

Attention Mechanism

Debapriyo Majumdar

debapriyo@isical.ac.in

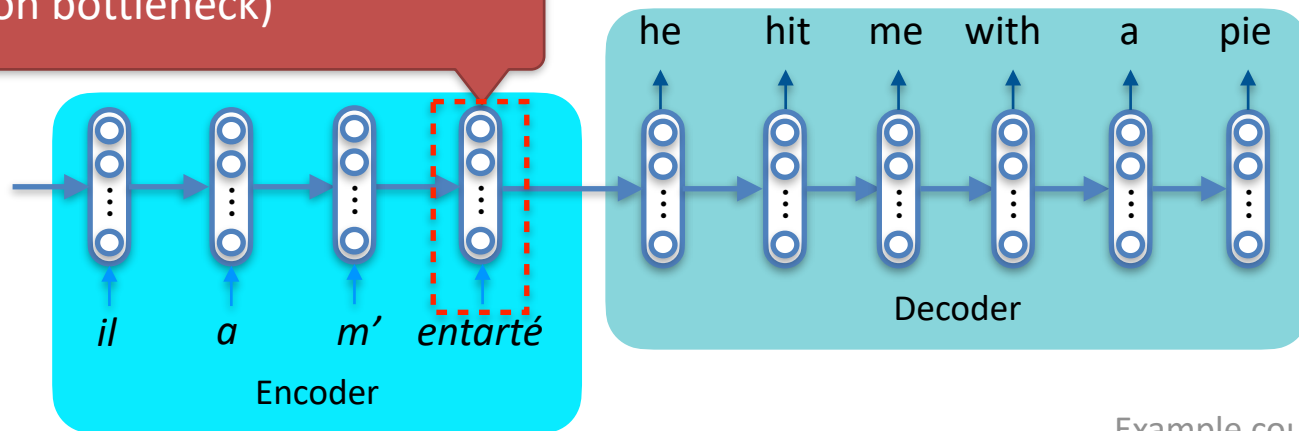
Attention: Motivation

2

- **Neural Machine Translation (NMT)** is the **flagship task** for NLP deep learning
- Basic architecture: an encoder RNN (or LSTM/GRU) producing an encoding of the sentence, followed by a decoder RNN (or LSTM/GRU)

Encoding of the source sentence: too much dependency on this vector (information bottleneck)

Attention: while decoding, focus on particular parts of the input



Example courtesy: Abigail See

Attention Mechanism

3

Amount of attention $y^{<1>}$
should pay on $a^{<i>}$

context vector
(attention output)

$$c^{<1>} = \sum_{i=1}^{T_x} \alpha^{<1,i>} a^{<i>}$$

attention $\alpha^{<1,i>}$
(probability
distribution)

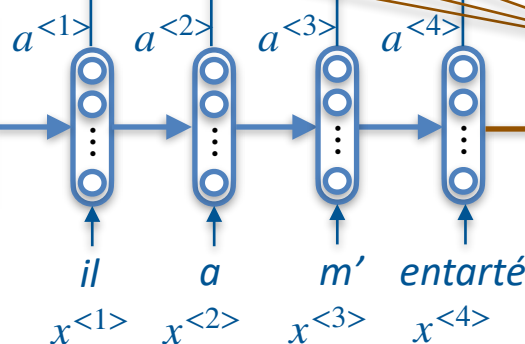
attention
scores $e^{<1,i>}$

softmax

$\hat{y}^{<1>}$
he

$c^{<1>}$

$s^{<0>}$



How do we
compute these?
Later.

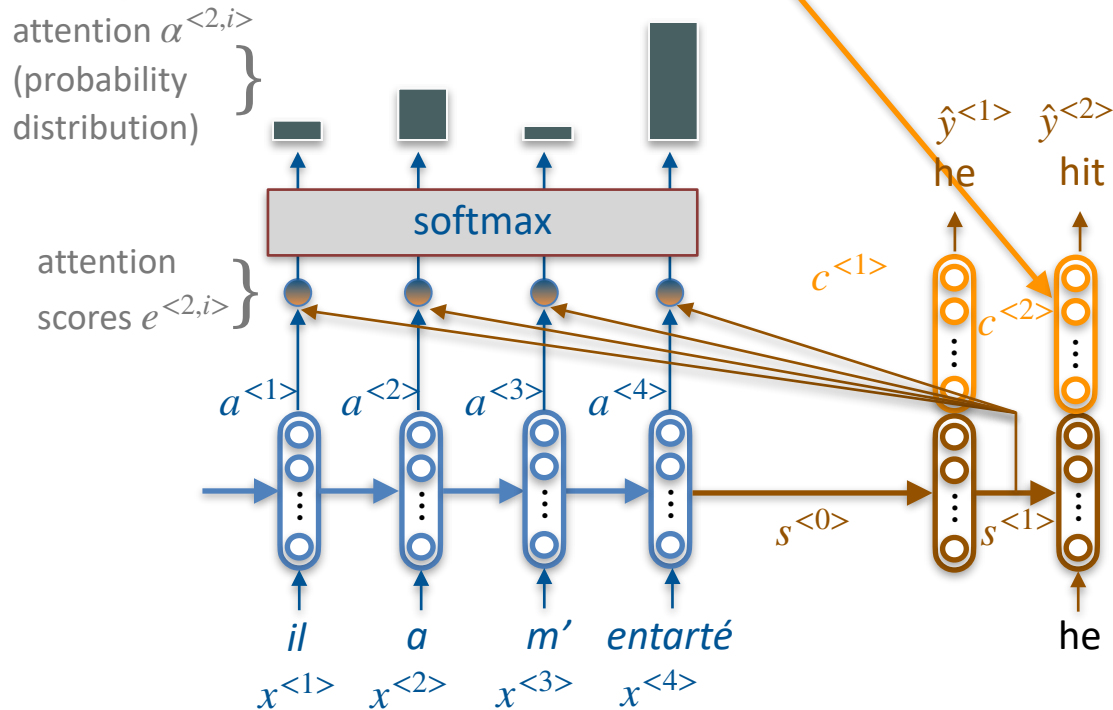
Attention Mechanism

4

Amount of attention $y^{<2>}$
should pay on $a^{<i>}$

context vector
(attention output)

$$c^{<2>} = \sum_{i=1}^{T_x} \alpha^{<2,i>} a^{<i>}$$



Attention Mechanism

5

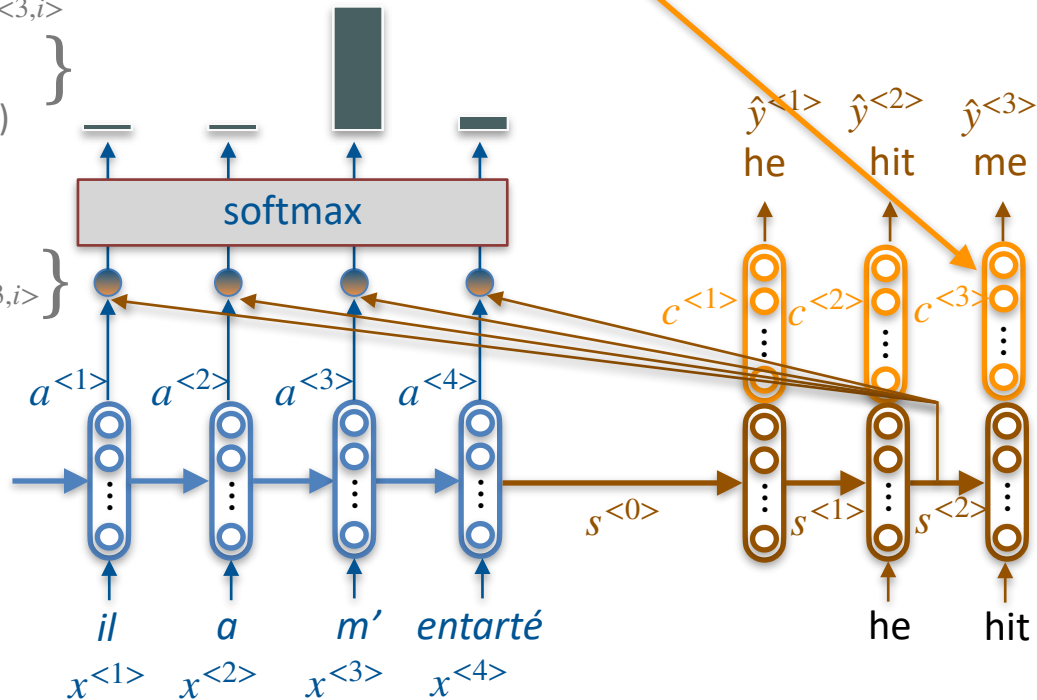
Amount of attention $y^{<3>}$
should pay on $a^{<i>}$

context vector
(attention output)

$$c^{<3>} = \sum_{i=1}^{T_x} \alpha^{<3,i>} a^{<i>}$$

attention $\alpha^{<3,i>}$
(probability
distribution)

attention
scores $e^{<3,i>}$



Attention Mechanism

6

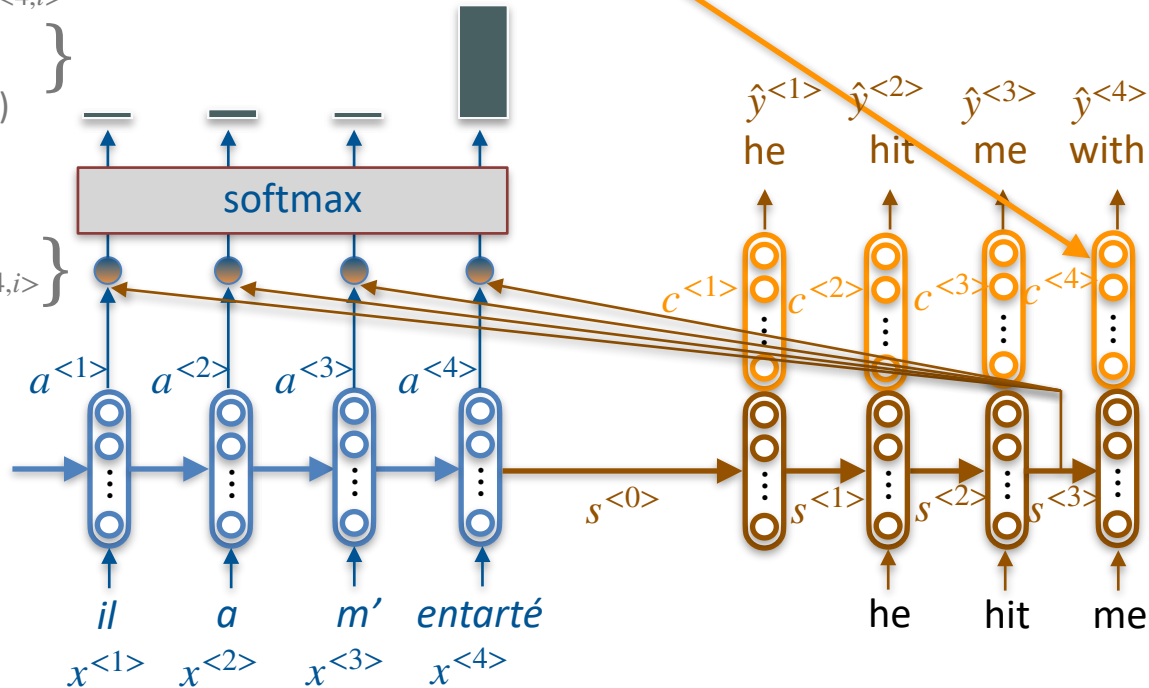
Amount of attention $y^{<4>}$
should pay on $a^{<i>}$

context vector
(attention output)

$$c^{<4>} = \sum_{i=1}^{T_x} \alpha^{<4,i>} a^{<i>}$$

attention $\alpha^{<4,i>}$
(probability
distribution)

attention
scores $e^{<4,i>}$



Attention Mechanism

7

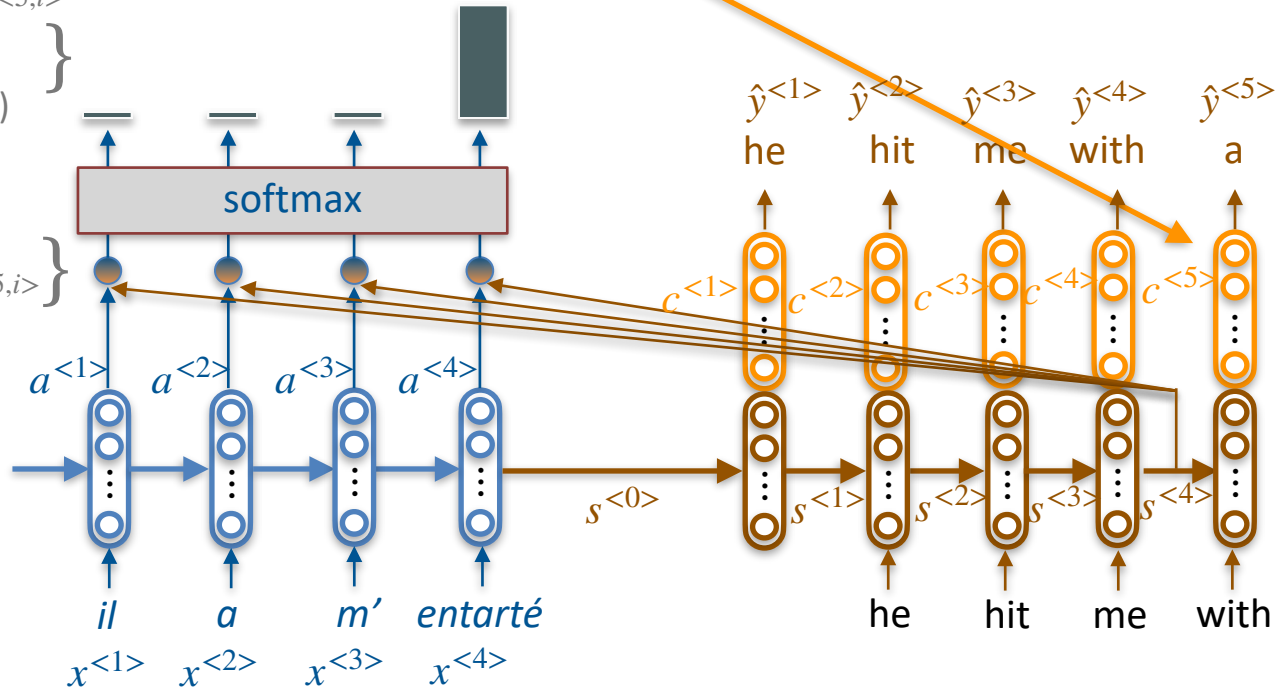
Amount of attention $y^{<5>}$
should pay on $a^{<i>}$

context vector
(attention output)

$$c^{<5>} = \sum_{i=1}^{T_x} \alpha^{<5,i>} a^{<i>}$$

attention $\alpha^{<5,i>}$
(probability
distribution)

attention
scores $e^{<5,i>}$



Attention Mechanism

Amount of attention $y^{<6>}$ should pay on $a^{<i>}$

context vector
(attention output)

$$c^{<6>} = \sum_{i=1}^{T_x} \alpha^{<6,i>} a^{<i>}$$

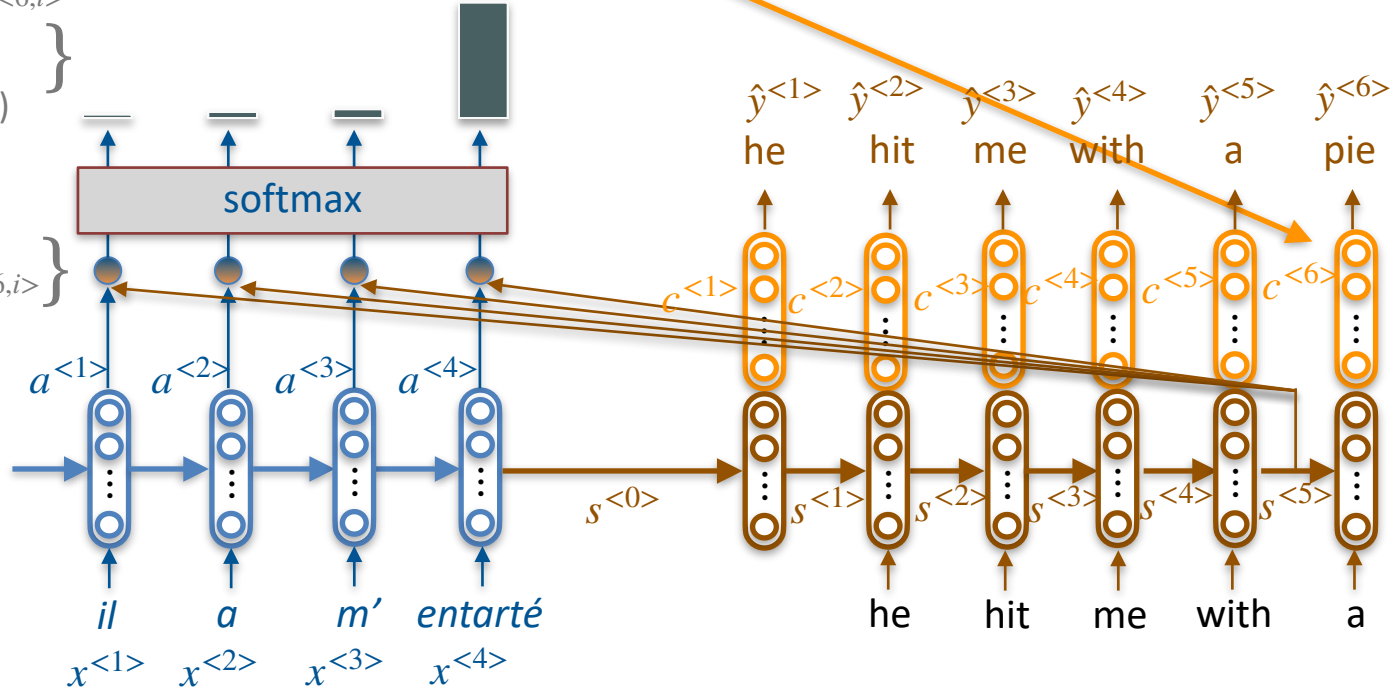
Attention matrix α
(transpose of)

	he	hit	me	with	a	pie
il						
a						
m'						
entarté						

Courtesy: Abigail See

attention $\alpha^{<6,i>}$
(probability distribution)

attention scores $e^{<6,i>}$



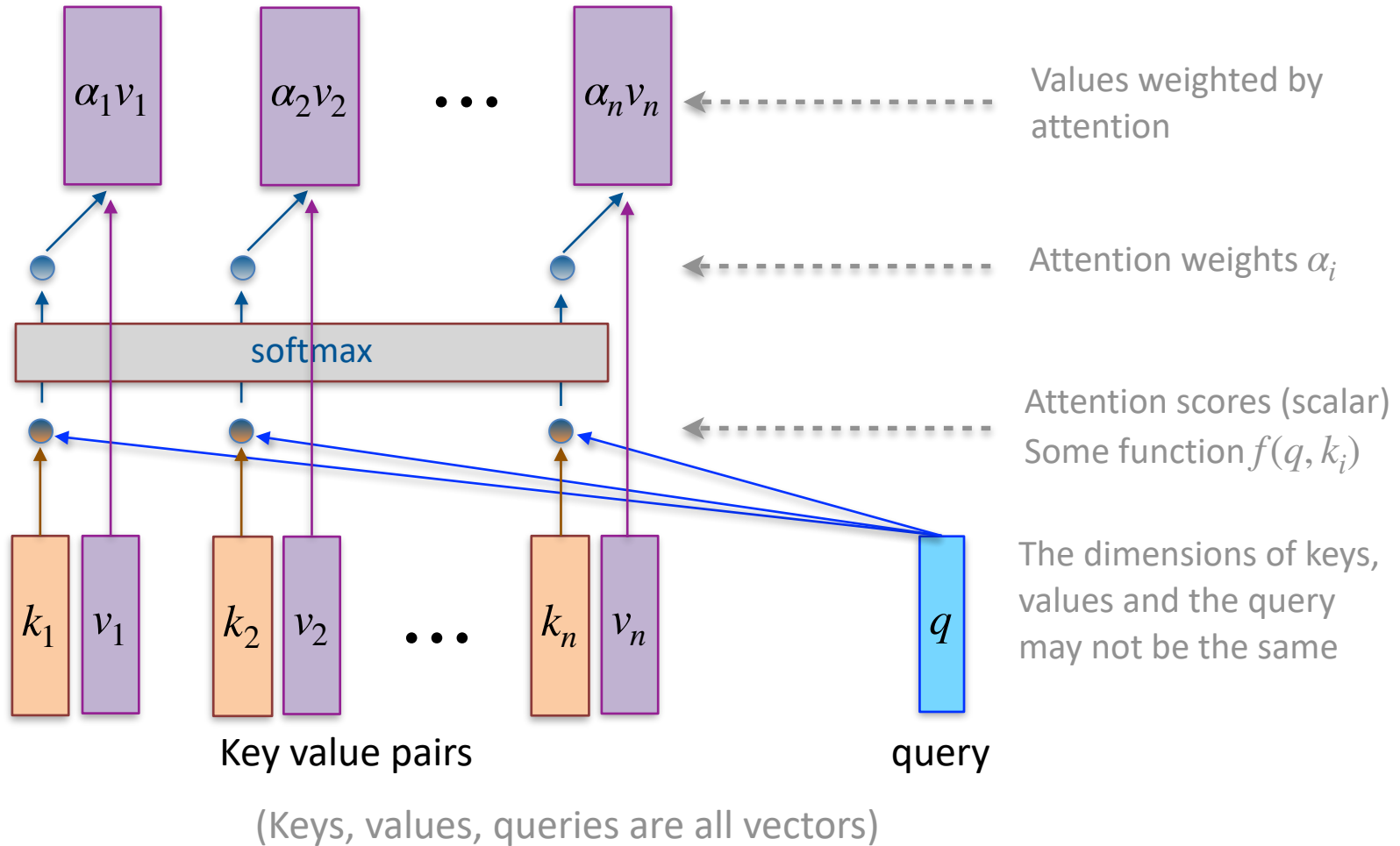
Computing attention scores

9

- Given: encoder hidden states $a^{<i>} \in \mathbb{R}^{d_1}$ and decoder hidden state $s^{<j>} \in \mathbb{R}^{d_2}$
 - Learn the attention scores (called *alignment model* by the original authors) using a neural network with $a^{<i>}$ and $s^{<j>}$ forming the input layer
 - $e^{<j,i>} = v^T \tanh(W_1 a^{<i>} + W_2 s^{<j>})$
 - where $W_1 \in \mathbb{R}^{d_3 \times d_1}$ and $W_2 \in \mathbb{R}^{d_3 \times d_2}$ are weight matrices, $v \in \mathbb{R}^{d_3 \times 1}$ is a weight vector and d_3 is the attention dimension (a hyperparameter)
 - Simplified versions
 - Dot product: $e^{<j,i>} = a^{<i>}^T s^{<j>}$ (in this case we require $d_1 = d_2$)
 - Multiplication with a weight matrix: $e^{<i,j>} = a^{<i>}^T W s^{<j>}$ where W is a matrix that is learnt
- The next steps: computing $\alpha^{<j,i>}$ from $e^{<j,i>}$ by **softmax**
- Attention output (context vector) $c^{<j>} = \sum_{i=1}^{T_x} \alpha^{<j,i>} a^{<i>}$
- In general: attention is a way to compute weighted sum of a given set of vector **values** w.r.t. a **query**

Attention: Generalization

10



- Attention allows the decoder to focus on certain parts of the input
 - Improves performances for several NLP tasks
- Solves the information bottleneck problem
- Helps with the vanishing gradient problem as well
- Attention is interpretable
 - The attention matrix shows what the decoder focussed on
 - Automatically trained soft alignment

- Chris Manning, Abigail See and other TAs. *Natural Language Processing with Deep Learning*. Stanford University Course (CS224n), Winter 2019. web.stanford.edu/class/archive/cs/cs224n/cs224n.1194/, Lecture 8 (Abigail See): www.youtube.com/watch?v=XXtpJxZBa2c
- Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. "Neural machine translation by jointly learning to align and translate." *arXiv preprint arXiv:1409.0473* (2014).