

Sequence to Sequence Modelling with RNNs

(Encoder-Decoder Architecture, Attention Mechanism)



Tirtharaj Dash

Dept. of CS & IS and APPCAIR
BITS Pilani, Goa Campus

November 11, 2021

Machine Translation

In the last lecture :

- ① RNN "cell"
- ② Vanishing gradients
 - GRU
 - LSTM

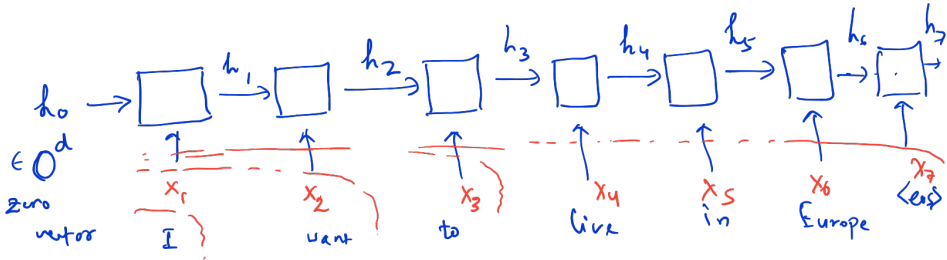
Seq-to-Seq Modelling: (Ref: TBI, Andrew Ng's course)

Machine Translation

English \rightarrow French

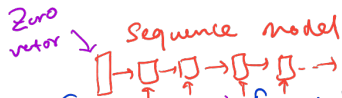
"I want to live in Europe" \rightarrow "Je veux vivre en Europe"

\hookrightarrow Summary \rightarrow (French)



Steps:

- ① Construct a RNN arch. to process the English sentence (creating a summary for the input)

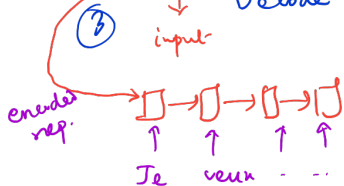


② Encoded representation for the input

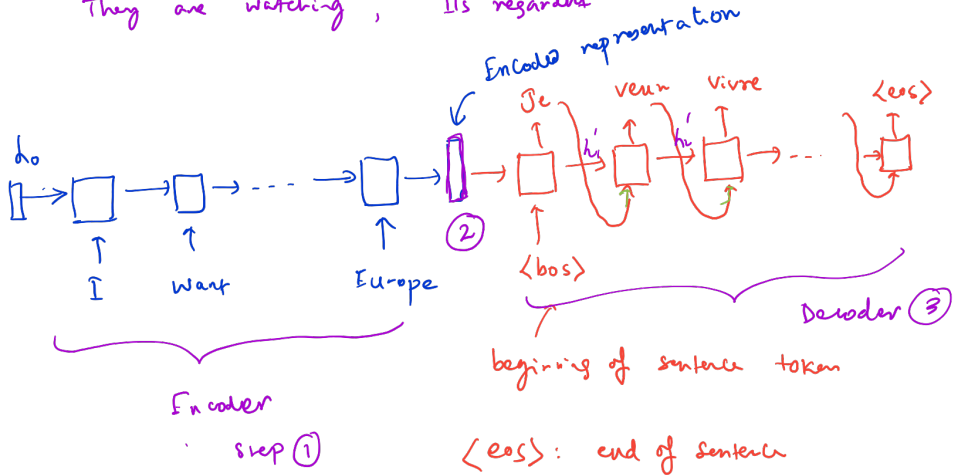


input

Decode the encoded representation into the target (output) language.

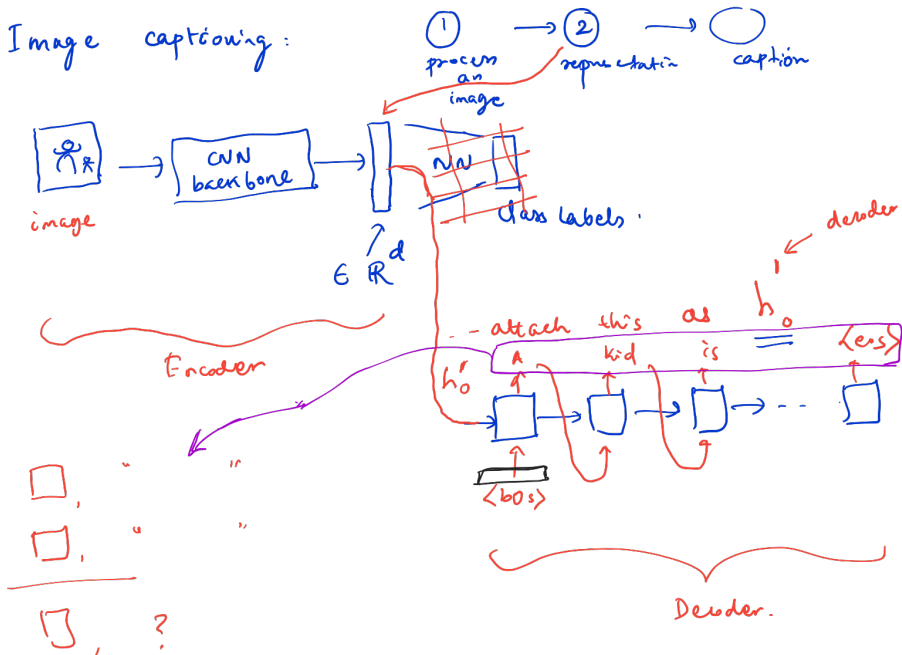


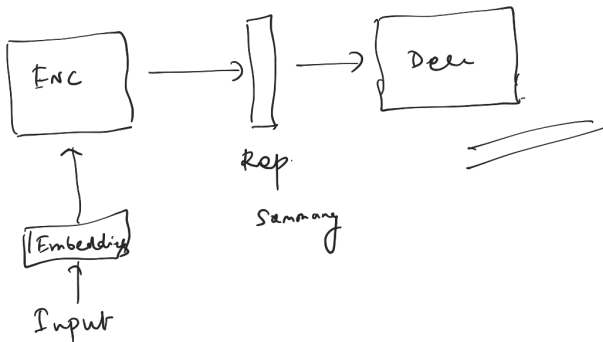
"I want to live in Europe", "Je veux vivre en Europe"
 "They are watching", "Ils regardent"



<eos>: end of sentence

Image captioning:





Machine Translation

as a Conditional Language Model.

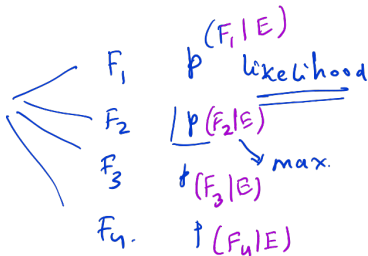


↓ τ

$$P(x_1 x_2 \dots x_T) = \prod_{t=1}^T P(x_t | x_1 \dots x_{t-1})$$

$P(\text{French sentence} \mid \text{English sentence})$

Given any input E
(in English)



$$\arg \max_{y_1, y_2, \dots, y_{T_y}} p(y_1, y_2, \dots, y_{T_y} \mid \underline{x})$$

English sentence
 x_1, x_2, \dots, x_{T_x}

$$\arg \max_i p(F_i \mid \underline{x})$$

→ Idea: most likely : Greedy search

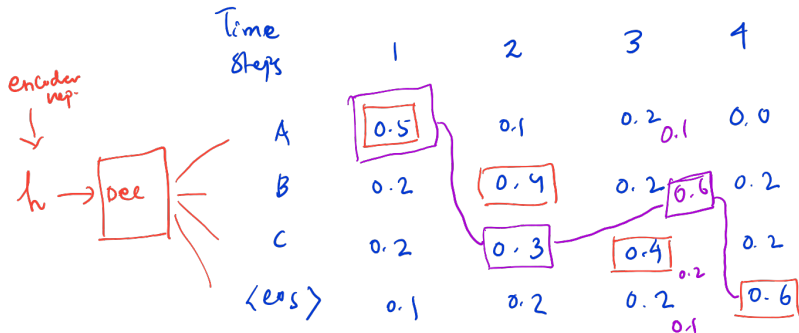
F_1 — 0.45 (optimal) ?

F_3 — 0.42

optimal ^{"Translation"} sequence is not guaranteed.

Let we have 4 output tokens

"A" "B" "C" "<eos>"



Greedy search: $P(ABC<eos> | h) = 0.5 \times 0.4 \times 0.4 \times 0.6 = 0.048$

$P(ACB<eos> | h) = 0.5 \times 0.4 \times 0.6 \times 0.6 = 0.072$

Solution:

Exhaustive search. X

Greedy
search

$$\frac{10000 \times 10}{}$$

$$O(|Y| T_y)$$

→ F_1

→ F_2

→ F_3

→ \vdots

\vdots

\vdots

\vdots

→ F_M



pick the one with best score

$|Y|$ # of output tokens = 10000

Max. length of the translation = 10

T_y

$$10000^{10}$$

$$O(|Y|^{T_y})$$

compute a
(Fast)



$K \downarrow$



$K \uparrow$

Accurate



greedy
search

Beam search.

Exhaustive
search

$$O(|Y| T_y)$$

$$O(\underset{\substack{\uparrow \\ \text{beam width}}}{K} |Y| T_y)$$

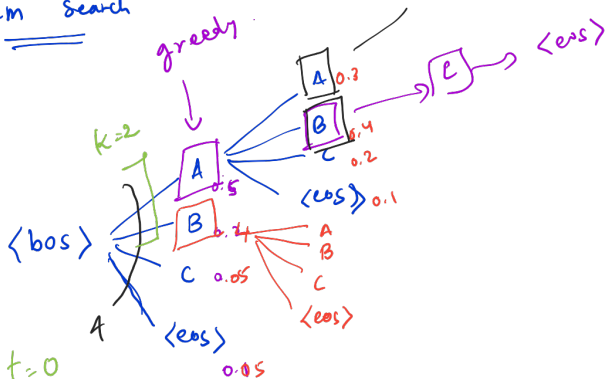
beam width.

$$O(|Y|^{T_y})$$

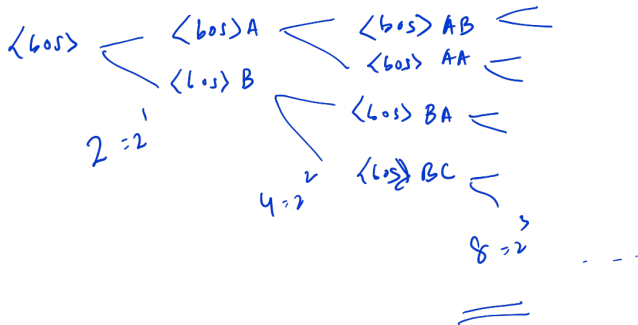
homework

Beam Search

"ABC <eos>"



$k=2$



Is it $O(k |Y| T_y)$?

Homework

Read about Beam search & How
the architectures of the Decoder will

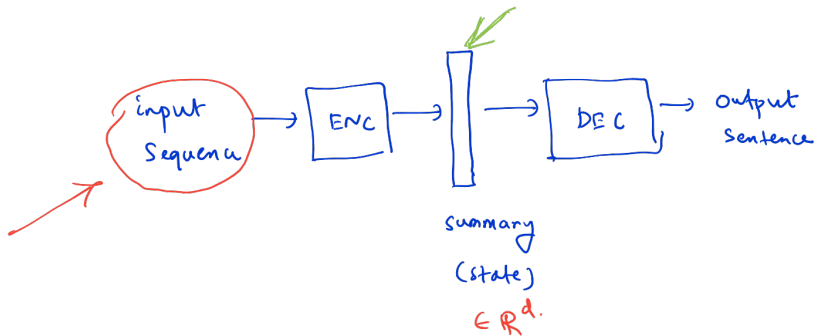
change to implement
Beam search ??

ATTENTION MECHANISM

Seq 2 seq learning.

2014

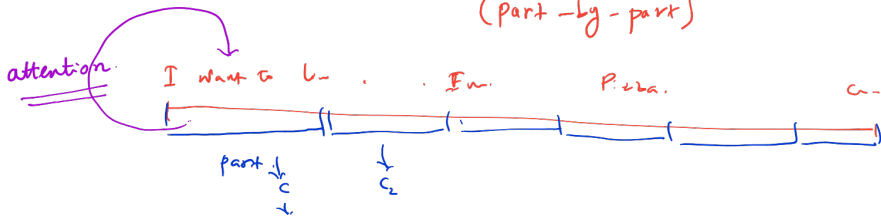
Machine translation.



Issues: Long sentences.

Solution: ① Bidirectional RNN $\begin{cases} \text{ENC} \\ \text{DEC} \end{cases}$

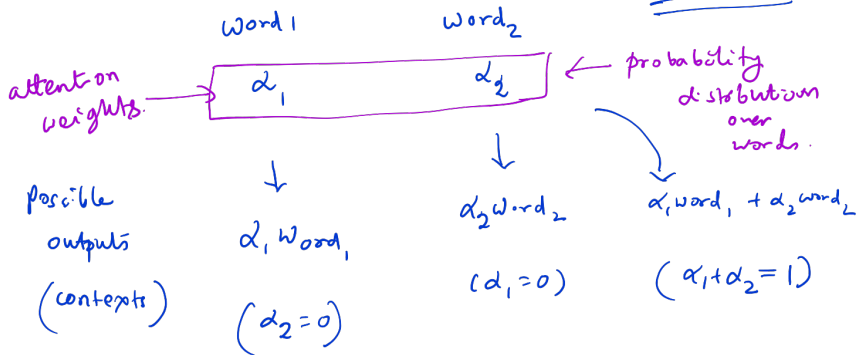
② Read \rightarrow summarise \rightarrow the input
(part-by-part)

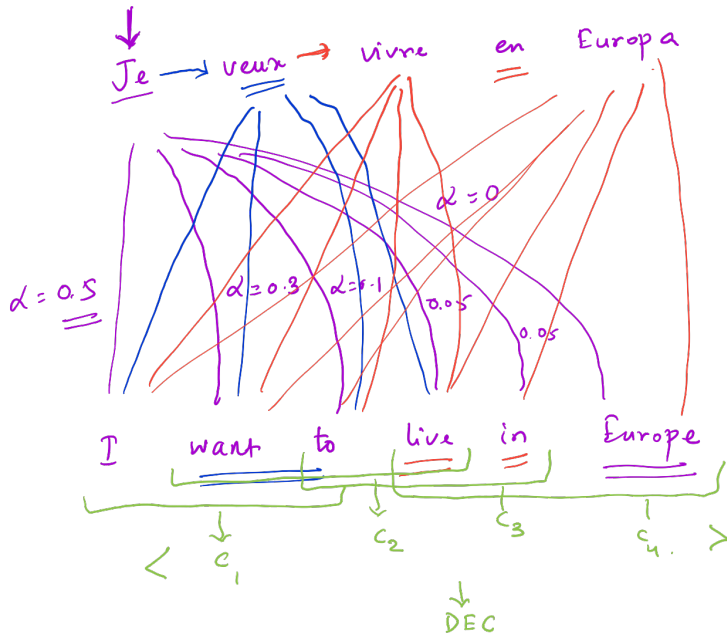


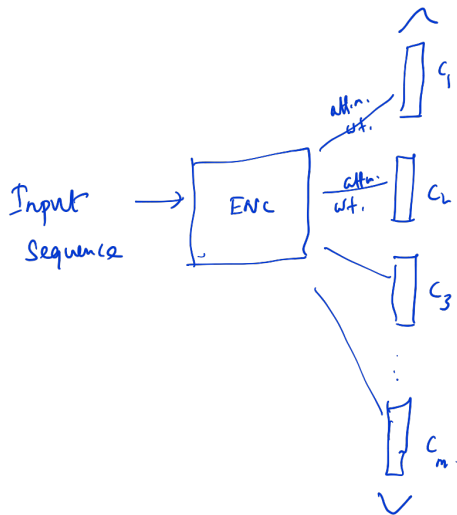
attention score : core of attention mechanism.

consider 2 words

$$\underline{\underline{\alpha_1 + \alpha_2 = 1}}$$







→ Input for the
Decoder
=.

Doubt

Bidirectional RNN



inputs
 \uparrow
 $[\vec{h}, \overleftarrow{h}]$

