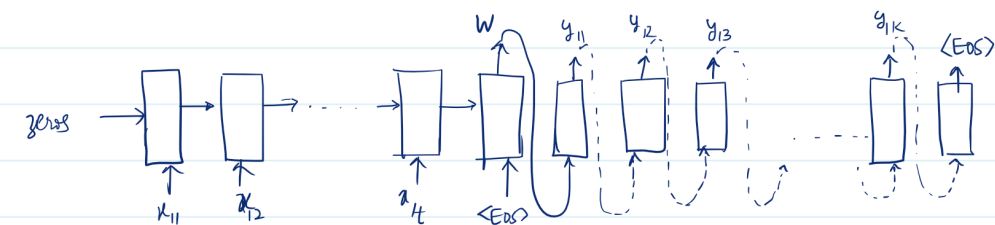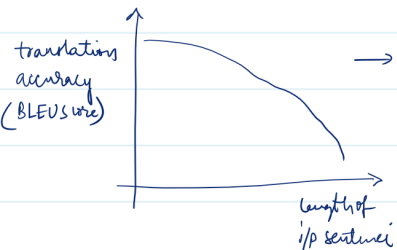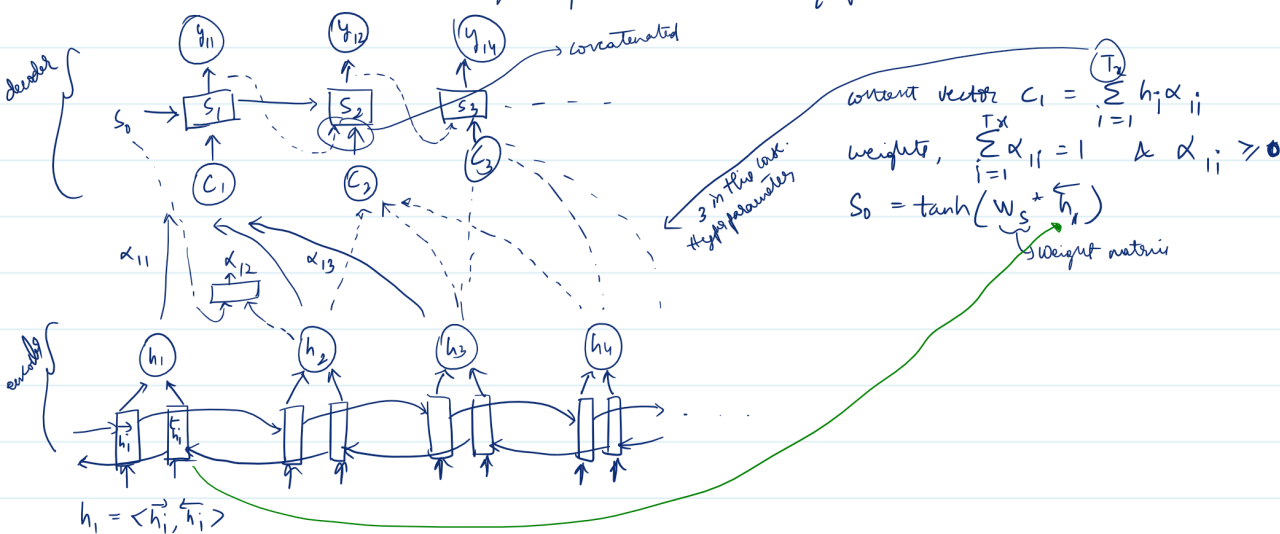# Attention Based Models :-



→ W on it's own is not able to capture the essence of the entire sentence if it's long.



→ for simple encoder - decoder models.

→ We use bidirectional RNN because it might depend on words before / after a word in encoder. Decoder is unidirectional RNN



content vector $C_i = \sum_{i=1}^{T_x} h_j \alpha_{ji}$

weights, $\sum_{i=1}^{T_x} \alpha_{1i} = 1$ & $\alpha_{1i} \geqslant 0$

$S_0 = \tanh \left( W_s * \overleftarrow{h_1} \right)$

↳ weight matrix

3 in this context. Hyperparameter

$h_1 = \langle \overrightarrow{h_1}, \overleftarrow{h_1} \rangle$

How to compute $\alpha_{ij}$'s :-

always the /o.

$\alpha_{ij} = \dfrac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$

denom will be same so sum = 1

$e_{ij} = a\left( S_{i-1}, h_j \right)$

prev s & current h

ex :- $e_{12} = a\left( S_0, h_2 \right)$

Attention model (how much attention needs to be given to encoder). We can use a simple feed forward NN for this.
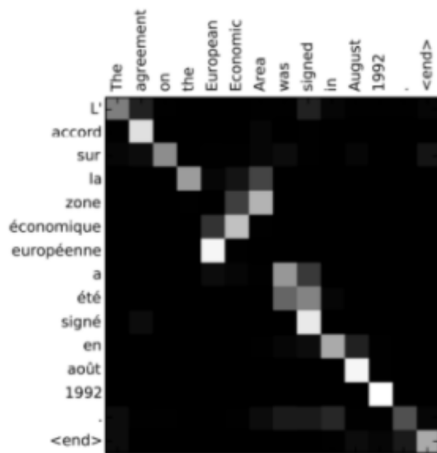
→ Train using Backprop through time using adam & let it converge.
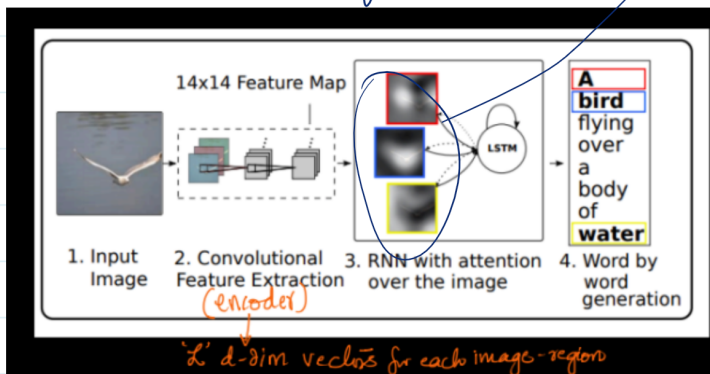
RNNSearch = Attention Model.



→ $T_x$

Drawback :- Time Complexity $= O\left( K_1 \cdot K_2 \right)$

length of O/p.

length of I/p

→ We can visualize $a_{ij}$'s like this to debug. Accuracy Measurement can be done by using metrics like BLEU score.



→ Image caption generation using Visual Attention :—

→ every region is basically like a word. if region is white ⟹ $\alpha$ of that region is large.



1. Input Image
2. Convolutional Feature Extraction (encoder)
3. RNN with attention over the image
4. Word by word generation

'L' d-dim vectors for each image-region

Visualizing $a_{ij}$'s:



(Can also be used as localization)