

Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>



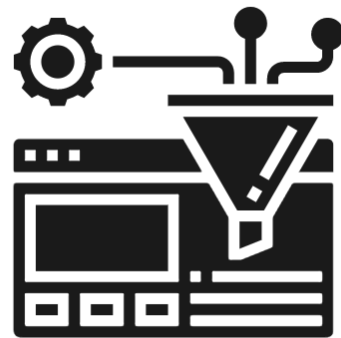
deeplearning.ai

Seq2Seq model for NMT



Outline

- Introduction to Neural Machine Translation
- Seq2Seq model and its shortcomings
- Solution for the information bottleneck



Neural Machine Translation

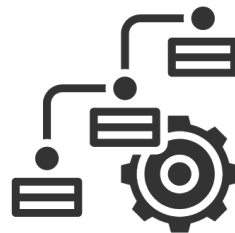
It's time for tea



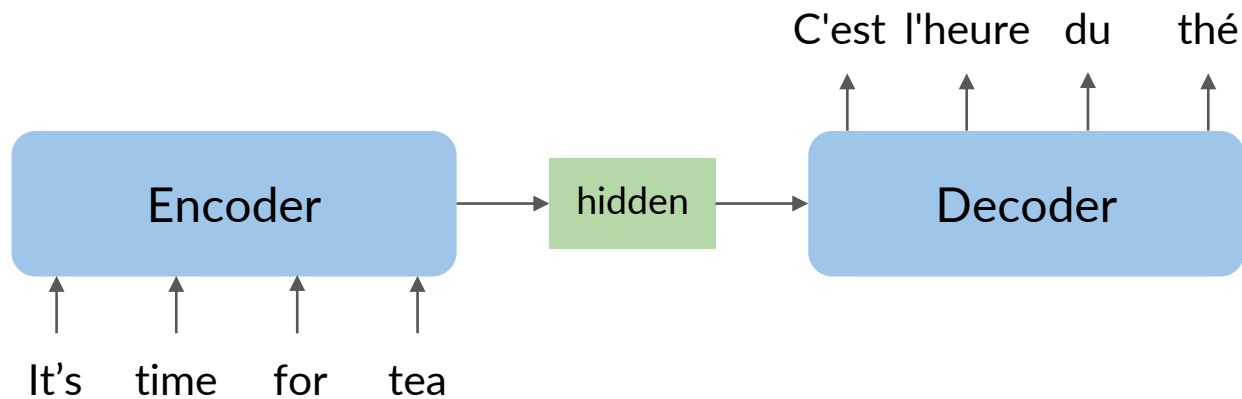
C'est l'heure du thé

Seq2Seq model

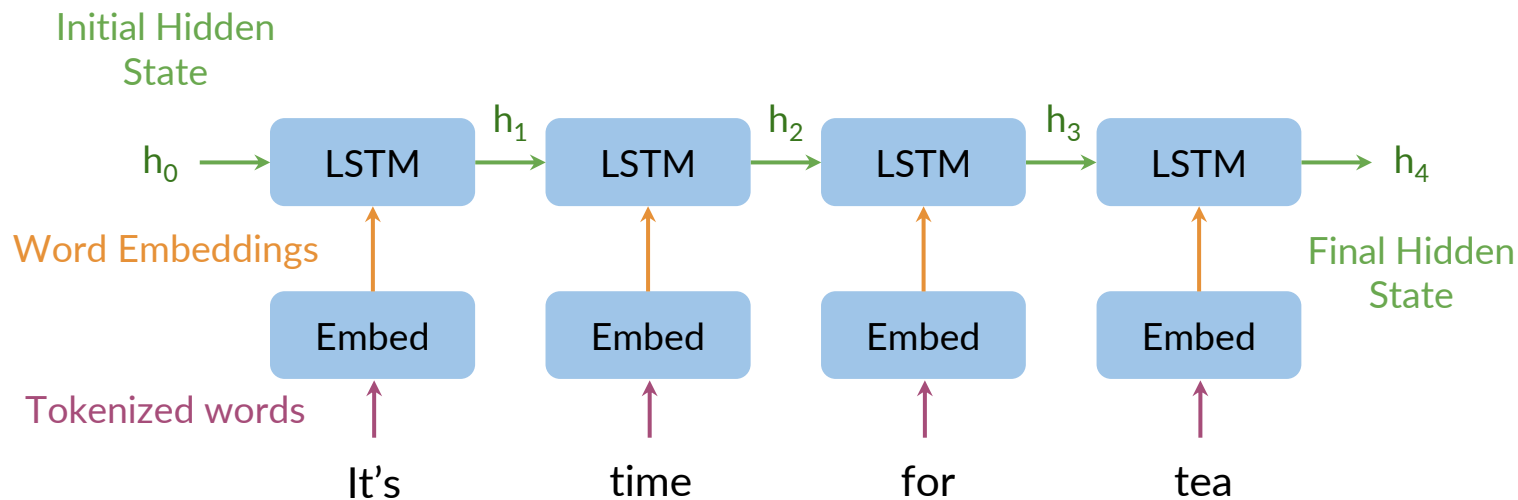
- Introduced by Google in 2014
- Maps variable-length sequences to fixed-length memory
- Inputs and outputs can have different lengths
- LSTMs and GRUs to avoid vanishing and exploding gradient problems



Seq2Seq model

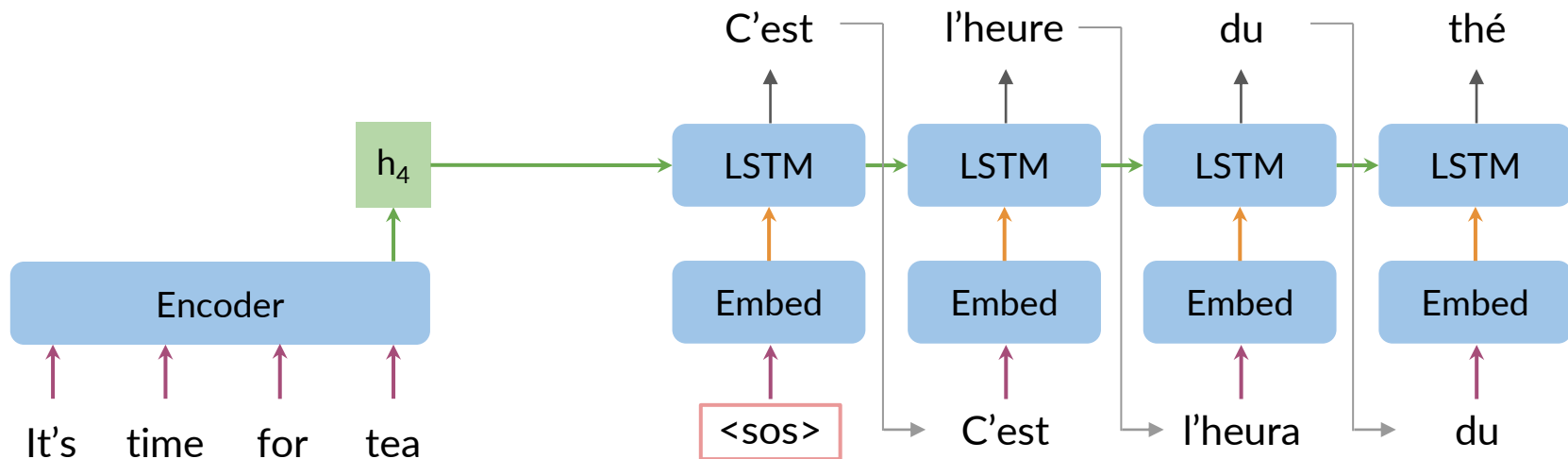


Seq2Seq encoder

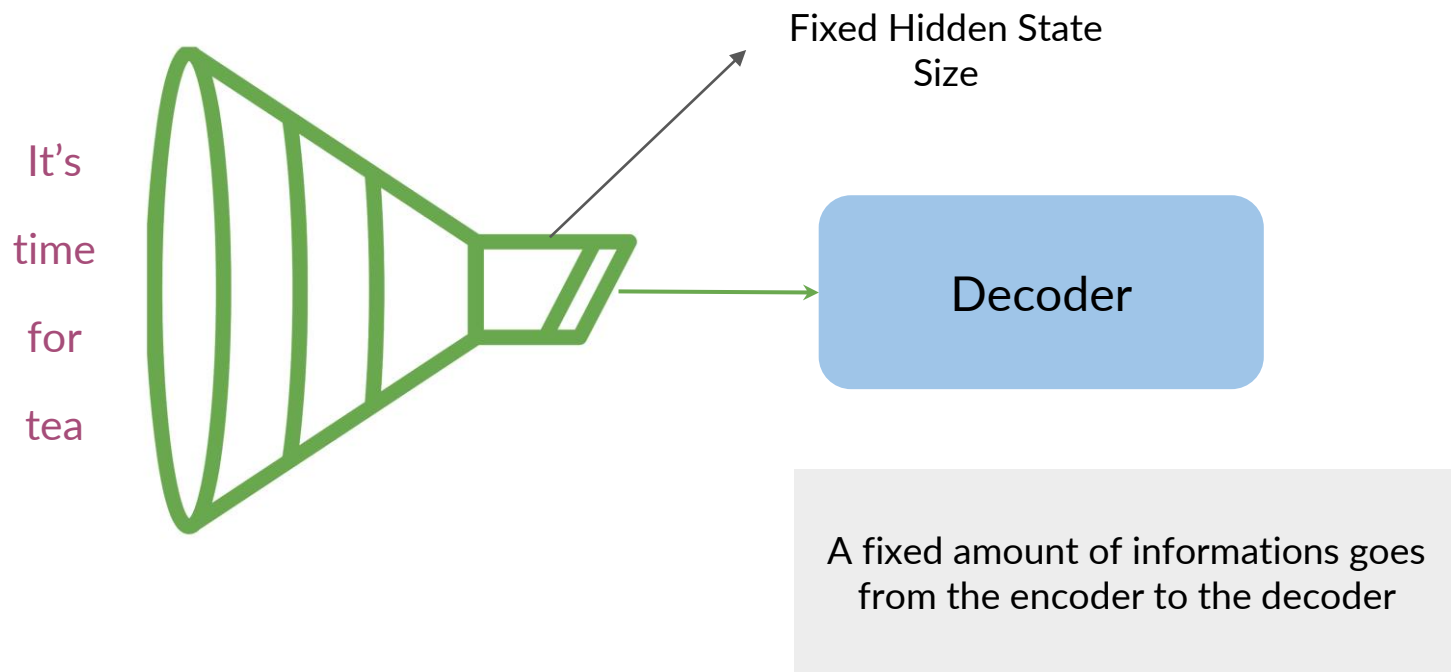


Encodes the overall meaning of the sentence

Seq2Seq decoder



The information bottleneck



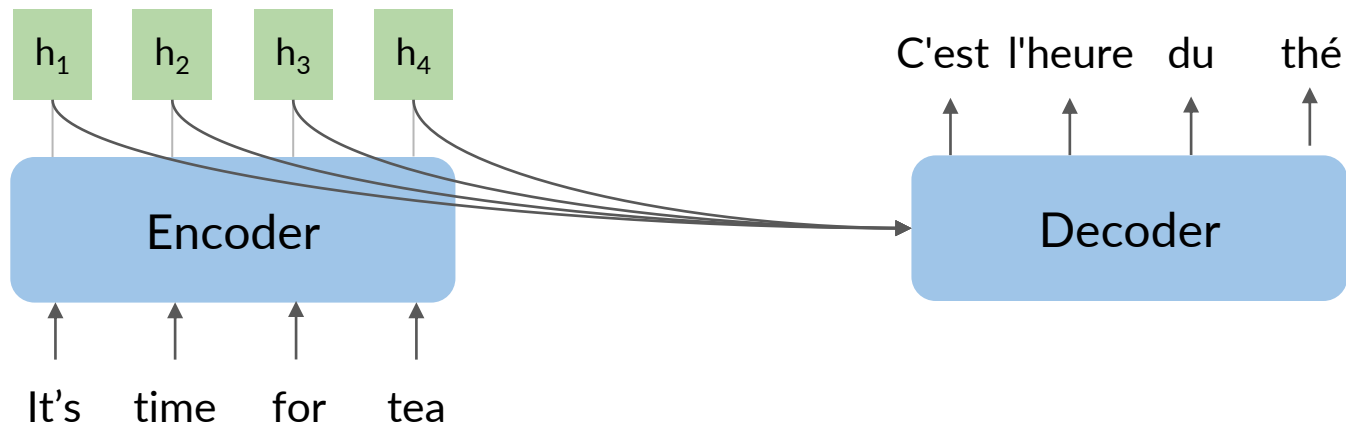
Seq2Seq shortcomings

- Variable-length sentences + fixed-length memory =

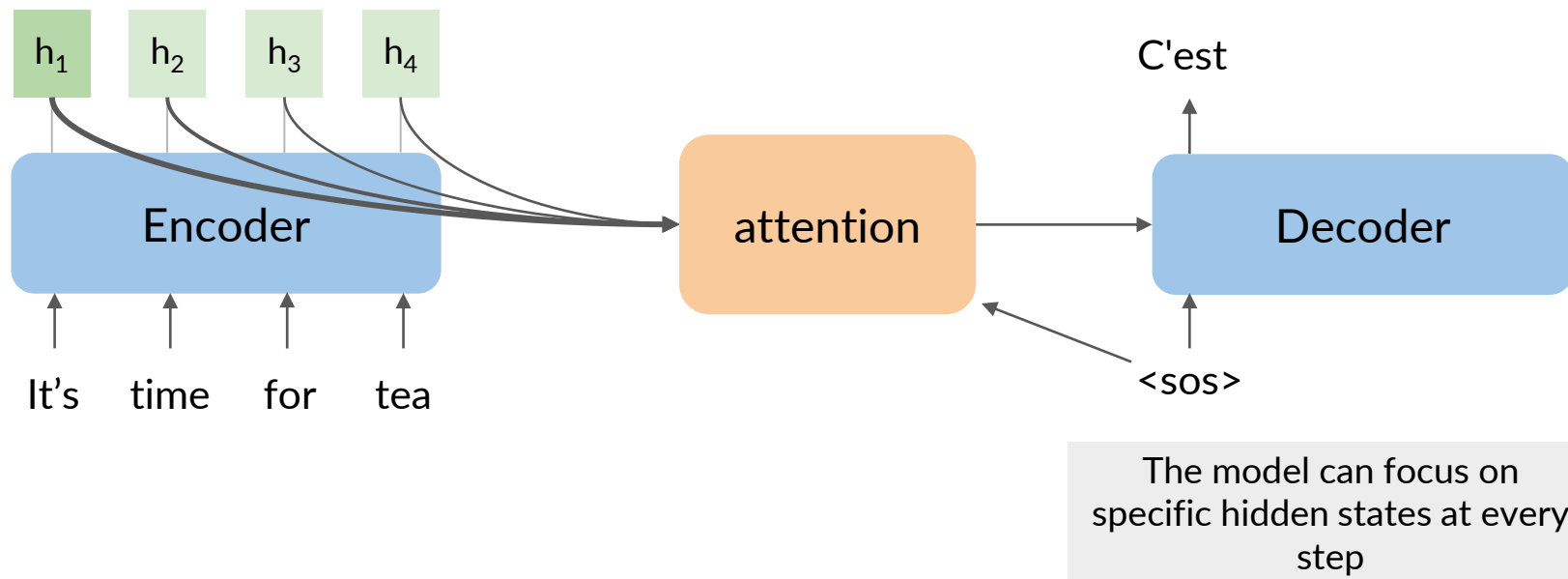


- As sequence size increases, model performance decreases

Use all the encoder hidden states?



Solution: focus attention in the right place





deeplearning.ai

Seq2Seq model with attention

NEURAL MACHINE TRANSLATION BY JOINTLY LEARNING TO ALIGN AND TRANSLATE

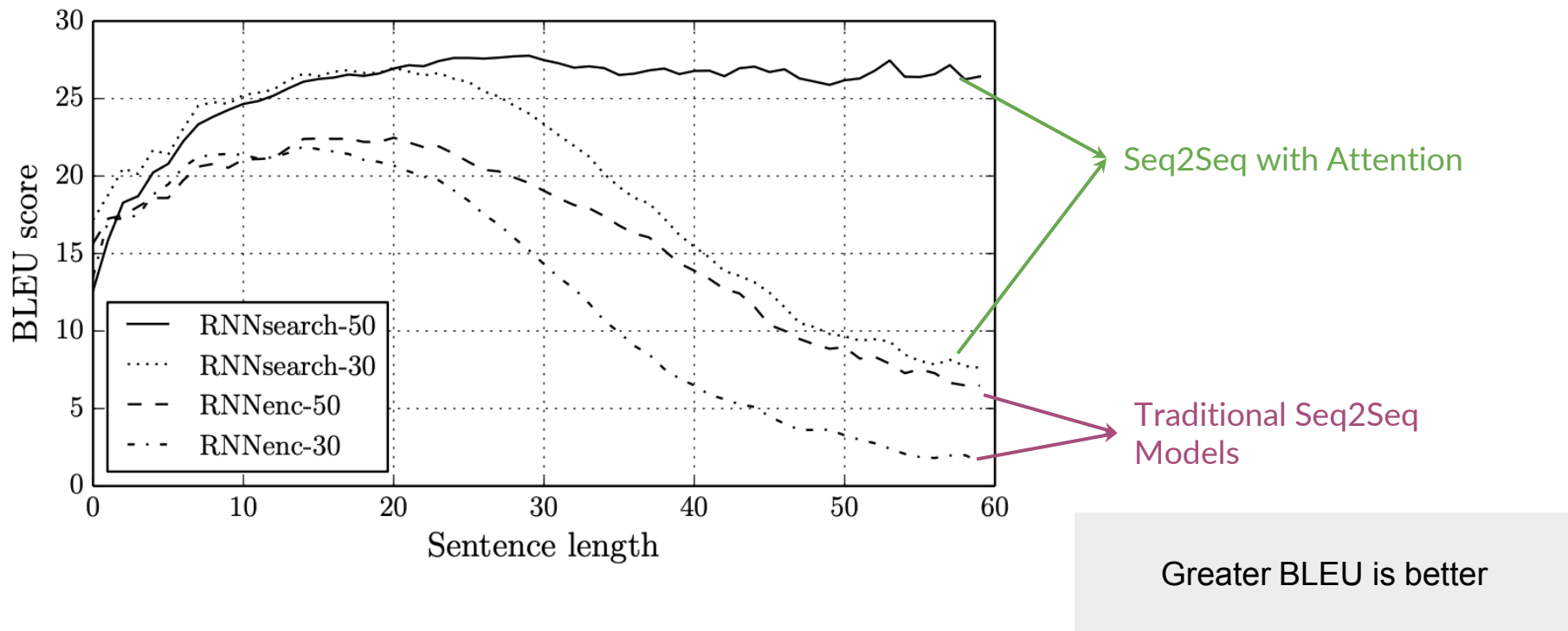
Dzmitry Bahdanau

Jacobs University Bremen, Germany

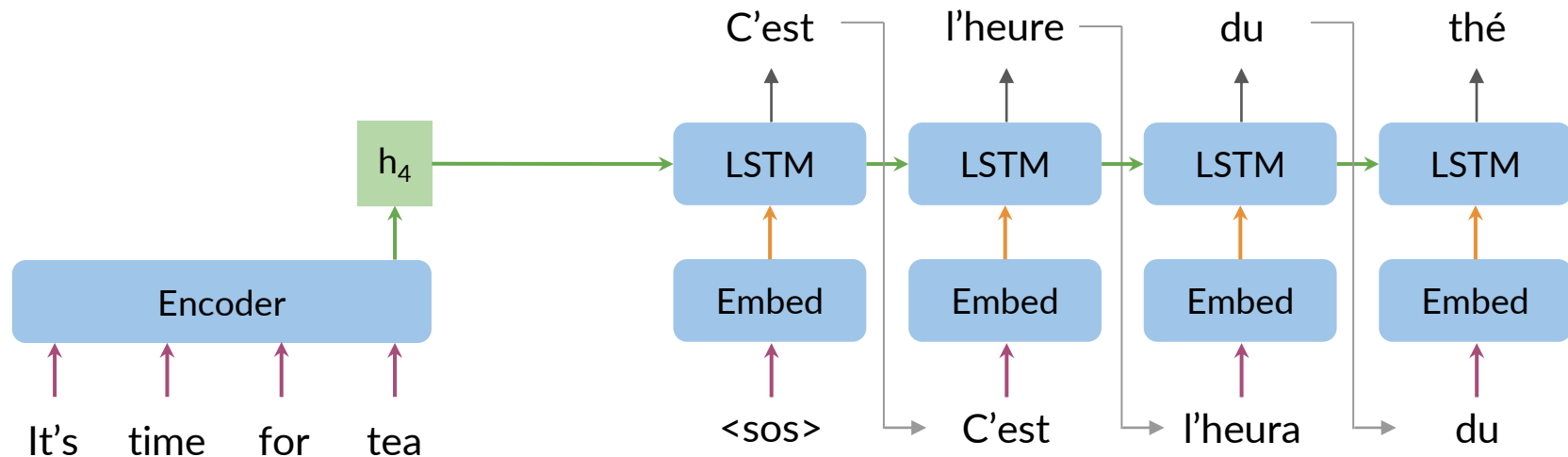
KyungHyun Cho Yoshua Bengio*

Université de Montréal

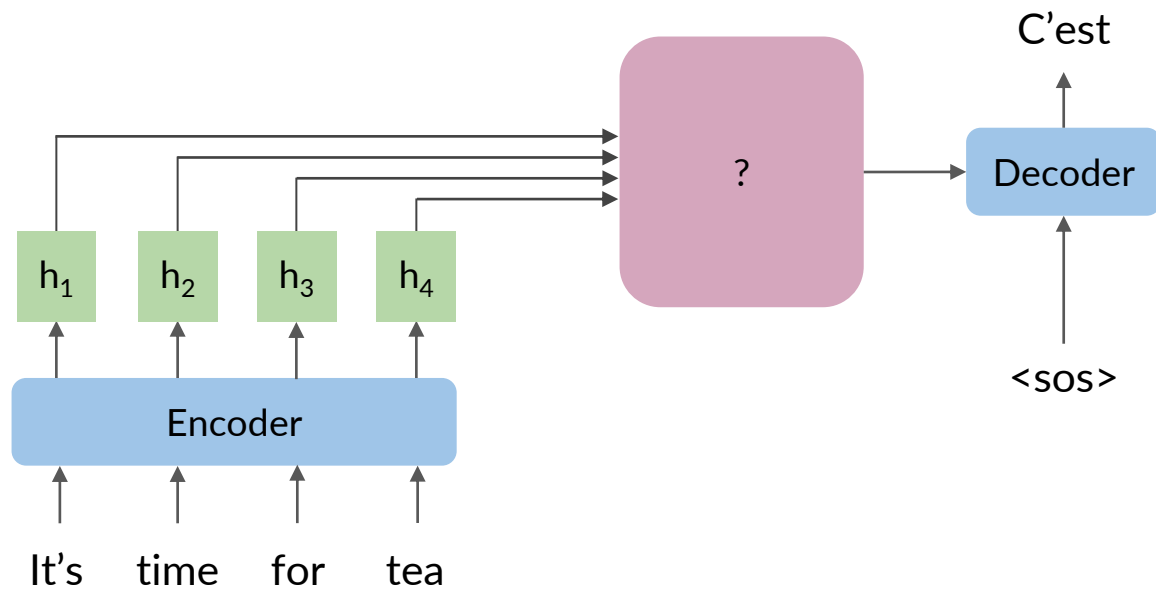
Performance



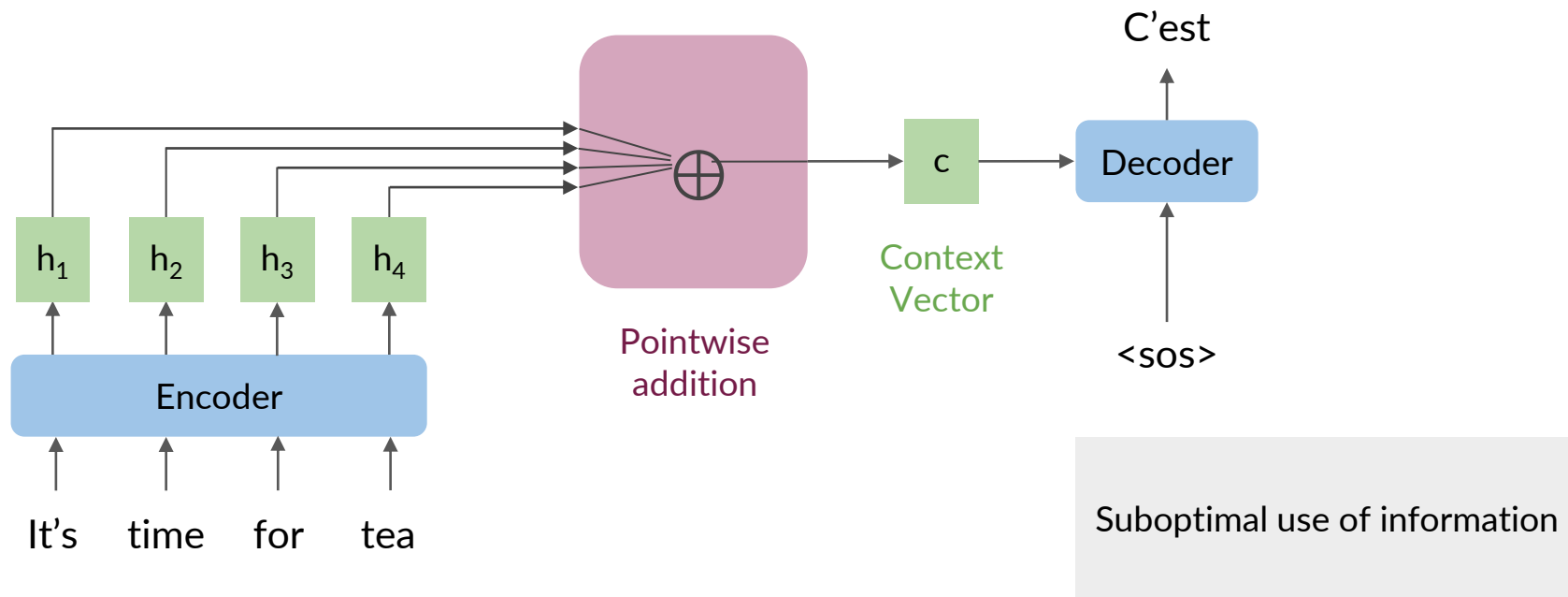
Traditional seq2seq models



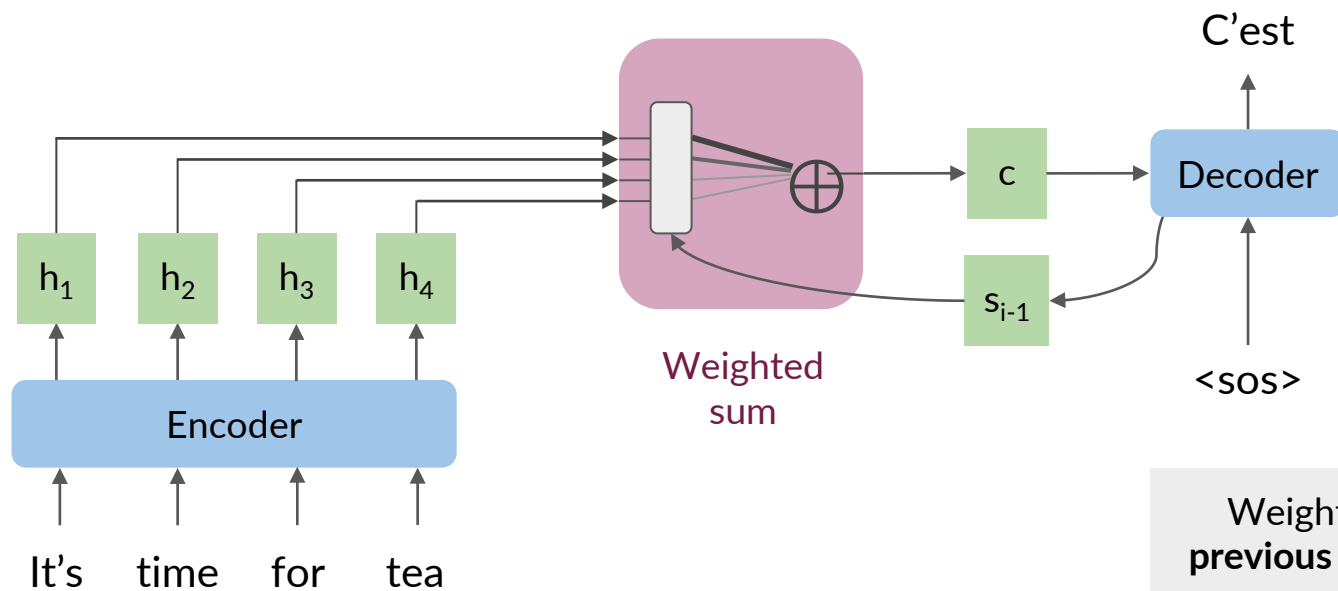
How to use all the hidden states?



How to use all the hidden states?

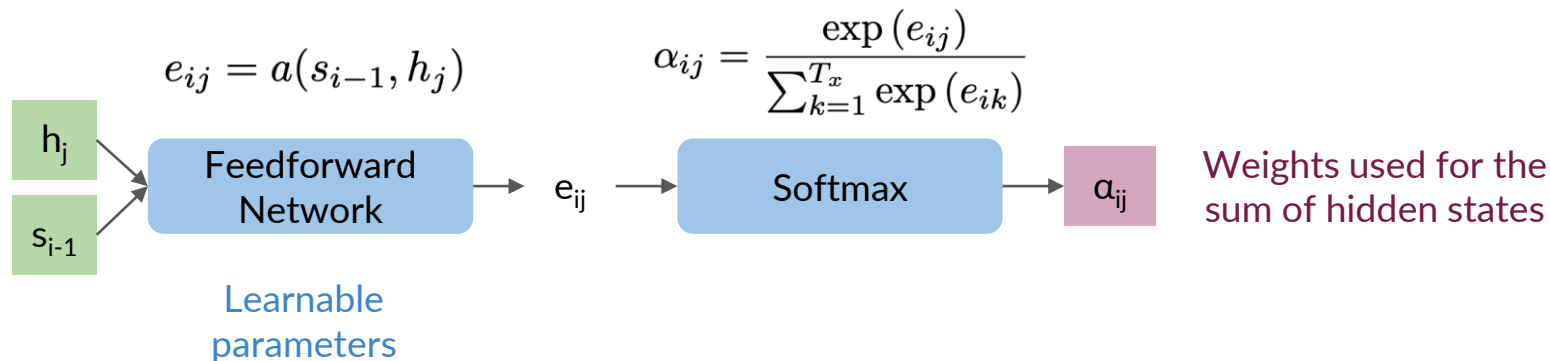


How to use all the hidden states?



Weights depend on the
previous hidden state in the
decoder

The attention layer in more depth



$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$$
$$\alpha_{i1} h_1 + \alpha_{i2} h_2 + \alpha_{i3} h_3 + \dots + \alpha_{iM} h_M \rightarrow c_i$$

Context Vector is an expected value

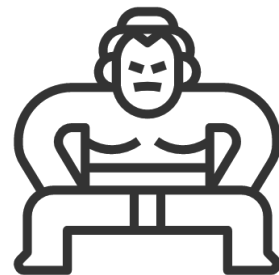


deeplearning.ai

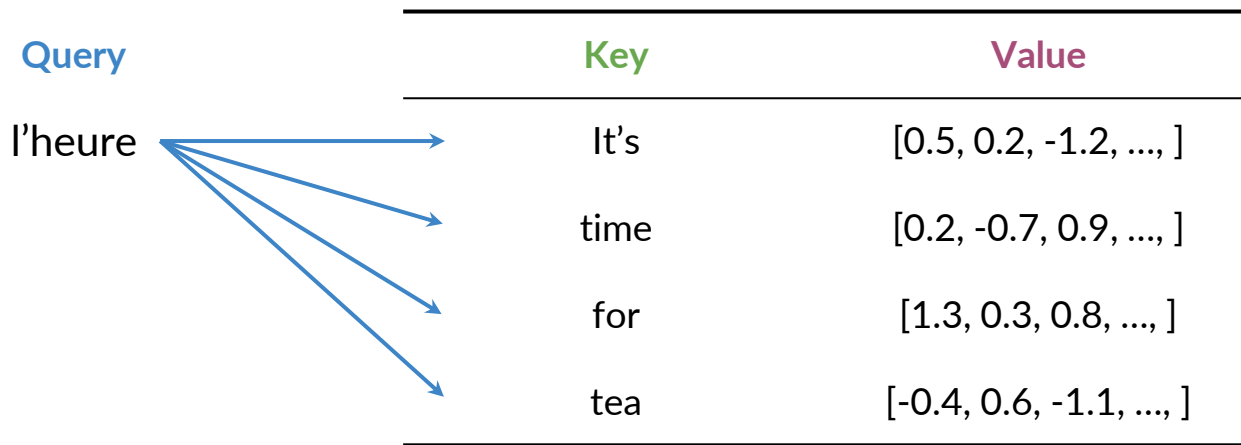
Queries, Keys, Values and Attention

Outline

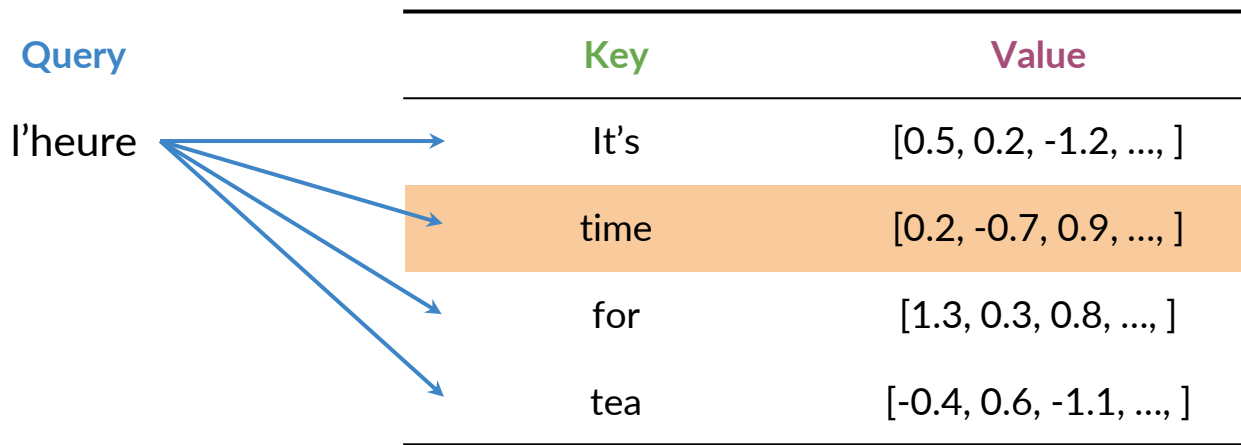
- Queries, Keys, and Values
- Alignment



Queries, Keys, Values



Queries, Keys, Values

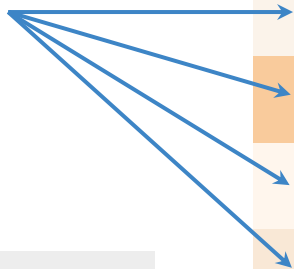


	Key	Value
l'heure	It's	[0.5, 0.2, -1.2, ...,]
	time	[0.2, -0.7, 0.9, ...,]
	for	[1.3, 0.3, 0.8, ...,]
	tea	[-0.4, 0.6, -1.1, ...,]

Queries, Keys, Values

Query

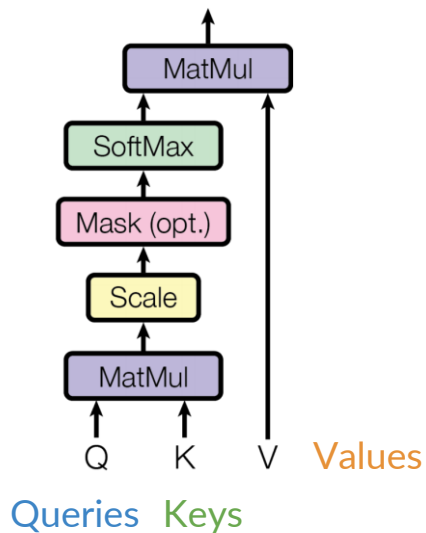
l'heure



Similarity is used in for
weighted sum

Key	Value
It's	[0.5, 0.2, -1.2, ...,]
time	[0.2, -0.7, 0.9, ...,]
for	[1.3, 0.3, 0.8, ...,]
tea	[-0.4, 0.6, -1.1, ...,]

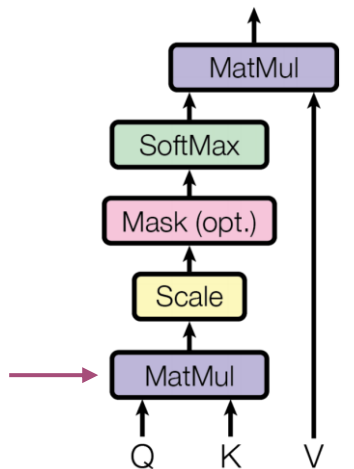
Scaled dot-product attention



(Vaswani et al., 2017)

$$\text{softmax} \left(\frac{QK^{\top}}{\sqrt{d_k}} \right) V$$

Scaled dot-product attention

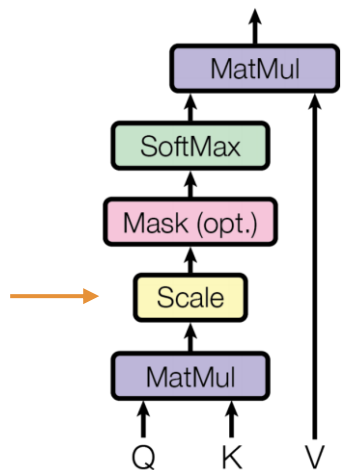


(Vaswani et al., 2017)

Similarity Between
Q and K

$$\text{softmax} \left(\frac{QK^{\top}}{\sqrt{d_k}} \right) V$$

Scaled dot-product attention

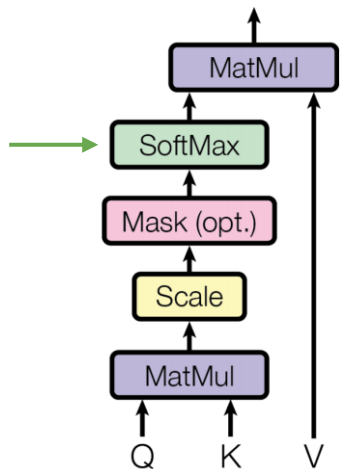


(Vaswani et al., 2017)

$$\text{softmax} \left(\frac{QK^{\top}}{\sqrt{d_k}} \right) V$$

Scale using the root
of the key vector
size

Scaled dot-product attention

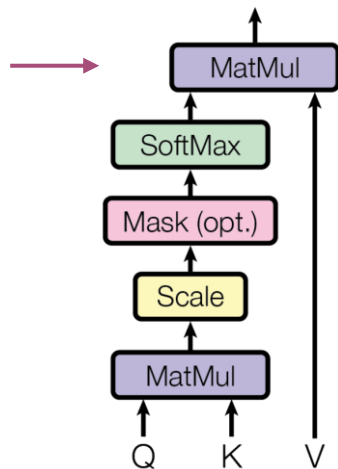


(Vaswani et al., 2017)

$$\text{softmax} \left(\frac{QK^{\top}}{\sqrt{d_k}} \right) V$$

Weights for the
weighted sum

Scaled dot-product attention



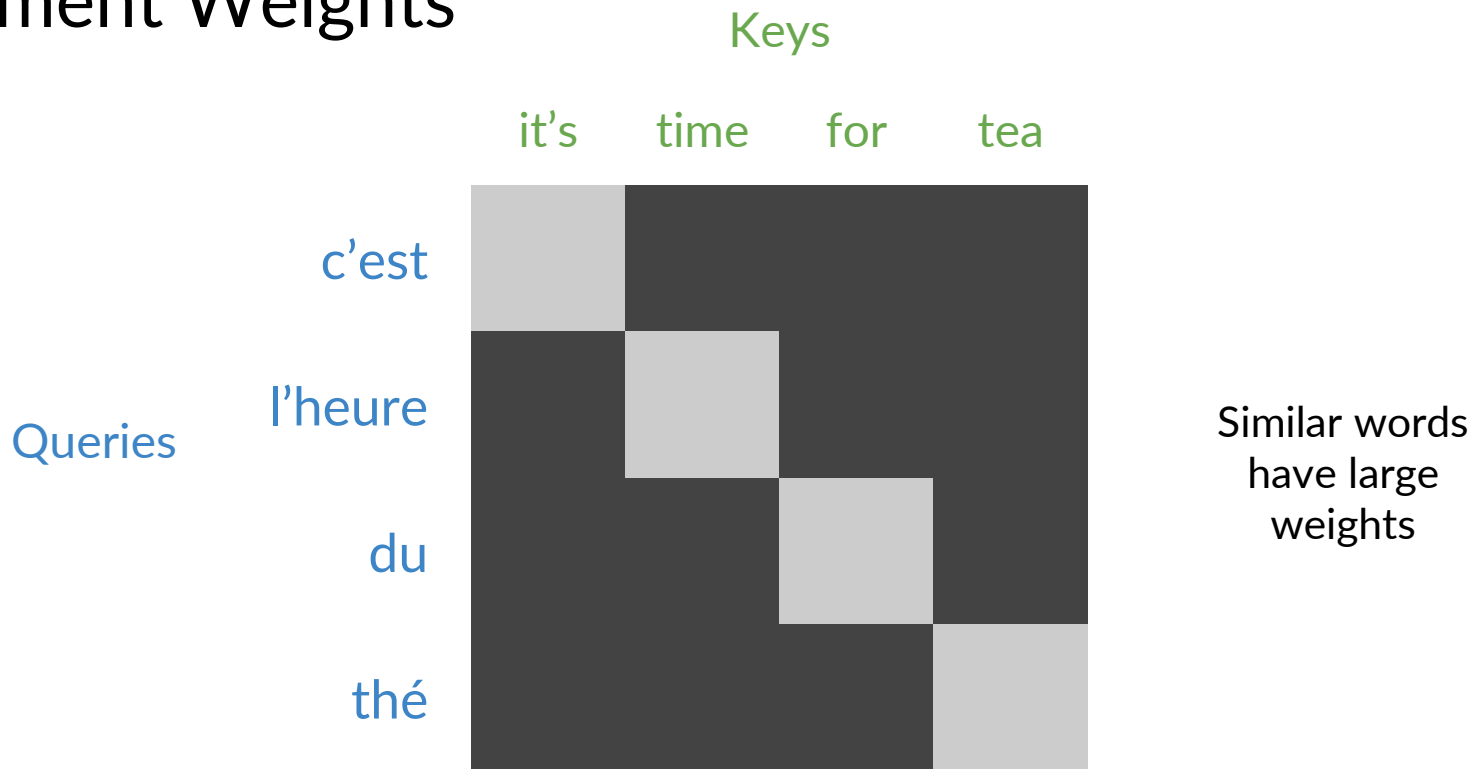
(Vaswani et al., 2017)

$$\text{softmax} \left(\frac{QK^{\top}}{\sqrt{d_k}} \right) V$$

Weighted sum of values V

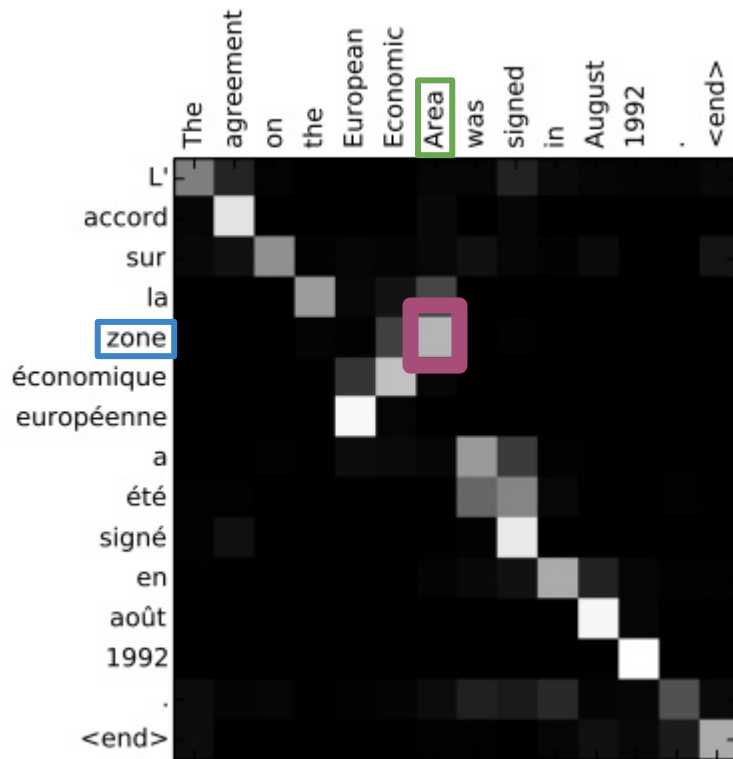
Just two matrix multiplications
and a Softmax!

Alignment Weights



Flexible attention

Works for languages with different grammar structures!



[Bahdanau et al., 2015](#)

Summary

- Attention is a layer that lets a model focus on what's important
- Queries, Values, and Keys are used for information retrieval inside the Attention layer
- Works for languages with very different grammatical structures





deeplearning.ai

Setup for machine translation

Data in machine translation

English	French
I am hungry!	J'ai faim!
...	...
I watched the soccer game.	J'ai regardé le match de football.

Attention! (pun intended) Assignment dataset is not as squeaky-clean as this example and contains some Spanish translations.

Machine translation setup

- Use pre-trained vector embeddings
- Otherwise, initially represent words with a one-hot vectors
- Keep track of index mappings with word2ind and ind2word dictionaries
- Add end of sequence tokens: **<EOS>**
- Pad the token vectors with zeros

Preparing to Translate to English

ENGLISH SENTENCE:

Both the ballpoint and the mechanical pencil in the series are equipped with a special mechanism: when the twist mechanism is activated, the lead is pushed forward.

TOKENIZED VERSION OF THE ENGLISH SENTENCE:

The diagram illustrates padding in a sequence. It shows three rows of numbers:

- Row 1: [4546, 4, 11358, 362, 8, 4, 23326, 20104, 1745, 8210, 9641, 5, 6]
- Row 2: 4 3103, 31 2767, 30 13 914 4797, 64 196 4, 22474, 5 4797, 16
- Row 3: 24864, 86, 2, 4, 1060, 16, 6413, 1138, 3, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0

A green box highlights the first 10 elements of Row 3. A red arrow points from the 10th element (16) to the text "<EOS>". The remaining elements (zeros) are labeled "Padding".

English to French

FRENCH TRANSLATION:

Le stylo à bille et le porte-mine de la série sont équipés d'un mécanisme spécial: lorsque le mécanisme de torsion est activé, le plomb est poussé vers l'avant.

TOKENIZED VERSION OF THE FRENCH TRANSLATION:

7	29587	9	18240	8	7	420	5	3440	2	6	156	39	7941	14	19
5548	2648	562	7	5548	2	23194	18	20114	1	7	5695	18	8865	149	
12	137	1	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0													

Padding

<EOS>

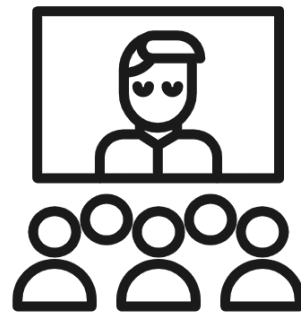


deeplearning.ai

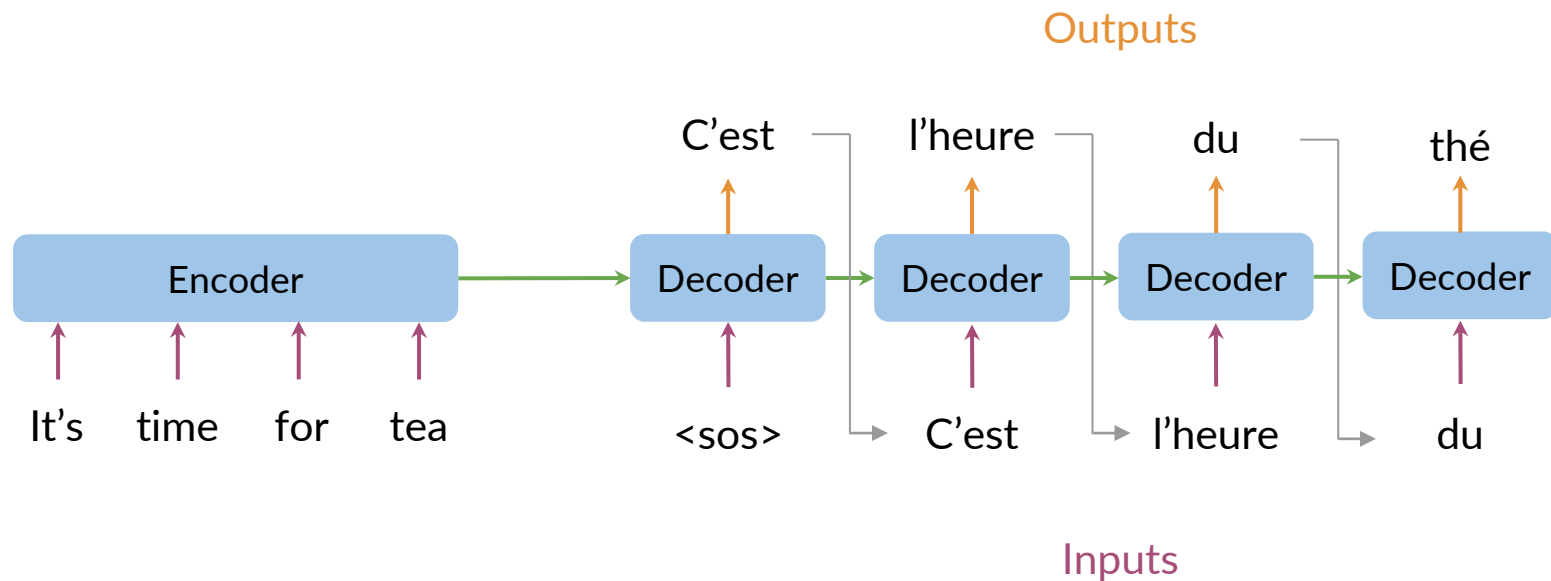
Teacher Forcing

Outline

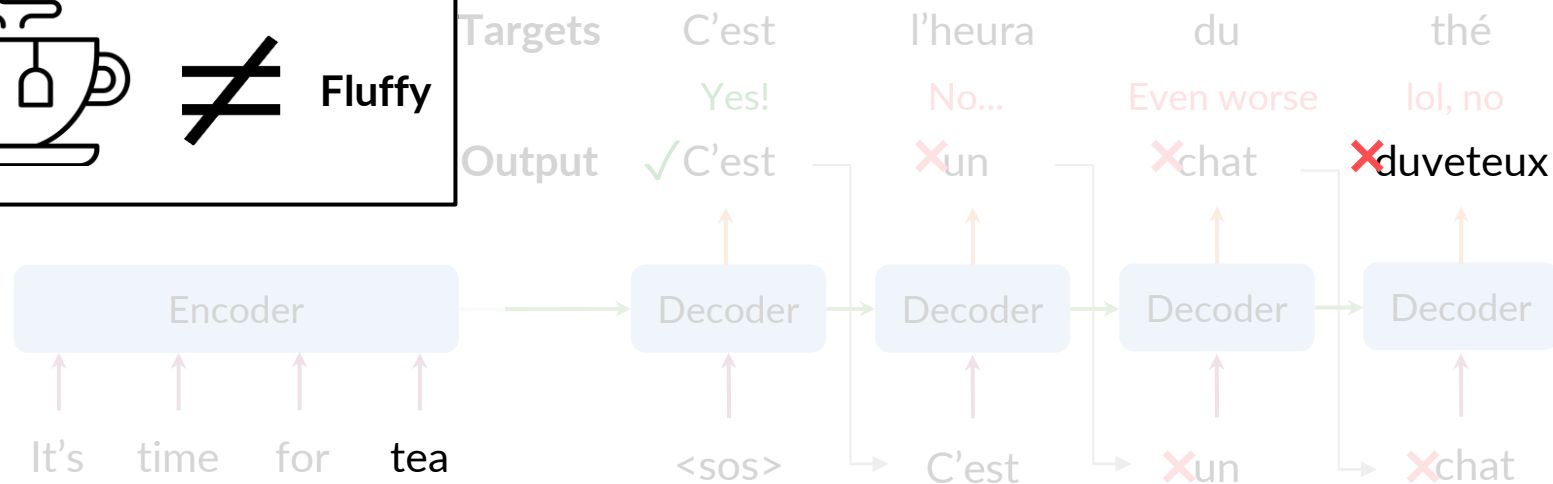
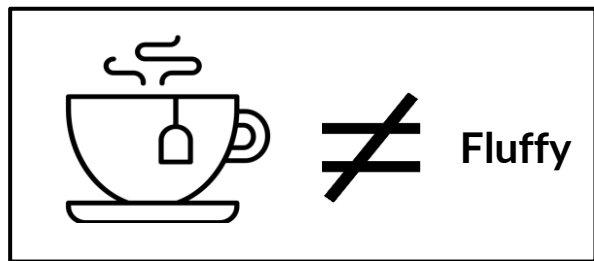
- Training for NMT
- Teacher forcing



Traditional seq2seq models

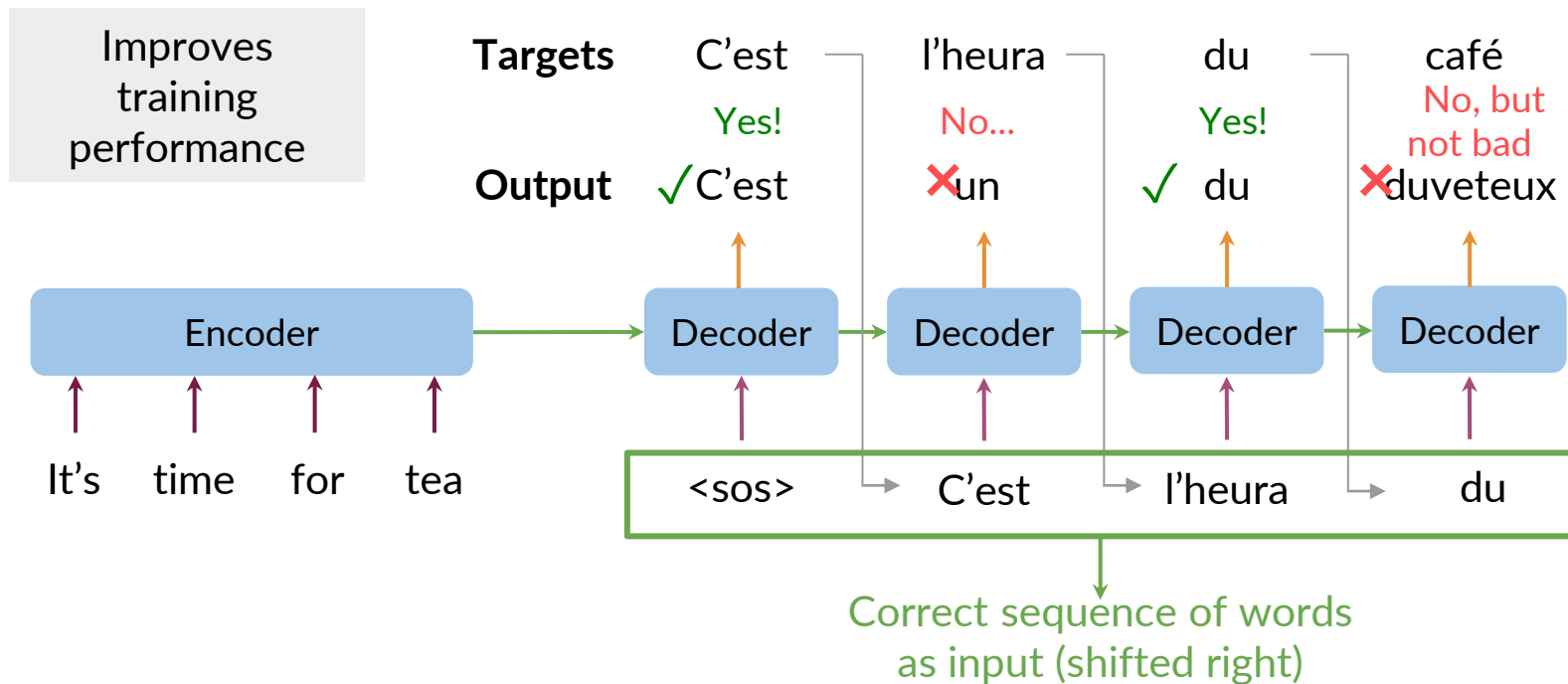


Training seq2seq models



Errors from early steps propagate

Teacher Forcing



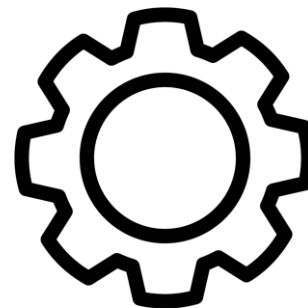


deeplearning.ai

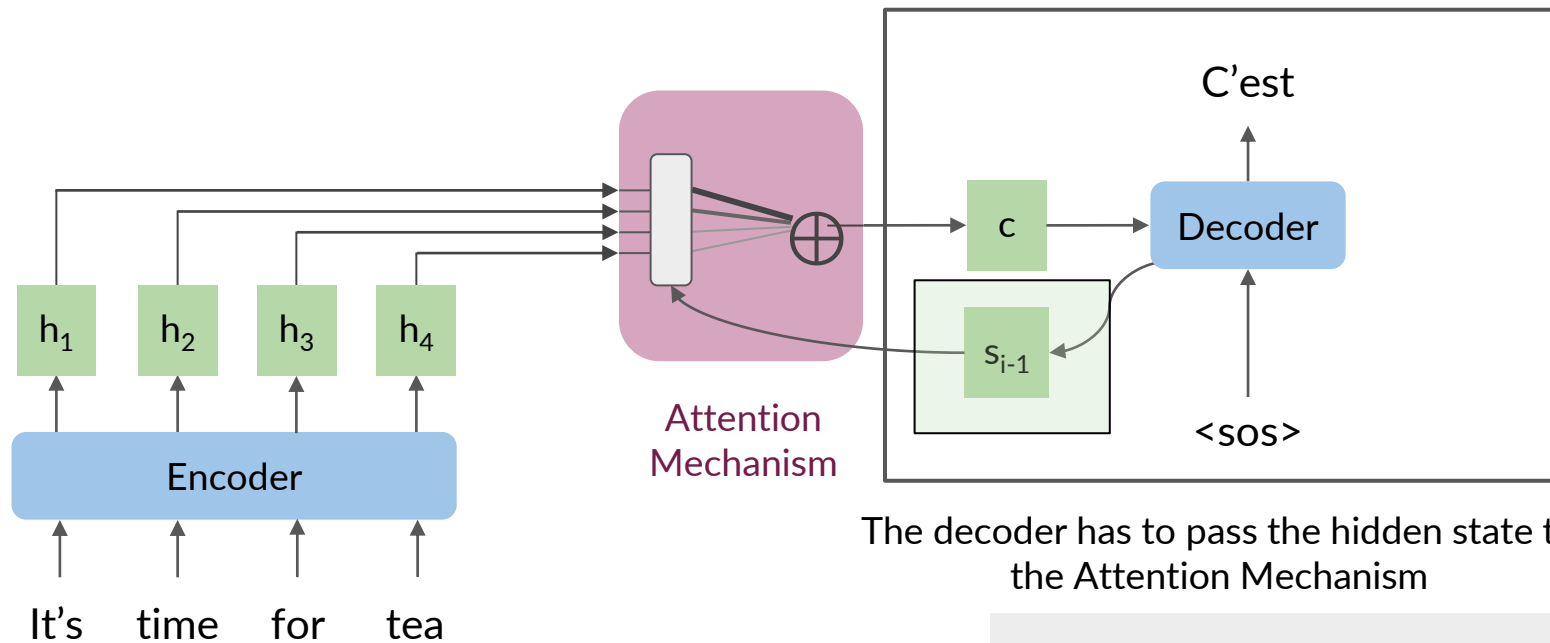
NMT Model with Attention

Outline

- How everything fits together
- NMT model in detail



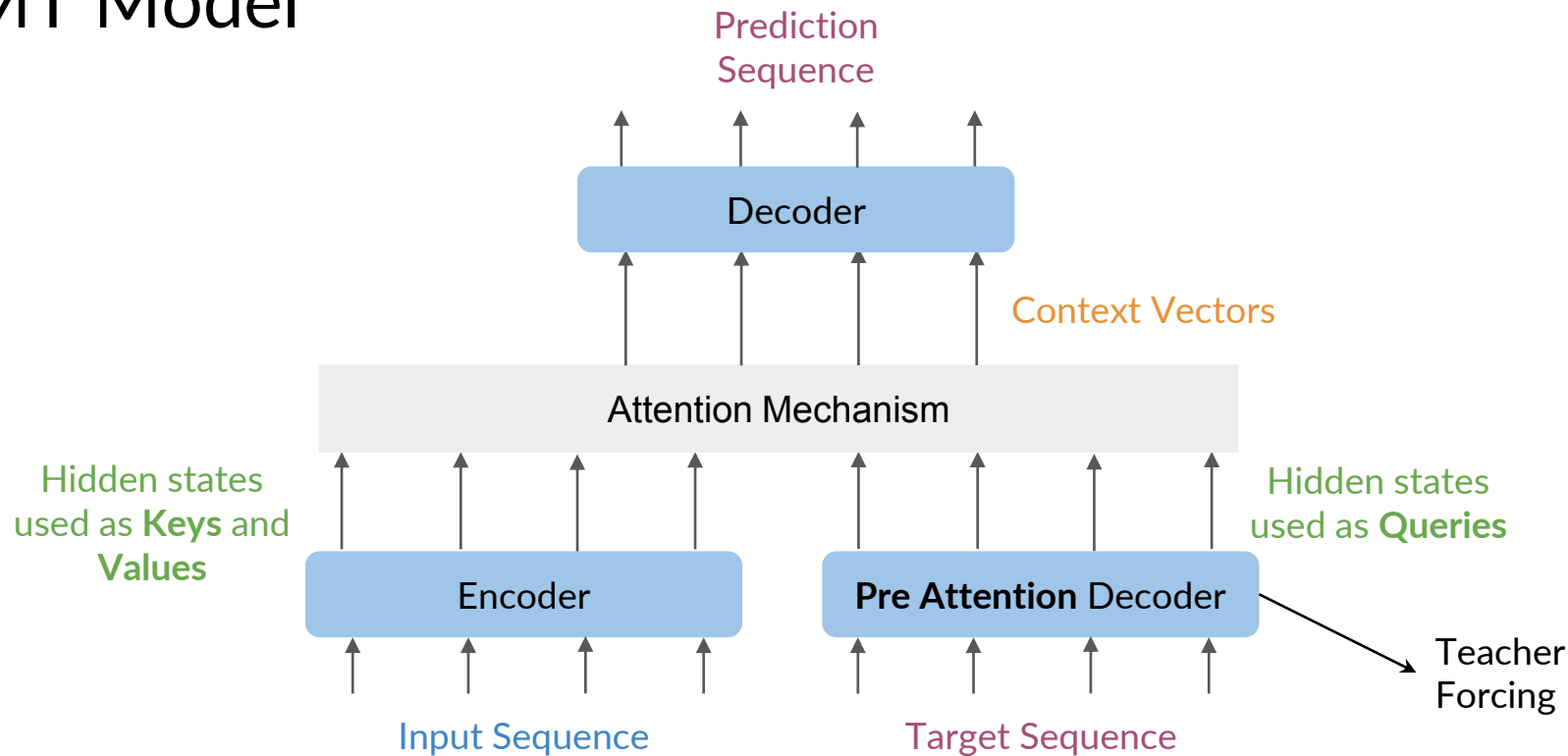
NMT Model



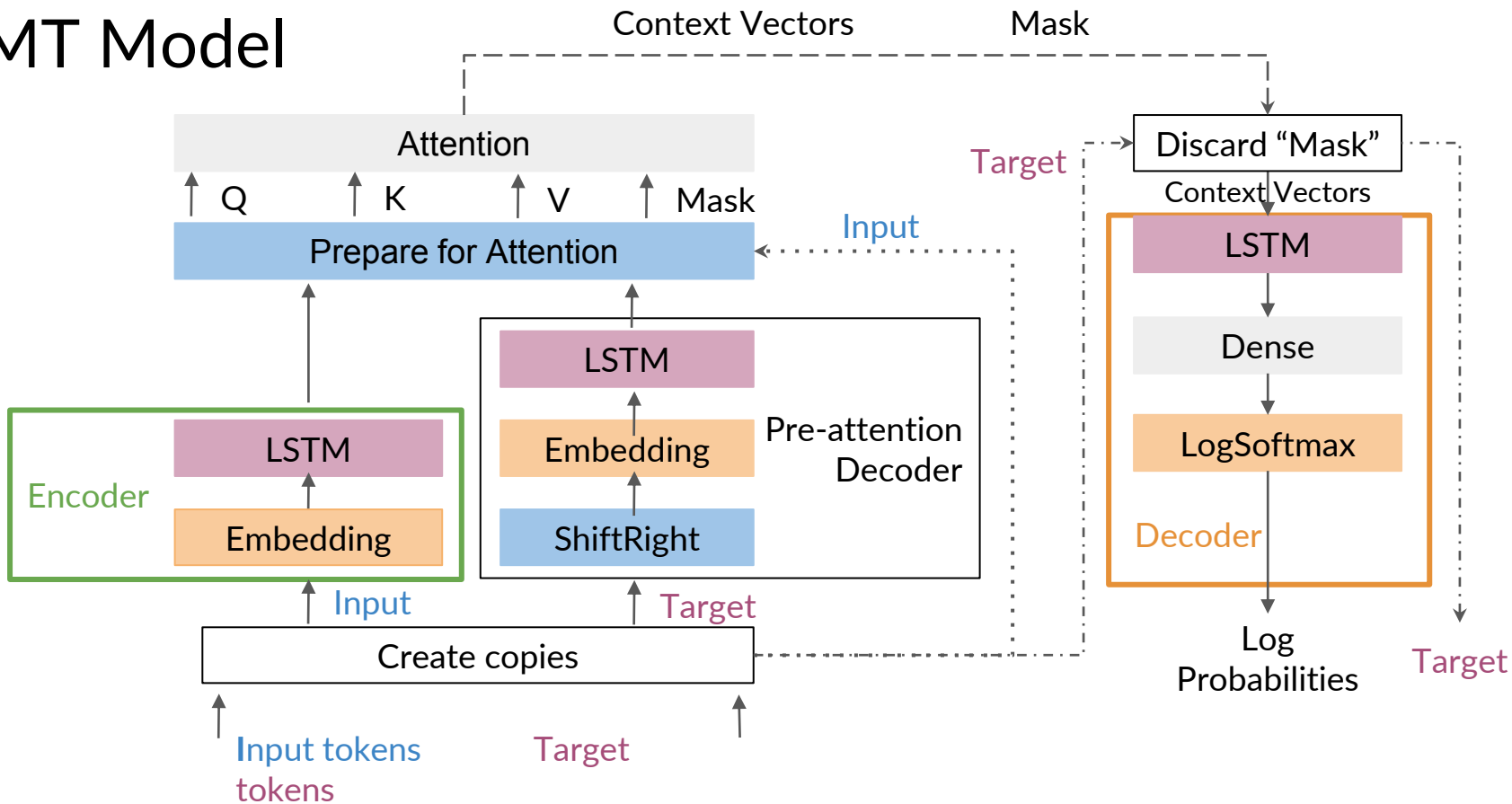
The decoder has to pass the hidden state to the Attention Mechanism

Difficult to implement, so a **pre-attention decoder** is introduced.

NMT Model



NMT Model





deeplearning.ai

BLEU Score

BLEU Score

BiLingual EValuation Understudy

Compares candidate translations to reference (human) translations

The closer to **1**, the better



BLEU Score

Candidate	I	I	am	I	
Reference 1	Younes	said	I	am	hungry
Reference 2	He	said	I	am	hungry

How many words from the **candidate** appear in the **reference** translations?

BLEU Score

Candidate	I	I	am	I	
Reference 1	Younes	said	<u>I</u>	<u>am</u>	hungry
Reference 2	He	said	<u>I</u>	<u>am</u>	hungry

Count: $\frac{1+1+1+1}{4} = 1$

A model that always
outputs common
words will do great!

BLEU Score (Modified)

Candidate	I	I	am	I	
Reference 1	Younes	said			hungry
Reference 2	He	said			hungry

Count: $\frac{1+1}{4} = 0.5$

Better than the
previous
implementation
version!

BLEU score is great, but...

Consider the following:

- BLEU doesn't consider semantic meaning
- BLEU doesn't consider sentence structure:

“Ate I was hungry because!”





deeplearning.ai

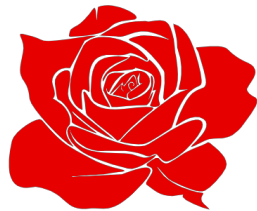
ROUGE-N Score

ROUGE

Recall-Oriented Understudy for Gisting Evaluation

Compares candidates with reference (human) translations

Multiple versions for this metric



ROUGE-N

Candidate	I	I	am	I	
Reference 1	Younes	said	I	am	hungry
Reference 2	He	said	I	am	hungry

How many words from the **reference** appear in the **candidate** translations?

ROUGE-N

Candidate	I	I	am	I	
Reference 1	Younes	said	I	am	hungry
Reference 2	He	said	I	am	hungry

$$\text{Count 1: } \frac{1+1}{5} = 0.4$$

$$\text{Count 2: } \frac{1+1}{5} = 0.4$$

ROUGE-N, BLEU and F1 score

Candidate	I	I	am	I	am
Reference 1	Younes	said	I	am	hungry
Reference 2	He	said	I	am	hungry

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \longrightarrow F1 = 2 \times \frac{\text{BLEU} \times \text{ROUGE-N}}{\text{BLEU} + \text{ROUGE-N}}$$

$$F1 = 2 \times \frac{0.5 \times 0.4}{0.5 + 0.4} = \frac{4}{9} \approx 0.44$$



deeplearning.ai

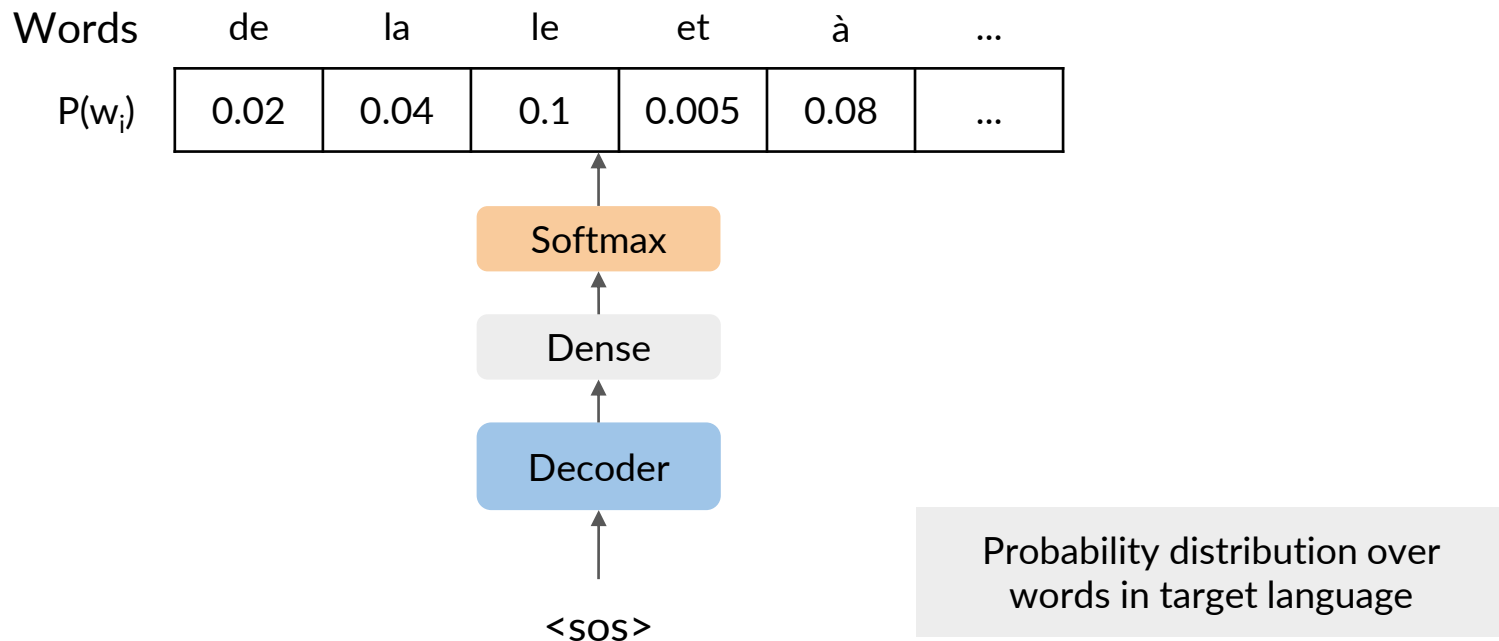
Sampling and Decoding

Outline

- Random sampling
- Temperature in sampling
- Greedy decoding



Seq2Seq model



Greedy decoding

Selects the most probable word at each step

But the best word at each step may not be the best for longer sequences...

Can be fine for shorter sequences, but limited by inability to look further down the sequence

J'ai faim.

I am hungry.

I am, am, am, am...

Random sampling

am	full	hungry	I	the
0.05	0.3	0.15	0.25	0.25

Often a little too random for accurate translation!

Solution: Assign more weight to more probable words, and less weight to less probable words.

Temperature

Can control for more or less randomness in predictions

Lower temperature setting : More confident, conservative network

Higher temperature setting : More excited, random network





deeplearning.ai

Beam Search

Beam search decoding

Most probable translation is **not** the one with the most probable word at each step



Solution

Calculate probability of multiple possible sequences



Beam search

Beam search decoding

Probability of multiple possible sequences at each step

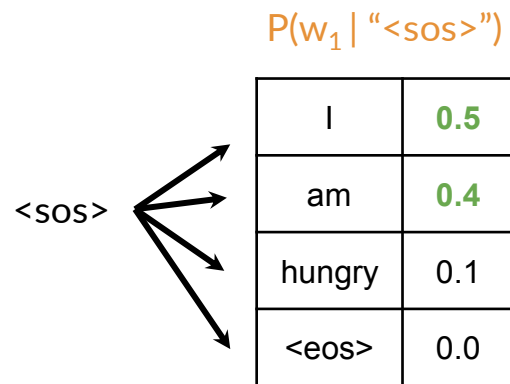
Beam width B determines number of sequences you keep

Until all B most probable sequences end with $\langle \text{EOS} \rangle$

Beam search with $B=1$
is **greedy decoding**.

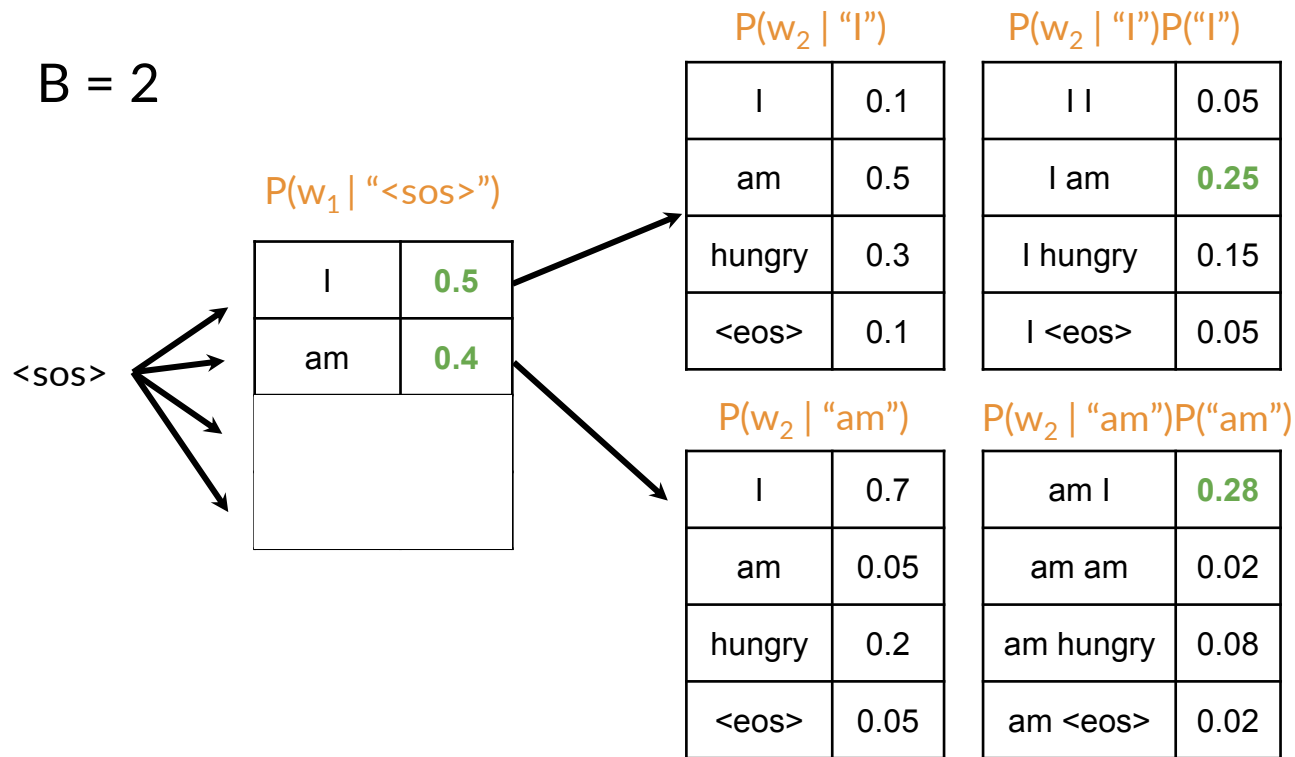
Beam search example

$B = 2$



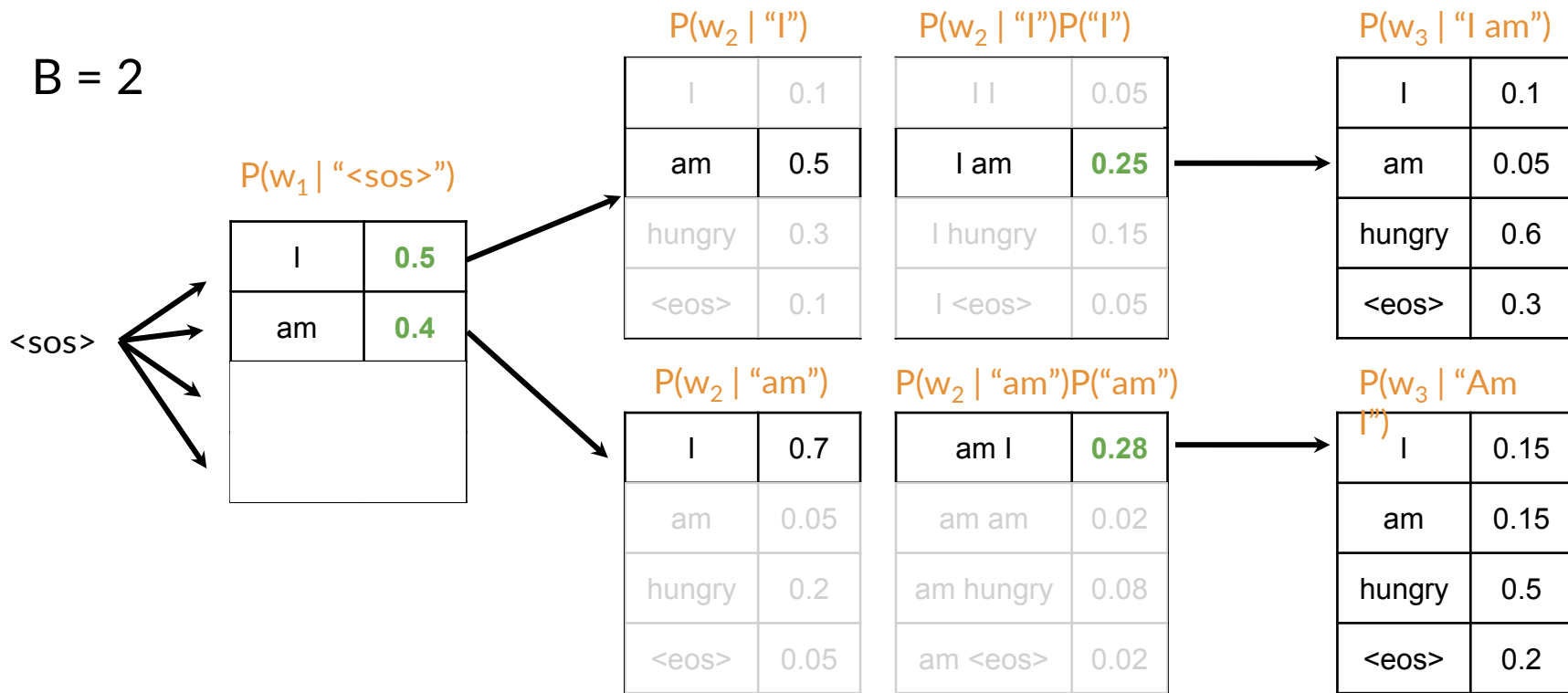
Beam search example

$B = 2$

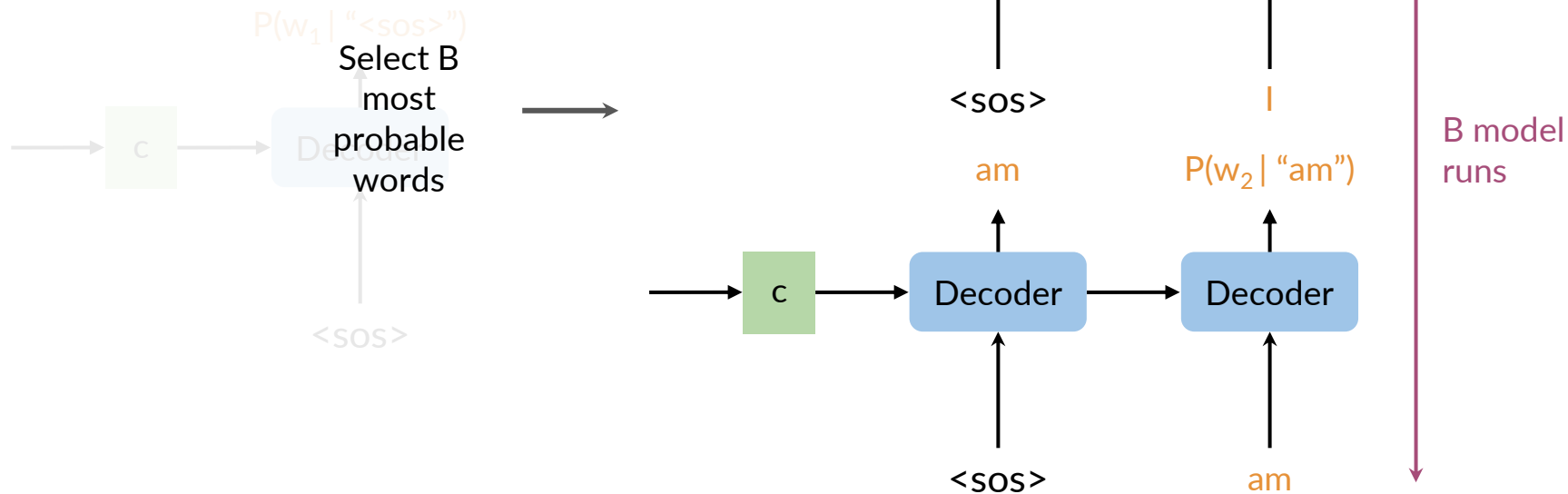


Beam search example

$B = 2$



Beam search decoding



Problems with beam search

Penalizes long sequences, so you should normalize by the sentence length

Computationally expensive and consumes a lot of memory



deeplearning.ai

Minimum Bayes Risk

Minimum Bayes Risk (MBR)

- Generate several candidate translations
- Assign a similarity to every pair using a similarity score (such as ROUGE!)
- Select the sample with the highest average similarity

Minimum Bayes Risk (MBR)

$$\arg \max_E \frac{1}{n} \sum_{E'} \text{ROUGE}(E, E')$$

Find the candidate translation that maximizes

Compare with every other candidate

ROUGE score between pair of candidates

The diagram illustrates the Minimum Bayes Risk (MBR) formula. The formula is presented as $\arg \max_E \frac{1}{n} \sum_{E'} \text{ROUGE}(E, E')$. The terms are annotated with colored boxes and arrows: a green box around $\arg \max_E$ points to the text 'Find the candidate translation that maximizes'; a purple box around $\sum_{E'}$ points to the text 'Compare with every other candidate'; and an orange box around $\text{ROUGE}(E, E')$ points to the text 'ROUGE score between pair of candidates'.

Example: MBR Sampling

ROUGE(C_1, C_2)

ROUGE(C_1, C_3)

ROUGE(C_1, C_4)

Compute average ROUGE

$$R_1 = \frac{1}{3} \sum_{i \neq 1} \text{ROUGE}(C_1, C_i)$$

Repeat for every candidate

Select the candidate with the highest average

R_1

R_2

R_3

R_4

Summary

- Compare several candidate translations
- Choose candidate with highest average similarity
- **Better performance** than random sampling and greedy decoding

