tasks, from content recommendation to identifying protein functions. However, most existing approaches require that all nodes in the graph are present during training of the embeddings; these previous approaches are inherently *transductive* and do not naturally generalize to unseen nodes. Here we



Hamilton, Will, Zhitao Ying, and Jure Leskovec. "Inductive representation learning on large graphs." *Advances in neural information processing systems*. 2017.

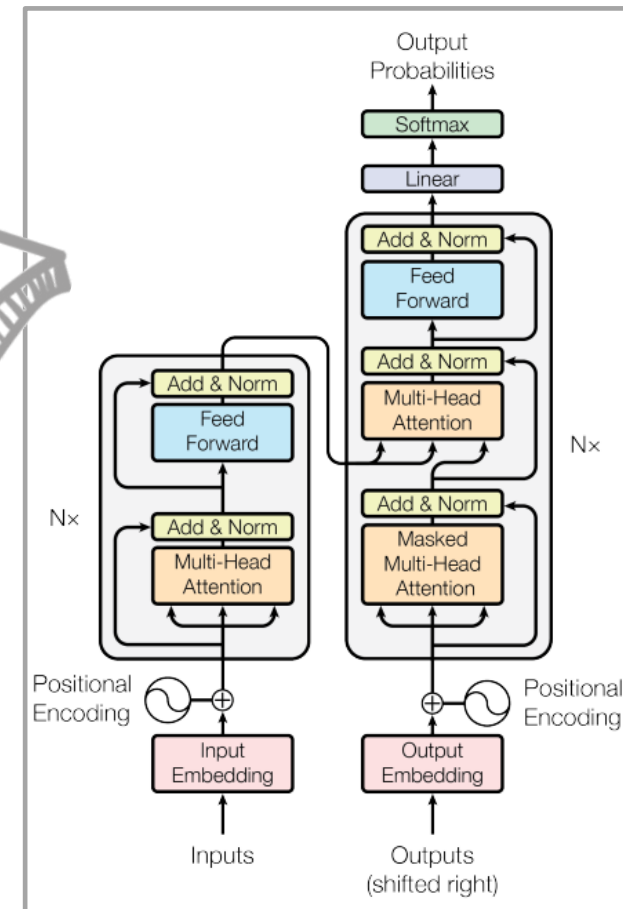
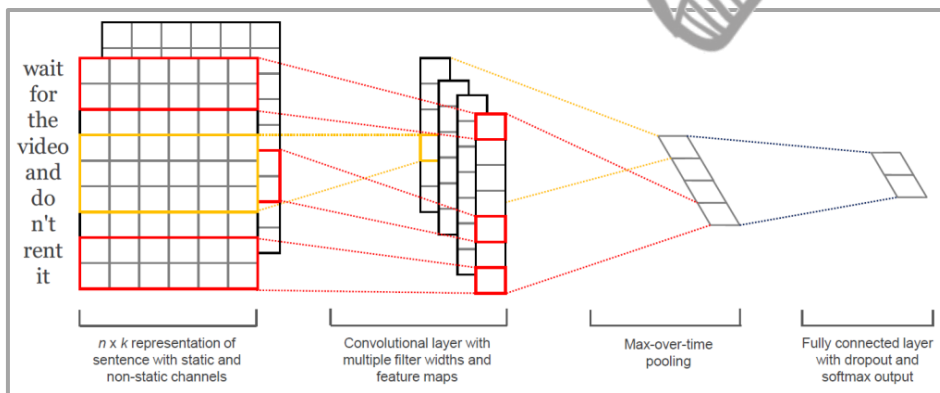
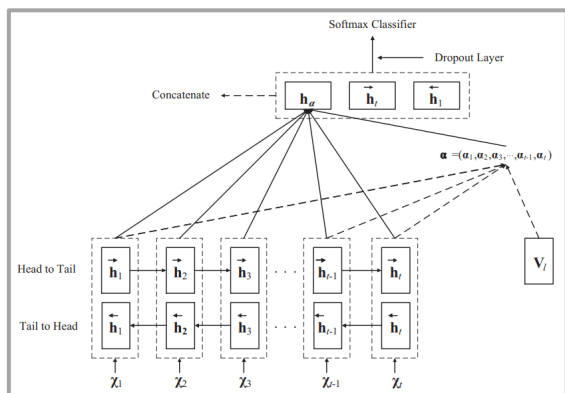
03 | Graph Neural Networks

■ GNN Variants

• Challenges

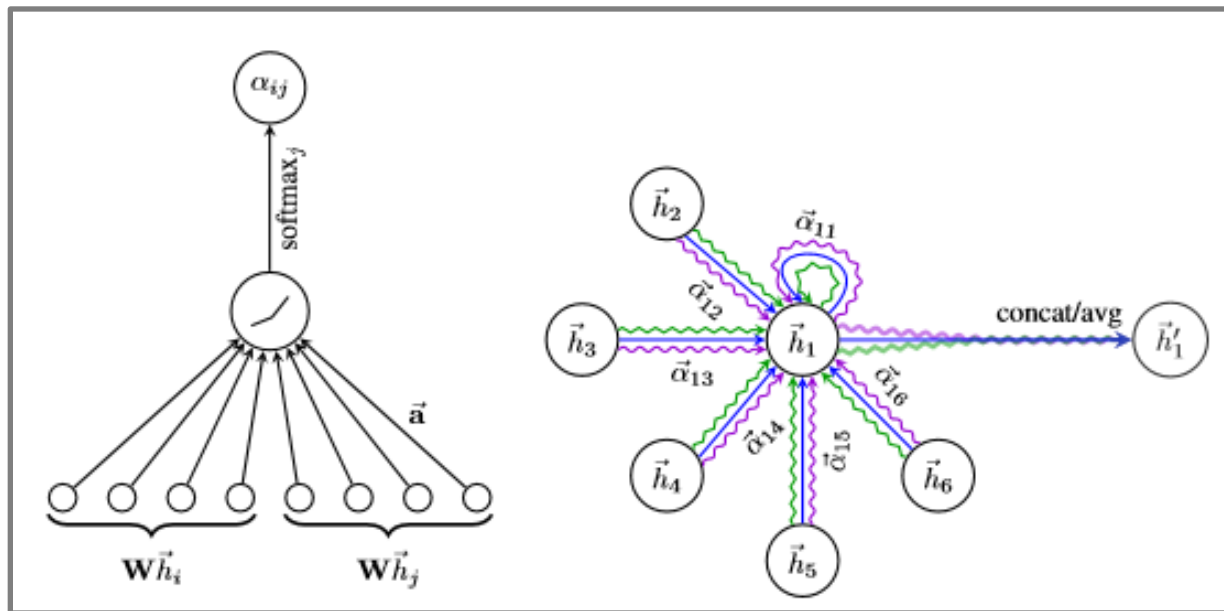
- ✓ Self-loop (Vanilla GNN)
- ✓ Node degrees (GCN)
- ✓ Node importance (GraphSAGE)
- ✓ Overfitting / Vanishing gradient (GGNN)

Attention is All You Need!



04 | Graph Attention Networks

- Graph Attention Networks: Paper
 - ICLR 2018
 - Cambridge University
 - Veličković, Petar, et al.
 - Recite: 1951



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

ABSTRACT

We present graph attention networks (GATs), novel neural network architectures that operate on graph-structured data, leveraging masked self-attentional layers to address the shortcomings of prior methods based on graph convolutions or their approximations. By stacking layers in which nodes are able to attend over their neighborhoods' features, we enable (implicitly) specifying different weights to different nodes in a neighborhood, without requiring any kind of costly matrix operation (such as inversion) or depending on knowing the graph structure upfront. In this way, we address several key challenges of spectral-based graph neural networks simultaneously, and make our model readily applicable to inductive as well as transductive problems. Our GAT models have achieved or matched state-of-the-art results across four established transductive and inductive graph benchmarks: the *Cora*, *Citeseer* and *Pubmed* citation network datasets, as well as a *protein-protein interaction* dataset (wherein test graphs remain unseen during training).

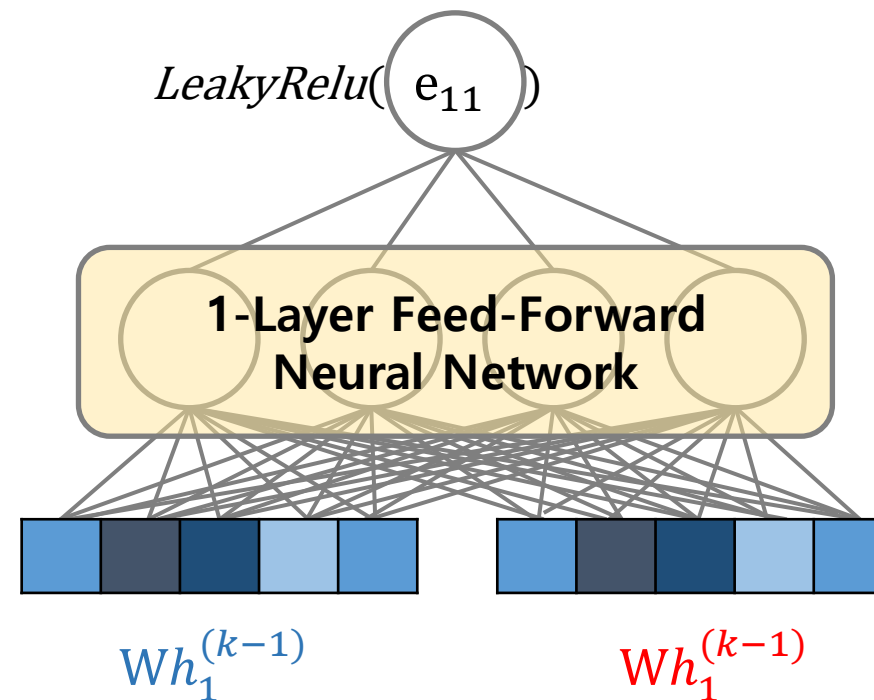
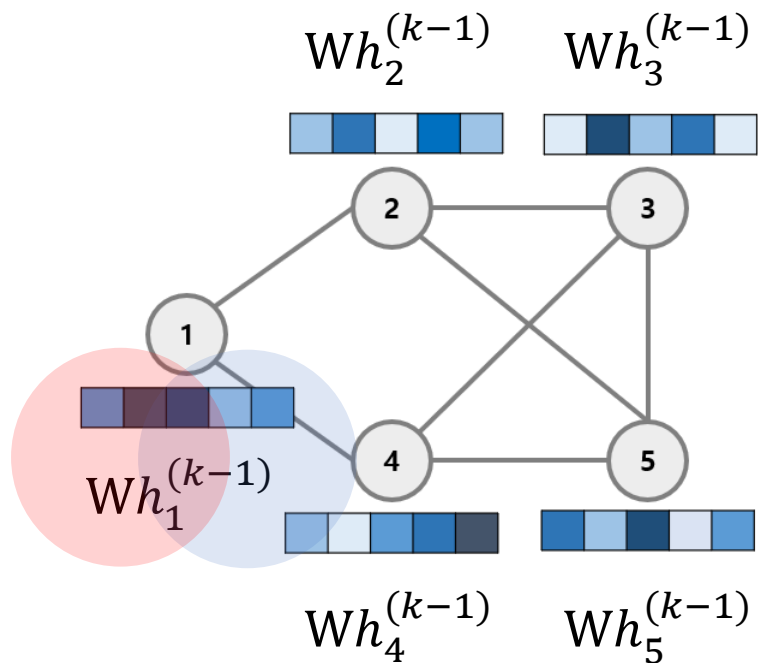
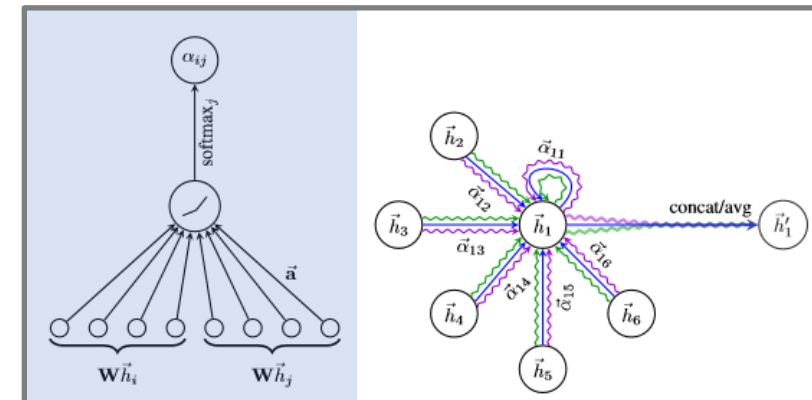
04 | Graph Attention Networks

Model Architecture

- Aggregate

- ✓ $a_v^{(k-1)} = \text{Attention}(\{h_u^{(k-1)}, u \in N(v) \cup \{v\}\})$
- ✓ **Key = Query = Value** = $h_u^{(k-1)}$
- ✓ **Similarity Function** (f) = 1 - layer Feed - Forward NN

$$A(q, K, V) = \sum_i \text{softmax}(f(K, q)) V$$



Veličković, Petar, et al. "Graph attention networks." *arXiv preprint arXiv:1710.10903* (2017).

04 | Graph Attention Networks

- Model Architecture

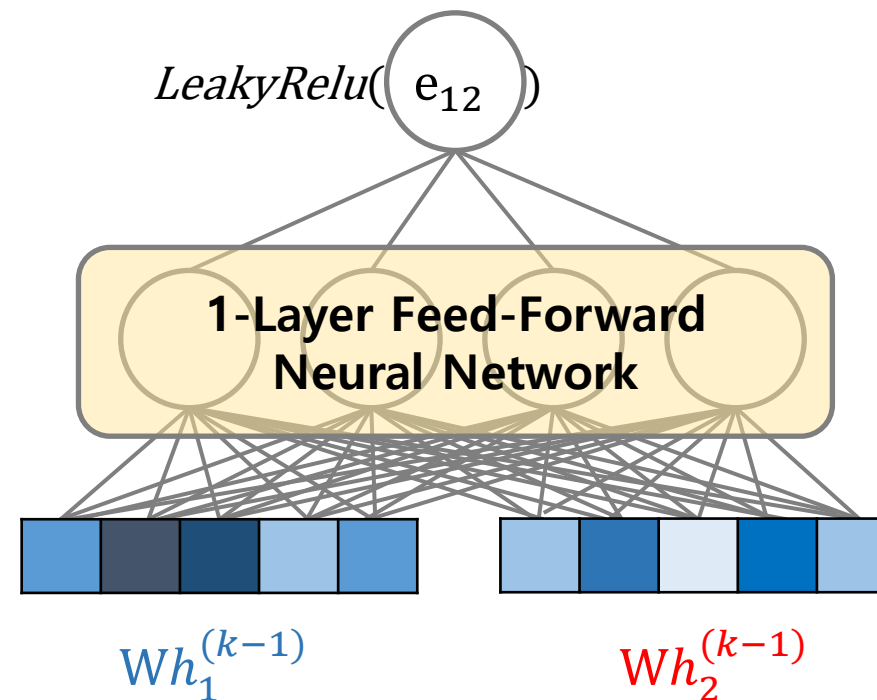
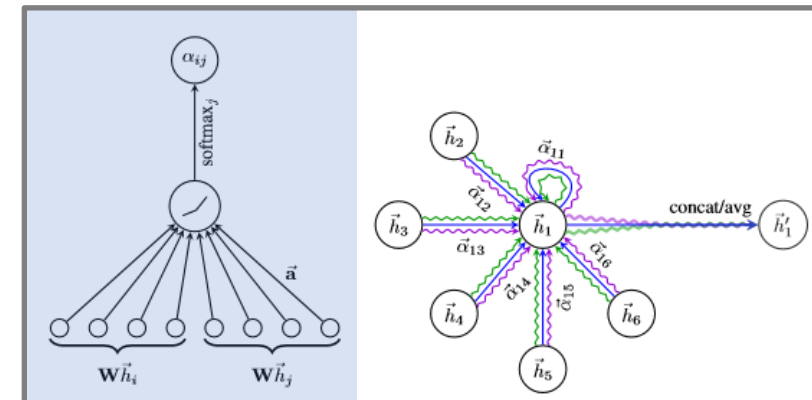
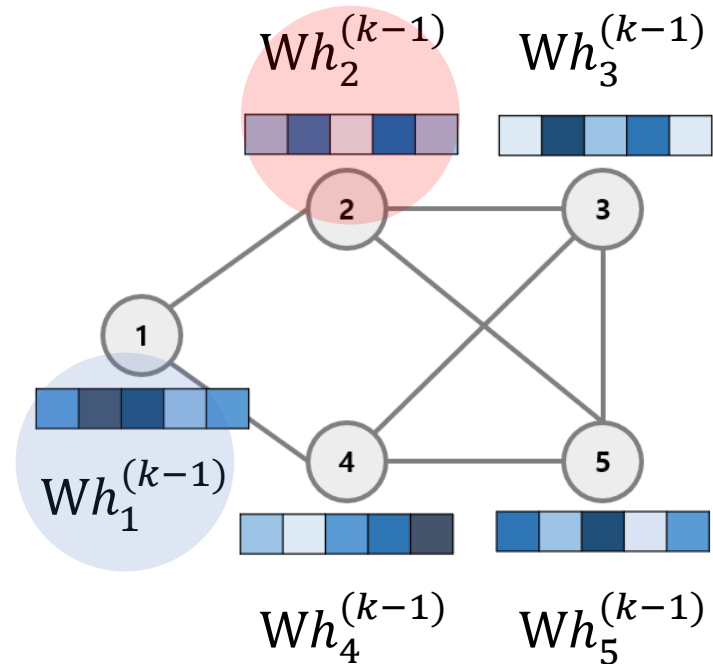
- Aggregate

- ✓ $a_v^{(k-1)} = \text{Attention}(\{h_u^{(k-1)}, u \in N(v) \cup \{v\}\})$

- ✓ **Key = Query = Value** = $h_u^{(k-1)}$

- ✓ **Similarity Function**(f) = 1 - layer Feed - Forward NN

$$A(q, K, V) = \sum_i \text{softmax}(f(K, q)) V$$



Veličković, Petar, et al. "Graph attention networks." *arXiv preprint arXiv:1710.10903* (2017).

04 | Graph Attention Networks

Model Architecture

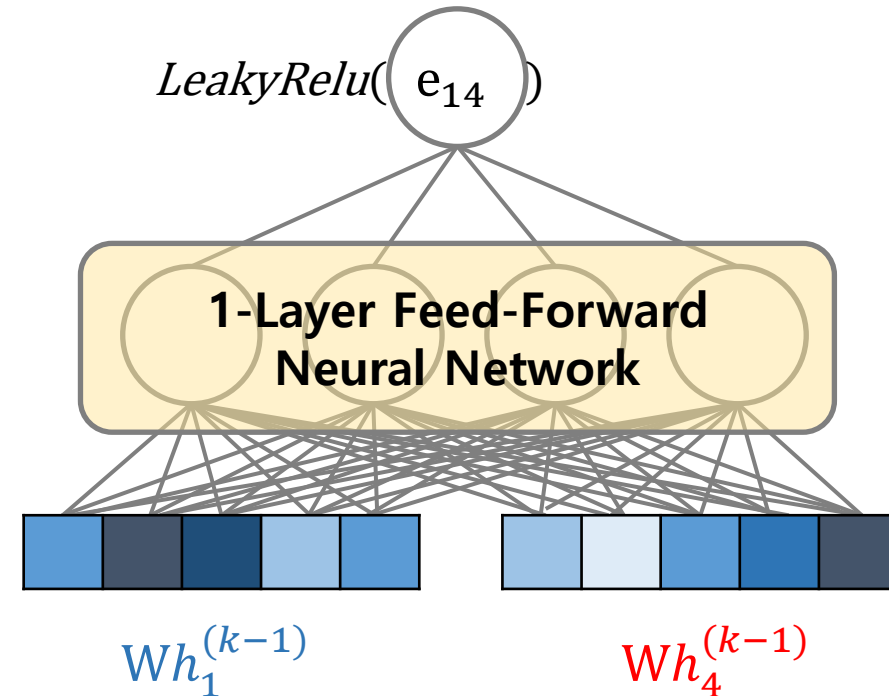
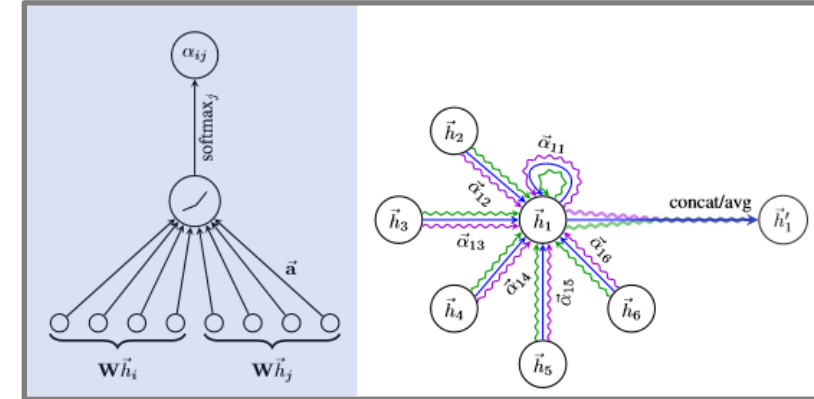
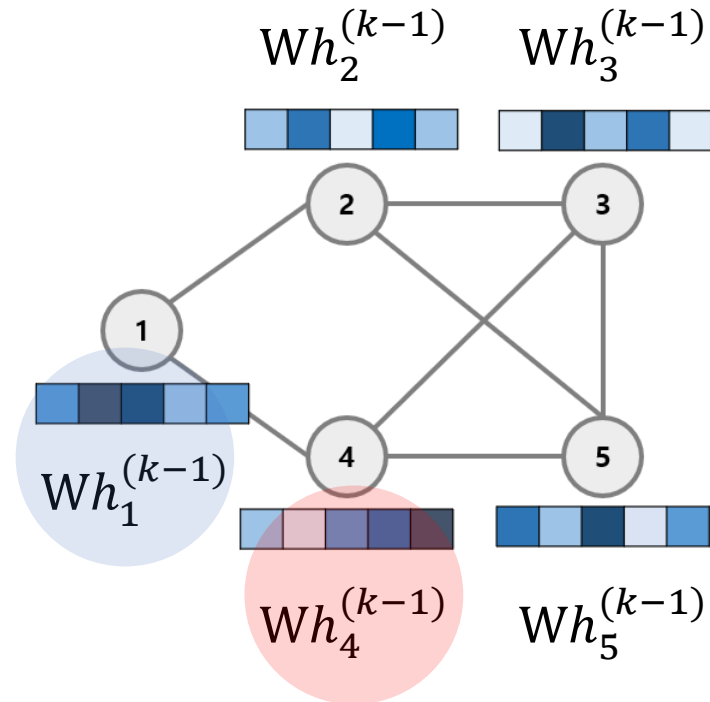
- Aggregate

- ✓ $a_v^{(k-1)} = \text{Attention}(\{h_u^{(k-1)}, u \in N(v) \cup \{v\}\})$

- ✓ **Key = Query = Value** = $h_u^{(k-1)}$

- ✓ **Similarity Function** (f) = 1 - layer Feed - Forward NN

$$A(q, K, V) = \sum_i \text{softmax}(f(K, q))V$$



Veličković, Petar, et al. "Graph attention networks." *arXiv preprint arXiv:1710.10903* (2017).

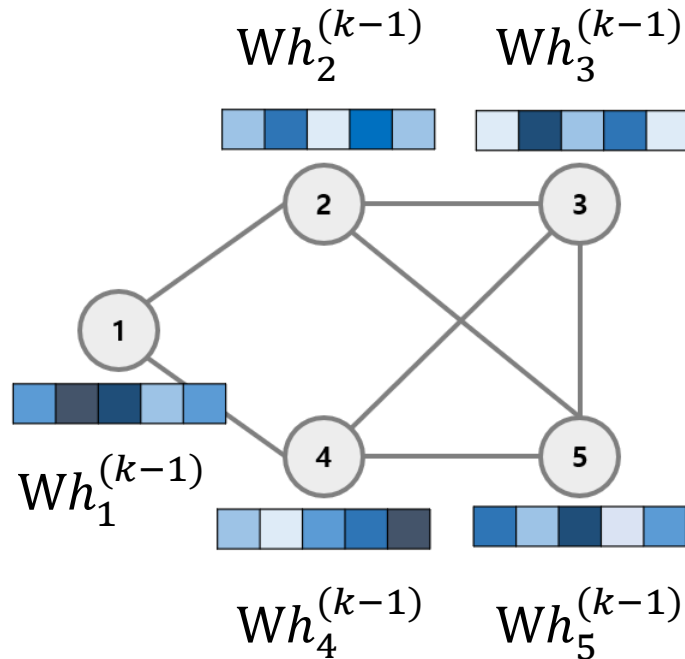
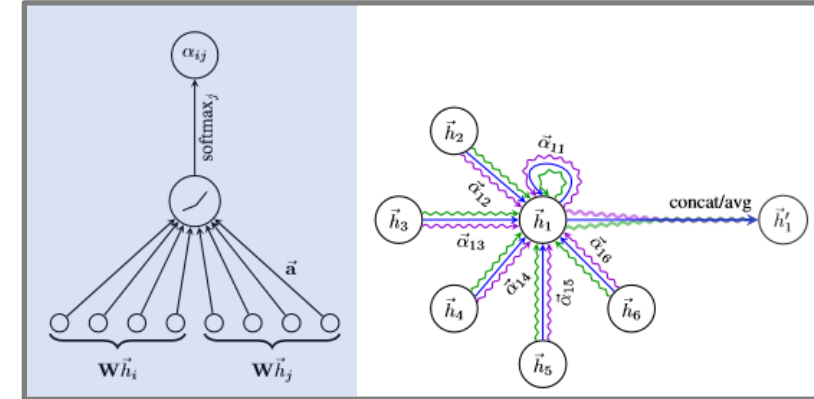
04 | Graph Attention Networks

- Model Architecture

- Aggregate

✓ $a_v^{(k-1)} = \text{Attention}(\{h_u^{(k-1)}, u \in N(v) \cup \{v\}\})$

$$A(q, K, V) = \sum_i \text{softmax}(f(K, q))V$$



e_{11}	e_{12}		e_{14}	
e_{21}	e_{22}	e_{23}		e_{25}
	e_{32}	e_{33}	e_{34}	e_{35}
e_{41}		e_{43}	e_{44}	e_{45}
	e_{25}	e_{53}	e_{54}	e_{55}

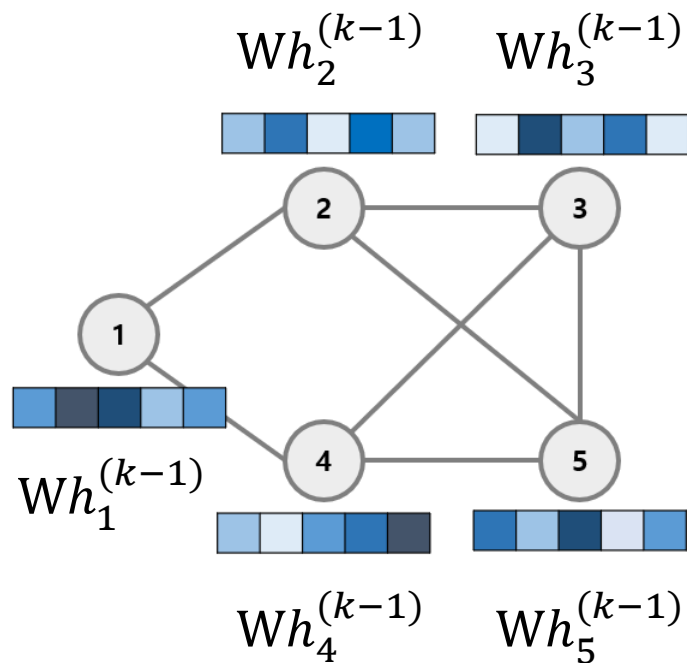
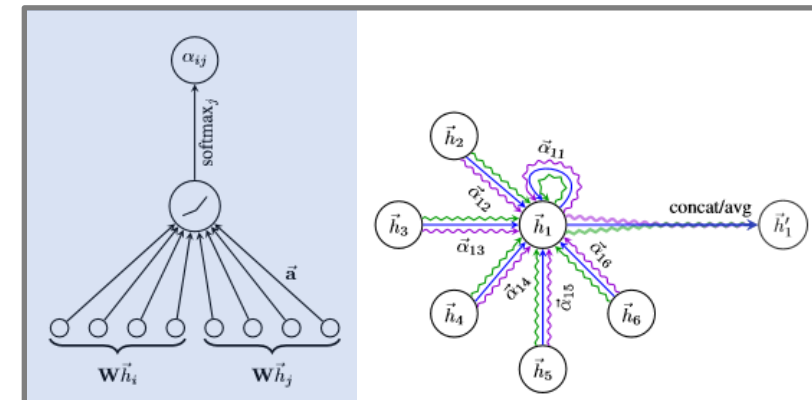
04 | Graph Attention Networks

- Model Architecture

- Aggregate

✓ $a_v^{(k-1)} = \text{Attention}(\{h_u^{(k-1)}, u \in N(v) \cup \{v\}\})$

$$A(q, K, V) = \sum_i \text{softmax}(f(K, q))V$$



e ₁₁	e ₁₂	Softmax	e ₁₄		
e ₂₁	e ₂₂	Softmax	e ₂₃	e ₂₅	
	e ₃₂	Softmax	e ₃₃	e ₃₄	e ₃₅
e ₄₁		Softmax	e ₄₃	e ₄₄	e ₄₅
	e ₅₂	Softmax	e ₅₃	e ₅₄	e ₅₅

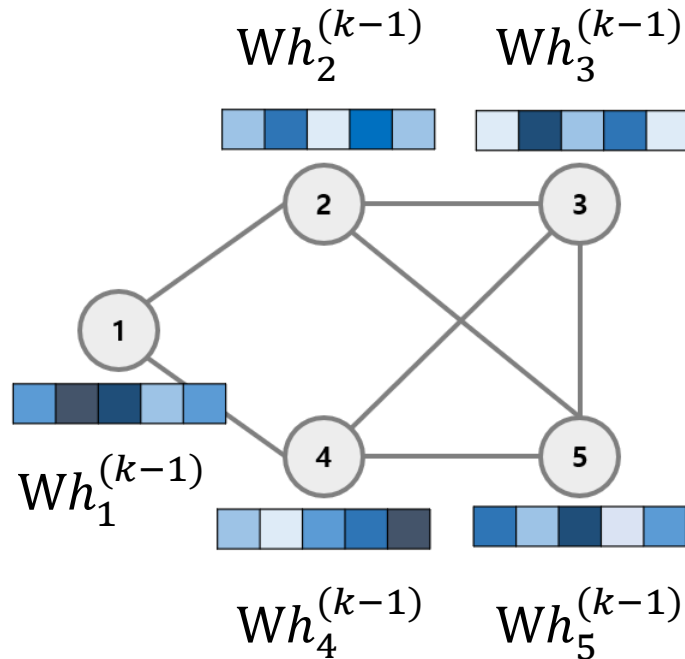
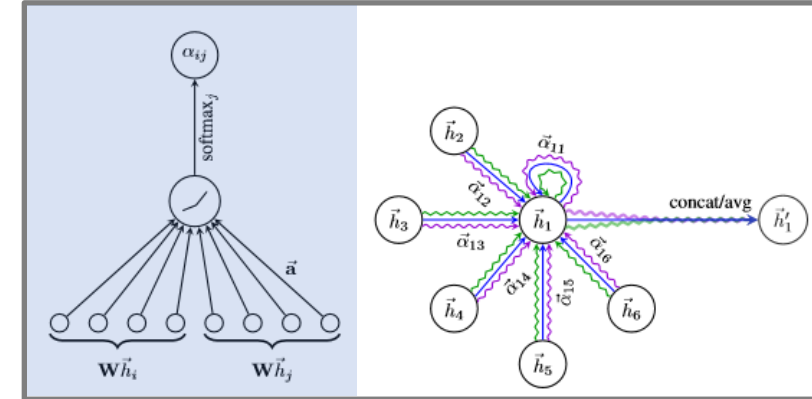
04 | Graph Attention Networks

- Model Architecture

- Aggregate

$$a_v^{(k-1)} = \text{Attention}(\{h_u^{(k-1)}, u \in N(v) \cup \{v\}\})$$

$$A(q, K, V) = \sum_i \text{softmax}(f(K, q))V$$



a_{11}	a_{12}		a_{14}	
a_{21}	a_{22}	a_{23}		a_{25}
	a_{32}	a_{33}	a_{34}	a_{35}
a_{41}		a_{43}	a_{44}	a_{45}
	a_{25}	a_{53}	a_{54}	a_{55}

04 | Graph Attention Networks

Model Architecture

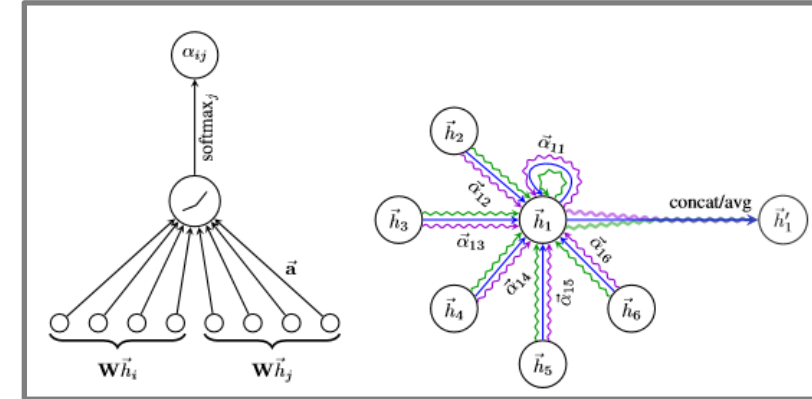
- Aggregate

$$a_v^{(k-1)} = \text{Attention}(\{h_u^{(k-1)}, u \in N(v) \cup \{v\}\})$$

- Masked Attention

Using adjacency matrix, masked value with negative values before softmax.

$$A(q, K, V) = \sum_i \text{softmax}(f(K, q))V$$



e_{11}	e_{12}	e_{13}	e_{14}	e_{15}
e_{21}	e_{22}	e_{23}	e_{24}	e_{25}
e_{31}	e_{32}	e_{33}	e_{34}	e_{35}
e_{41}	e_{42}	e_{43}	e_{44}	e_{45}
e_{51}	e_{52}	e_{53}	e_{54}	e_{55}

Similarity

e_{11}	e_{12}	$-9e^{10}$	e_{14}	$-9e^{10}$
e_{21}	e_{22}	e_{23}	$-9e^{10}$	e_{25}
$-9e^{10}$	e_{32}	e_{33}	e_{34}	e_{35}
e_{41}	$-9e^{10}$	e_{43}	e_{44}	e_{45}
$-9e^{10}$	e_{52}	e_{53}	e_{54}	e_{55}

Masking

a_{11}	a_{12}	0	a_{14}	0
a_{21}	a_{22}	a_{23}	0	a_{25}
0	a_{32}	a_{33}	a_{34}	a_{35}
a_{41}	0	a_{43}	a_{44}	a_{45}
0	a_{52}	a_{53}	a_{54}	a_{55}

Softmax



04 | Graph Attention Networks

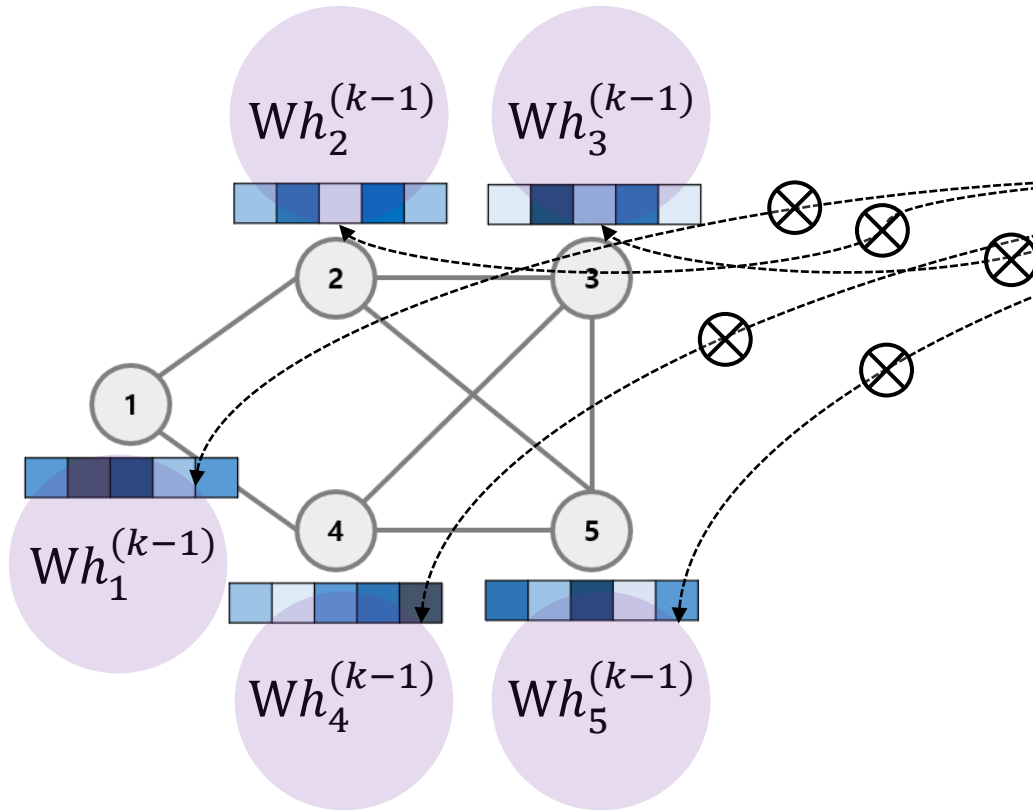
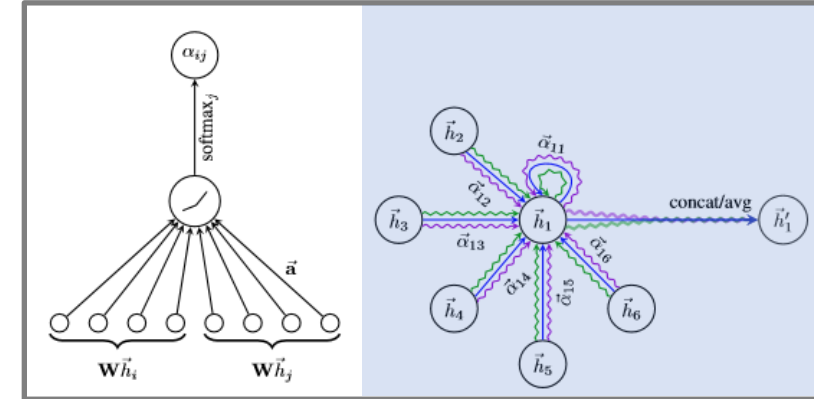
Model Architecture

- Combine

- $$h_v^{(k)} = \sum_{u \in N(v) \cup \{v\}} a_u^{(k-1)} h_u^{(k-1)}$$

- Weight-sum of aggregated information & k-1 hidden state (Value)

$$A(q, K, V) = \sum_i \text{softmax}(f(K, q))V$$



a_{11}	a_{12}	0	a_{14}	0
a_{21}	a_{22}	a_{23}	0	a_{25}
0	a_{32}	a_{33}	a_{34}	a_{35}
a_{41}	0	a_{43}	a_{44}	a_{45}
0	a_{25}	a_{53}	a_{54}	a_{55}

04 | Graph Attention Networks

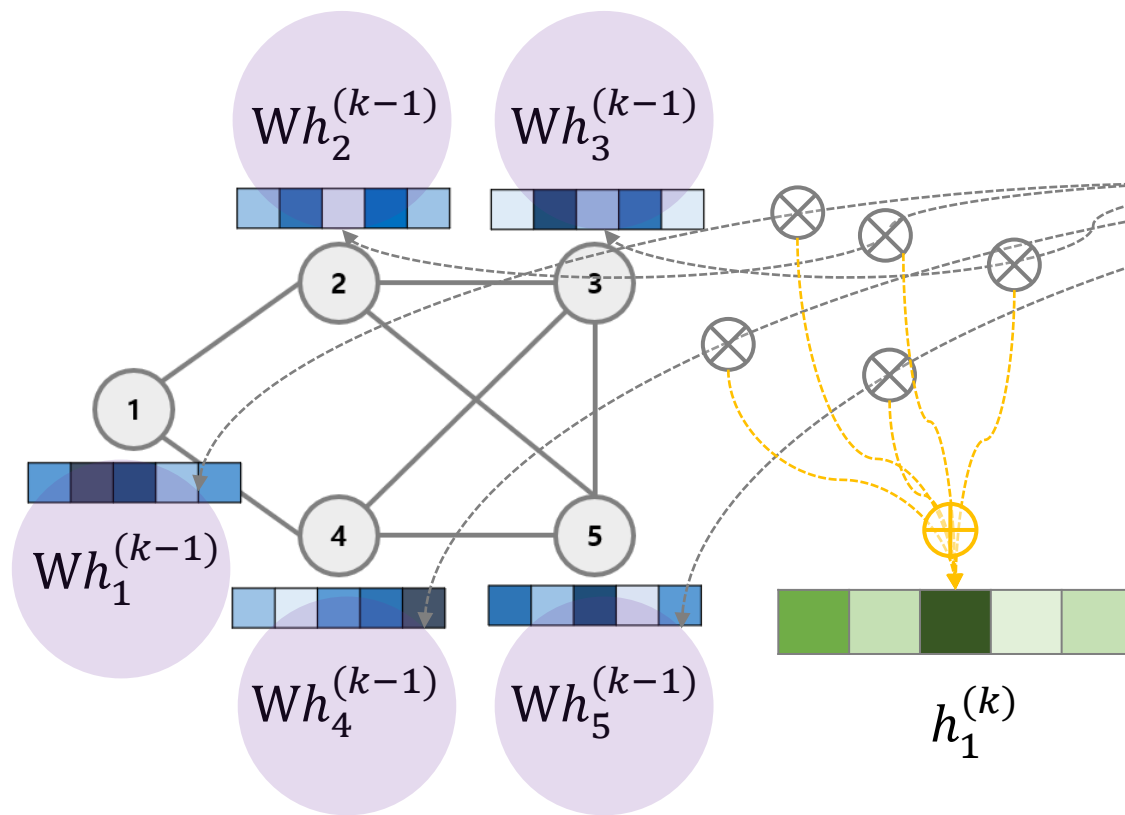
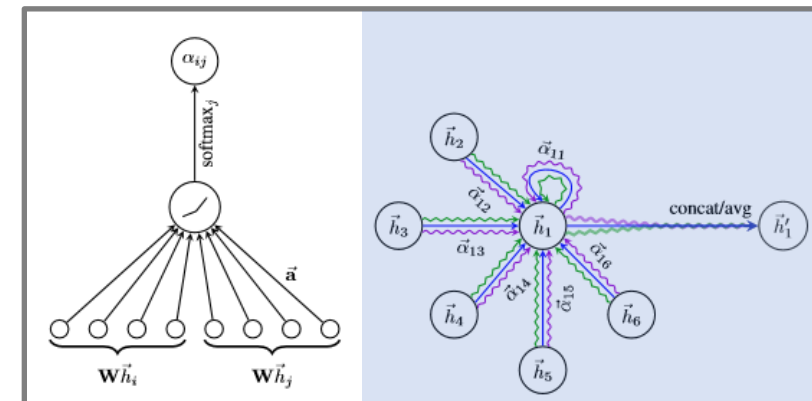
Model Architecture

- Combine

- $$h_v^{(k)} = \sum_{u \in N(v) \cup \{v\}} a_u^{(k-1)} h_u^{(k-1)}$$

- Weight-sum of aggregated information & k-1 hidden state (Value)

$$A(q, K, V) = \sum_i \text{softmax}(f(K, q))V$$



a_{11}	a_{12}	0	a_{14}	0
a_{21}	a_{22}	a_{23}	0	a_{25}
0	a_{32}	a_{33}	a_{34}	a_{35}
a_{41}	0	a_{43}	a_{44}	a_{45}
0	a_{25}	a_{53}	a_{54}	a_{55}

04 | Graph Attention Networks

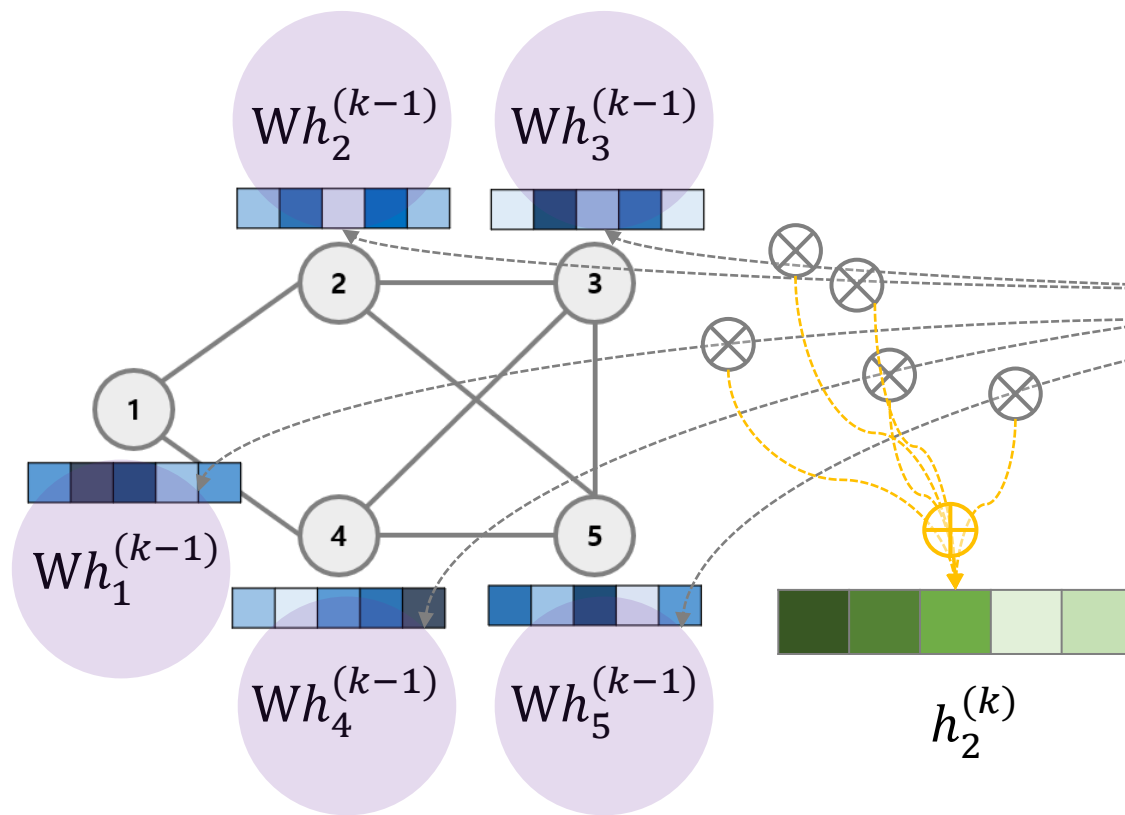
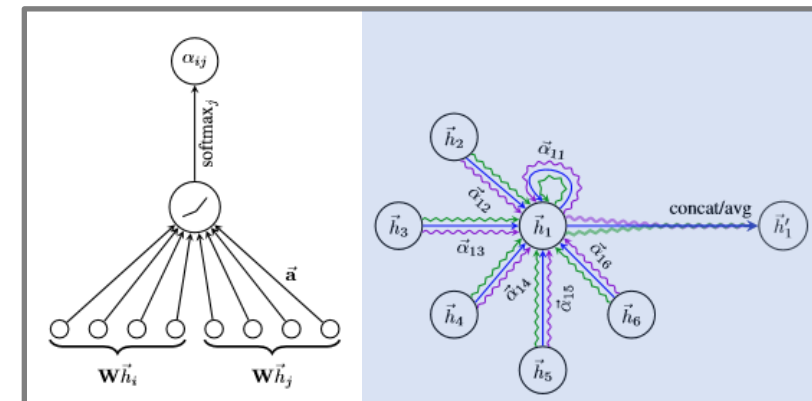
Model Architecture

- Combine

✓ $h_v^{(k)} = \sum_{u \in N(v) \cup \{v\}} a_u^{(k-1)} h_u^{(k-1)}$

- ✓ Weight-sum of aggregated information & k-1 hidden state (Value)

$$A(q, K, V) = \sum_i \text{softmax}(f(K, q))V$$



a_{11}	a_{12}	0	a_{14}	0
a_{21}	a_{22}	a_{23}	0	a_{25}
0	a_{32}	a_{33}	a_{34}	a_{35}
a_{41}	0	a_{43}	a_{44}	a_{45}
0	a_{25}	a_{53}	a_{54}	a_{55}

04 | Graph Attention Networks

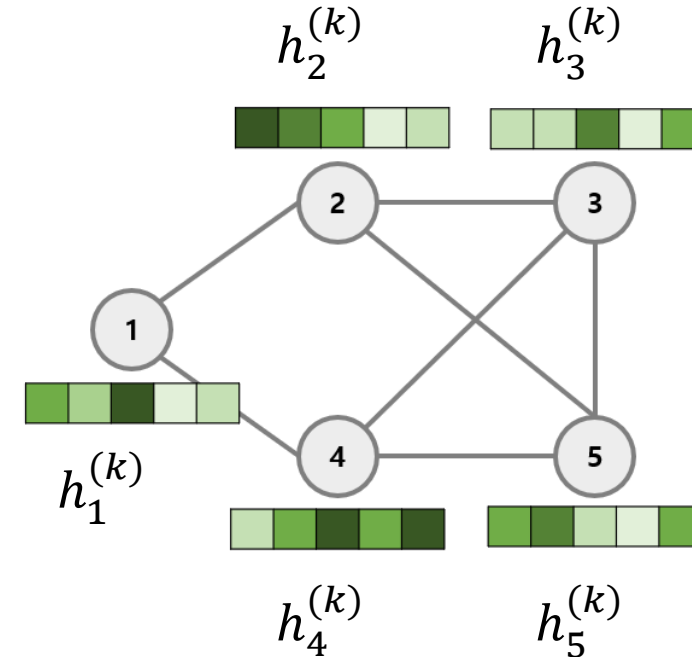
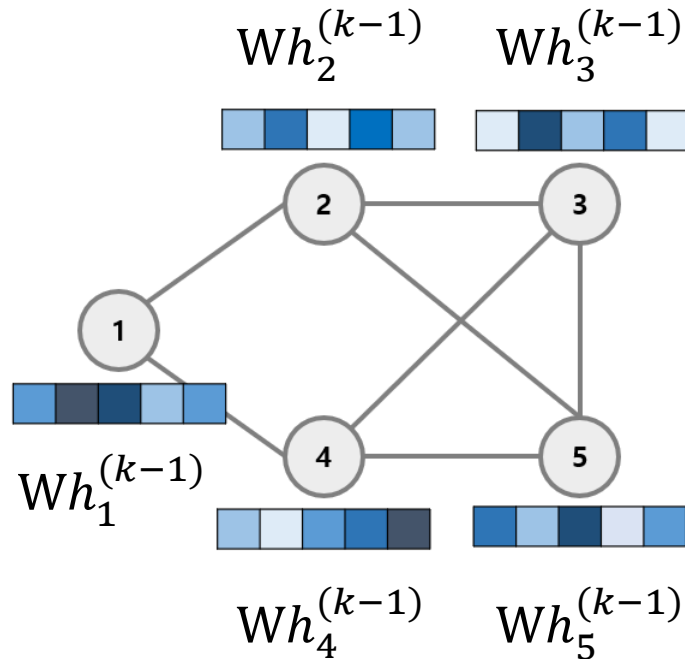
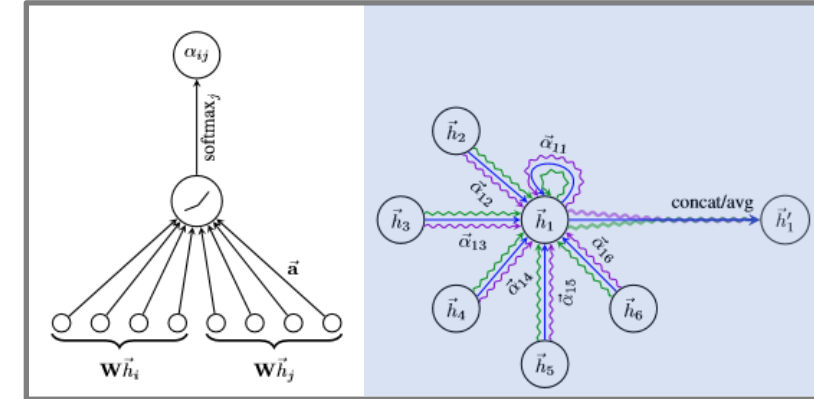
- Model Architecture

- Combine

- ✓ $h_v^{(k)} = \sum_{u \in N(v) \cup \{v\}} a_u^{(k-1)} h_u^{(k-1)}$

- ✓ Weight-sum of aggregated information & k-1 hidden state (Value)

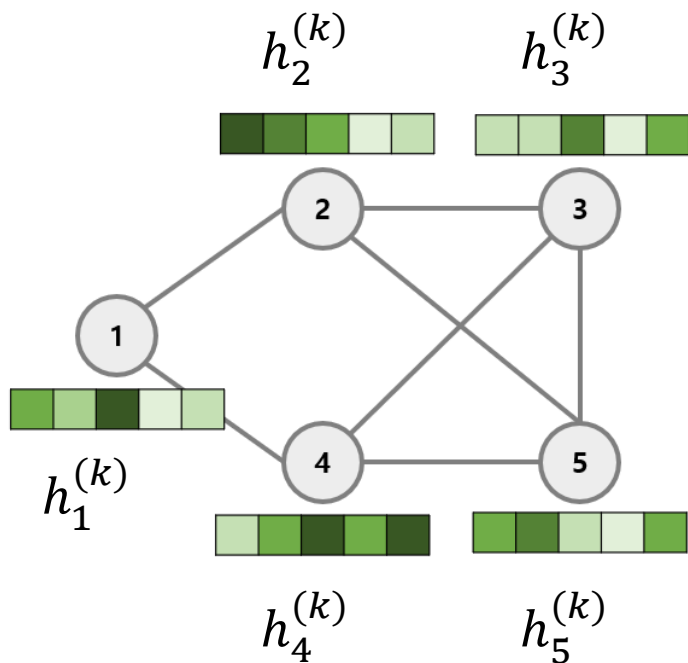
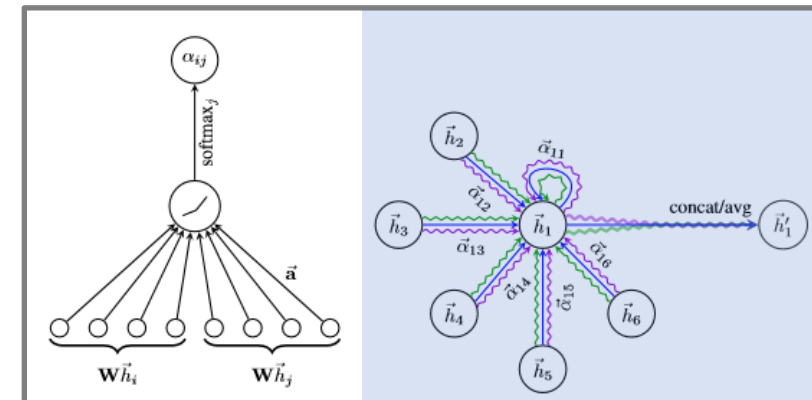
$$A(q, K, V) = \sum_i \text{softmax}(f(K, q))V$$



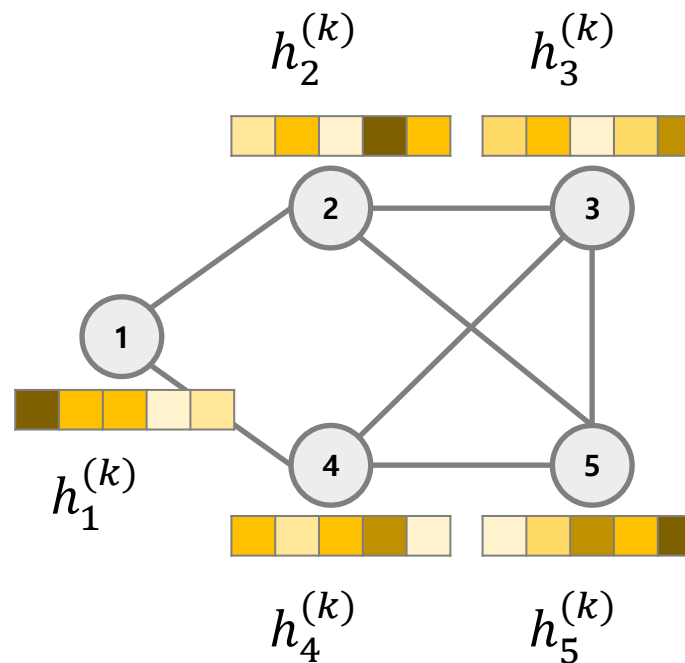
04 | Graph Attention Networks

- Model Architecture
 - Multi-head Attention
 - Concatenate or average

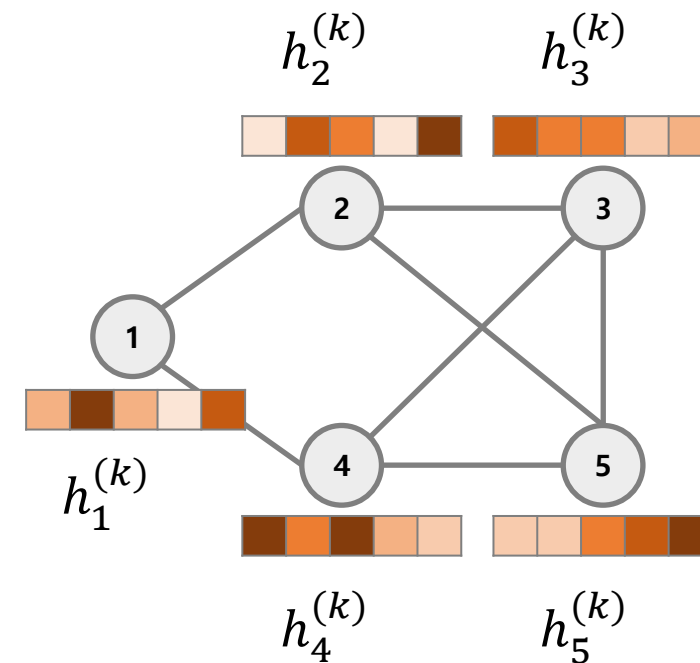
$$A(q, K, V) = \sum_i \text{softmax}(f(K, q))V$$



Head_1



Head_2



Head_3



04 | Graph Attention Networks

- Result

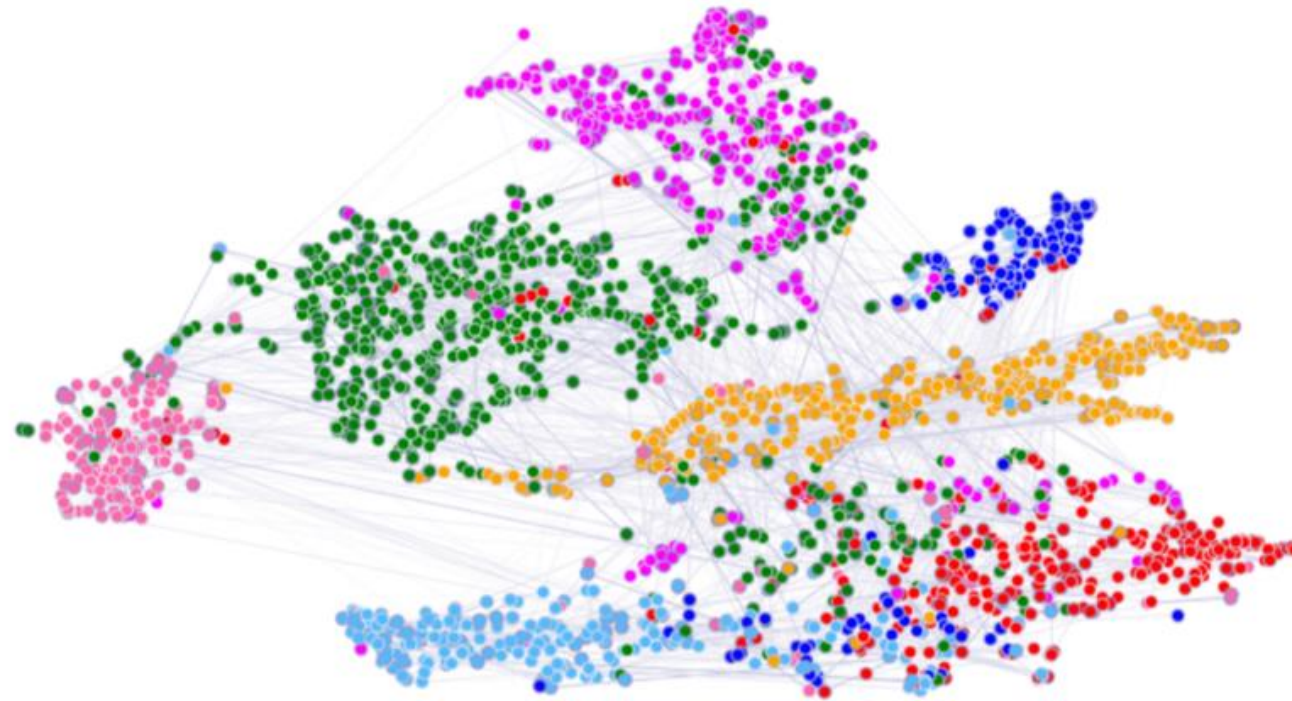
<i>Inductive</i>	
Method	PPI
Random	0.396
MLP	0.422
GraphSAGE-GCN (Hamilton et al., 2017)	0.500
GraphSAGE-mean (Hamilton et al., 2017)	0.598
GraphSAGE-LSTM (Hamilton et al., 2017)	0.612
GraphSAGE-pool (Hamilton et al., 2017)	0.600
GraphSAGE*	0.768
Const-GAT (ours)	0.934 ± 0.006
GAT (ours)	0.973 ± 0.002

<i>Transductive</i>			
Method	Cora	Citeseer	Pubmed
MLP	55.1%	46.5%	71.4%
ManiReg (Belkin et al., 2006)	59.5%	60.1%	70.7%
SemiEmb (Weston et al., 2012)	59.0%	59.6%	71.7%
LP (Zhu et al., 2003)	68.0%	45.3%	63.0%
DeepWalk (Perozzi et al., 2014)	67.2%	43.2%	65.3%
ICA (Lu & Getoor, 2003)	75.1%	69.1%	73.9%
Planetoid (Yang et al., 2016)	75.7%	64.7%	77.2%
Chebyshev (Defferrard et al., 2016)	81.2%	69.8%	74.4%
GCN (Kipf & Welling, 2017)	81.5%	70.3%	79.0%
MoNet (Monti et al., 2016)	81.7 ± 0.5%	—	78.8 ± 0.3%
GCN-64*	81.4 ± 0.5%	70.9 ± 0.5%	79.0 ± 0.3%
GAT (ours)	83.0 ± 0.7%	72.5 ± 0.7%	79.0 ± 0.3%

04 | Graph Attention Networks

■ Attention Score Visualization

- Color = Class of node
- Edge thickness = average of multi-head attention score



■ Conclusion

- Attention
 - ✓ Attention: Query와 key의 유사도를 계산한 후 value의 가중합을 계산하는 과정
- Graph Neural Network
 - ✓ Graph 구조를 학습하는 딥러닝 모델
 - ✓ Text, Image 데이터를 그래프로 표현해서 학습하는 모델들도 연구가 많이 되고 있음
- Graph Attention Network
 - ✓ Attention 개념을 GNN에 적용하여 explainability + model performance

- Pytorch-Geometric
 - Message Passing(Aggregate + Combine)

```
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 γ_{Θ} and ϕ_{Θ} denote differentiable functions such as MLPs. See [here](#) for the accompanying tutorial.

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")
- node_dim** (*int, optional*) – The axis along which to propagate. (default: -2)

```
CLASS GCNConv ( in_channels: int, out_channels: int, improved: bool = False, cached: bool = False, add_self_loops: bool = True, normalize: bool = True, bias: bool = True, **kwargs ) [source]
```

The graph convolutional operator from the "Semi-supervised Classification with Graph Convolutional Networks" paper

$$\mathbf{X}' = \hat{\mathbf{D}}^{-1/2} \hat{\mathbf{A}} \hat{\mathbf{D}}^{-1/2} \mathbf{X} \Theta,$$

where $\hat{\mathbf{A}} = \mathbf{A} + \mathbf{I}$ denotes the adjacency matrix with inserted self-loops and $\hat{D}_{ii} = \sum_{j=0} \hat{A}_{ij}$ its diagonal degree matrix.

PARAMETERS

- in_channels** (*int*) – Size of each input sample.
- out_channels** (*int*) – Size of each output sample.

```
CLASS GATConv ( in_channels: Union[int, Tuple[int, int]], out_channels: int, heads: int = 1, concat: bool = True, negative_slope: float = 0.2, dropout: float = 0.0, add_self_loops: bool = True, bias: bool = True, **kwargs ) [source]
```

The graph attentional operator from the "Graph Attention Networks" paper

$$\mathbf{x}'_i = \alpha_{i,i} \Theta \mathbf{x}_i + \sum_{j \in \mathcal{N}(i)} \alpha_{i,j} \Theta \mathbf{x}_j,$$

where the attention coefficients $\alpha_{i,j}$ are computed as

$$\alpha_{i,j} = \frac{\exp\left(\text{LeakyReLU}\left(\mathbf{a}^T [\Theta \mathbf{x}_i \parallel \Theta \mathbf{x}_j]\right)\right)}{\sum_{k \in \mathcal{N}(i) \cup \{i\}} \exp\left(\text{LeakyReLU}\left(\mathbf{a}^T [\Theta \mathbf{x}_i \parallel \Theta \mathbf{x}_k]\right)\right)}$$

감사합니다

