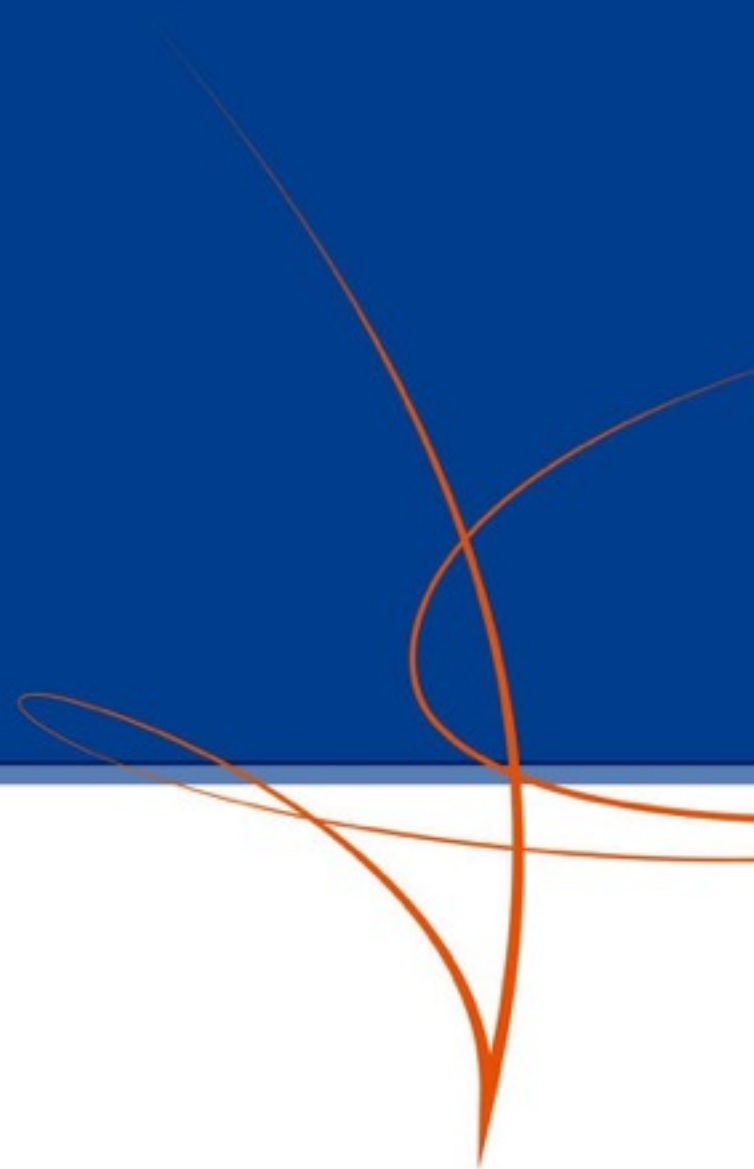


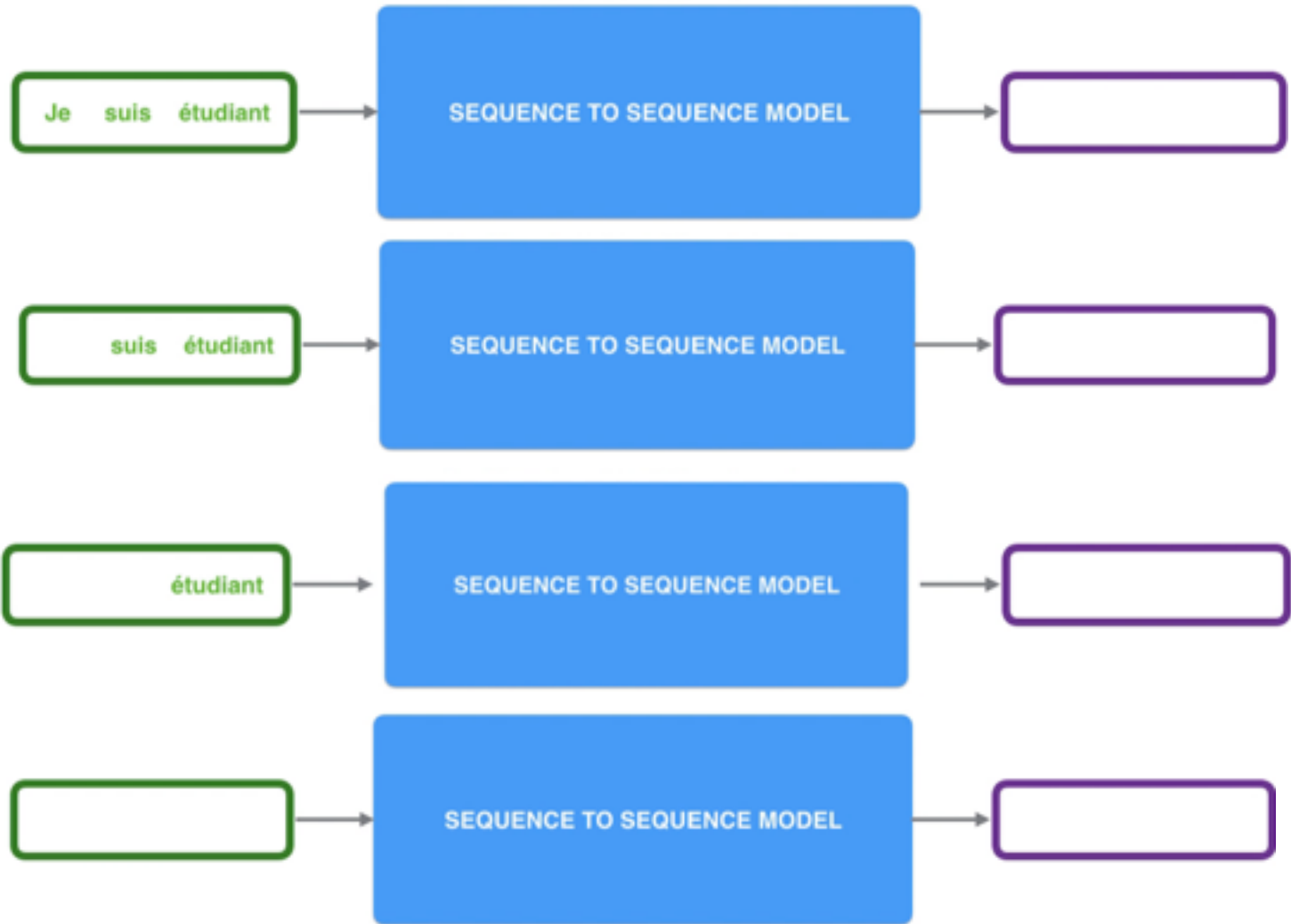
Sequence to Sequence Model



- Sequence to Sequence Model Overview
- Encoder-Decoder
- Attention Mechanism

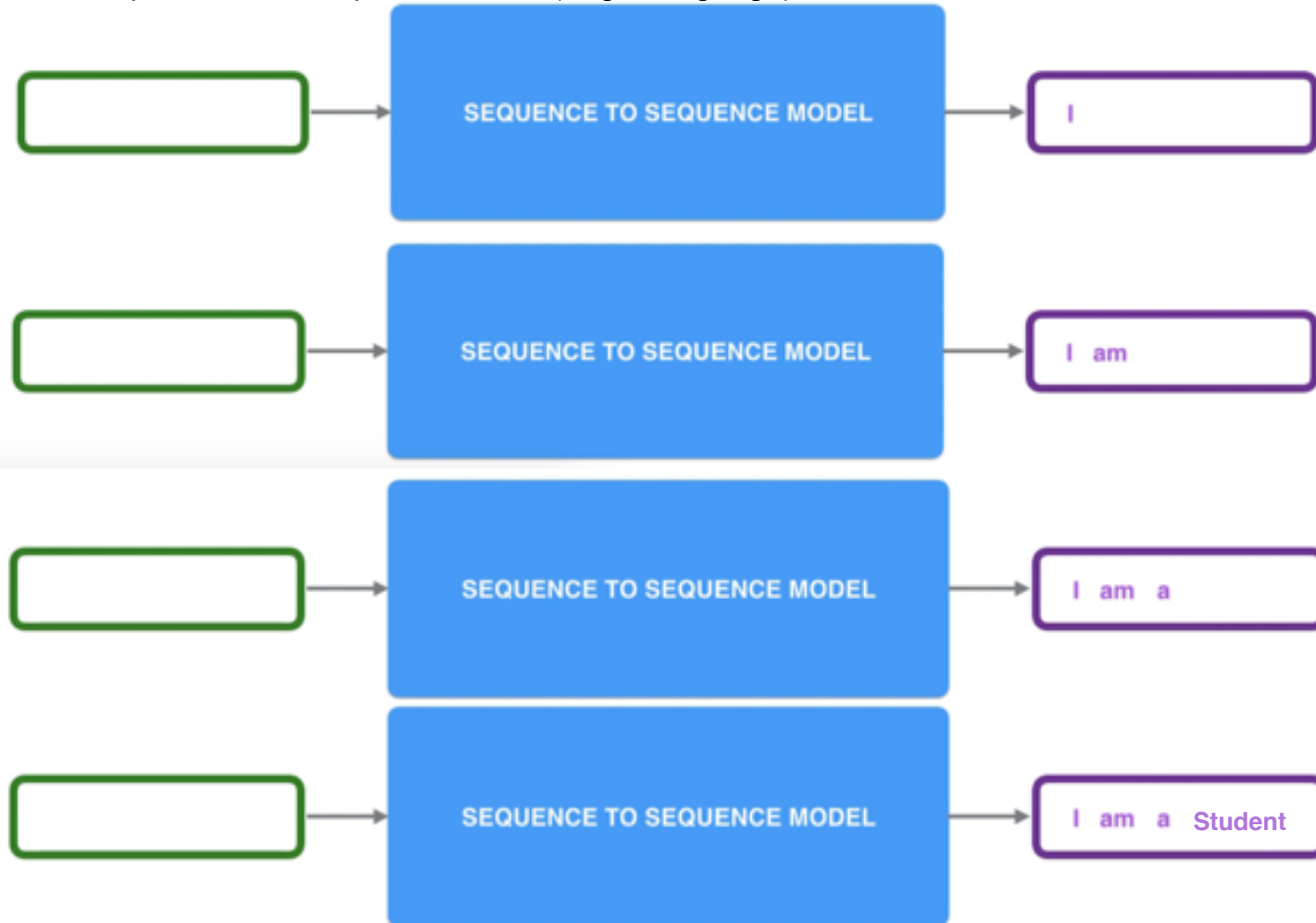
Sequence to Sequence Model Overview

Sequence and Sequence model is an interesting NLP model that can do **mapping** between input text and output text. Specifically, the input can be a **sequence of text** (source language).



Sequence to Sequence Model Overview

The output is also a sequence of text (target language) in machine translation use case.

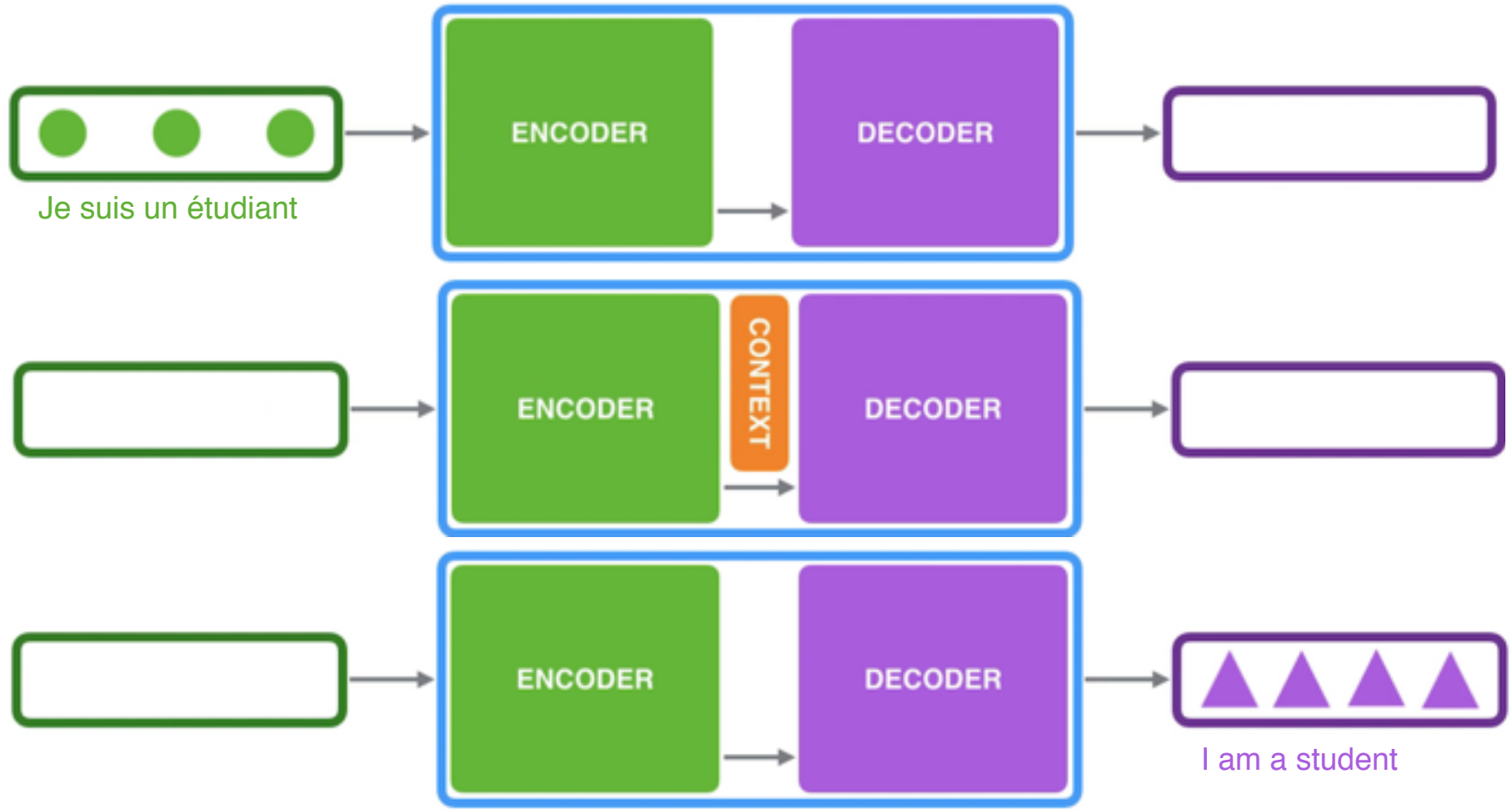


Sequence to Sequence Model Structure: Encoder-Decoder

Encoder: Handles **each element** of text sequence (a word) can transform them into **Context Vector**

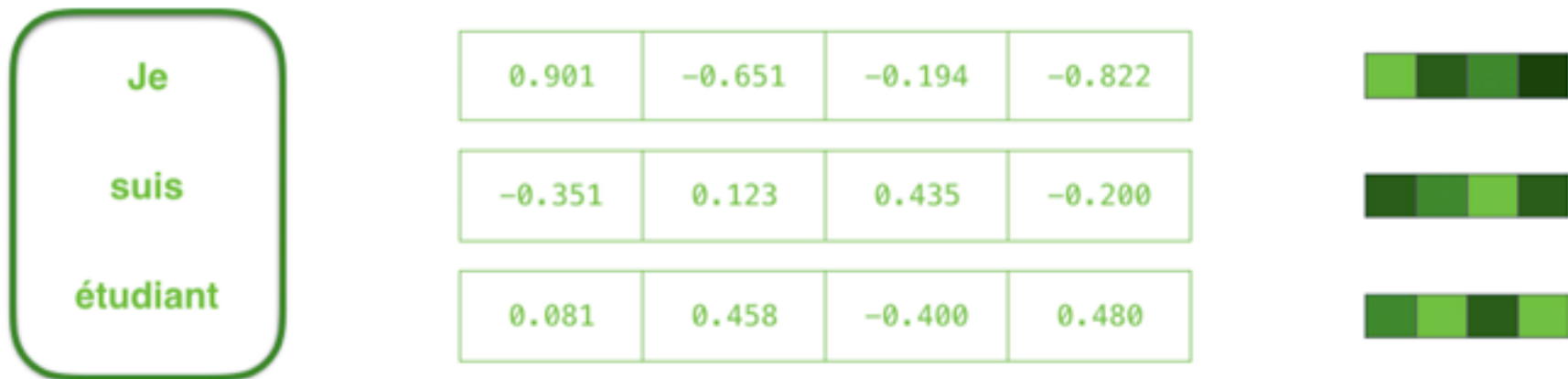
Decoder: Generates **output text sequence** based on Context Vector

SEQUENCE TO SEQUENCE MODEL



Input: Word Embedding

Word Embedding: Normally, the input text sequence is encoded as a **dense vector**



Context Vector

Context Vector: Generated by **Encoder** and will be feed into **Decoder**

CONTEXT

0.11
0.03
0.81
-0.62

0.11
0.03
0.81
-0.62

Implementation of Encoder and Decoder: Recurrent Neural Network

Recurrent Neural Network: The **input vector** of each timestamp will be processed with the **previous hidden state** and then generate the corresponding **hidden state** and output.

Time step #1:

An RNN takes two input vectors:



hidden state #0



input vector #1

Processes them

Then produces two output vectors:



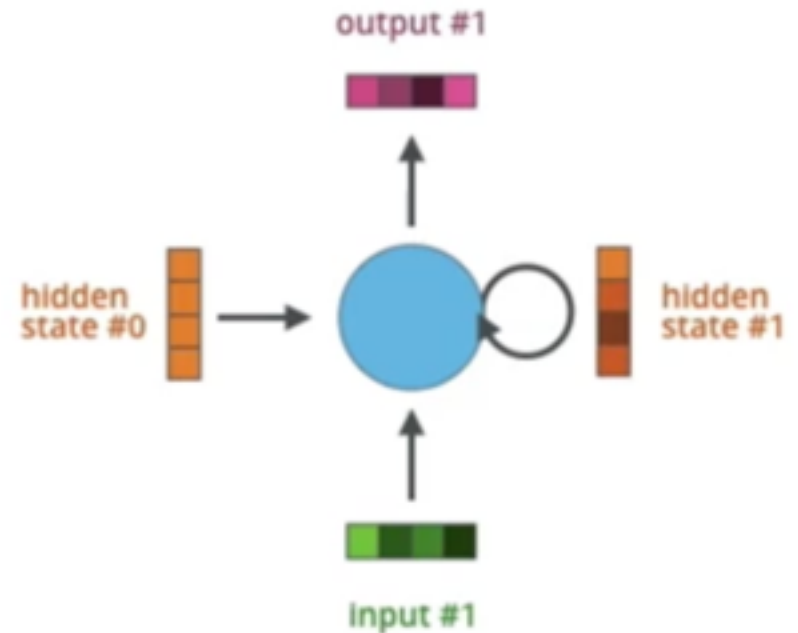
hidden state #1



output vector #1

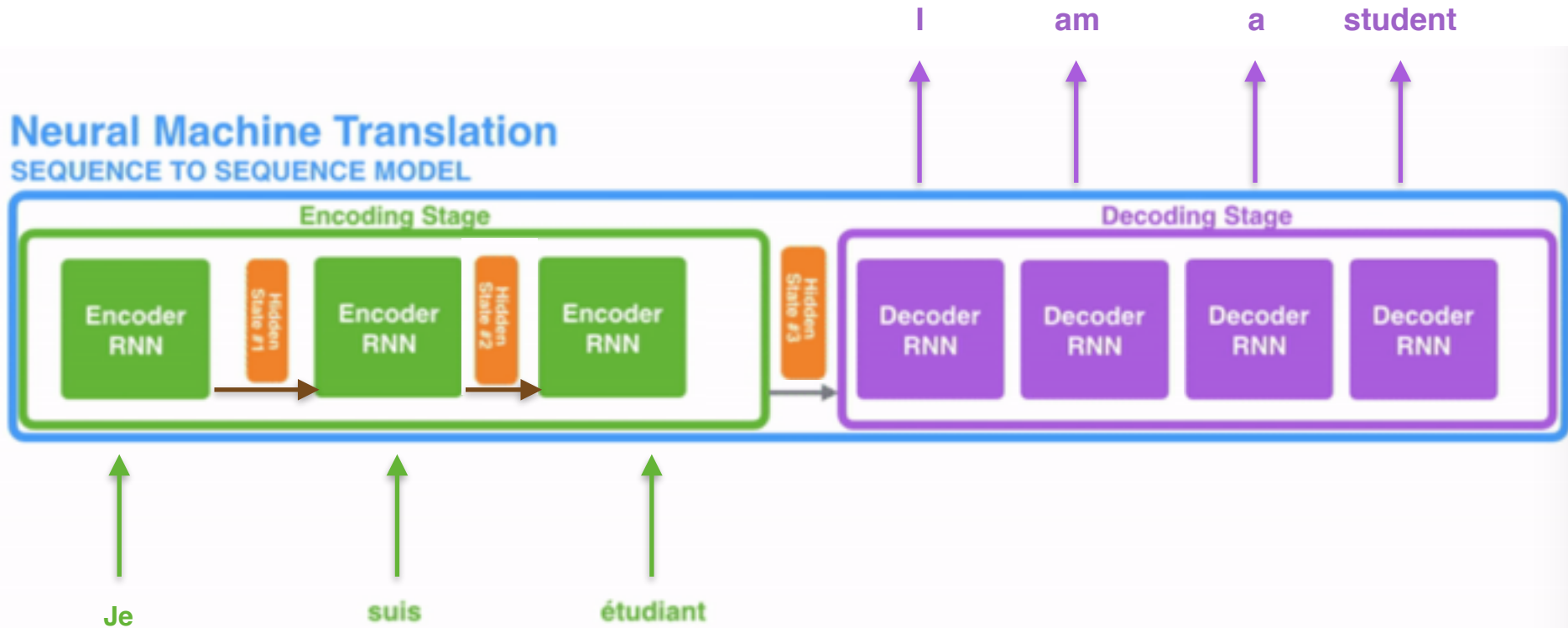
$$S_t = f(UX_t + WS_{t-1})$$

$$O_t = \text{softmax}(VS_t)$$

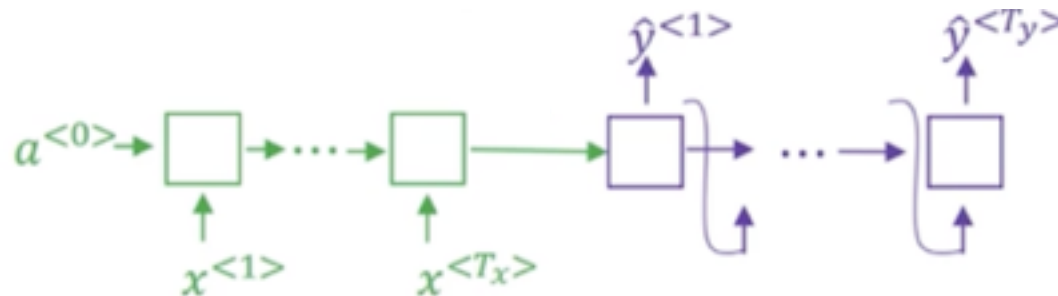


RNN in Encoder and Decoder

Neural Machine Translation SEQUENCE TO SEQUENCE MODEL



Problems of long sequences

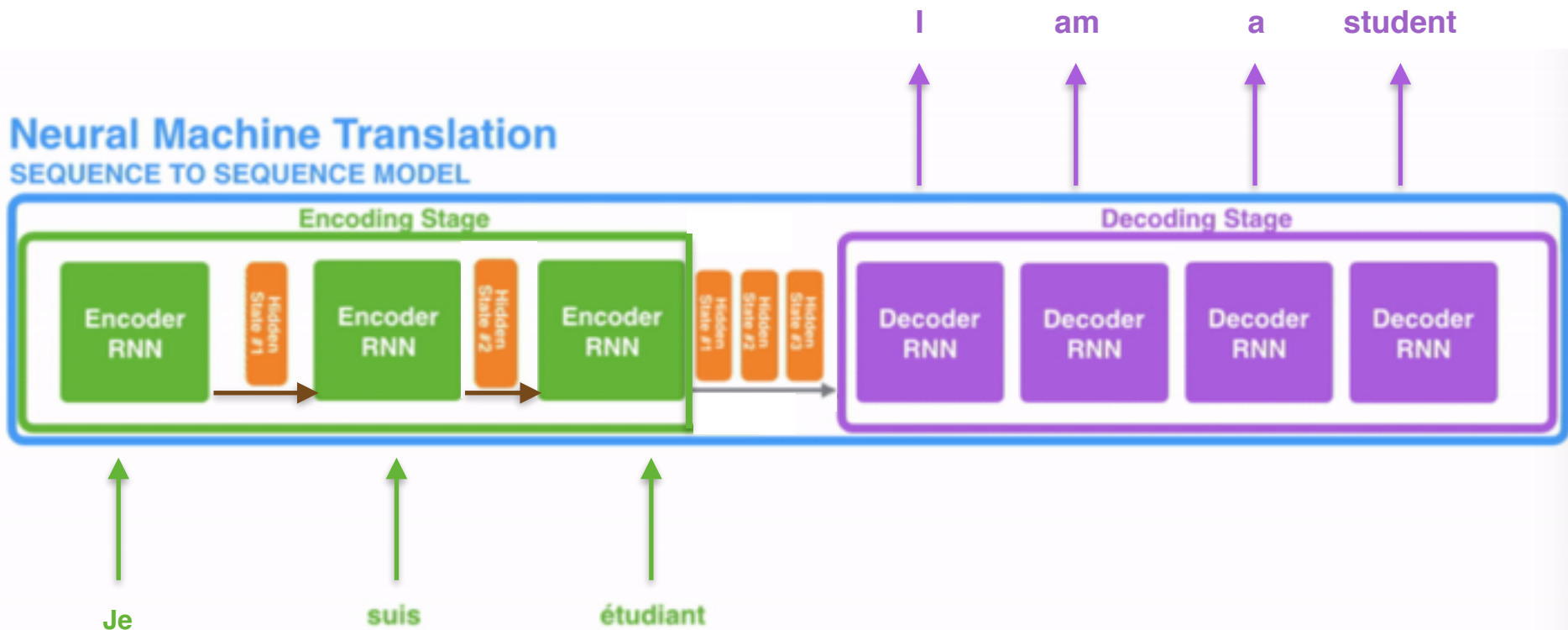


Jane s'est rendue en Afrique en septembre dernier, a apprécié la culture et a rencontré beaucoup de gens merveilleux; elle est revenue en parlant comment son voyage était merveilleux, et elle me tente d'y aller aussi.

Jane went to Africa last September, and enjoyed the culture and met many wonderful people; she came back raving about how wonderful her trip was, and is tempting me to go too.

Solution: Attention in Encoder and Decoder

To learn the alignment between the source language and the targeted language, **Attention** is applied in Sequence and Sequence model instead of using the context information that is put into a single context vector.

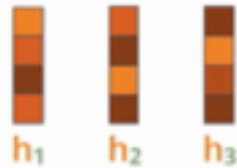


Attention Mechanism

Attention mechanism can be understood as a **weighted approach (softmax)** based on the **importance** of input. Specifically, it sets **different weights** for the **hidden states** generated by **Encoder**.

Attention at time step 4

1. Prepare inputs

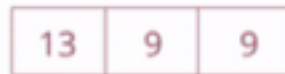


Encoder
hidden
states



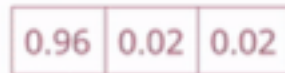
Decoder hidden
state at time step 4

2. Score each hidden state



scores
Attention weights for
decoder time step #4

3. Softmax the scores



softmax scores

4. Multiply each vector by
its softmaxed score



=

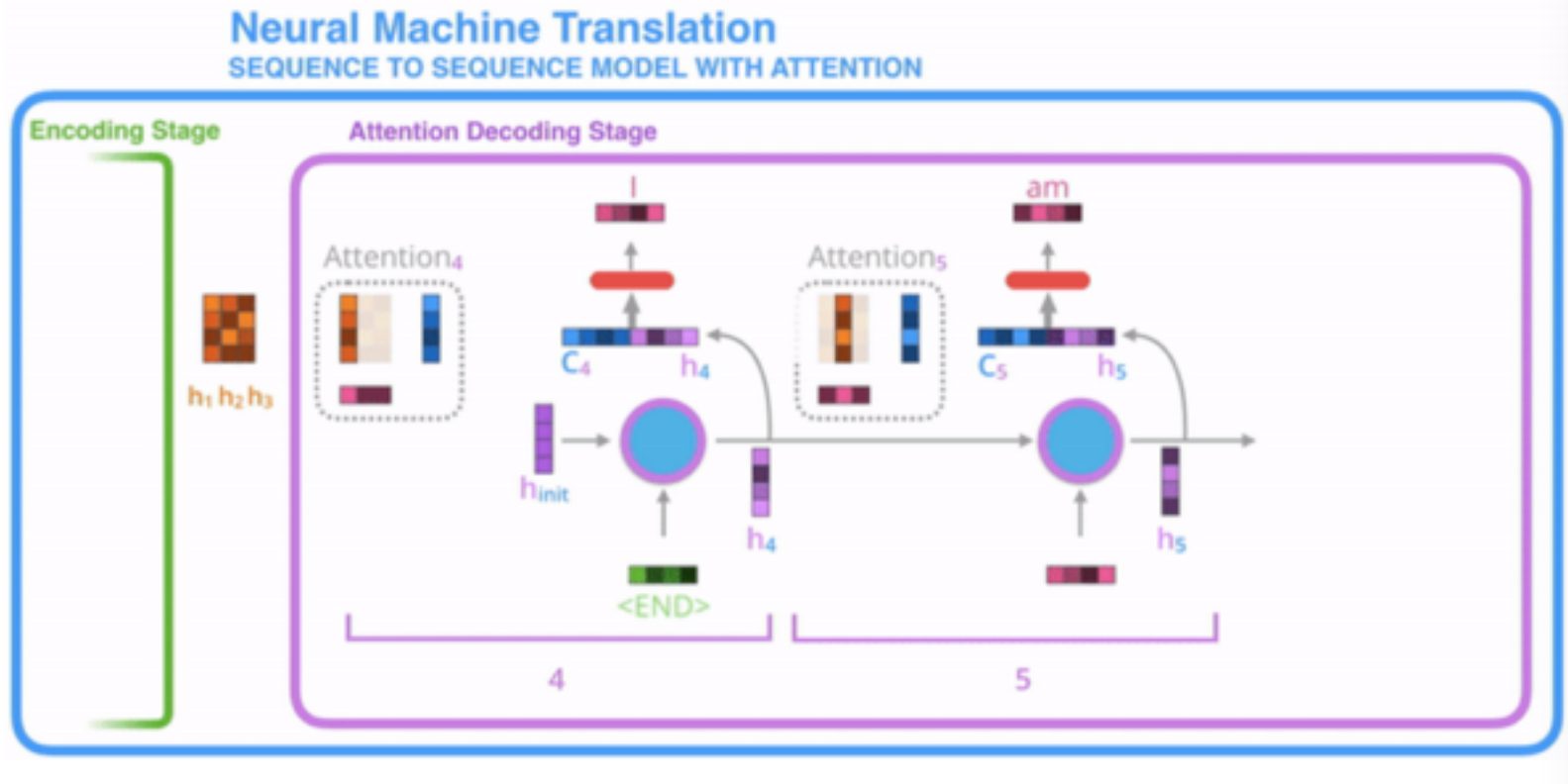
5. Sum up the weighted
vectors



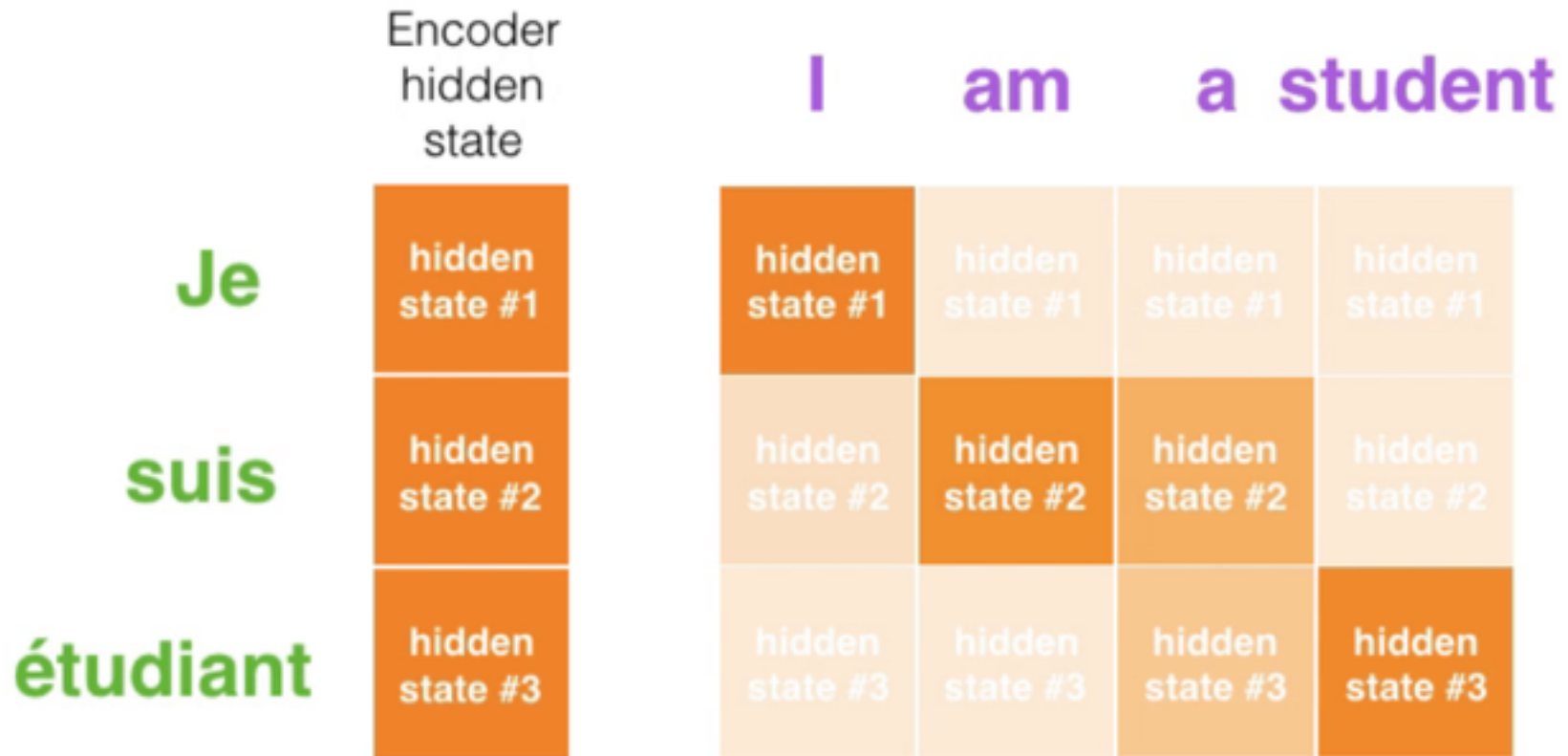
Context vector for
decoder time step #4

Attention Mechanism

1. RNN Decoder with attention receives **word embedding** and a **init hidden state** as **input**.
2. RNN cell handles input, generates output(discard) and a **new hidden state** (h_4).
3. Then **attention** mechanism is applied to **calculate context vector (c_4)** using **Encoder's hidden states and the new hidden state(h_4)**.
4. **Concatenates** context vector(c_4) with the new hidden state (h_4).
5. Finished decoder step with the concatenated vector as input using **fully connected layer (softmax)**.
6. **Repeats** the above steps for each RNN cell

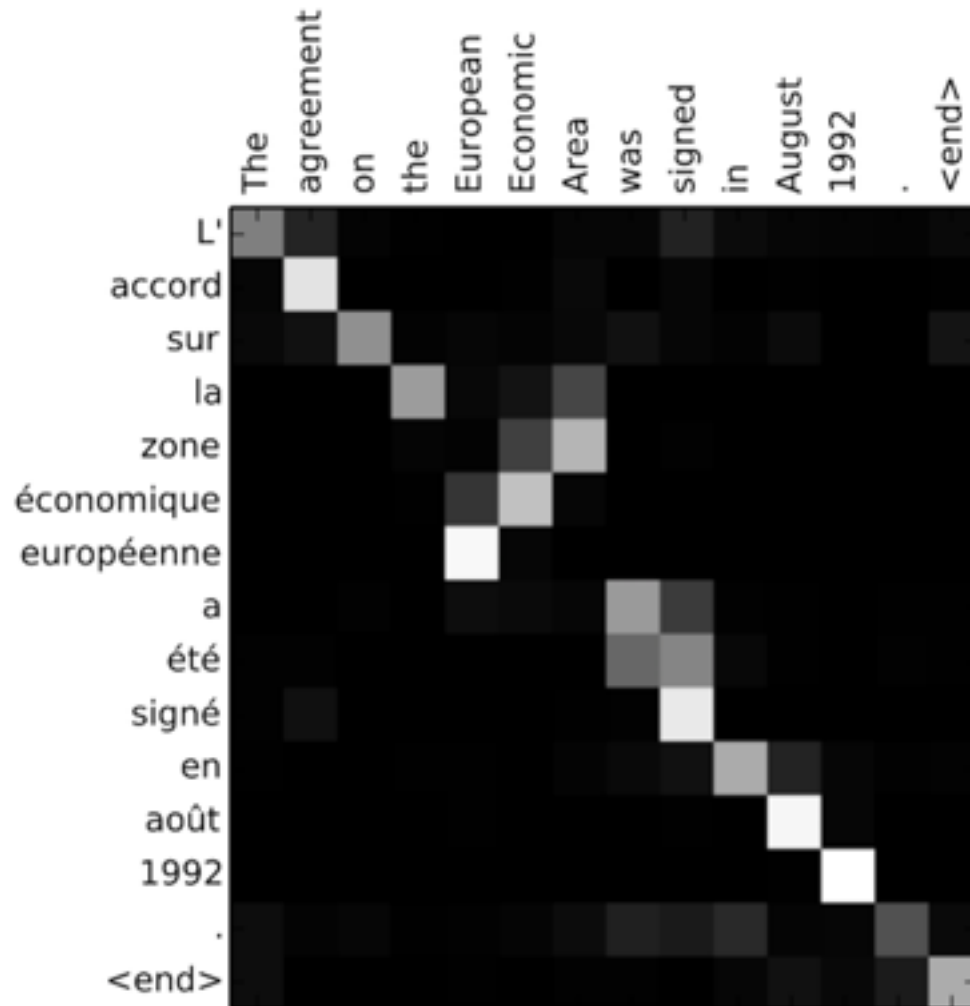


Attention Mechanism



Attention Mechanism

Attention mechanism can be used to learn the alignment between the source language and the targeted language



Attention Mechanism Details

Input: $x=(x_1,...,x_T)$

Output: $y=(y_1,...,y_T)$

(1) $h_t=RNN_{enc}(x_t,h_{t-1})$: each RNN cell in Encoder receives **word embedding** of each input word and the **hidden state generated by the previous cell**. The output is the **hidden state of this current cell**.

(2) $s_t=RNN_{dec}(y_{t-1},s_{t-1})$: each RNN cell in Decoder receives **word embedding** of each targeted word and **hidden state generated by the previous cell**.

(3) $c_t = \sum_{j=1}^{T_x} \alpha_{tj} h_j$: **Context vector** is a **weighted** average of hidden states generated by Encoder.

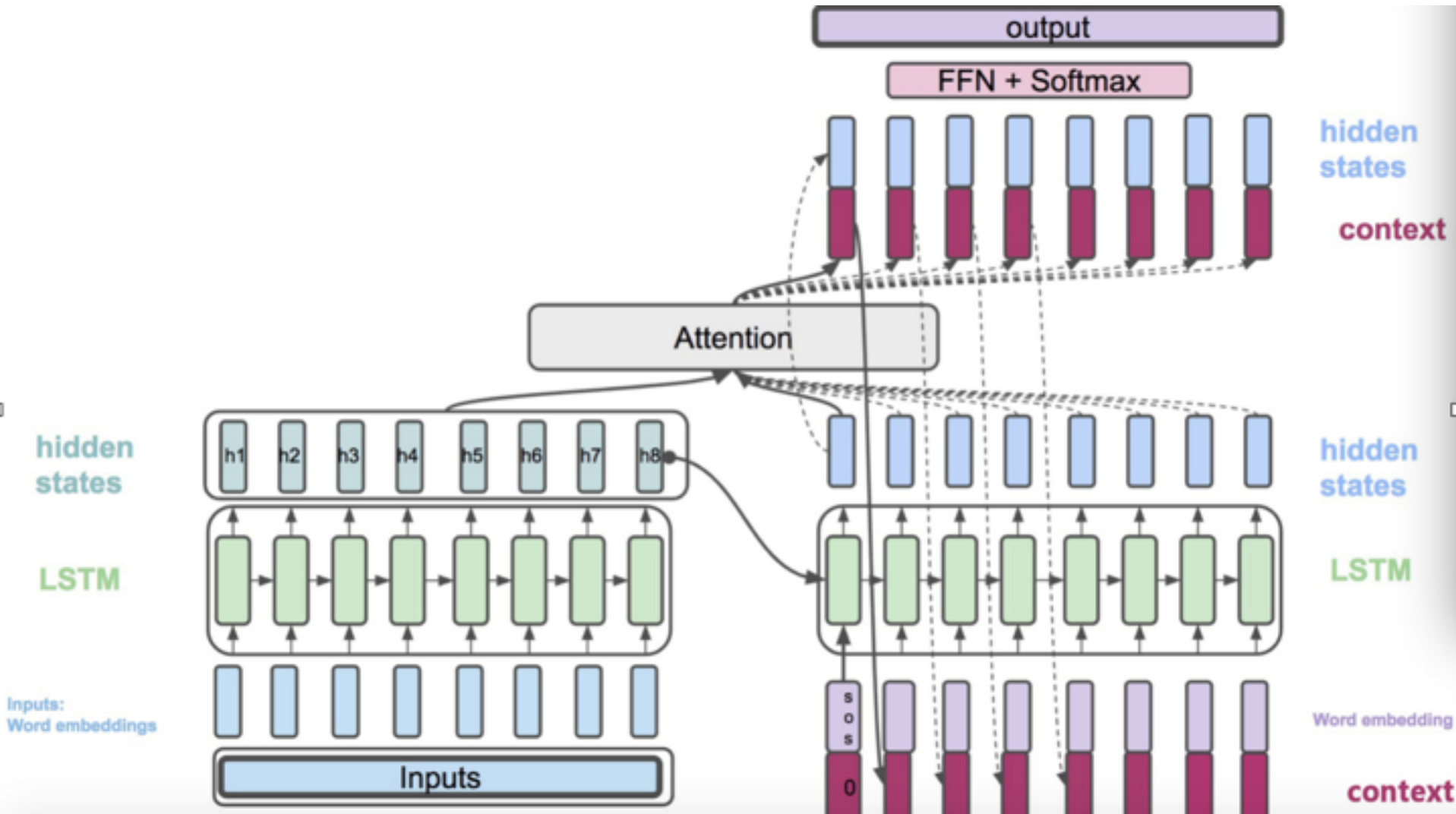
(4) $\alpha_{tj} = \frac{\exp(e_{tj})}{\sum_{k=1}^{T_x} \exp(e_{tk})}$: The **weight** of each hidden state generated by Encoder.

(5) $e_{tj} = \text{score}(s_t, h_j)$: The **score** is calculated by **hidden states generated by both Encoder and Decoder**. It is used for the weight calculation in step (4).

(6) $\hat{s}_t = \tanh(W_c[c_t; s_t])$: **Concatenates** context vector with the hidden state generated by Decoder.

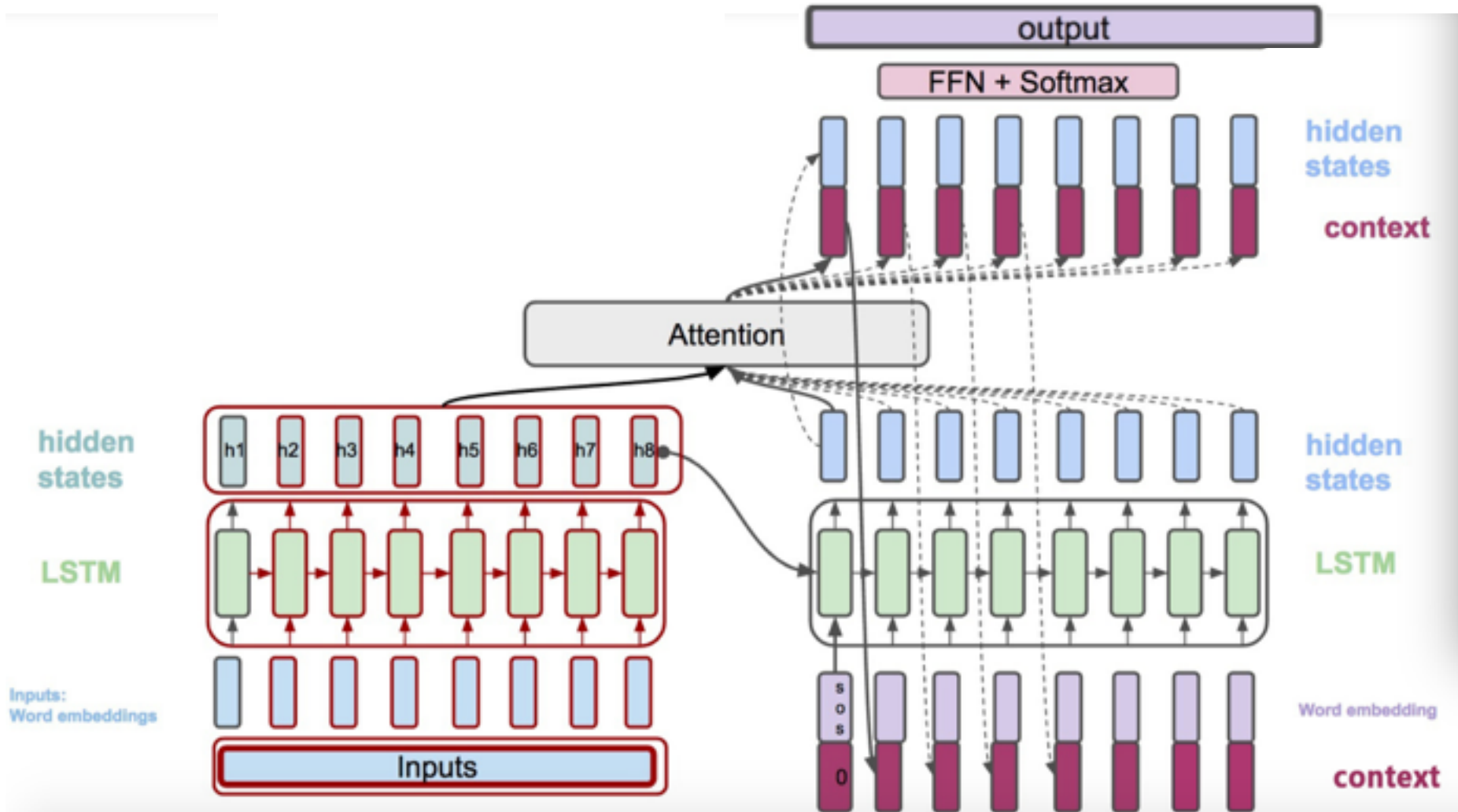
(7) $p(y_t|y_{<t}, x) = \text{softmax}(W_s \hat{s}_t)$: Calculates the final output probability of each word in dictionary.

Attention Mechanism Details



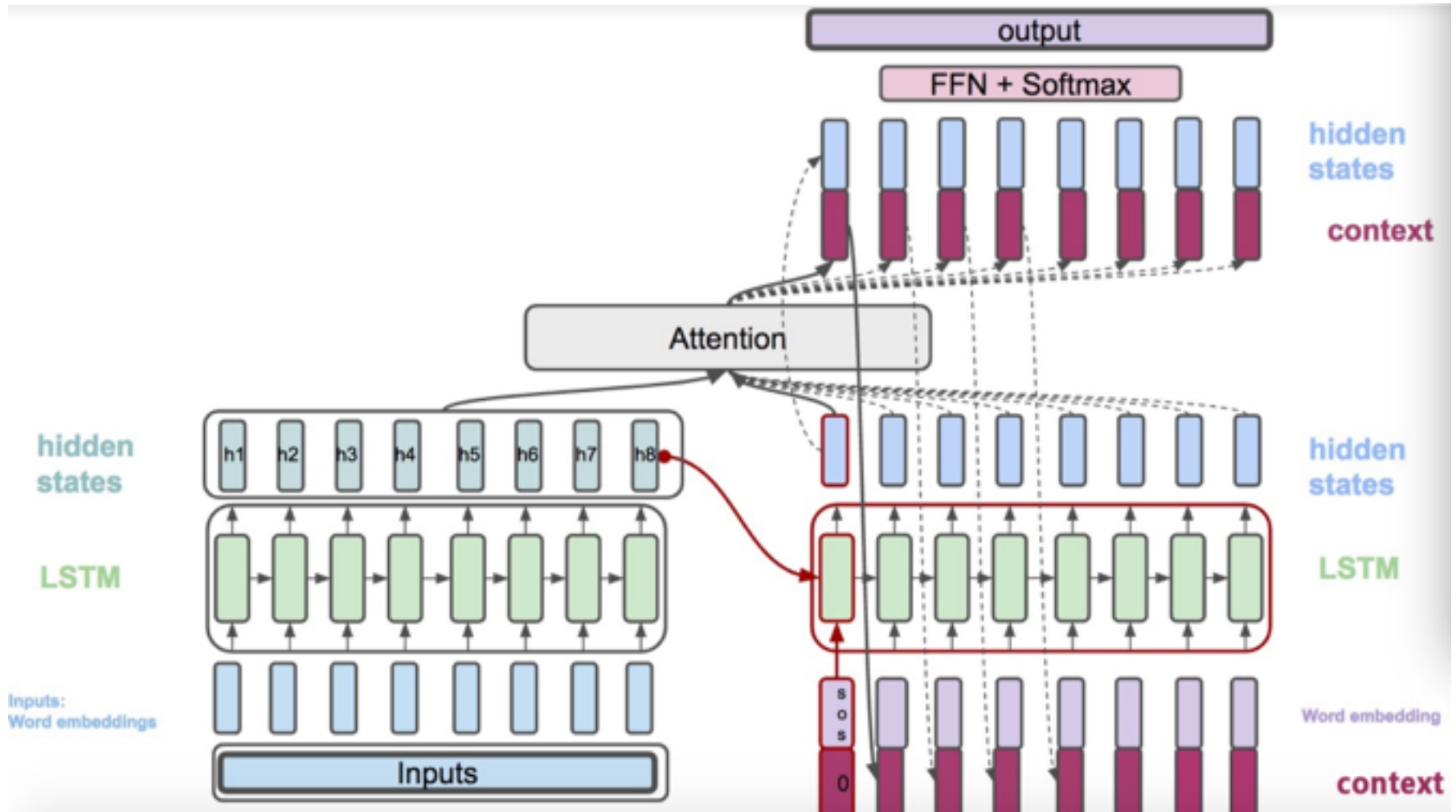
Attention Mechanism Details

1) $h_t = RNN_{enc}(x_t, h_{t-1})$: each RNN cell in Encoder receives **word embedding** of each input word and the **hidden state generated by the previous cell**. The output is the **hidden state of this current cell**.



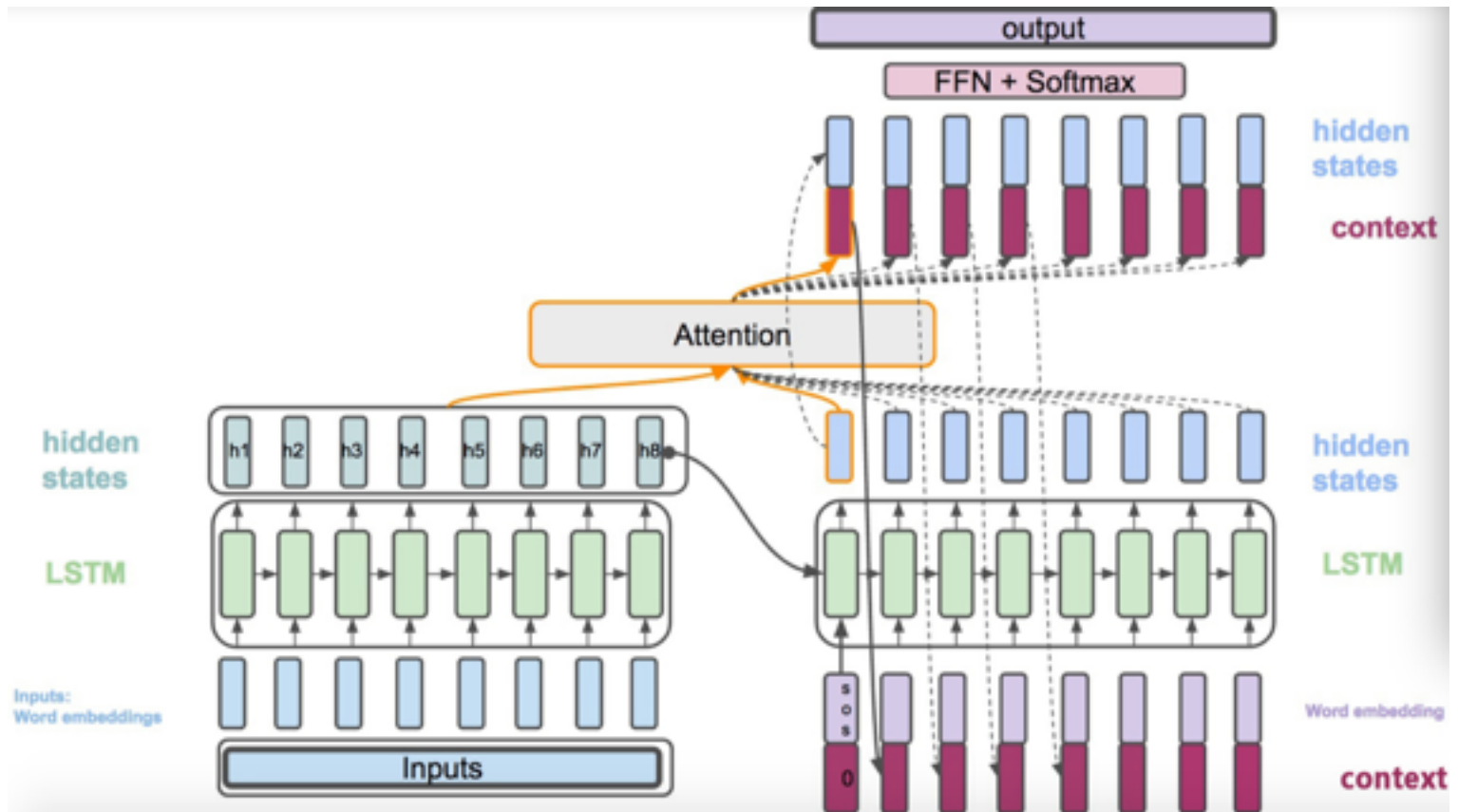
Attention Mechanism Details

(2) $st = RNN_{dec}(y_{t-1}^{\wedge}, st-1)$: each RNN cell in Decoder receives **word embedding** of each targeted word and the **hidden state generated by the previous cell**.



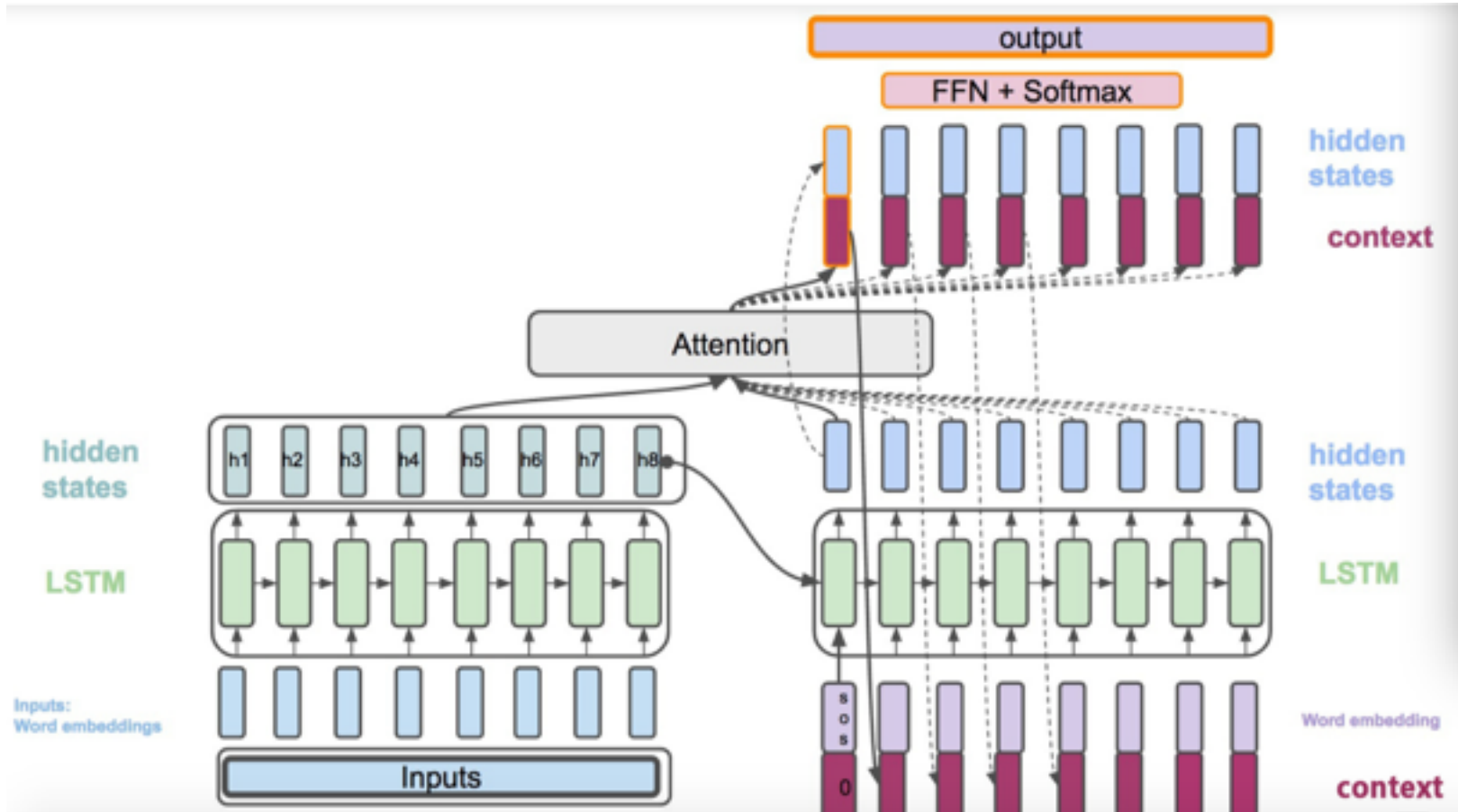
Attention Mechanism Details

- (3) $c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$: **Context vector** is a **weighted** average of hidden states generated by Encoder.
- (4) $\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$: The **weight** of each hidden state generated by Encoder.
- (5) $e_{ij} = \text{score}(s_i, h_j)$: The **score** is calculated by **hidden states generated by both Encoder and Decoder**. It is used for the weight calculation in step (4).



Attention Mechanism Details

- (6) $\hat{s}_t = \tanh(W_c[c_t; s_t])$: **Concatenates** context vector with the hidden state generated by Decoder.
- (7) $p(y_t|y_{<t}, x) = \text{softmax}(W_s \hat{s}_t)$: Calculates the final output probability of each word in dictionary.



$$\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s) = \begin{cases} \mathbf{h}_t^\top \bar{\mathbf{h}}_s & \text{dot} \\ \mathbf{h}_t^\top \mathbf{W}_a \bar{\mathbf{h}}_s & \text{general} \\ \mathbf{W}_a[\mathbf{h}_t; \bar{\mathbf{h}}_s] & \text{concat} \end{cases}$$

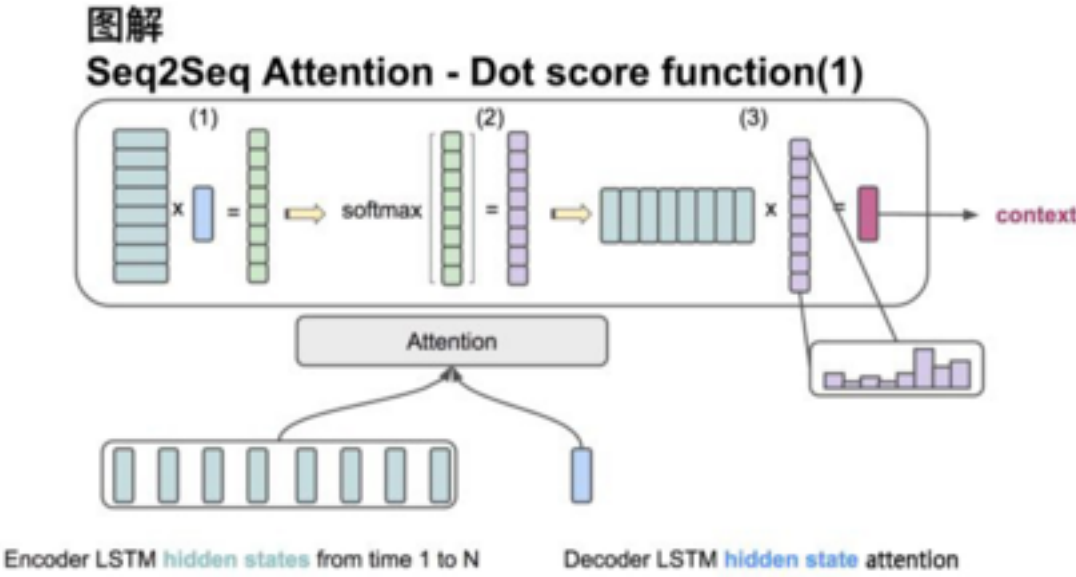
Attention Mechanism: Score Calculation

Approach 1

Input: The dimension of matrix represents **all hidden states generated by Encoder (H)**: (hidden dimension, input sequence length)

The dimension of **hidden state generated by decoder at any time point (S)**: (hidden dimension, 1)

- 1) Get score matrix (sequence length, 1) by transposing H to (sequence length, hidden dimension) and conducting **dot** with S
- 2) Applying **softmax** to these scores to normalize the sum of weighted scores = 1
- 3) Generating Context Vector by conducting Dot to H and the weighted scores of step (2)



Sequence to Sequence Model Use Cases

- Chatbot
- Machine Translation
- Auto Summarization