

## cs224n 스터디 1팀

2024.05.07

발표자 : 유금현, 조다영

## 스터디원 소개 및 만남 인증



스터디원 1 : 송경준

스터디원 2 : 유금현

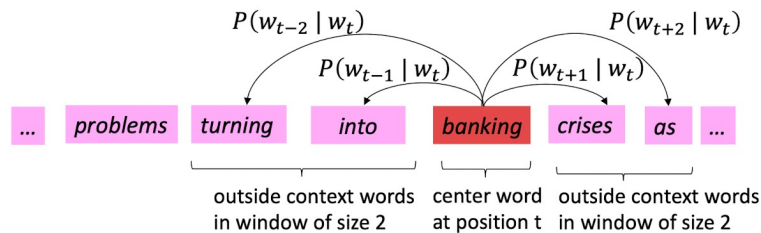
스터디원 3 : 조다영

# Word Embedding

: 단어를 벡터화하는 방식

## Word2vec

- 단어 간의 유사도를 반영
- 텍스트 전체의 정보를 반영하지 못한다는 단점

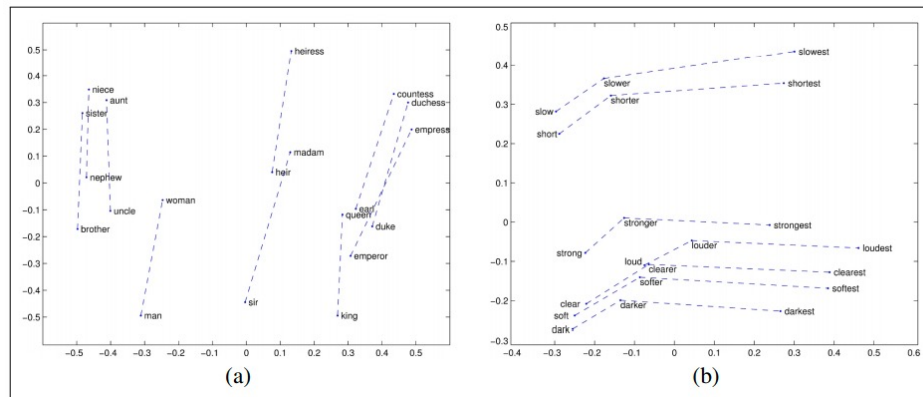


$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(w_{t+j} | w_t; \theta)$$

# Word Embedding

GloVe

$$w_i \cdot w_j = \log P(i|j)$$



## Word Embedding

Model	Dev	Test	ACE	MUC7
Discrete	91.0	85.4	77.4	73.4
SVD	90.8	85.7	77.3	73.7
SVD-S	91.0	85.5	77.6	74.3
SVD-L	90.5	84.8	73.6	71.5
HPCA	92.6	<b>88.7</b>	81.7	80.7
HSMN	90.5	85.7	78.7	74.7
CW	92.2	87.4	81.7	80.2
CBOW	93.1	88.2	82.2	81.1
GloVe	<b>93.2</b>	88.3	<b>82.9</b>	<b>82.2</b>

Model	Size	WS353	MC	RG	SCWS	RW
SVD	6B	35.3	35.1	42.5	38.3	25.6
SVD-S	6B	56.5	71.5	71.0	53.6	34.7
SVD-L	6B	65.7	<u>72.7</u>	75.1	56.5	37.0
CBOW <sup>†</sup>	6B	57.2	65.6	68.2	57.0	32.5
SG <sup>†</sup>	6B	62.8	65.2	69.7	<u>58.1</u>	37.2
GloVe	6B	<u>65.8</u>	<u>72.7</u>	<u>77.8</u>	53.9	<u>38.1</u>
SVD-L	42B	74.0	76.4	74.1	58.3	39.9
GloVe	42B	<b><u>75.9</u></b>	<b><u>83.6</u></b>	<b><u>82.9</u></b>	<b><u>59.6</u></b>	<b><u>47.8</u></b>
CBOW*	100B	68.4	79.6	75.4	59.4	45.5

Model	Dim.	Size	Sem.	Syn.	Tot.
SVD	300	6B	6.3	8.1	7.3
SVD-S	300	6B	36.7	46.6	42.1
SVD-L	300	6B	56.6	63.0	60.1
CBOW <sup>†</sup>	300	6B	63.6	<u>67.4</u>	65.7
SG <sup>†</sup>	300	6B	73.0	66.0	69.1
GloVe	300	6B	<u>77.4</u>	67.0	<u>71.7</u>

## Dependency parsing

San Jose cops kill man with knife

Text

Paper

Translate

Listen

Close

# San Jose cops kill man with knife

경찰이 칼을 든 남자를 죽였다

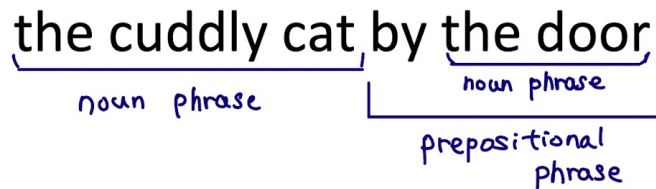
VS

경찰이 칼로 남자를 죽였다

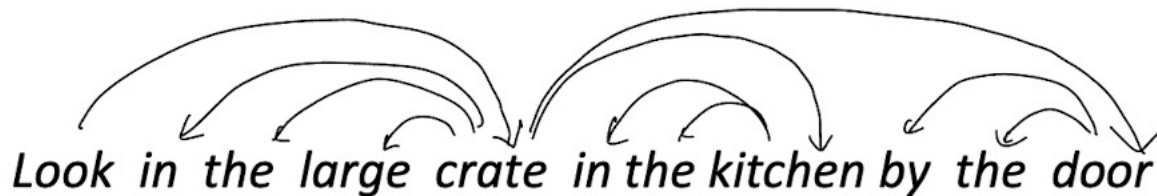
=> 모호성 문제

## Dependency parsing

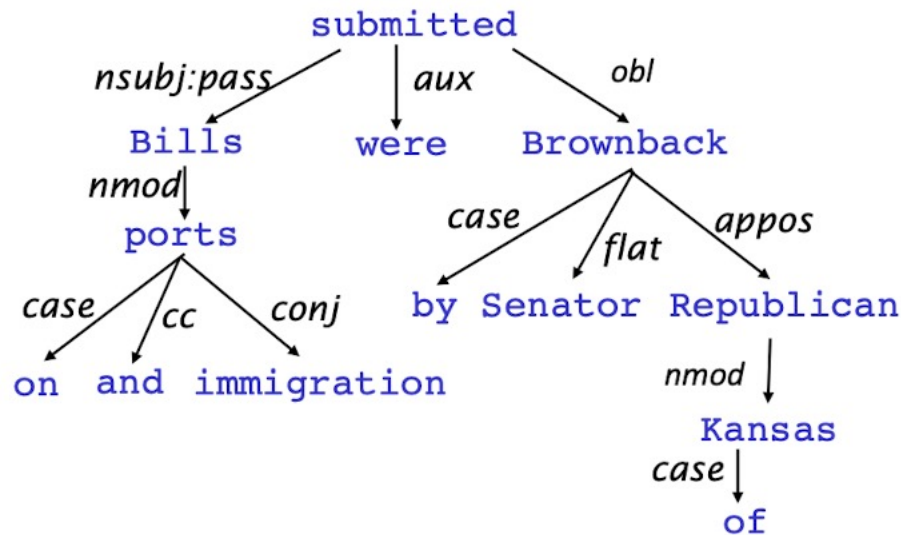
### 1) Phrase Structure



### 2) Dependency Structure



## Dependency parsing

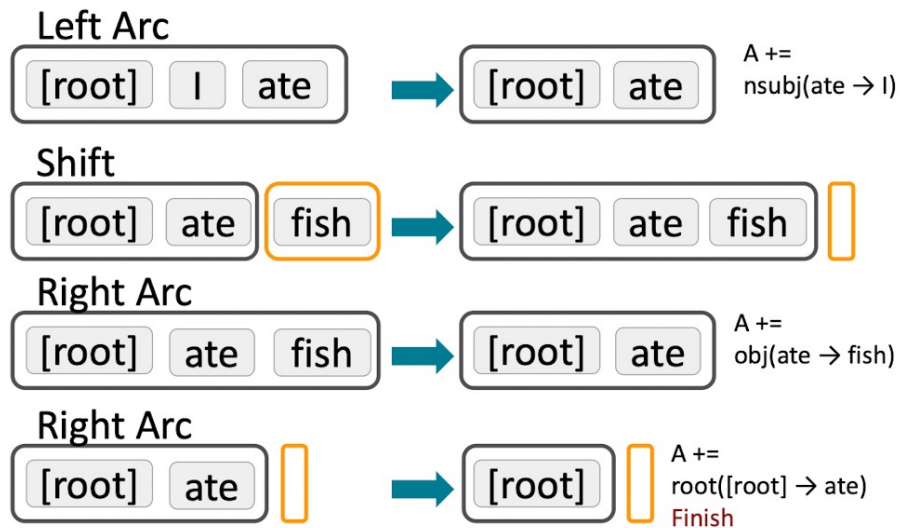


- Dynamic programming
- Graph algorithms
- Constraint satisfaction
- Transition-based parsing



## Dependency parsing

### Transition-based parsing



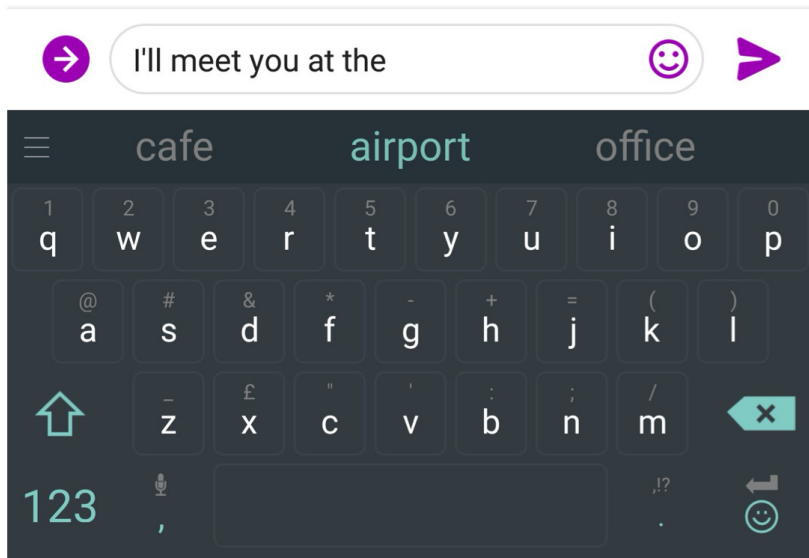
## Dependency parsing

Parser	UAS	LAS	sent. / s
MaltParser	89.8	87.2	469
MSTParser	91.4	88.1	10
TurboParser	<b>92.3</b>	89.6	8
C & M 2014	92.0	<b>89.7</b>	<b>654</b>

→ Neural Dependency Parser

## Language Model

Language Model : 문장 다음에 올 단어를 예측하는 일



## N-gram Language Model

~~as the proctor started the clock, the~~ *students opened their* \_\_\_\_\_  
 discard { condition on this

$$P(\mathbf{w} | \text{students opened their}) = \frac{\text{count}(\text{students opened their } \mathbf{w})}{\text{count}(\text{students opened their})}$$

## N-gram Language Model

Problems with n-gram Language Model

Sparsity Problems

$$P(\mathbf{w} | \text{students opened their}) = \frac{\text{count}(\text{students opened their } \mathbf{w})}{\text{count}(\text{students opened their})} = 0$$

## N-gram Language Model

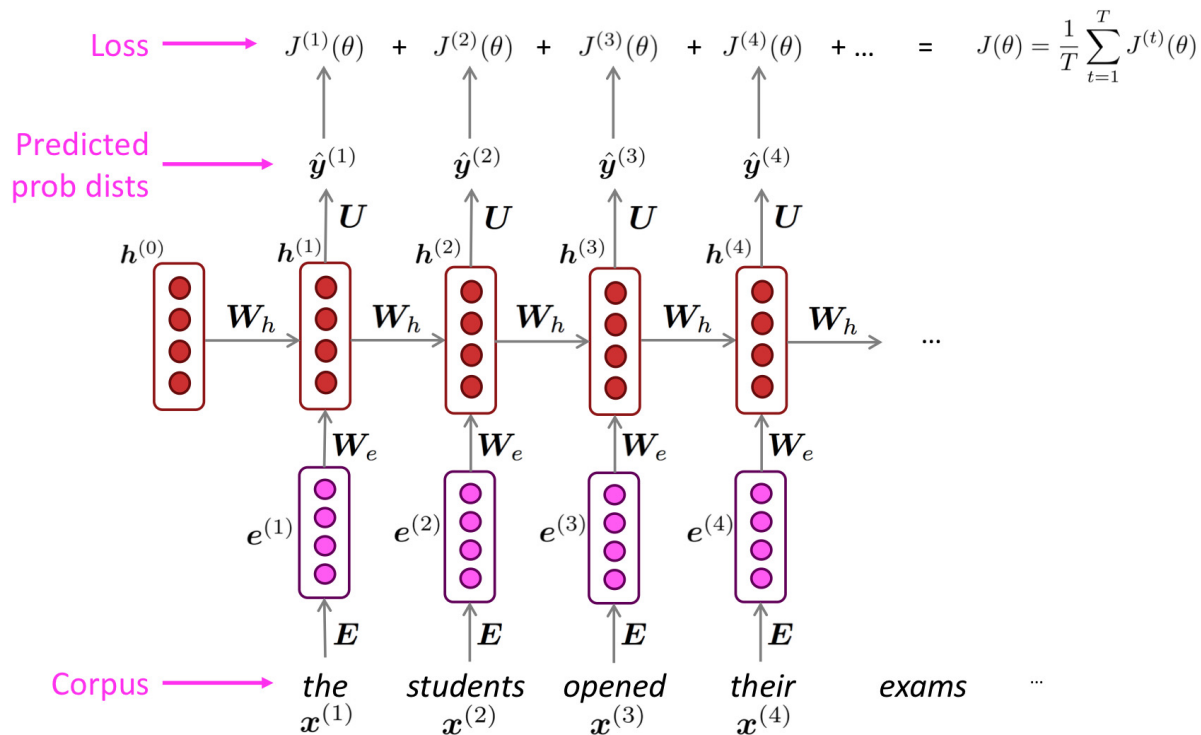
Problems with n-gram Language Model

Storage Problems

$$P(\mathbf{w}|\text{students opened their}) = \frac{\text{count}(\text{students opened their } \mathbf{w})}{\text{count}(\text{students opened their})}$$

## RNN

### Recurrent Neural Network(RNNs)



## RNN

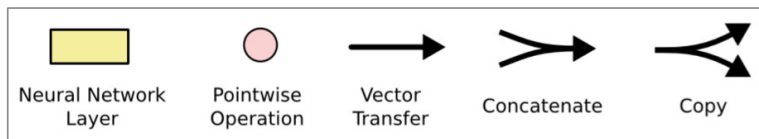
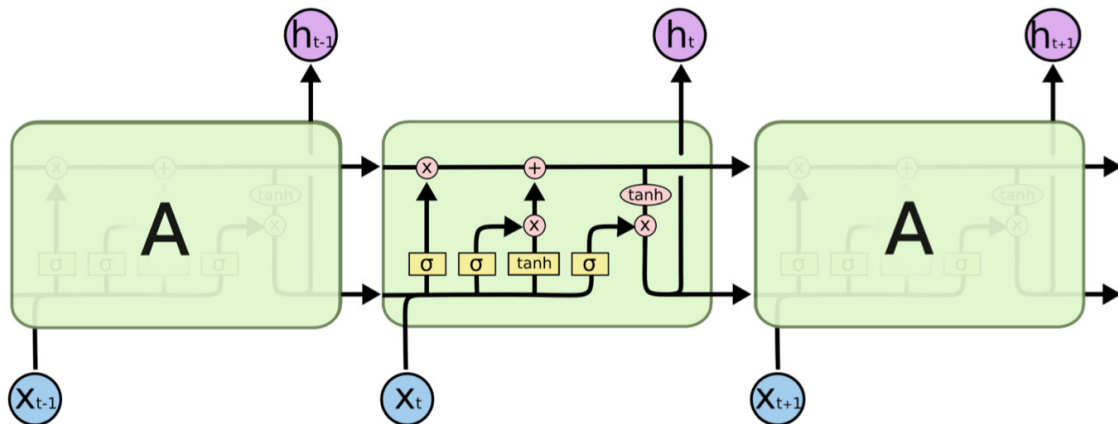
### Recurrent Neural Network(RNNs)

장점	단점
Input의 크기에 구애 받지 않음	계산이 느림
Input의 크기가 늘어나도 모델의 크기는 늘어나지 않음	
같은 W를 적용해 input 처리 과정이 대칭적임	아주 먼 단어의 정보를 쓰기가 어려움
아주 먼 단어의 정보도 사용할 수 있음	

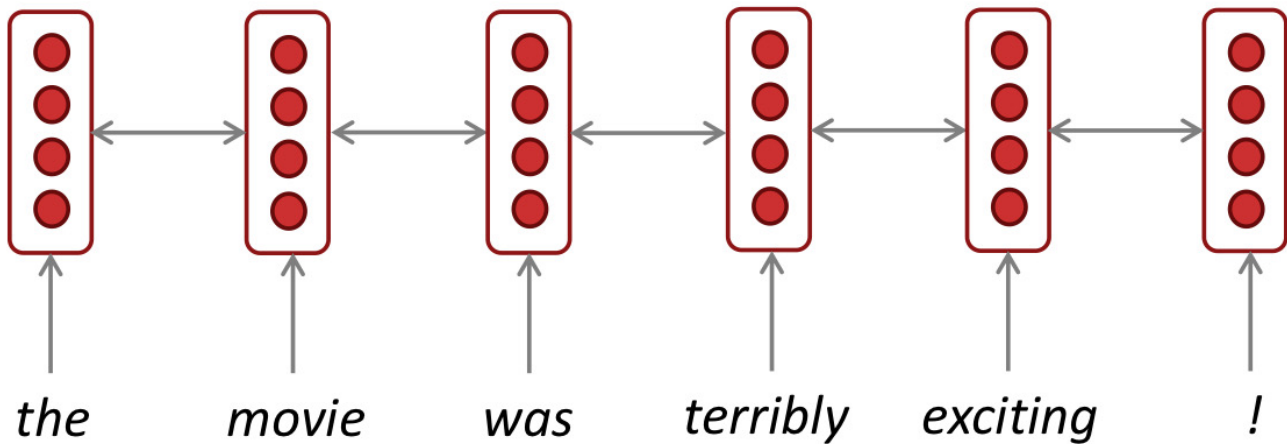


## LSTM

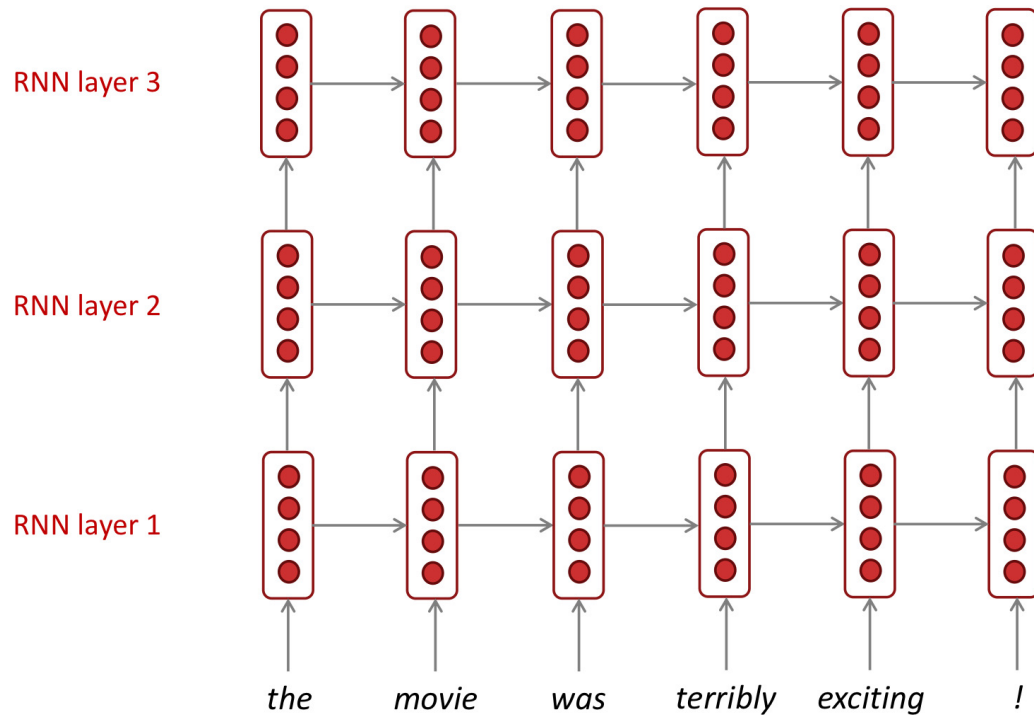
Long Short-Term Memory RNNs(LSTMs)



## Bidirectional RNNs



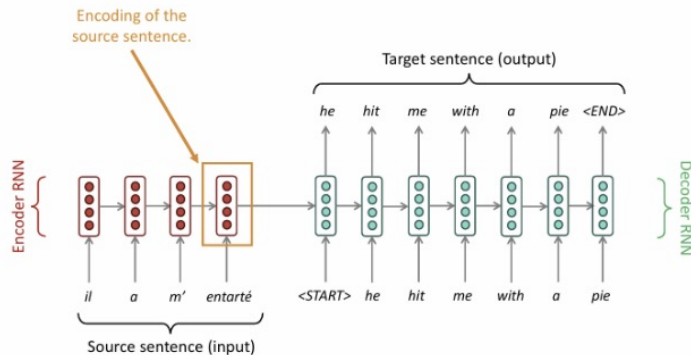
## Multi-layer RNNs



# 기존 RNN의 문제

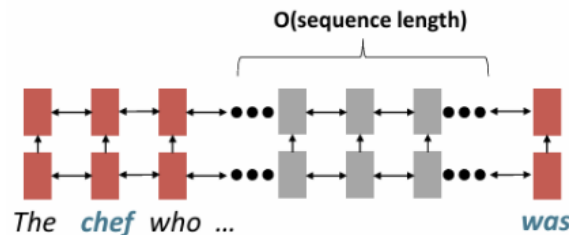
## 1. Bottleneck problem

- Encoder의 마지막 block은 input sentence의 모든 정보를 담아야 함



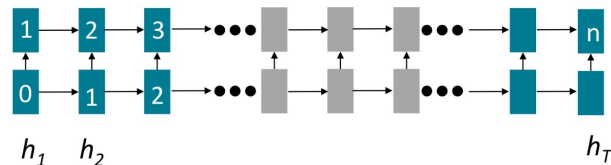
## 2. Long-distance dependency problem

- 서로 다른 단어가 상호작용하기 위해  $O(\text{sequence length})$  만큼의 step이 필요함
- LSTM으로도 위 문제를 완벽하게 극복하기 어려움



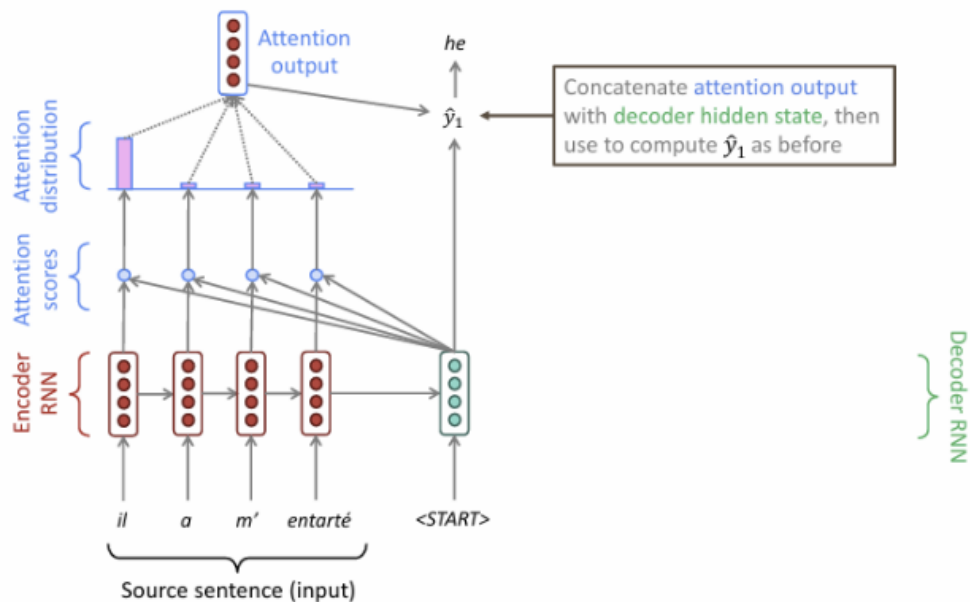
## 3. Lack of parallelizability

- 연산을 병렬로 수행할 수 없음
- 이전 hidden state가 계산 되기 전까지 그 다음 hidden state를 계산할 수 없기 때문



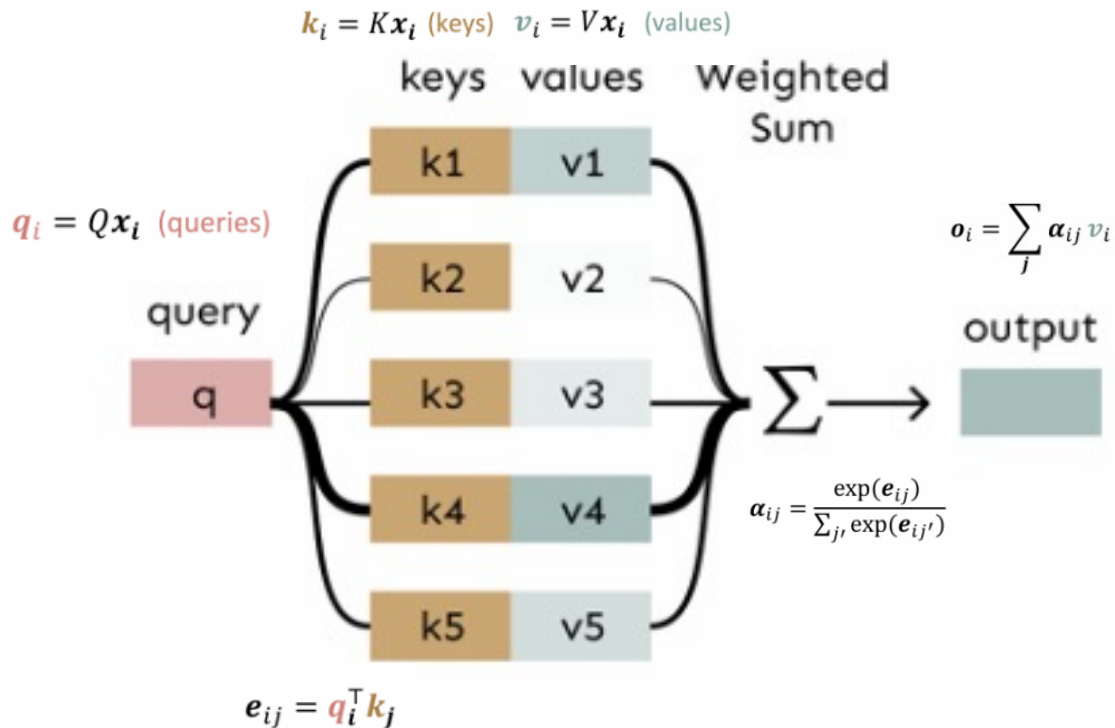
# Attention

Core idea : Decoder에서 encoder에 direct connection을 사용, input sentence의 특정 부분에 집중하자!



## Self-Attention

Attention과 차이점 : Query, Key, Value가 모두 동일한 embedding vector에서 도출



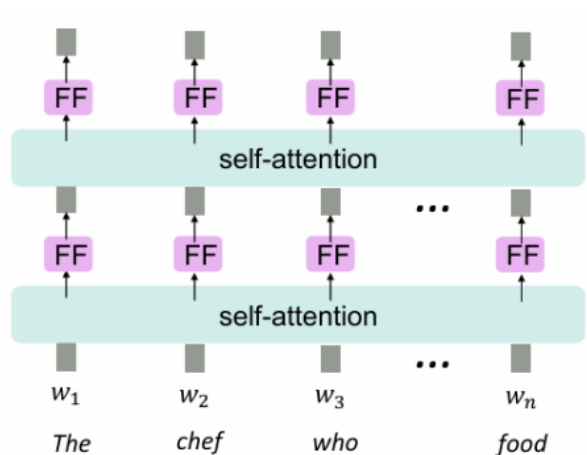
# Transformer

## 1. Positional Representation

$p_i \in \mathbb{R}^d$ , for  $i \in \{1, 2, \dots, n\}$  are position vectors

$$\tilde{x}_i = x_i + p_i$$

## 2. Feed-forward Network



## 3. Masking



## Transformer

