

Distinguished Engineering

Transformer 2/5 **- Attention, please**

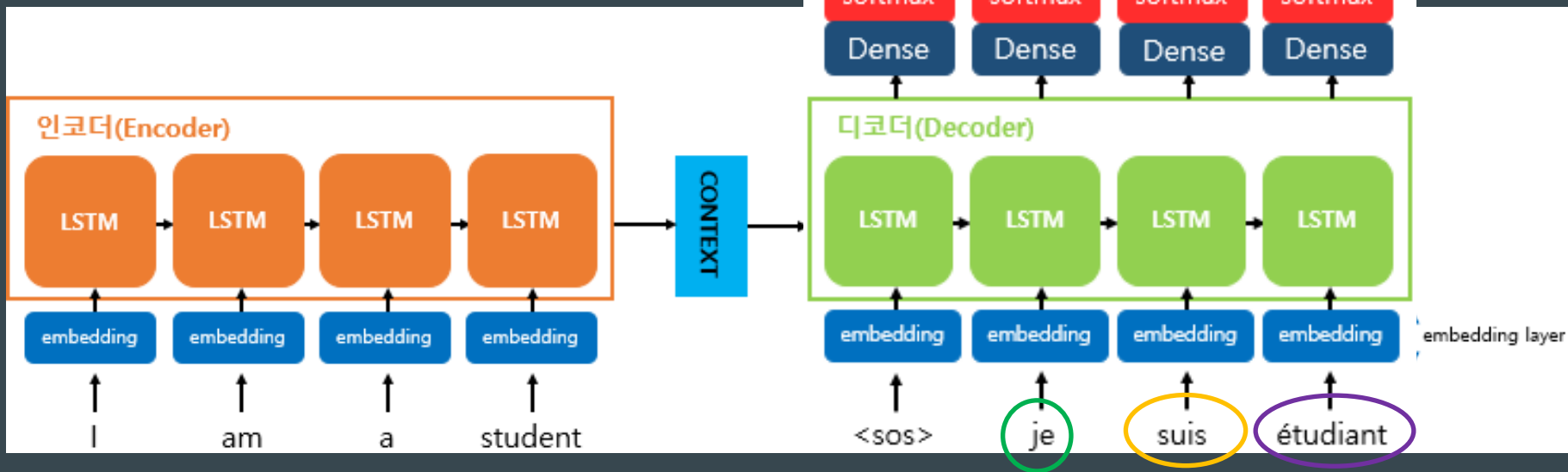
...

BW

Plan

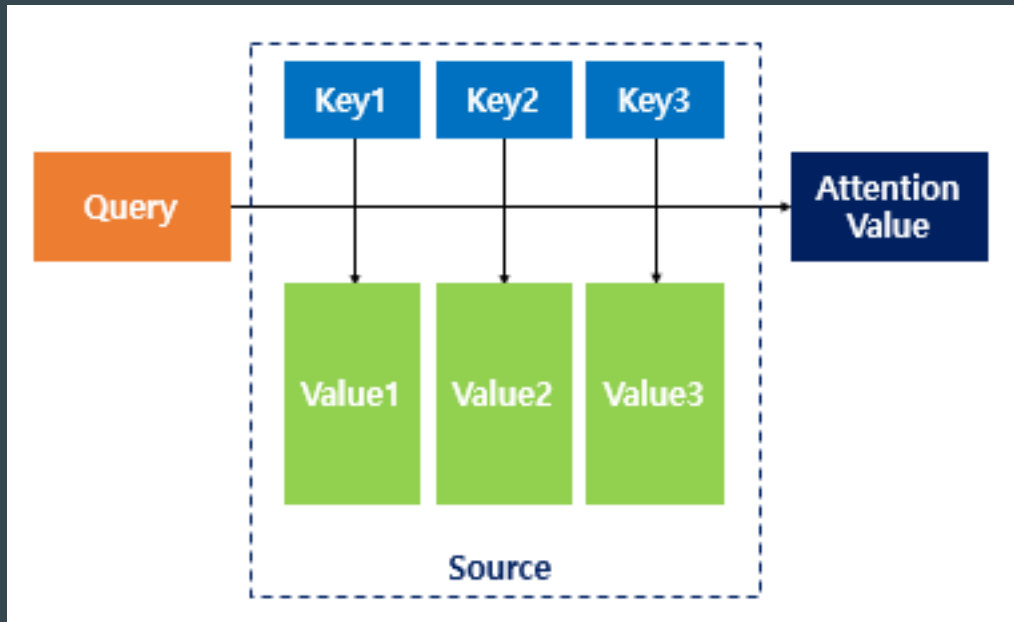
- Prologue, seq2seq
- **Attention, please**
- Transformer, a new hope
- Transformer, revenge of the fallen
- Transformer, vision

Sequence-to-Sequence



I		0.157	am		0.78	a		0.75	student		0.88
		-0.25			0.29			-0.81			-0.17
		0.478			-0.96			0.96			0.29
		-0.78			0.52			0.12			0.48

Attention



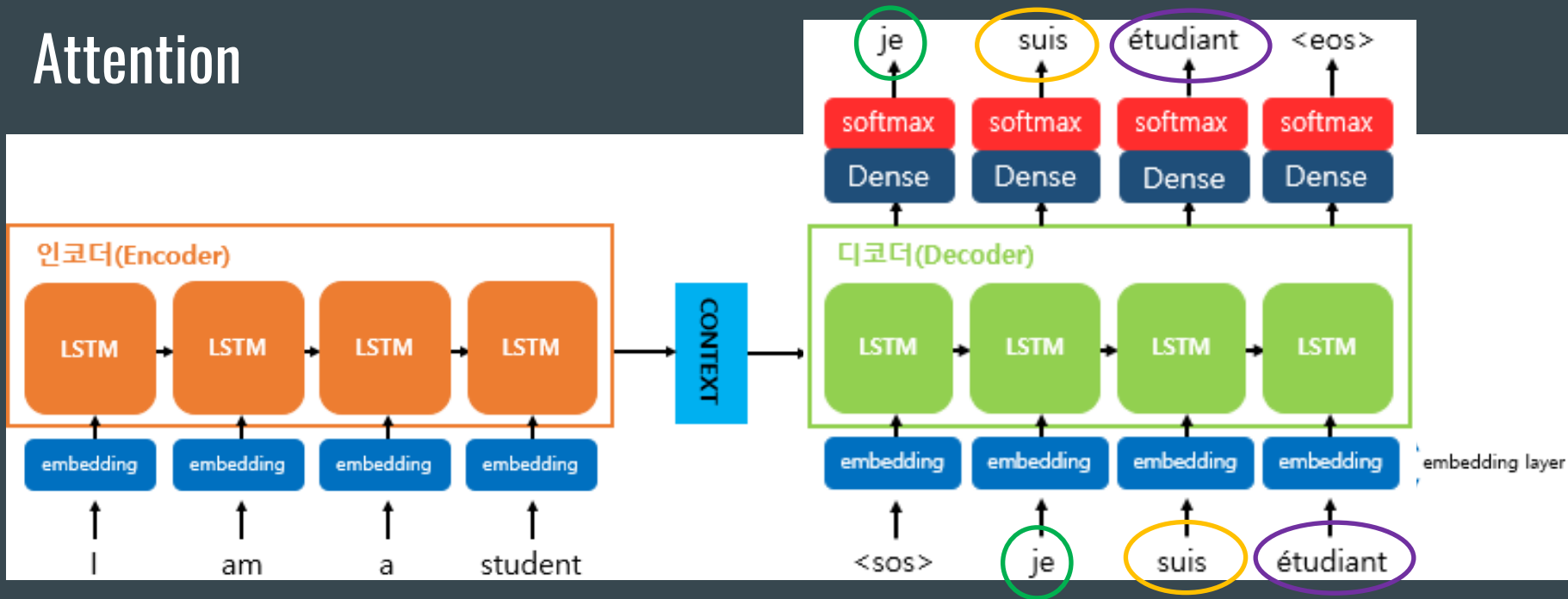
Attention(Q, K, V) = Attention Value

Q = Query : t 시점의 **Decoder** 셀에서의 은닉 상태

K = Keys : 모든 시점의 **Encoder** 셀의 은닉 상태들

V = Values : 모든 시점의 **Encoder** 셀의 은닉 상태들

Attention



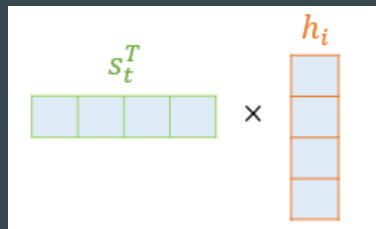
Attention(Q, K, V) = Attention Value

Q = Query : t 시점의 **Decoder** 셀에서의 은닉 상태
K = Keys : 모든 시점의 **Encoder** 셀의 은닉 상태들
V = Values : 모든 시점의 **Encoder** 셀의 은닉 상태들

Attention

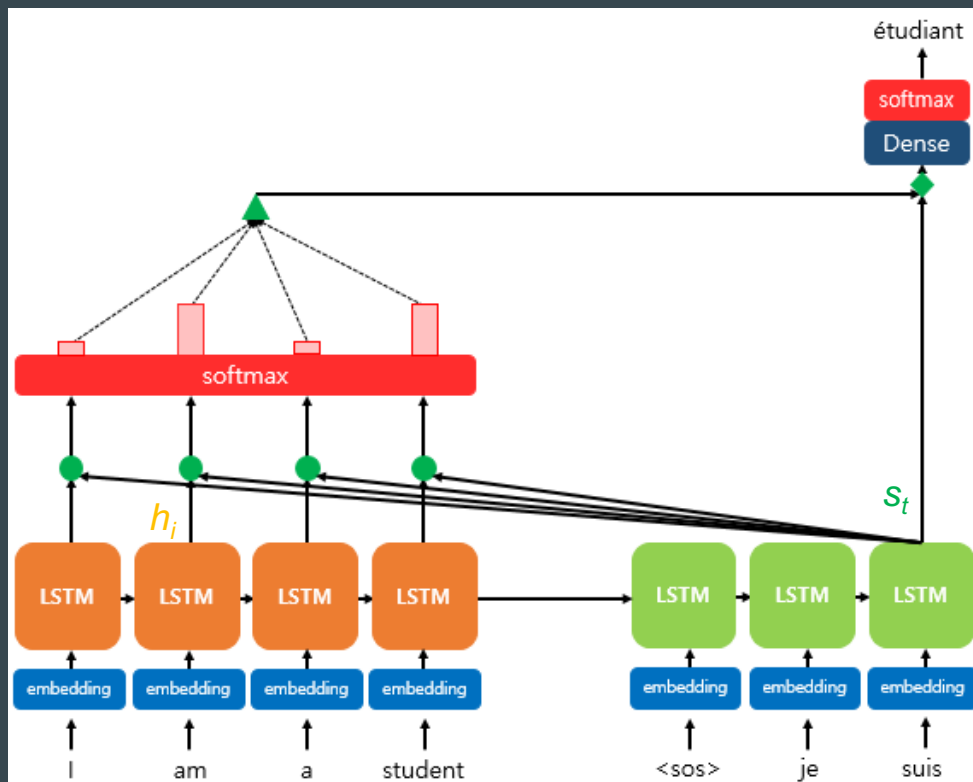
- Dot product attention

Attention score



$$\text{score}(s_t, h_i) = s_t^T h_i$$

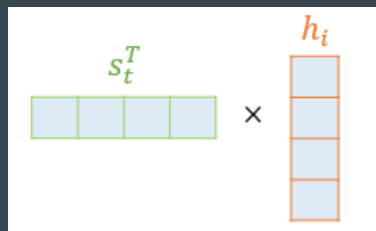
$$e^t = [s_t^T h_1, \dots, s_t^T h_N]$$



Attention

- Dot product attention

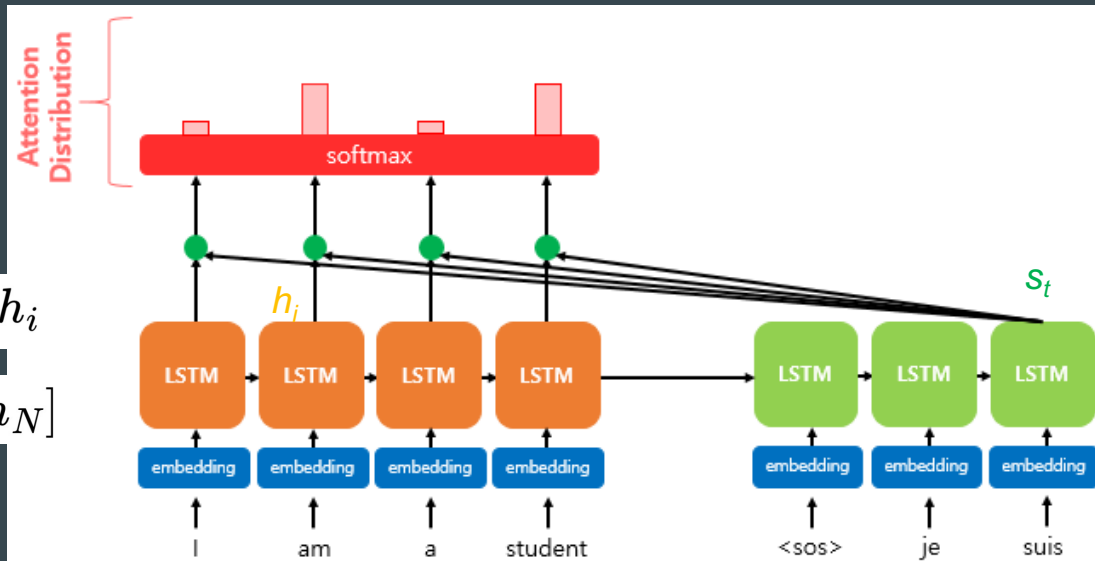
Attention score \rightarrow Softmax



$$\text{score}(s_t, h_i) = s_t^T h_i$$

$$e^t = [s_t^T h_1, \dots, s_t^T h_N]$$

$$\alpha^t = \text{softmax}(e^t)$$



Attention

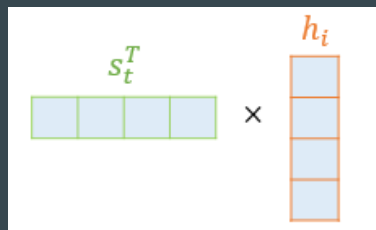
- Dot product attention

Attention value,
context vector

s_t

$$a_t = \sum_{i=1}^N \alpha_i^t h_i$$

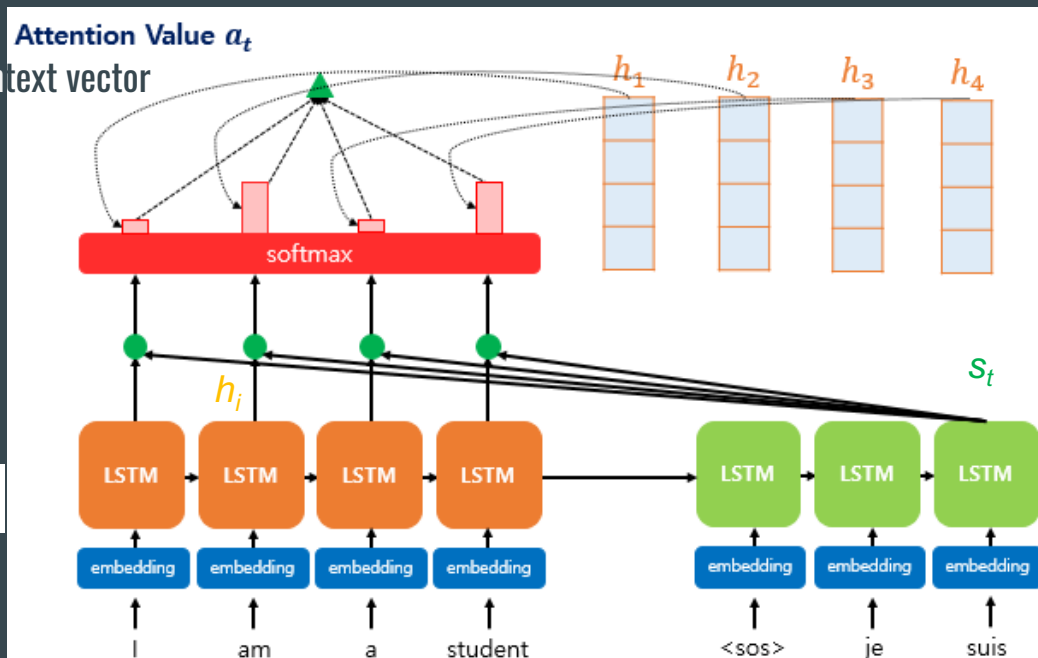
Attention score \rightarrow Softmax \rightarrow Context vector



$$\text{score}(s_t, h_i) = s_t^T h_i$$

$$e^t = [s_t^T h_1, \dots, s_t^T h_N]$$

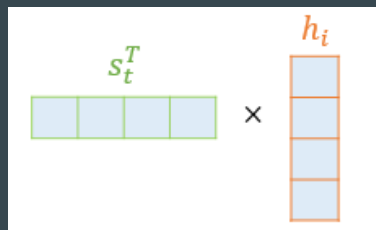
$$\alpha^t = \text{softmax}(e^t)$$



Attention

- Dot product attention

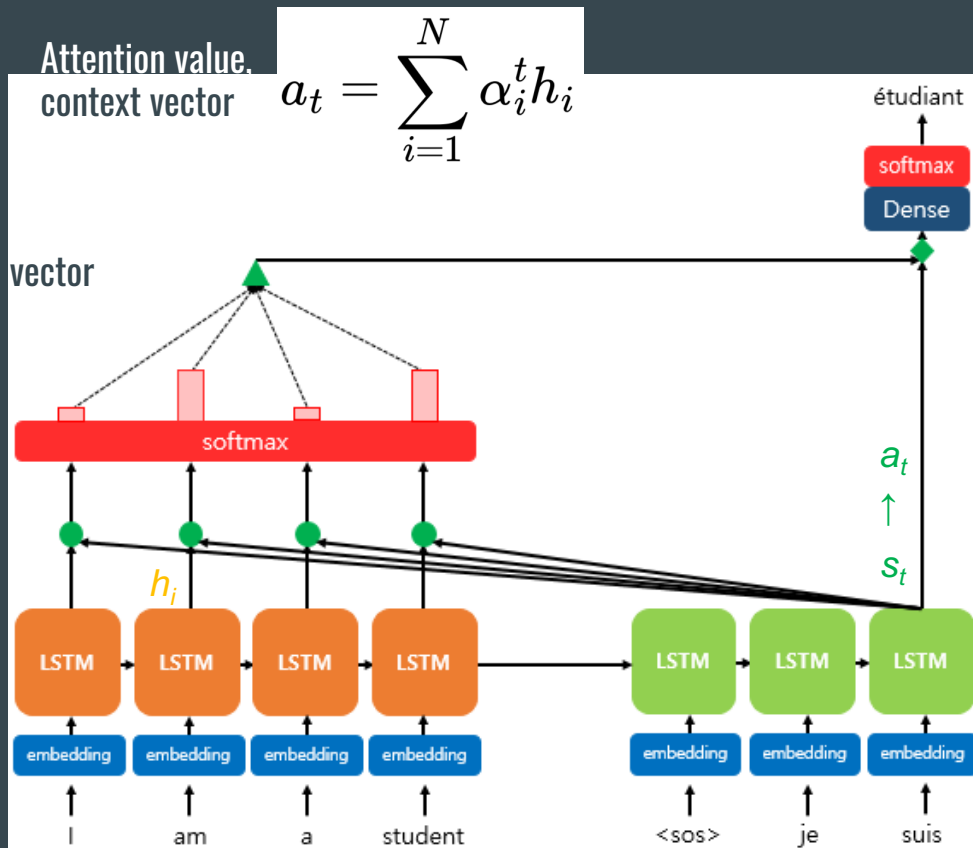
Attention score \rightarrow Softmax \rightarrow Context vector



$$\text{score}(s_t, h_i) = s_t^T h_i$$

$$e^t = [s_t^T h_1, \dots, s_t^T h_N]$$

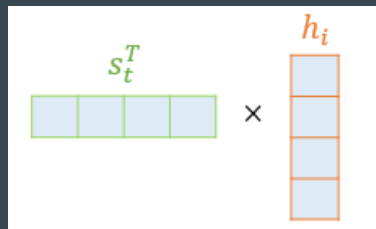
$$\alpha^t = \text{softmax}(e^t)$$



Attention

- Dot product attention

Attention score \rightarrow Softmax \rightarrow Context vector \rightarrow Predict

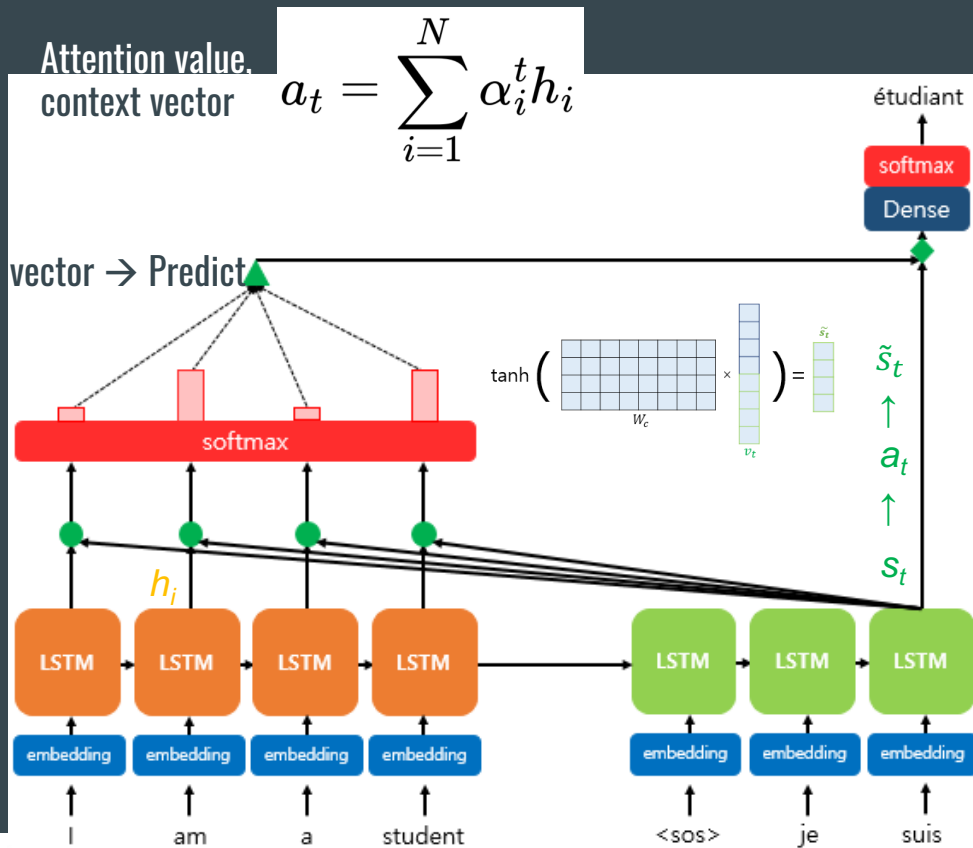


$$\text{score}(s_t, h_i) = s_t^T h_i$$

$$e^t = [s_t^T h_1, \dots, s_t^T h_N]$$

$$\alpha^t = \text{softmax}(e^t)$$

$$\tilde{s}_t = \tanh(\mathbf{W}_c [a_t; s_t] + b_c)$$



Everything can be found in

<https://wikidocs.net/book/2155>

<https://github.com/ukairia777/tensorflow-nlp-tutorial>