# Attention Models

## Sargur N. Srihari
srihari@buffalo.edu

This is part of lecture slides on Deep Learning:
http://www.cedar.buffalo.edu/~srihari/CSE676

# Topics in NLP

1. N-gram Models
2. Neural Language Models
3. High-Dimensional Outputs
4. Combining Neural Language Models with n-grams
5. Neural Machine Translation
6. Attention Models
7. Transformer Models
8. Historical Perspective

# Topics in Attention Models

- Attention Vector
- General and Self Attention
- Encoder-Decoder without Attention
- Types of Attention Models
- Encoder-Decoder with Attention
- Neural network for Self-Attention

# What is Attention?

- ## Can I have your Attention please!

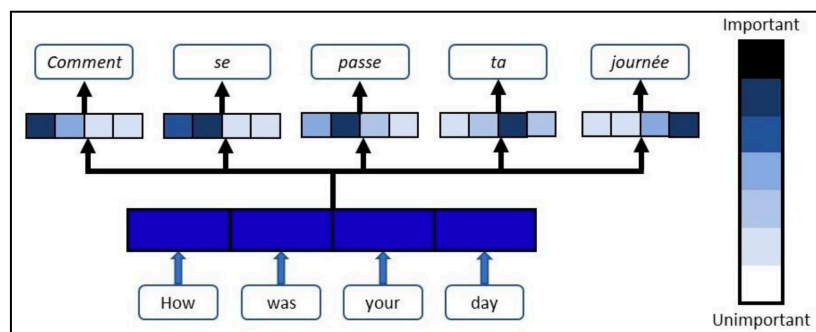| ITALIAN | GERMAN | FRENCH | KANNADA | HINDI | TAMIL |

L' attention

ಗಮನ
"gamana"

ध्यान
"dhyan"

கவனம்
"kavanam"

- ## It means to focus on something and take notice

- ## Paying greater attention to certain factors when processing the input

# General and Self-Attention

- The attention **component** of a network manages and quantifies the **interdependence**:
  - General Attention: Between input and output elements
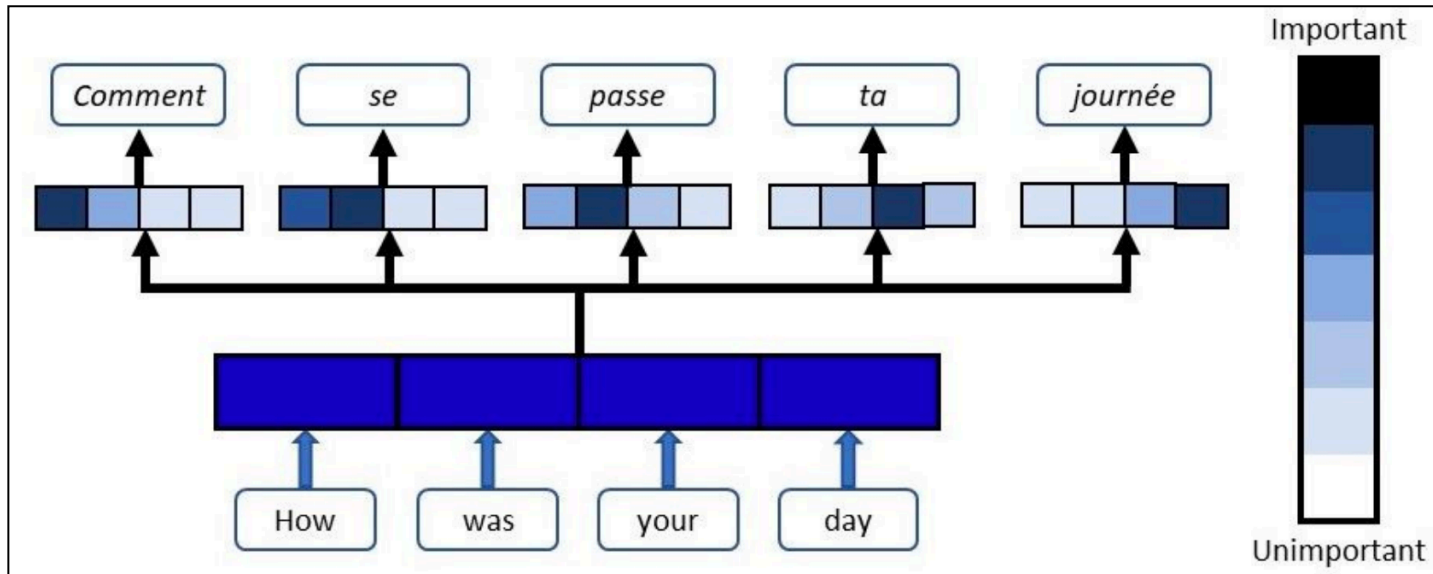


  - Self Attention: Within input elements



5

# Attention Vectors in NMT

- ## In Neural Machine Translation, for each output word
  Assign higher weights to relevant words in input
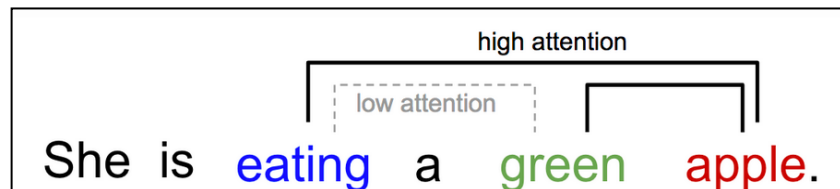- ## English to French Translation
  How was your day
  Comment se passe ta journee

6

# Self Attention

- ## How words correlate words in a sentence



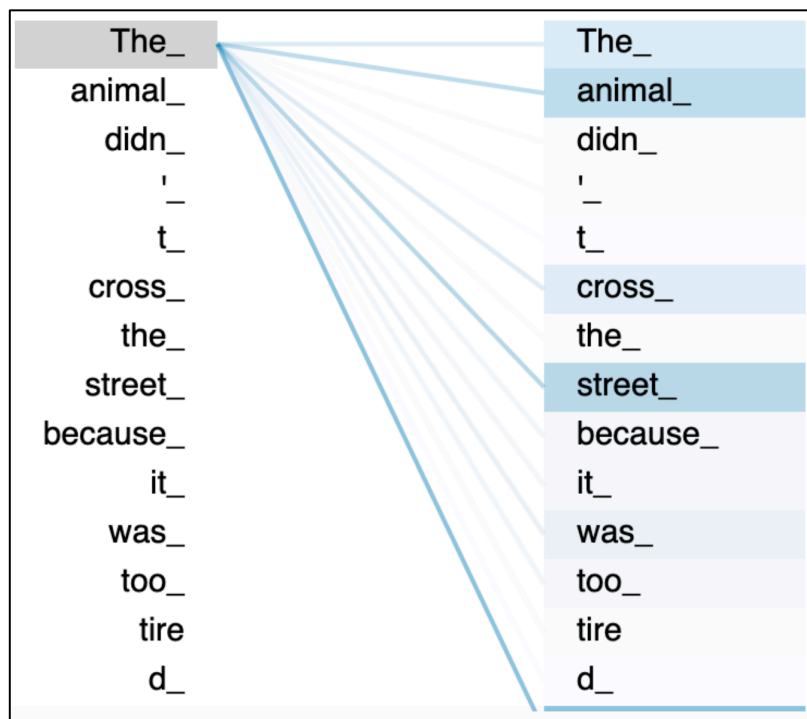https://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html

  – When we see "eating", we expect a food word soon
    - "Green" describes food, but next word depends more more  on "eating"
  – Attention is a vector of importance weights
    - To predict a word in a sentence, we estimate using the attention vector  how strongly it is correlated with (or "*attends to*") other elements

- ## This is a case of self-attention
  – Attention within different parts of same input

# Self Attention Vector

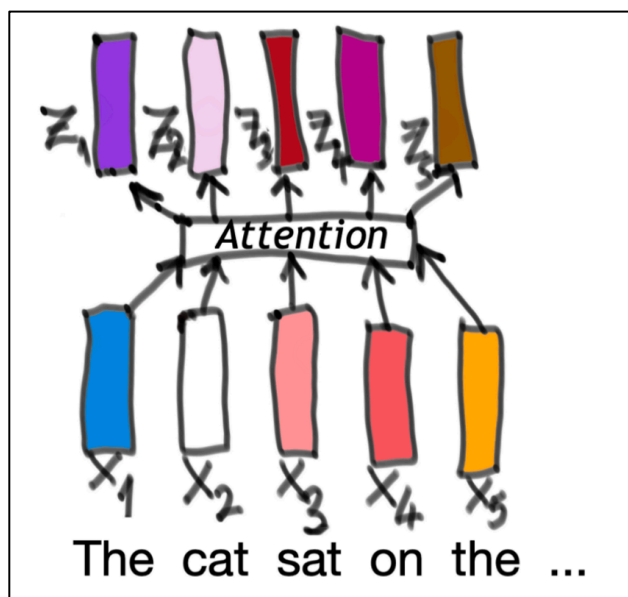- ## We want ML system to learn relationships between words in input sentence



Color coding shows that that there is some connection between *animal*,
*cross, street* and *the*
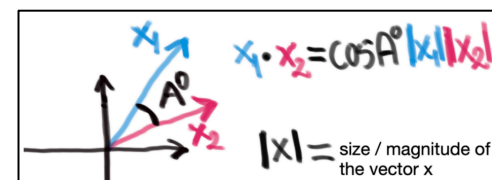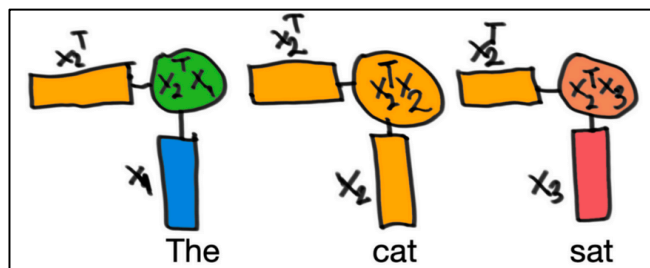because they're all related to *animal*,
the subject of the sentence.

# The Self-Attention Process

## Self-Attention is a process

A sequence of vectors $X$ is **encoded** into another sequence of vectors $Z$

Its corresponding $Z$ vector represents both the original word *and* its

**relationship** with the other words around it



Cosine similarity for $x_2$ considering only $x_1 x_2 x_3$
Multiply each by $x_2$



Raw weights :

$$x_2^T x_1 = w_{21}$$
$$x_2^T x_2 = w_{22}$$
$$x_2^T x_3 = w_{23}$$

Softmax weights

Compute $z_2$

https://towardsdatascience.com/self-attention-5b95ea164f61

# Multiple attention vectors

# We discuss five Attention Models

1. Basic encoder-decoder RNN model

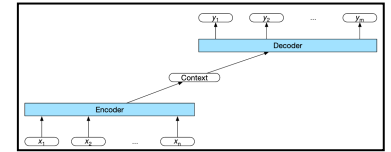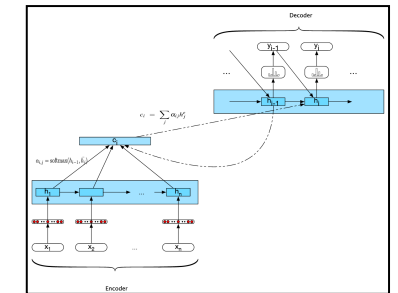   – Without attention

2. Encoder-decoder RNN with Attention

3. Transformer model with no recurrence

4. Generative pre-training of Attention model

5. BERT with no decoder, using bidirectionality

# Encoder-decoder without attention

# Roles of Encoder-decoder in NMT



**Encoder** compresses input sequence into a context vector

It is a "sentence embedding" or "thought" vector of a *fixed length*.

It is expected to be a good summary of the meaning of the *whole* sequence

**Decoder** is initialized with the context vector to emit the transformed output

# Basic RNN-based Encoder-Decoder



- Notation: At position $t$, $h_t^e$ and $h_t^d$ are the hidden states of the encoder and decoder RNN respectively

- Final encoder hidden state $h_n^e$ serves as context for $h_1^d$ in decoder

- Both the encoder and decoder are RNNs
  e.g.,  LSTM or GRU units

14

Source:  Jurafsky and Martin, Speech and Language Processing, 2020

# Sequence-to-sequence RNN model

- Encoder & decoder have the same structure (RNN)

  - Encoder encodes input sequence $\{x_1, x_2 \ldots x_n\}$ into fixed length vectors $\{h_1^e, h_2^e \ldots h_n^e\}$

  - Decoder takes single fixed length vector $\{h_n^e\}$ as input and generates output $\{y_1, y_2 \ldots y_m\}$, token by token

    - Uses Autoregression

      - Consumes previous output



15

# Summary of Encoder-Decoder



## Encoder

Simple RNNs, LSTMs, GRUs, CNNs, transformer networks.
Only a single network layer shown.
In a stacked architecture output states from
the top layer are taken as the final representation.
A widely used encoder design makes use of stacked Bi-LSTMs

## Decoder

Autoregressive generation is used to produce an output sequence,
an element at a time, until an end-of-sequence marker is generated.
This incremental process is guided by the context provided by the encoder
as well as any items generated for earlier states by the decoder.
Again, a typical approach is to use an LSTM or GRU-based RNN

Source: Jurafsky and Martin, Speech and Language Processing, 2020

# Deficiencies of traditional encoder-decoder

1. Encoder compresses all input into a single fixed length vector $h_n$ to be passed to decoder
   – Leads to loss of information

2. Unable to model alignment between input and output sequences
   – Which is essential for structured output tasks such as translation or summarization

# Intuition of what is needed

- In sequence-to-sequence tasks, each output token is expected to be more influenced by some specific parts of the input sequence

- However, decoder lacks any mechanism to selectively focus on relevant input tokens while generating each output token

# Types of Attention Models

1. Basic encoder-decoder RNN model
   – Without attention

2. Encoder-decoder RNN with Attention

3. Transformer model with no recurrence

4. Generative pre-training of Attention model

5. BERT with no decoder, using bidirectionality

# Encoder-Decoder without and with attention



Basic Encoder-Decoder

Final encoder hidden state $h_n^e$ serves as context for $h_1^d$ in decoder

Attention model

1. Accesses entire encoded input $\{h_1, h_2 \dots h_n\}$
2. Induces attention weights over input to prioritize positions for generating next output
3. Dynamically updates during decoding
4. Can be embodied in a fixed-size vector

$$c_i = \sum_j \alpha_{ij} h_j^e$$

$$\alpha_{i,j} = \text{softmax}(h_{i-1}, \bar{h}_i)$$

# Encoder-Decoder with Attention



*Decoder hidden state* $h_i^d$ is based on previous hidden state $h_{i-1}^d$, previous word generated $y_{i-1}$, and current context vector $c_i$

$$h_i^d = g(\hat{y}_{i-1}, h_{i-1}^d, c_i)$$

Context vector $c_i$
obtained by comparing the previous hidden state $h_{i-1}^d$ to all of the encoder hidden states $h_j^e$
*Score* is a measure of similarity of the decoder hidden state to each encoder hidden state

Alignment score functions

Dot Product

$$score(h_{i-1}^d, h_j^e) = h_{i-1}^d \cdot h_j^e$$

General Weight Matrix

$$score(h_{i-1}^d, h_j^e) = h_{t-1}^d W_s h_j^e$$

Neural network determines attention weights $\alpha_{ij}$

$$\alpha_{ij} = \text{softmax}(score(h_{i-1}^d, h_j^e) \ \forall j \in e)$$

$$= \frac{exp(score(h_{i-1}^d, h_j^e))}{\sum_k exp(score(h_{i-1}^d, h_k^e))}$$

$$c_i = \sum_j \alpha_{ij} h_j^e$$

$c_i$: Context vector for output $y_i$
$\alpha_{ij}$: How well two words $y_i$ and $x_j$ are aligned

# Bidirectional Encoder



Encoder is a bidirectional RNN with forward state $h_i^{\rightarrow}i$ and backward $h_i^{\leftarrow}$

A simple concatenation of two represents the encoder state.

Motivation:
Include both preceding and following words in annotating one word.
$h_i = [h_i^{\rightarrow \top}; h_i^{\leftarrow \top}]^{\top}$, $i = 1, \ldots, n$

The decoder network has hidden state $s_t = f(s_{t-1}, y_{t-1}, \mathbf{c}_t)$
for the output word at position t, $t = 1, \ldots, m$, where the context vector
$\mathbf{c}_t$ is a sum of hidden states of the input sequence, weighted by alignment scores:

$$\mathbf{c}_t = \sum_{i=1}^{n} \alpha_{t,i} \boldsymbol{h}_i \qquad\qquad ; \text{Context vector for output } y_t$$

$$\alpha_{t,i} = \text{align}(y_t, x_i) \qquad\qquad ; \text{How well two words } y_t \text{ and } x_i \text{ are aligned.}$$

$$= \frac{\exp(\text{score}(\boldsymbol{s}_{t-1}, \boldsymbol{h}_i))}{\sum_{i'=1}^{n} \exp(\text{score}(\boldsymbol{s}_{t-1}, \boldsymbol{h}_{i'}))} \qquad ; \text{Softmax of some predefined alignment score..}$$

22

# Example of Alignment Matrix $\alpha_{ij}$

- Matrix of alignment scores is a byproduct to explicitly show the correlation between source and target words.

- Alignment matrix of French input
  *"L'accord sur la zone économique europénne a été signé en août 1992"*

- Its English translation
  *"The agreement on the European Economic Area was signed in August 1992"*



23

# Alignment score functions

$$\alpha_{ij} = \text{softmax}(score(h_{i-1}^d, h_j^e) \ \forall j \in e)$$

$$= \frac{exp(score(h_{i-1}^d, h_j^e))}{\sum_k exp(score(h_{i-1}^d, h_k^e))}$$

| Name | Alignment score function | Citation |
|---|---|---|
| Content-base attention | $score(s_t, h_i) = \text{cosine}[s_t, h_i]$ | Graves2014 |
| Additive(*) | $score(s_t, h_i) = \mathbf{v}_a^\top \tanh(\mathbf{W}_a[s_t; h_i])$ | Bahdanau2015 |
| Location-Base | $\alpha_{t,i} = \text{softmax}(\mathbf{W}_a s_t)$ <br> Note: This simplifies the softmax alignment to only depend on the target position. | Luong2015 |
| General | $score(s_t, h_i) = s_t^\top \mathbf{W}_a h_i$ <br> where $\mathbf{W}_a$ is a trainable weight matrix in the attention layer. | Luong2015 |
| Dot-Product | $score(s_t, h_i) = s_t^\top h_i$ | Luong2015 |
| Scaled Dot-Product(^) | $score(s_t, h_i) = \frac{s_t^\top h_i}{\sqrt{n}}$ <br> Note: very similar to the dot-product attention except for a scaling factor; where n is the dimension of the source hidden state. | Vaswani2017 |

# Self-Attention

- ## Self-attention, also known as *intra-attention*

  - Relates different positions of a single sequence to compute a representation of the same sequence

  - As opposed to attention over entire input space



high attention

low attention

She is *eating* a *green* *apple*.

- Useful in machine reading, text summarization, image description generation

- ## Broad categories of Attention:

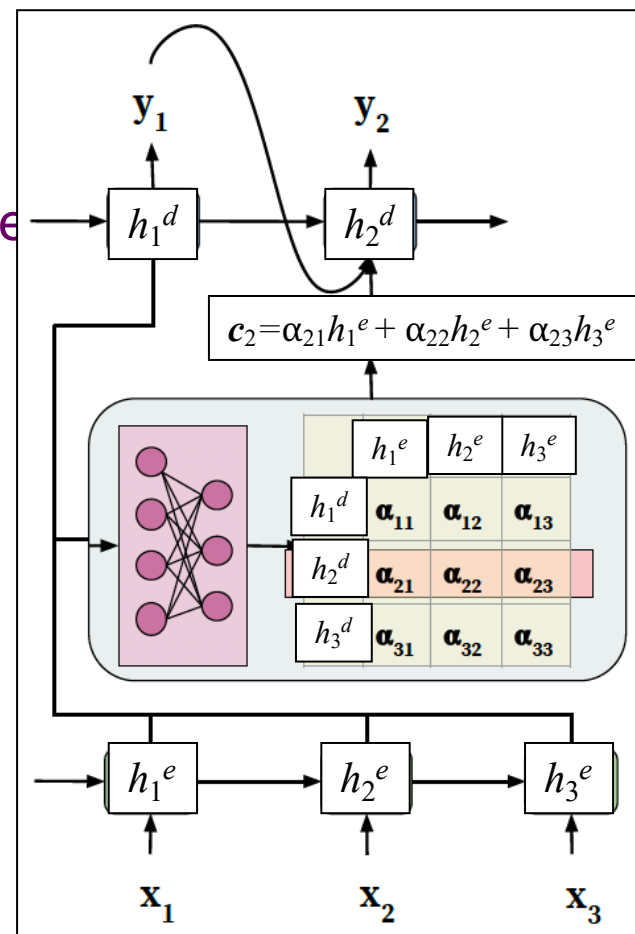| Name | Definition | Citation |
|---|---|---|
| Self-Attention(&) | Relating different positions of the same input sequence. Theoretically the self-attention can adopt any score functions above, but just replace the target sequence with the same input sequence. | Cheng2016 |
| Global/Soft | Attending to the entire input state space. | Xu2015 |
| Local/Hard | Attending to the part of input state space; i.e. a patch of the input image. | Xu2015; Luong2015 |

25

# Neural network for learning $\alpha_{ij}$

- $\alpha_{ij}$ capture relevance between
  - $h_i^d$ (decoder hidden state, or *query* state)
  - $h_j^e$ (encoder hidden state, or candidate state)

- $\alpha_{ij}$ used to build a context vector $c_i$
  at each decoding position $i$,
  - A weighted sum of encoder hidden states and attention weights

$$c_i = \sum_j \alpha_{ij} h_j^e$$

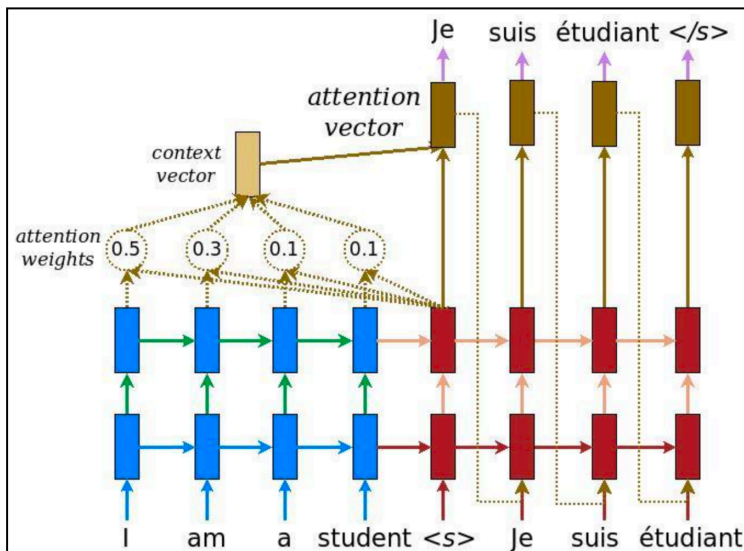- Context vectors $c_i$ determine next decoder hidden states

$$h_i^d = g(\hat{y}_{i-1}, h_{i-1}^d, c_i)$$

Context vector is how the decoder can access entire input sequence and also focus on its relevant positions



$y_1$ $y_2$

$h_1^d$ $h_2^d$

$c_2 = \alpha_{21} h_1^e + \alpha_{22} h_2^e + \alpha_{23} h_3^e$

|        | $h_1^e$ | $h_2^e$ | $h_3^e$ |
|--------|---------|---------|---------|
| $h_1^d$ | $\alpha_{11}$ | $\alpha_{12}$ | $\alpha_{13}$ |
| $h_2^d$ | $\alpha_{21}$ | $\alpha_{22}$ | $\alpha_{23}$ |
| $h_3^d$ | $\alpha_{31}$ | $\alpha_{32}$ | $\alpha_{33}$ |

$h_1^e$ $h_2^e$ $h_3^e$

$x_1$ $x_2$ $x_3$

# Tensorflow: NMT with Attention

https://www.tensorflow.org/tutorials/text/nmt_with_attention



$$\alpha_{ts} = \frac{\exp\left(\text{score}(\boldsymbol{h}_t, \bar{\boldsymbol{h}}_s)\right)}{\sum_{s'=1}^{S} \exp\left(\text{score}(\boldsymbol{h}_t, \bar{\boldsymbol{h}}_{s'})\right)} \qquad \text{[Attention weights]} \qquad (1)$$

$$\boldsymbol{c}_t = \sum_s \alpha_{ts} \bar{\boldsymbol{h}}_s \qquad \text{[Context vector]} \qquad (2)$$

$$\boldsymbol{a}_t = f(\boldsymbol{c}_t, \boldsymbol{h}_t) = \tanh(\boldsymbol{W_c}[\boldsymbol{c}_t; \boldsymbol{h}_t]) \qquad \text{[Attention vector]} \qquad (3)$$

$$\text{score}(\boldsymbol{h}_t, \bar{\boldsymbol{h}}_s) = \begin{cases} \boldsymbol{h}_t^{\top} \boldsymbol{W} \bar{\boldsymbol{h}}_s & \text{[Luong's multiplicative style]} \\ \boldsymbol{v}_a^{\top} \tanh\left(\boldsymbol{W_1} \boldsymbol{h}_t + \boldsymbol{W_2} \bar{\boldsymbol{h}}_s\right) & \text{[Bahdanau's additive style]} \end{cases} \qquad (4)$$

Plot of attention weights



Input: <start> hace mucho frio aqui . <end>
Predicted translation: it s very cold here . <end>