**ICT이노베이션스퀘어 AI복합교육 고급 언어과정**

# 자연어처리를 위한 Attention



**현청천**

2021.04.19
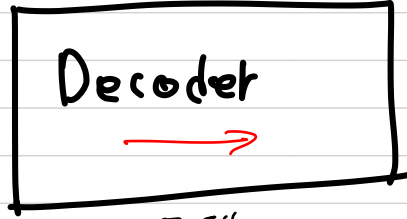
분류 문제에 좋음

Target (Next)

```
┌─────────────┐        ┌─────────────┐
│   Encoder   │ ────── │   Decoder   │
│     ←→      │        │      →      │
└─────────────┘        └─────────────┘
```

∘ feature 추출                ∘ 생성 모델

↑                             ↑
Soure                    Target (Previous)

∘ Transformer
∘ bert

∘ GPT ex) 소설. 대화. 오타 검색 (교정)

∘ google 에서 최근 점향

# What is Attention

주어진 **목적**에 맞는 **Source**에 **집중**

**ex) 감정분석**

할머니 만나는 부분에서 울었습니다. **감동적인 영화**입니다.

**긍정**

# What is Attention

주어진 **목적**에 맞는 **Source**에 **집중**
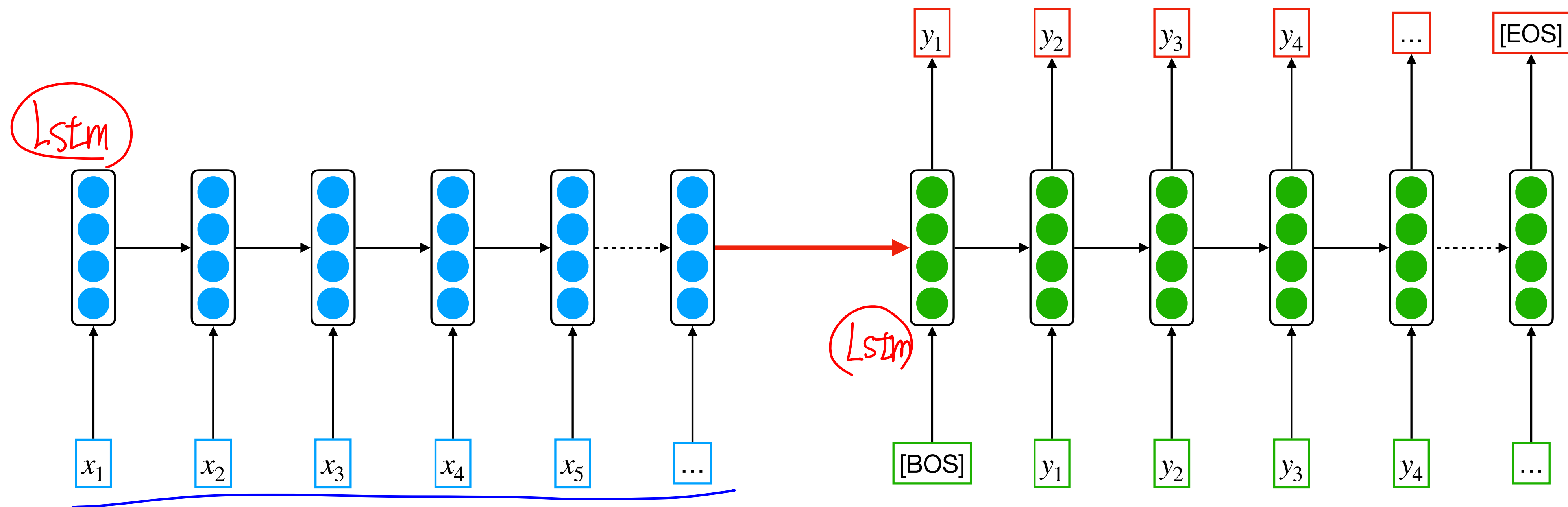
이미지분류



**고양이**

# What is Attention

N-M-T

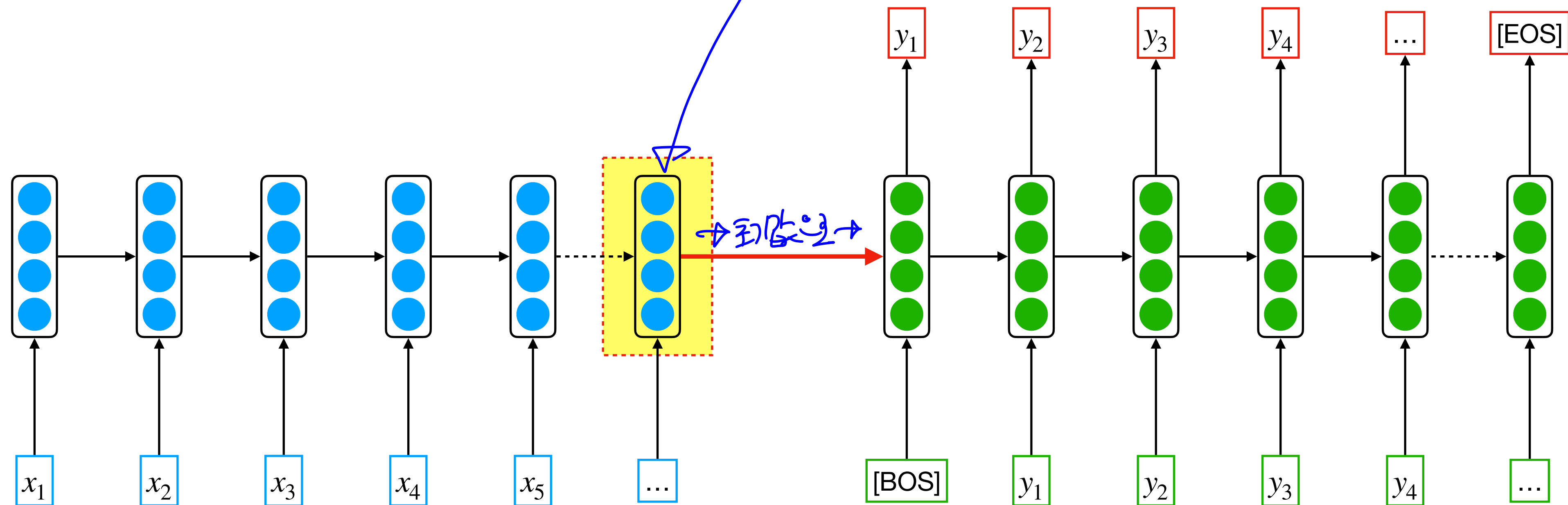나는 학생 입니다



Lstm

Lstm

Source

I - am Student

# What is Attention

Encoder Input 정보를 하나의 벡터로 저장

# What is Attention

**Encoder Input 정보를 하나의 벡터로 저장**

한계 ① 긴 문장을 하나의 벡터로 변환하기 어려움 **(Information bottleneck)**

⇒ 당연히 정보 손실이 생기지!

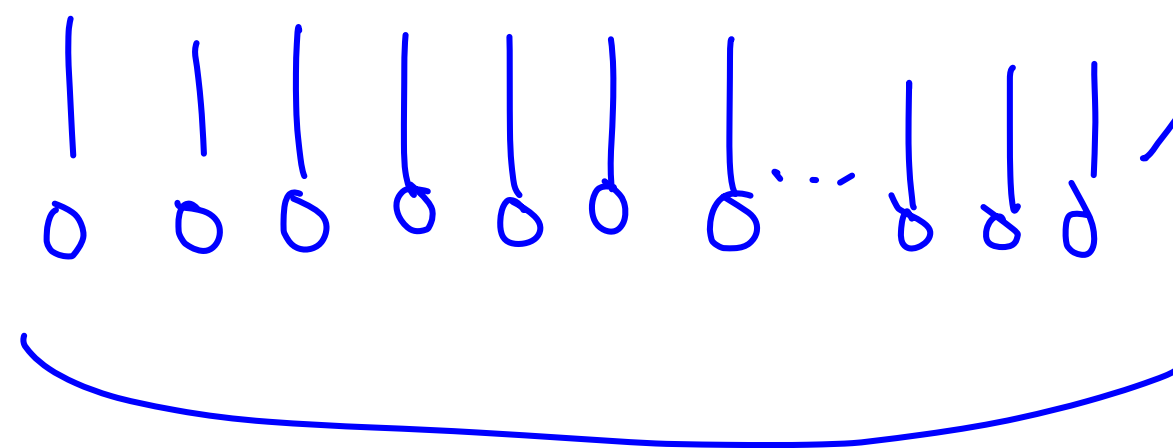if) token이 500개 있다면.

The dominant sequence transduction
models are based on complex

…

…

...

requiring significantly less time to train

500 개.

6

# What is Attention



**Encoder Input 정보를 하나의 벡터로 저장**

과거의 정보가 점점 사라짐 **(Vanishing gradient)**

# What is Attention

- 두가지 주요 문제 해결
  - 긴 문장을 하나의 벡터로 변환하면서 발생하는 Information bottleneck
  - 과거의 정보가 점점 사라지는 vanishing gradient

Encoder의   특정 벡터

**Source**의 **특정 부분을 집중하기 위해 Decoder가 Encoder의 정보를 직접 접근함**

# Attention Model



$o$ Encoder의 정보에
접근하고싶어요.

@2021 cchyun@gmail.com., Ltd All Rights Reserved.

9

# Attention Model

**Dot-product**
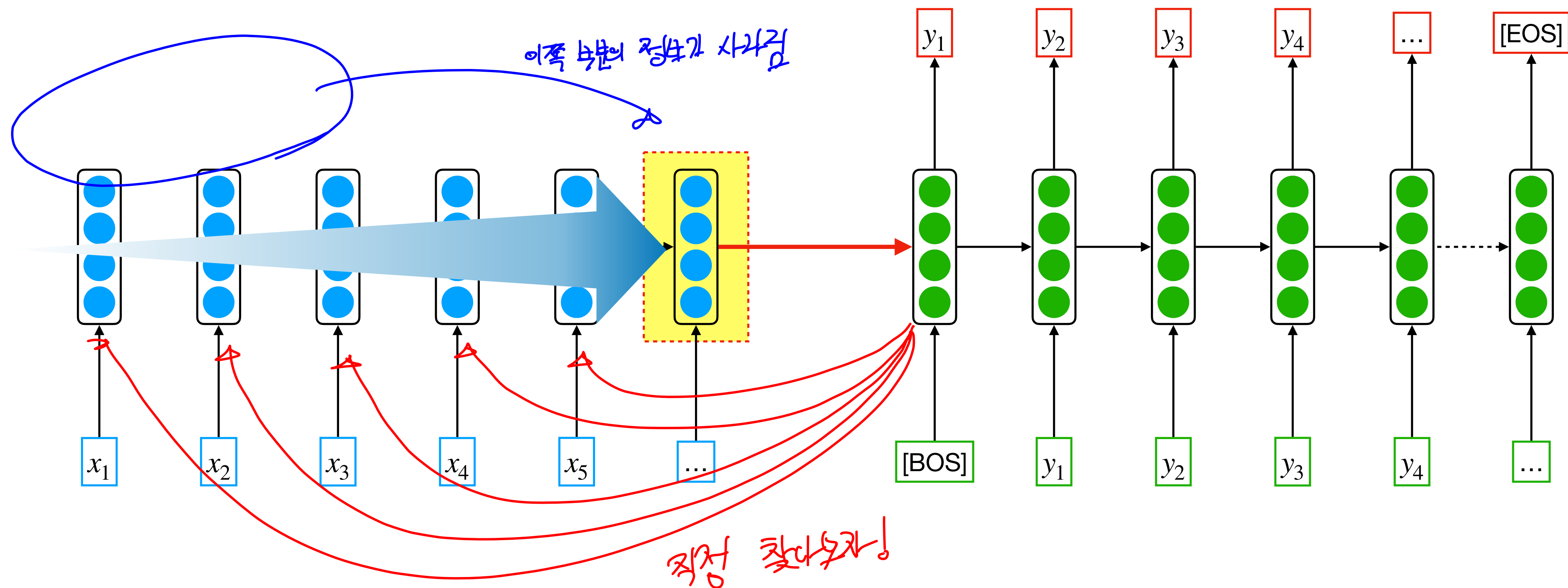
$$\boxed{s_1} \quad \boxed{h_1}$$

$$(h, 1) \qquad (h, 1)$$

$$\boxed{\frac{s_1^T}{(1, h)}} \cdot \boxed{\frac{h_1}{(h, 1)}} = \underline{1} = \underline{\text{스칼라 값}}$$

∘ Dot-Product

$$\vec{a} \circ \vec{b} = |a||b| \cos\theta$$

유사도값

$s_1 h_1^T$

# Attention Model

오shape이 감당되리만는 리나.

## Dot-product

# Attention Model

다 상을4까 "

$|s_1| \cdot |h_1| \cdot \cos\theta_1$

**Dot-product**

$|s_1| \cdot |h_1| \cos\theta_2$    $|s_1| |h_1| \cos\theta_3$

o $s_1$ 에대해 각 $h_i$ 들의 dot-product를 통해 유사도를 측정.

$s_1 h_1^T$    $s_1 h_2^T$    $s_1 h_3^T$    $s_1 h_4^T$    $s_1 h_5^T$    $s_1 h_n^T$

$s_1$

$h_1$    $h_2$    $h_3$    $h_4$    $h_5$    $h_n$

[BOS]

$x_1$    $x_2$    $x_3$    $x_4$    $x_5$    ...

# Attention Model



Softmax

$s_1$과 $\{h_1, \ldots, h_n\}$의 유사도를 확률분포로 변환

이 상음값을 softmax로 취하는거지!

13

# Attention Model



**Weighted Sum**
**Encoder hidden state의 weighted sum**

$$a_1 = \sum_{i=1}^{n} \alpha_{1i} h_i$$

attention output

○위사도를 통해 a를 계산
○ 그럼 때시 h 와 곱해서 합 ⇒ 앞시 a

$0.6 h_1 + 0.07 h_2 + 0.03 h_3 + \cdots 0.001 \cdot h_n$

$= h_1$은 60% 영향.
$= h_2$는 7% 영향.
⇒ · · · .

# Attention Model



encoder-weight-sum ⇒ encoder의 정보만 있는거지
S1은 단순히 Score를 나타내거니까
(⇒ 유다 3)

$a_1$

$y_1$

$\hat{y}_1$

**Concatenate attention output with decoder hidden state $[s_1; a_1]$ and compute output**

$\alpha_{11}$

$\alpha_{12}$ $\alpha_{13}$ $\alpha_{14}$ $\alpha_{15}$ $\alpha_{1n}$ ) 확률인간

$s_1 h_1^T$ $s_1 h_2^T$ $s_1 h_3^T$ $s_1 h_4^T$ $s_1 h_5^T$ $s_1 h_n^T$ ) Score

$h_1$ $h_2$ $h_3$ $h_4$ $h_5$ $h_n$

$s_1$

① Encoder의 정보만 있것아!
그래니까 Si랑 다시
Concat해야지!

concat  Si  a  a₁

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$ …

[BOS]

Q Sore를 계산한거지
그걸 통해서 학률관고
그런 h₁ 끓해서 각기 가중치를 주거짆다
그니까 Encoder 정부만있는거지!
그래서 decoder 정보를 합치기 위해서
Concat 하네함

# Attention Model

16

# Attention Model

# Attention Model (Equation)



Encoder hidden state: $h_1, \ldots, h_n \in \mathbb{R}^h$

h 차원수

h dimension

18

# Attention Model (Equation)

$$a_1$$

$$\alpha_{11}$$

$$\alpha_{12} \quad \alpha_{13} \quad \alpha_{14}$$

$$h_1, \ldots, h_n \in \mathbb{R}^h$$

$$y_1$$

Decoder hidden state on timestamp t: $s_t \in \mathbb{R}^h$

$$\hat{y}_1$$

이것도 $h$여러개랑

이거하고 dot produce 되니까.

$$s_1 h_1^T \quad s_1 h_2^T \quad s_1 h_3^T \quad s_1 h_4^T \quad s_1 h_5^T \quad s_1 h_n^T$$

$$s_1$$

$$h_1 \quad h_2 \quad h_3 \quad h_4 \quad h_5 \quad h_n$$

[BOS]

$$x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5 \quad \ldots$$

# Attention Model (Equation)



$$h_1, \ldots, h_n \in \mathbb{R}^h$$

$$s_t \in \mathbb{R}^h$$

Attention score: $e^t = [s_t^T h_1, \ldots, s_t^T h_n] \in \mathbb{R}^n$

$(1.h) \circ (h,1)$

$= (1.1)$

$\sim$ scalar

scalar가 n개!

$\circ$ 얼마나 맞는지 유사도해서

$cos\theta = \dfrac{|a||b|}{|a \cdot b|}$ 이게 유사도?

# Attention Model (Equation)



$$h_1, \ldots, h_n \in \mathbb{R}^h$$

$$s_t \in \mathbb{R}^h$$

$$e^t = [s_t^T h_1, \ldots, s_t^T h_n] \in \mathbb{R}^n$$

Attention probability distribution: $\alpha^t = softmax(e^t) \in \mathbb{R}^n$

# **Attention Model (Equation)**



$$h_1, \dots, h_n \in \mathbb{R}^h$$

$$s_t \in \mathbb{R}^h$$

$$e^t = [s_t^T h_1, \dots, s_t^T h_n] \in \mathbb{R}^n$$

$$\alpha^t = softmax(e^t) \in \mathbb{R}^n$$

Weighted sum of encoder hidden state: $\quad a_t = \sum_{i=1}^{n} \alpha_i^t h_i \in \mathbb{R}^h$

$$(1,n) \cdot (n,h)$$
$$= (1,h)$$

# Attention Model (Equation)



$$h_1, \ldots, h_n \in \mathbb{R}^h$$

$$s_t \in \mathbb{R}^h$$

$$e^t = [s_t^T h_1, \ldots, s_t^T h_n] \in \mathbb{R}^n$$

$$\alpha^t = softmax(e^t) \in \mathbb{R}^n$$

$$a_t = \sum_{i=1}^{n} \alpha_i^t h_i \in \mathbb{R}^h$$

concat

Concatenate attention output and decoder hidden state:   $[a_t; s_t] \in \mathbb{R}^2 h$

# Attention Model (Equation)

Encoder hidden state: $h_1, \ldots, h_n \in \mathbb{R}^h$
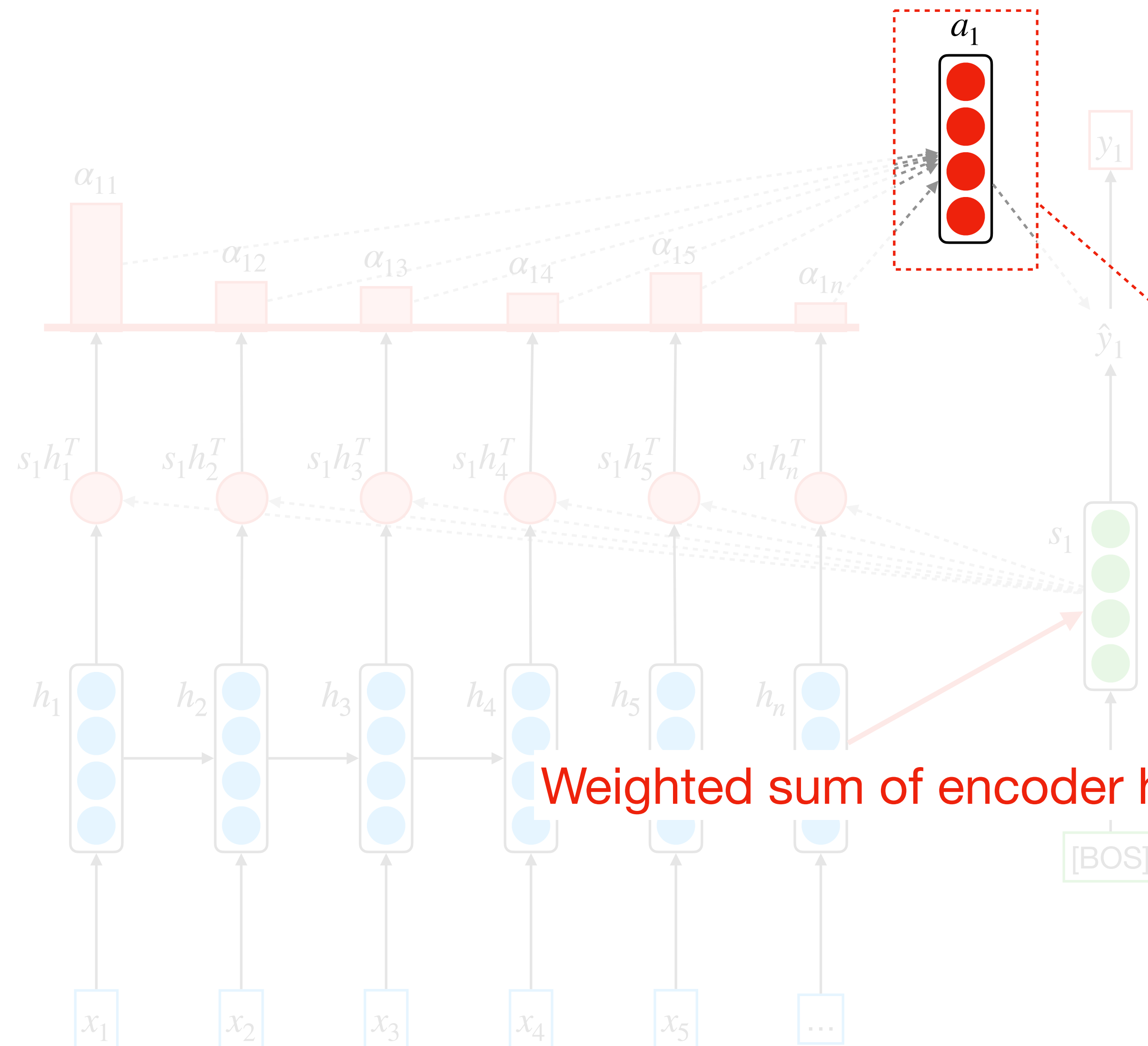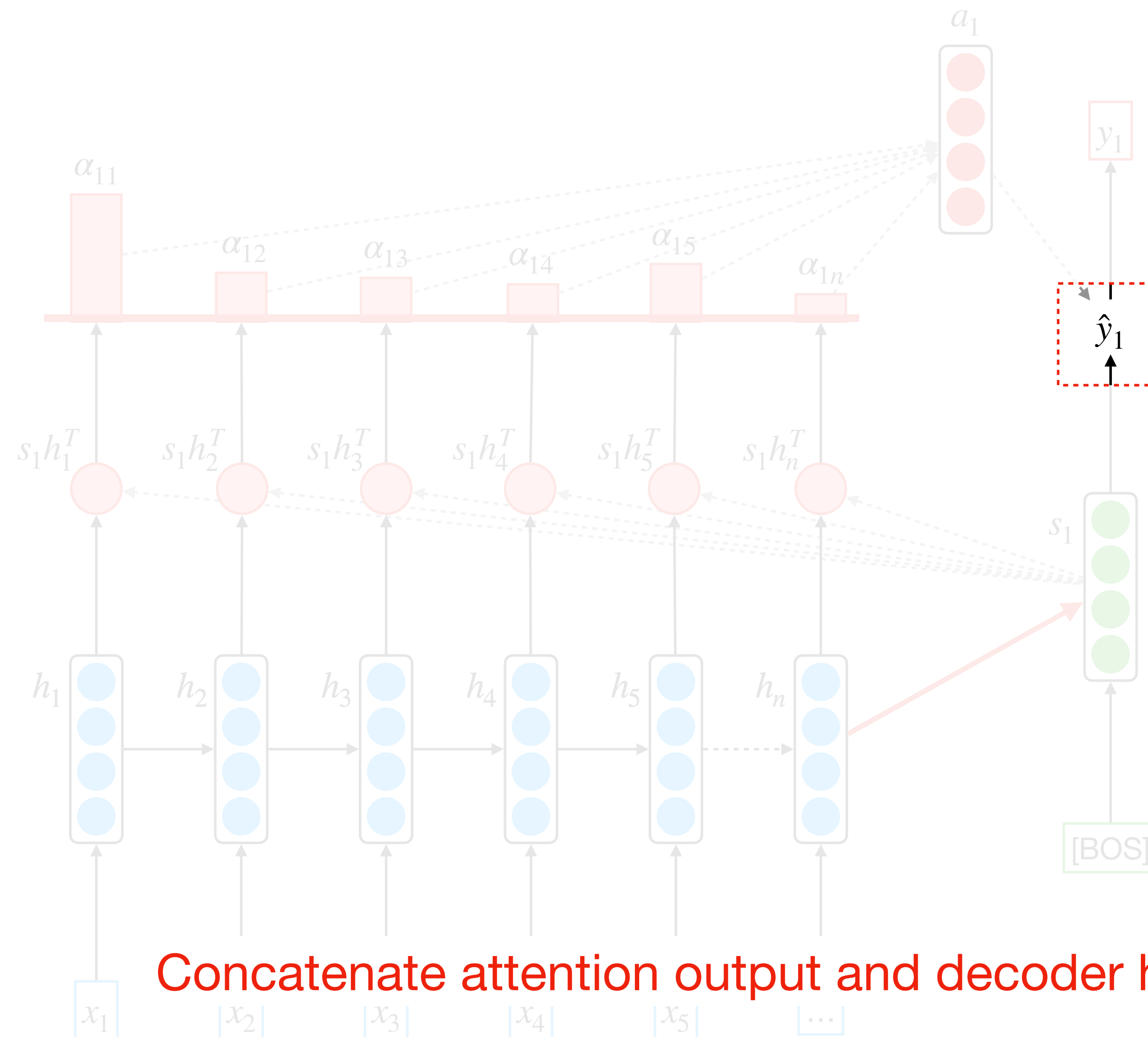
Decoder hidden state on timestamp t: $s_t \in \mathbb{R}^h$

Attention score: $e^t = [s_t^T h_1, \ldots, s_t^T h_n] \in \mathbb{R}^n$

Attention probability distribution: $\alpha^t = softmax(e^t) \in \mathbb{R}^n$

Weighted sum of encoder hidden state: $a_t = \sum_{i=1}^{n} \alpha_i^t h_i \in \mathbb{R}^h$

concat.

Concatenate attention output and decoder hidden state: $[a_t; s_t] \in \mathbb{R}^2 h$

# Attention Model (Advantage)

- Attention을 이용해 NMT의 성능이 많이 좋아짐 → 가중치를 준거임.
  - Decoder가 source의 특정 부분에 집중하도록 한 것이 매우 효과적임
- 한계1 해결 · Information bottleneck 문제를 해결 함
  - Decoder가 source에 직접 접근하도록 함
- 한계2 해결 · Vanishing gradient 문제를 해결 함
  - 거리가 먼 source의 정보를 접근 할 수 있음
- Attention이 alignment를 학습함

feature map을 보고는 해석이 어려워지지만,
attention map을 보고 해석이 가능해.

| | Education | is | most | powerful | weapon |
|---|---|---|---|---|---|
| 교육은 | ■ | | | | |
| 가장 | | | ■ | | |
| 강력한 | | | | ■ | |
| 무기 | | | | | ■ |
| 입니다 | | ■ | | | |

# Attention Model (Variants)

Values $\longrightarrow$ Encoder hidden $\quad$ $\boxed{(h_1), \ldots, h_n \in \mathbb{R}^h}$ $\qquad$ decoder hidden $\quad \boxed{s_t \in \mathbb{R}^h}$

○ $S \cdot h$의 dot-product를 했지.

Attention score $\longrightarrow$ $S_t^\mathsf{T} \cdot h_i = \boxed{e^t \in \mathbb{R}^n}$

**attention score를 계산하는 다양한 방법이 있음**

Attention probability $\longrightarrow$ $\boxed{\alpha^t = softmax(e^t) \in \mathbb{R}^n}$

○ $h$가 2개 동위

Attention output $\longrightarrow$ $\boxed{a_t = \sum_{i=1}^{n} \alpha_i^t (h_i) \in \mathbb{R}^h}$

# Attention Model (Variants)

*(handwritten, top right)* bi-lstm ⇒ unit= $\frac{d\text{-model}}{2}$ 야지  
out을 맞춰야하는거고.

- **Dot-product attention**
  - $e_i^t = s_t^T h_i \in \mathbb{R}$
    
    *(handwritten)* 내적.

- Multiplicative attention   *(handwritten)* weight shape이 다르다!
  - $e_i^t = s_t^T W h_i \in \mathbb{R}$
  - where $W \in \mathbb{R}^{d_s \times d_h}$   *(handwritten)* 곱하기서 ch로 바꿔줌!

$$e_t = [e_{t1}, \ldots, e_{tn}] \in \mathbb{R}^n$$

- Additive attention
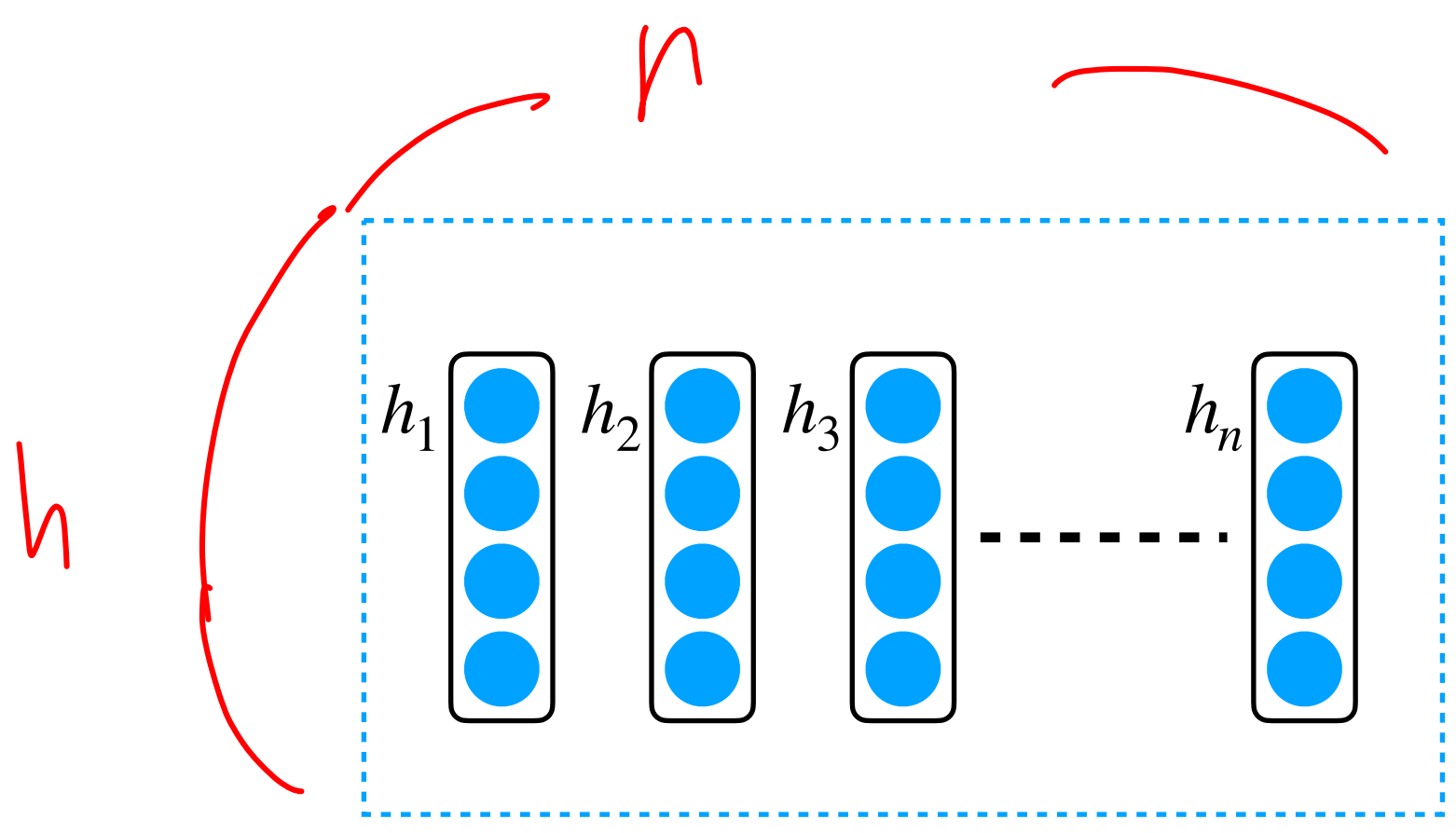  - $e_i^t = v^T \tanh(W_h h_i + W_s s_t) \in \mathbb{R}$
  - where $W_h \in \mathbb{R}^{d_v \times d_h}, W_s \in \mathbb{R}^{d_v \times d_s}, v \in \mathbb{R}^{d_v}$

# **Attention Tutorial**
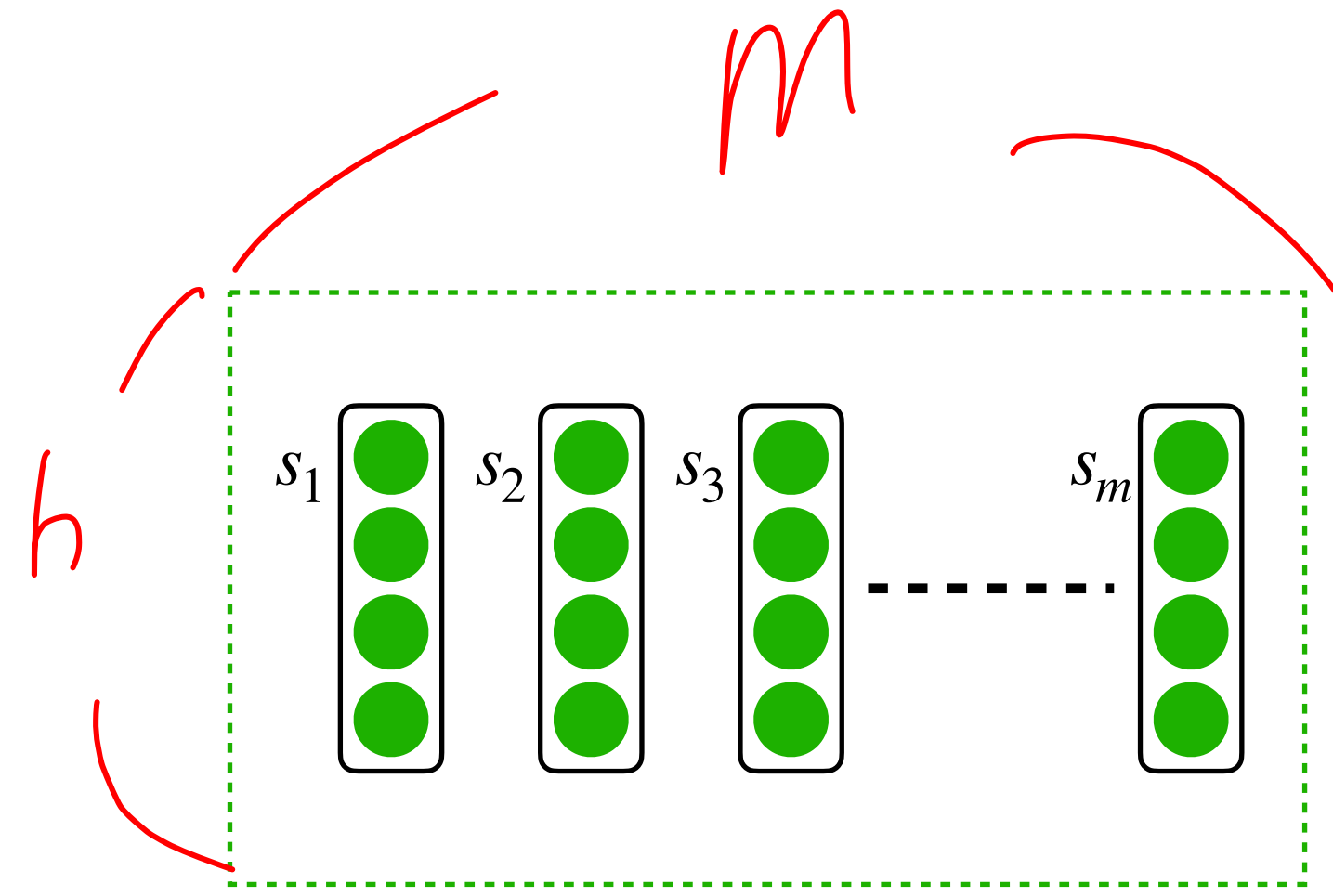
# Attention Tutorial (inputs)

$h, m$    $n \cdot h$

$h \times n$        $h \cdot m$

$(n \times h) \times (h \cdot m)$
$= (n \cdot m)$

$n$

$h$

$h_1$ $h_2$ $h_3$ $h_n$

Encoder hidden state: $h \in \mathbb{R}^{h \times n}$

$(h \times n)$

$m$

$h$

$s_1$ $s_2$ $s_3$ $s_m$

Decoder hidden state: $s \in \mathbb{R}^{h \times m}$

$(h \times m)$

$(h \times m) \times (n \times h)$

# Attention Tutorial (score)

$\circ$ $h^T$ = enc_hidden

$\circ$ $s^T$ = dec_hidden

방향을대면 score가 행렬에 펼쳐진다!

decoder - hidden

$s^T$

$$e = \boxed{s^T h} \in \mathbb{R}^{m \times n}$$

encoder- hidden $T$

$(m \times h) \times (h \times n)$

$= (m \times n)$

$h$

$m$

| | $h_1$ | $h_2$ | $h_3$ | ... | $h_n$ |
|---|---|---|---|---|---|
| $s_1^T$ | $s_1^T h_1$ | $s_1^T h_2$ | $s_1^T h_3$ | ... | $s_1^T h_n$ |
| $s_2^T$ | $s_2^T h_1$ | $s_2^T h_2$ | $s_2^T h_3$ | ... | $s_2^T h_n$ |
| $s_3^T$ | $s_3^T h_1$ | $s_3^T h_2$ | $s_3^T h_3$ | ... | $s_3^T h_n$ |
| ... | ... | ... | ... | ... | ... |
| $s_m^T$ | $s_m^T h_1$ | $s_m^T h_2$ | $s_m^T h_3$ | ... | $s_m^T h_n$ |

이 score를

$s_1 \cdot h_n$의 Score

$h_1$ $h_2$ $h_3$ $h_n$

Encoder hidden state: $h \in \mathbb{R}^{h \times n}$

Attention score: $e \in \mathbb{R}^{m \times n}$

Decoder hidden state: $s^T \in \mathbb{R}^{m \times h}$
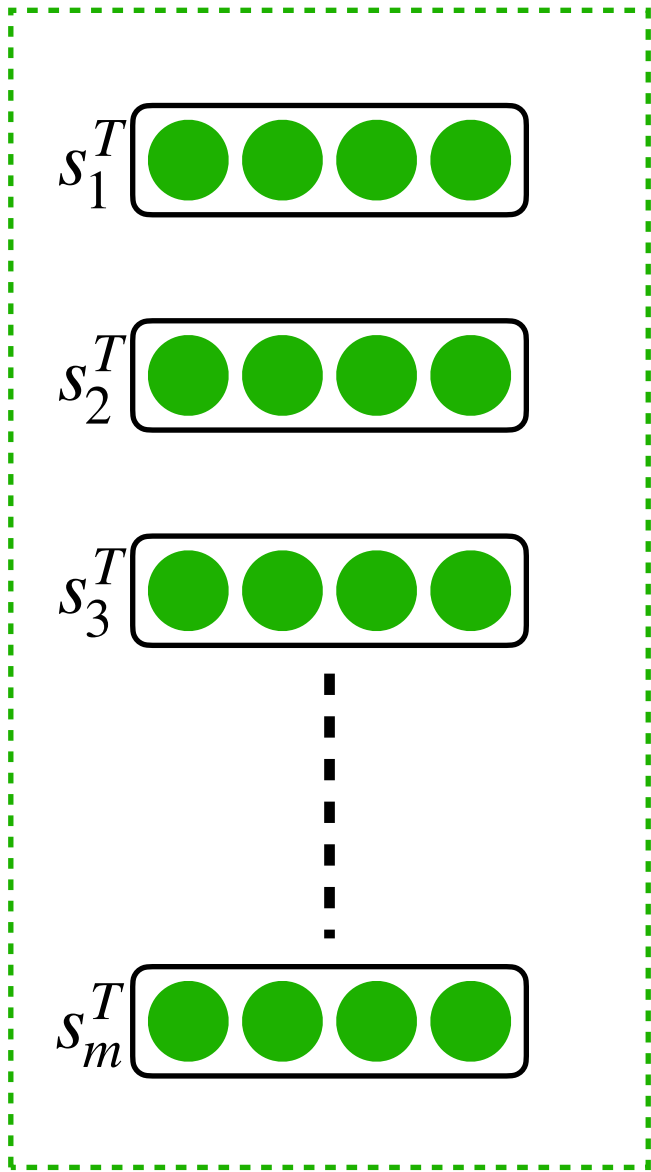
$\begin{pmatrix} a & b \\ c & d \end{pmatrix} \begin{pmatrix} e & f \\ g & h \end{pmatrix} = \begin{pmatrix} ae+bg & \cdots \end{pmatrix}$

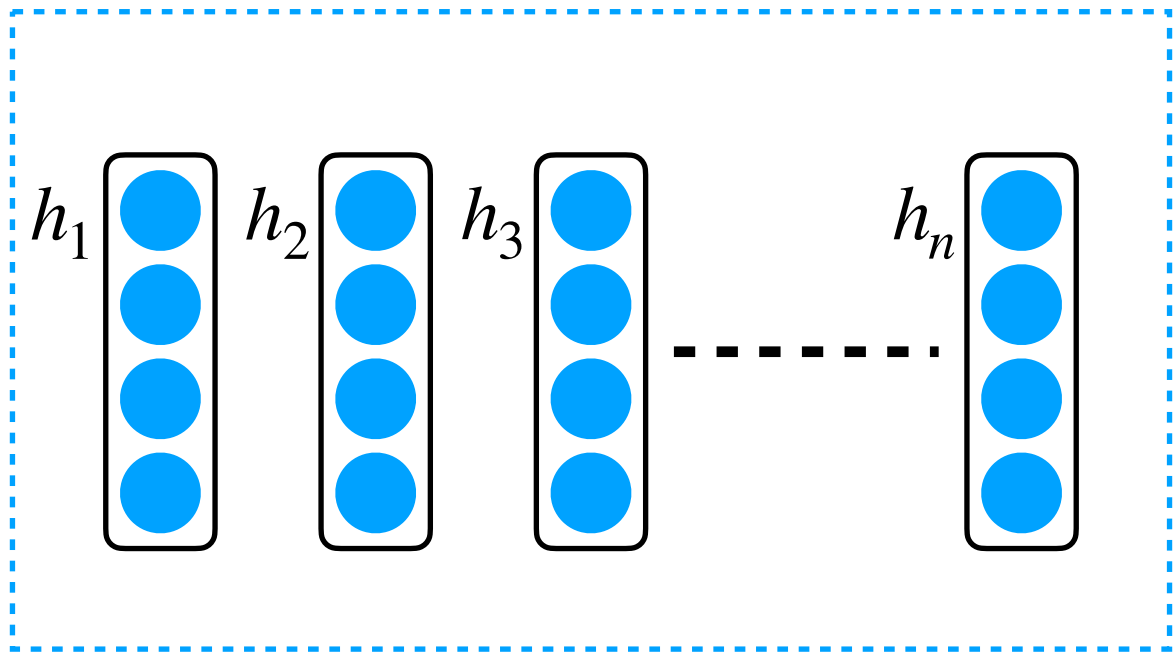$e_j^i = s_i^T h_j \in \mathbb{R}$

30

# Attention Tutorial (score)

$$e = s^T h \in \mathbb{R}^{m \times n}$$

그림에서 (분홍에서) (행벡터)    소수.

실제 코드 (열벡터)    numpy . tf.



Decoder hidden state: $s^T \in \mathbb{R}^{m \times h}$

Encoder hidden state: $h \in \mathbb{R}^{h \times n}$

Attention score: $e \in \mathbb{R}^{m \times n}$

$$e_j^i = s_i^T h_j \in \mathbb{R}$$

# Attention Tutorial (prob)

$$\alpha = softmax(e) \in \mathbb{R}^{m \times n}$$

$softmax($

|  | $h_1$ | $h_2$ | $h_3$ | ... | $h_n$ |
|---|---|---|---|---|---|
| $s_1^T$ | $s_1^T h_1$ | $s_1^T h_2$ | $s_1^T h_3$ | ... | $s_1^T h_n$ |
| $s_2^T$ | $s_2^T h_1$ | $s_2^T h_2$ | $s_2^T h_3$ | ... | $s_2^T h_n$ |
| $s_3^T$ | $s_3^T h_1$ | $s_3^T h_2$ | $s_3^T h_3$ | ... | $s_3^T h_n$ |
| ... | ... | ... | ... | ... | ... |
| $s_m^T$ | $s_m^T h_1$ | $s_m^T h_2$ | $s_m^T h_3$ | ... | $s_m^T h_n$ |

$)$

$=$

|  | $h_1$ | $h_2$ | $h_3$ | ... | $h_n$ |
|---|---|---|---|---|---|
| $s_1^T$ | $\alpha_1^1$ | $\alpha_2^1$ | $\alpha_3^1$ | ... | $\alpha_n^1$ |
| $s_2^T$ | $\alpha_1^2$ | $\alpha_2^2$ | $\alpha_3^2$ | ... | $\alpha_n^2$ |
| $s_3^T$ | $\alpha_1^3$ | $\alpha_2^3$ | $\alpha_3^3$ | ... | $\alpha_n^3$ |
| ... | ... | ... | ... | ... | ... |
| $s_m^T$ | $\alpha_1^m$ | $\alpha_2^m$ | $\alpha_3^m$ | ... | $\alpha_n^m$ |

$S_1$에대한 확률 분포 $\Rightarrow$ $\sum = 1$

$S_2$에 대한 확률분포 $\Rightarrow$ $\sum = 1$

Attention score: $e \in \mathbb{R}^{m \times n}$

Attention prob: $\alpha \in \mathbb{R}^{m \times n}$

행 단위 softmax

# Attention Tutorial (output)
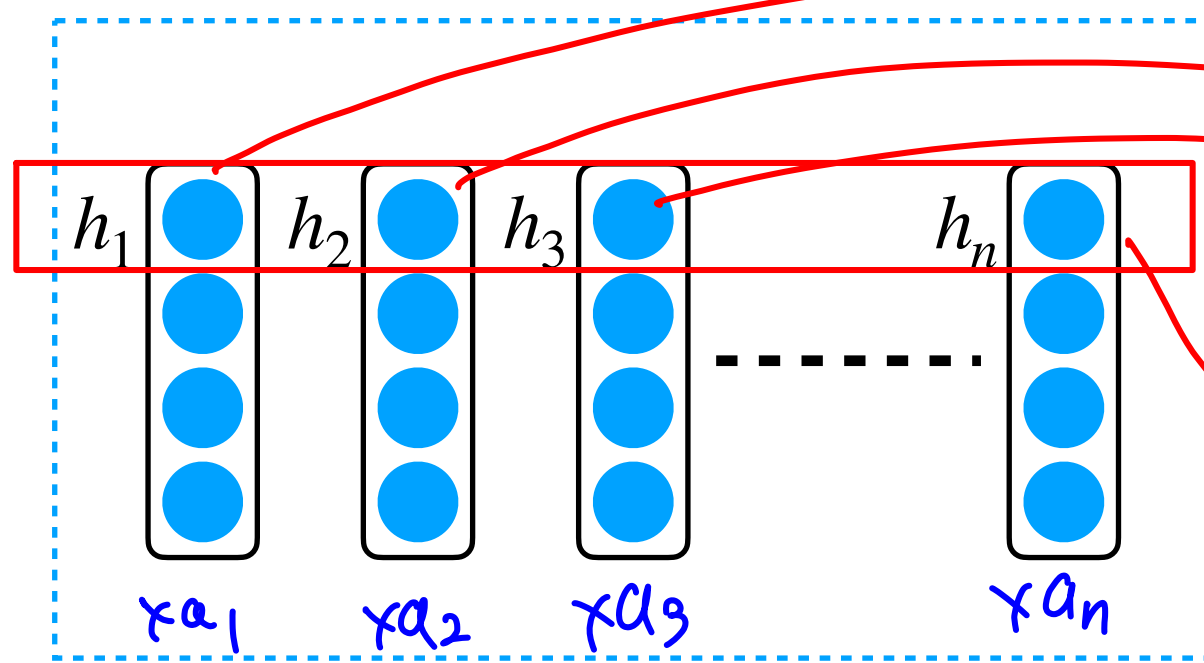
$$(A \cdot B)^T = B^T \cdot A^T$$

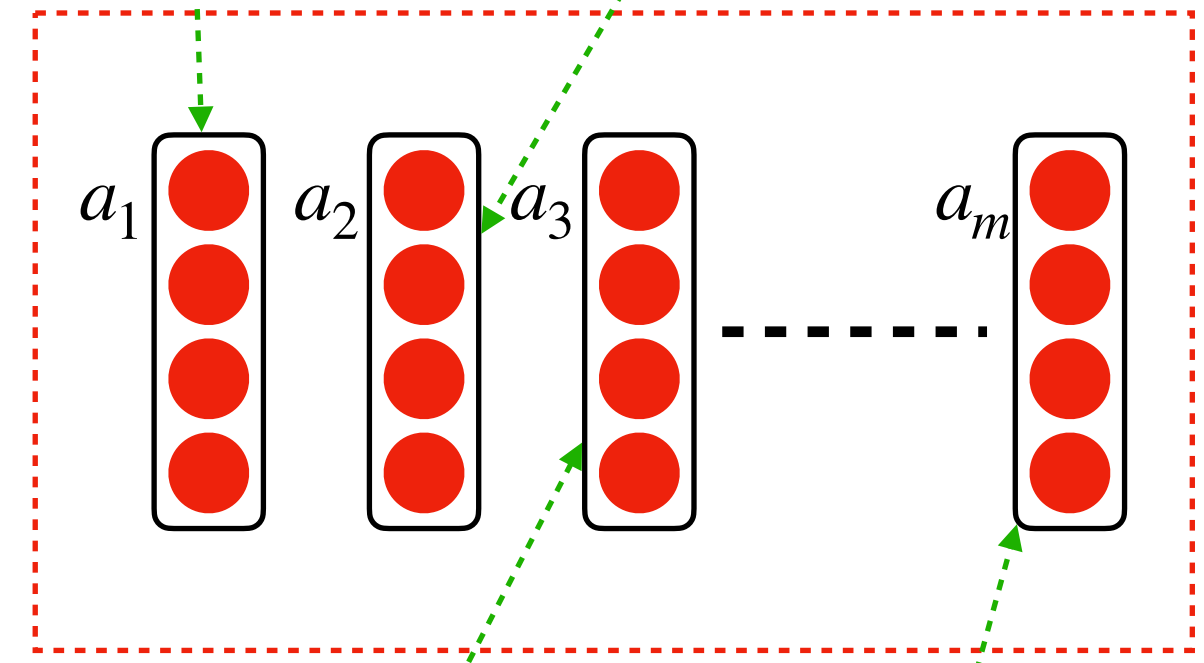$$a = h\alpha^T \in \mathbb{R}^{h \times m}$$

$$\alpha^T = \alpha \cdot h^T$$

$$\alpha_1^1 h_1 + \alpha_2^1 h_2 + \alpha_3^1 h_3 + \ldots + \alpha_n^1 h_n$$

$$\alpha_1^2 h_1 + \alpha_2^2 h_2 + \alpha_3^2 h_3 + \ldots + \alpha_n^2 h_n$$

$h_1 \alpha$

| | $s_1$ | $s_2$ | $s_3$ | $\ldots$ | $s_m$ |
|---|---|---|---|---|---|
| $h_1^T$ | $\alpha_1^1$ | $\alpha_1^2$ | $\alpha_1^3$ | $\ldots$ | $\alpha_1^m$ |
| $h_2^T$ | $\alpha_2^1$ | $\alpha_2^2$ | $\alpha_2^3$ | $\ldots$ | $\alpha_2^m$ |
| $h_3^T$ | $\alpha_3^1$ | $\alpha_3^2$ | $\alpha_3^3$ | $\ldots$ | $\alpha_3^m$ |
| $\ldots$ | $\ldots$ | $\ldots$ | $\ldots$ | $\ldots$ | $\ldots$ |
| $h_n^T$ | $\alpha_n^1$ | $\alpha_n^2$ | $\alpha_n^3$ | $\ldots$ | $\alpha_n^m$ |

0.6

$\times$

$h_1$  $h_2$  $h_3$  $h_n$  $\times a_1$  $\times a_2$  $\times a_3$  $\times a_n$

$=$

$a_1$  $a_2$  $a_3$  $a_m$

**Encoder hidden state:** $h \in \mathbb{R}^{h \times n}$

scalar $a_1^i$ $\begin{bmatrix} h_1 \end{bmatrix}$ + $a_2^i$ $\begin{bmatrix} h_2 \end{bmatrix}$ + $a_3^i$ $\begin{bmatrix} h_3 \end{bmatrix}$ + $\ldots$ + $a_n^i$ $\begin{bmatrix} h_n \end{bmatrix}$

0.6  0.1  0.01  0.00

$(h \times n)$

$= (h \times m)$

**Attention prob:** $\alpha^T \in \mathbb{R}^{n \times m}$

Transpose 된 상태

$(n \times m)$

**Attention output:** $a \in \mathbb{R}^{h \times m}$

$$\alpha_1^m h_1 + \alpha_2^m h_2 + \alpha_3^m h_3 + \ldots + \alpha_n^m h_n$$

$$\alpha_1^3 h_1 + \alpha_2^3 h_2 + \alpha_3^3 h_3 + \ldots + \alpha_n^3 h_n$$

O attention은 weight를 늘린게 아니라. 계산량을 늘린거지

O 지표도 보통 blue를.. count 기반

extraction → square

생성모델 →

# 감사합니다.