

Attention

- * Attention is based on Encoder-Decoder network structure or sequence-to-sequence (seq2seq) modelling.
- Before delving into seq2seq, it is a good idea to briefly talk about word embedding and word2vec.

- ① In NLP, if we want to plug in words into a deep network we need a way to turn words into numbers.
- * We can use random numbers or one-hot vectors to represent words → But this way doesn't capture relationship between words: good ↔ bad, perfect ↔ awesome, and etc.
- * Assigning a long vector to a number helps NN to capture different meaning of a word in different contexts:

This phone is great! ←→ my phone is broken, great!

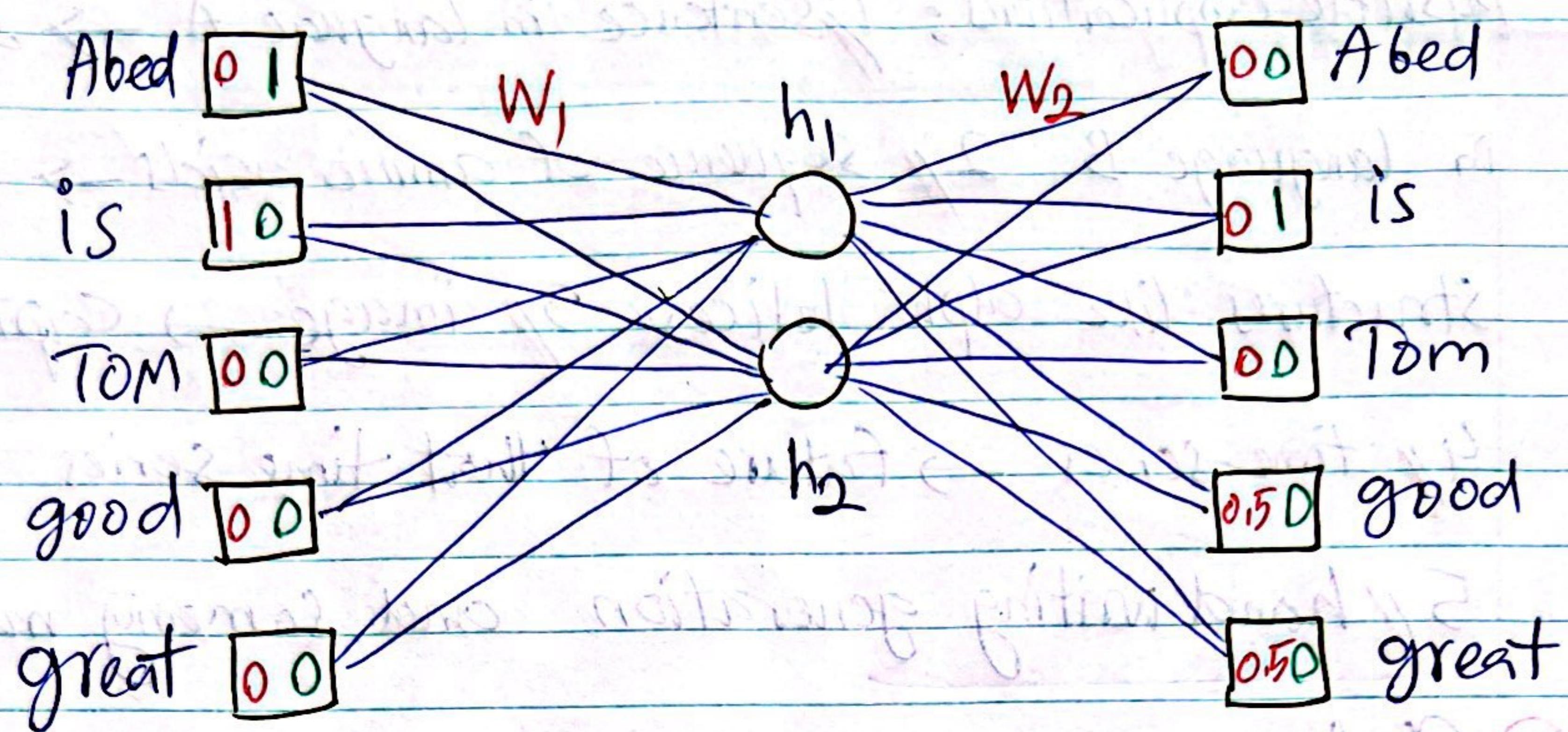
- * consider that we have this dataset:

1. Abed is great

2. Tom is good

We want to adjust weights w_1 and w_2 in a

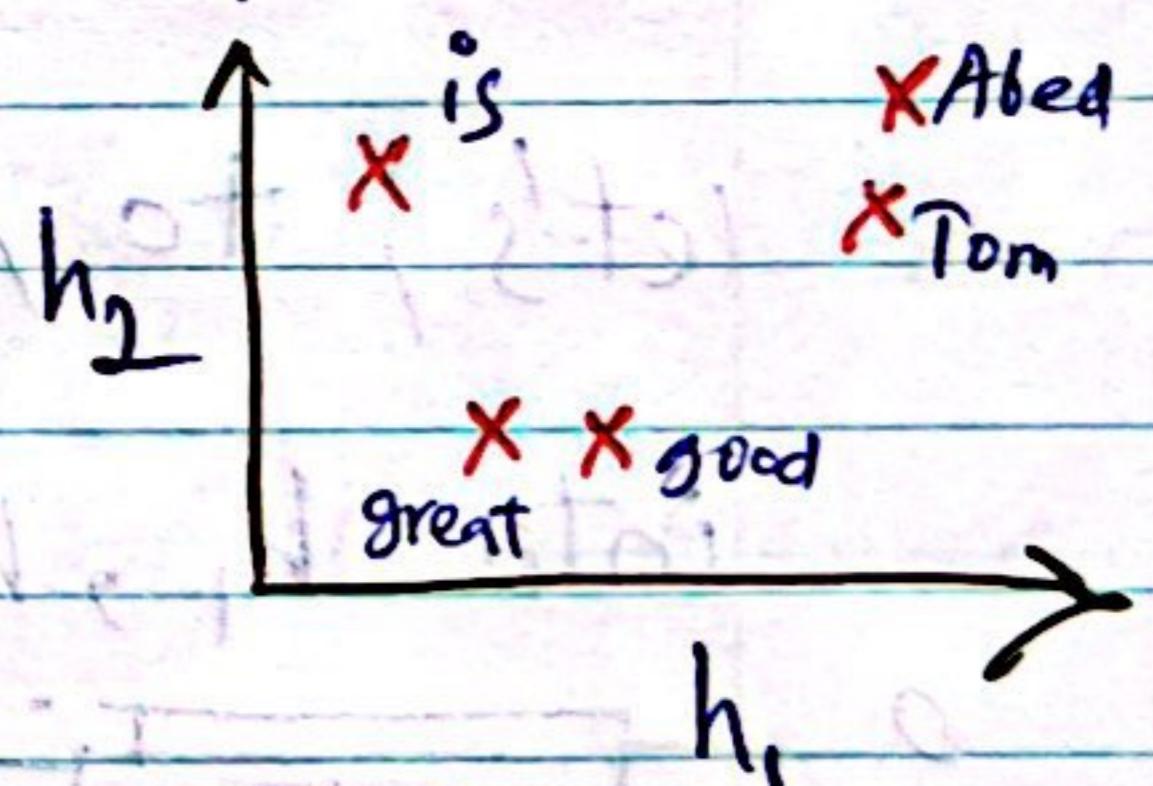
in



way that for each input word, the network predicts

the next word. Using backpropagation, we optimize W_1 and W_2

and the result seems like this.



Although, similar words are embedded

near to each other, this is not a perfect way to embed

words. We have these methods,

Continuous Bag of Words improves the context in vector h

by using the surrounding words to predict what occurs in the middle

Skip Gram improves the context in h by using the word in the middle to predict its surroundings

- Now, we will move to seq2seq & these are possible applications:
- 1 // Sentence in language A \rightarrow sentence in language B
 - 2 // sequence of amino acids \rightarrow 3D structures like alpha helices.
 - 3 // image \rightarrow caption
 - 4 // time-series \rightarrow future of that time series
 - 5 // hand writing generation and so many more!

* Challenge 3:

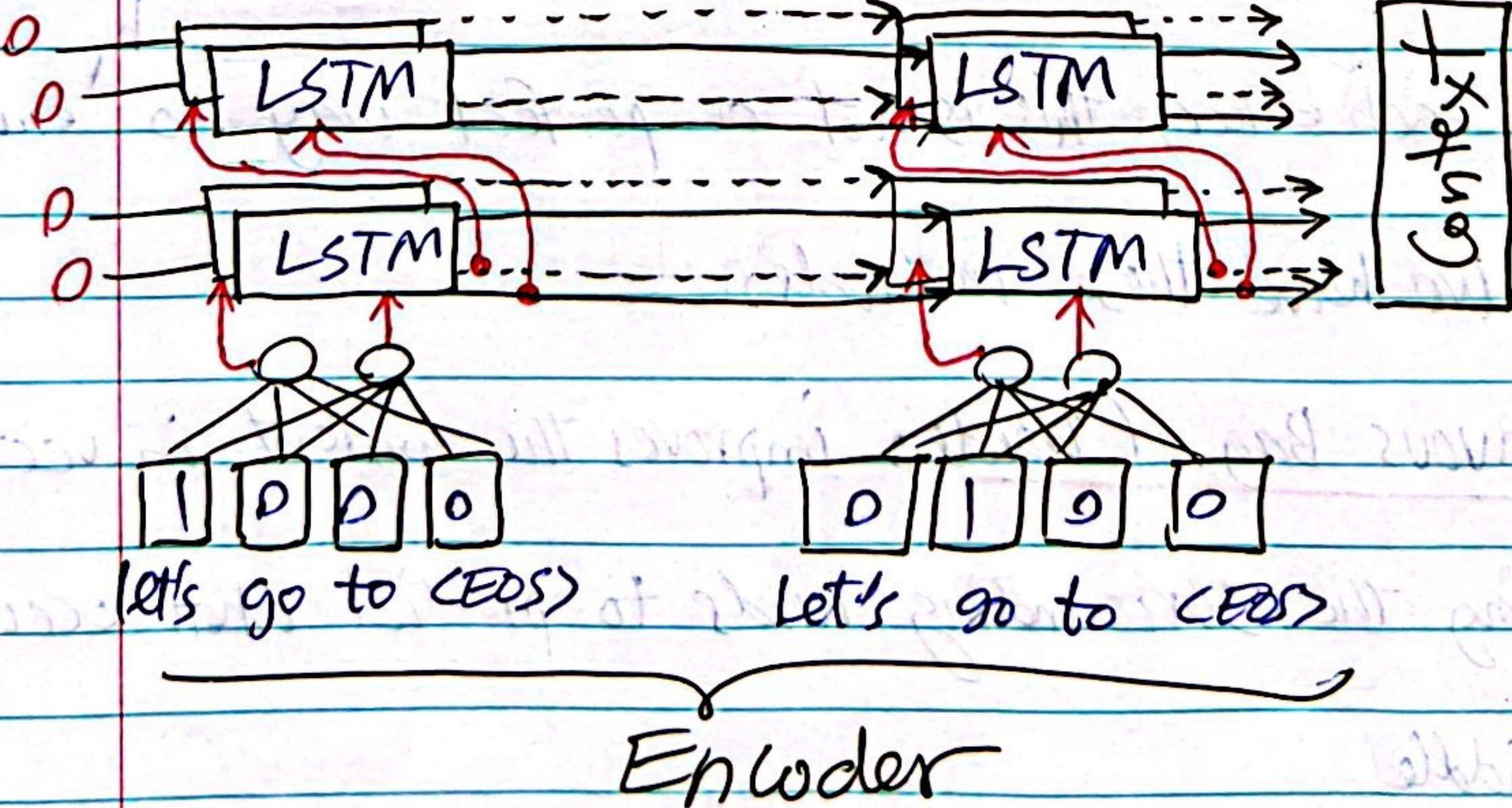
Different length of inputs/outputs \rightarrow RNN/LSTM

Can handle this! Let's say, we have 4 tokens

let's, to, go, <EOS> and we want to embedd

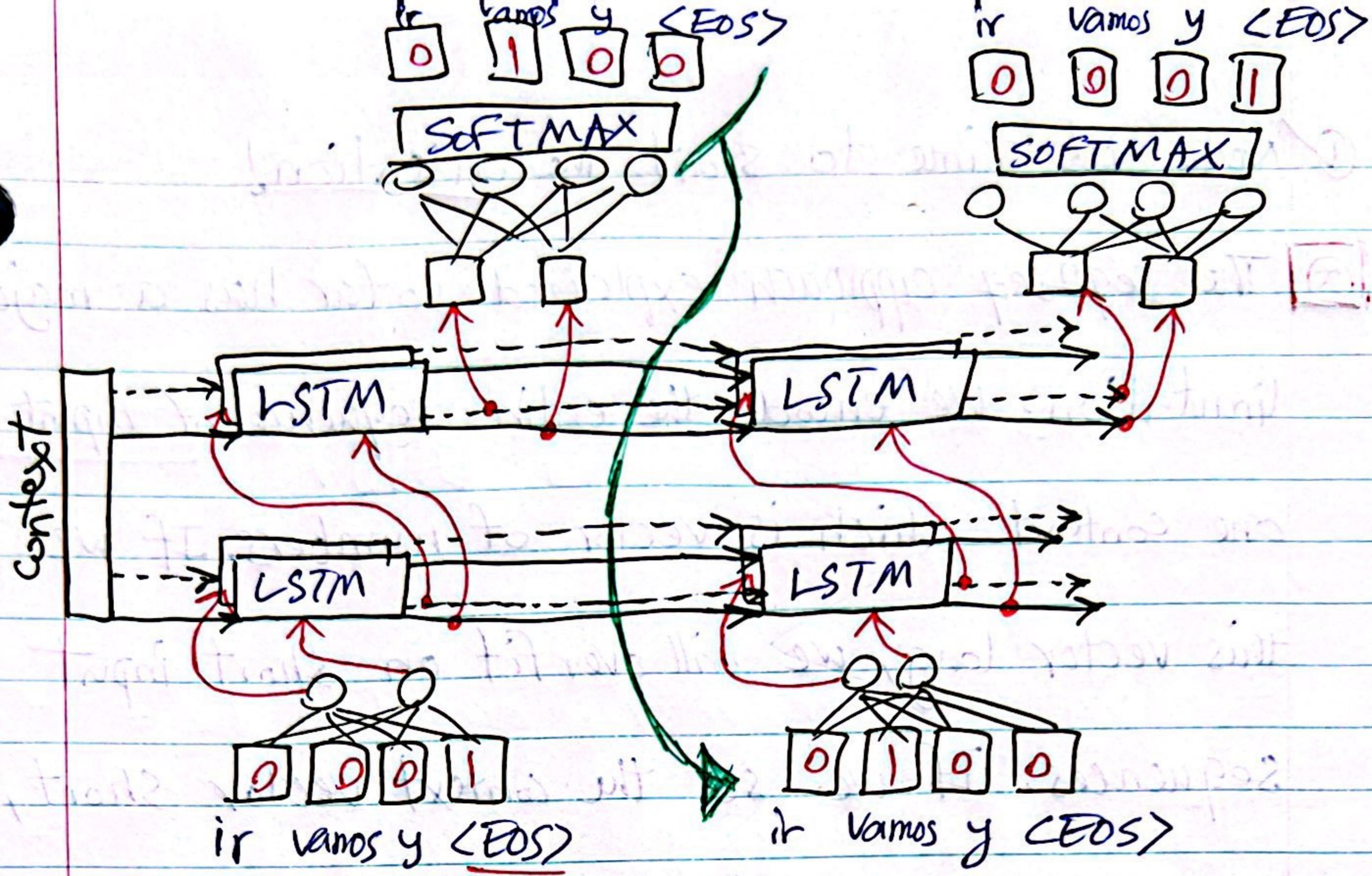
into h_1, h_2 !

layers of LSTM



Now we feed context to decoders. Decoder has

the structure of LSTMs with different weights.



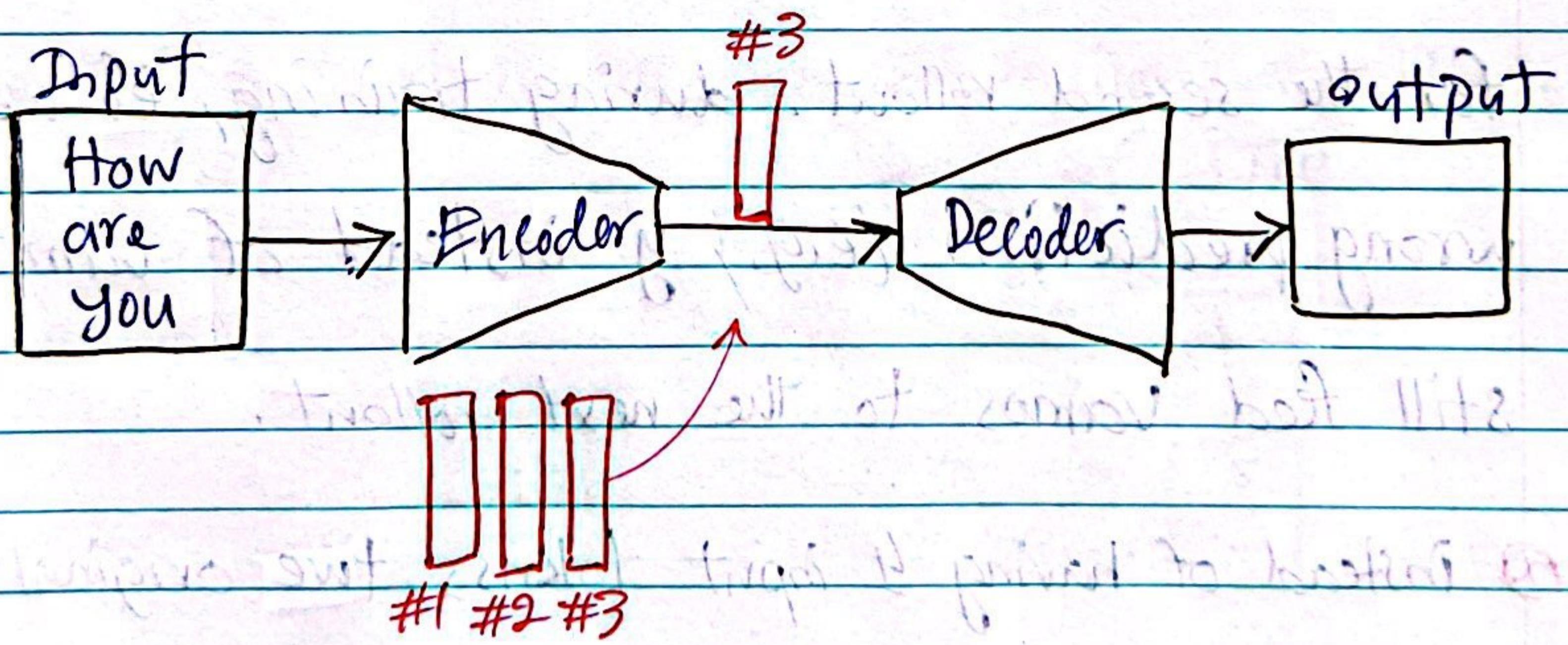
always we start with
EoS to start the decoder

④ As you see, the predicted word vamos is used as input for the second rollout. during training, if we have a wrong prediction (e.g., y instead of vamos) we still feed vamos to the next rollout.

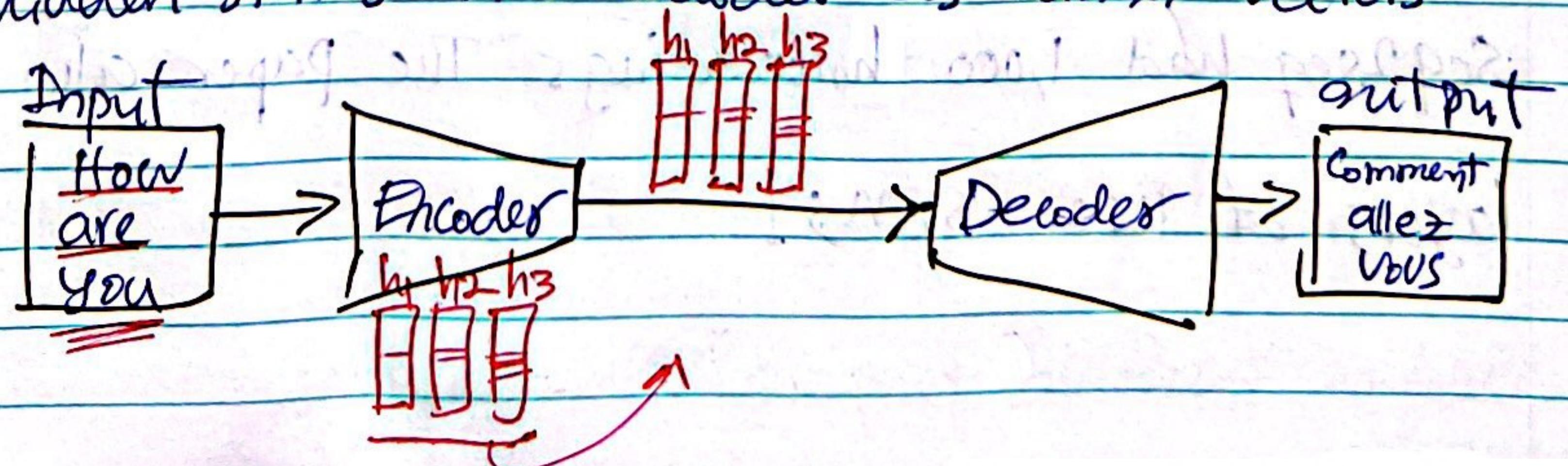
④ Instead of having 4 input tokens, the original paper has 160,000 tokens! The output vocabulary had 80,000 tokens. instead of having two embedding values, the original seq2seq had 1,000 embeddings. The paper also had 4 layers of 1000 LSTMs!

① Now, it's time to start the attention!

② The seq2seq approach explained so far has a major limitation. We encode the entire sequence of input into one context which is vector of numbers. If we set this vector long, we will overfit on short input sequences. If we set the context vector short, we cannot encode large sentences and we may forget concepts at the start of the input when we reach the end of sentence!



However in attention-based models we feed all of the hidden states into decoder as context vectors



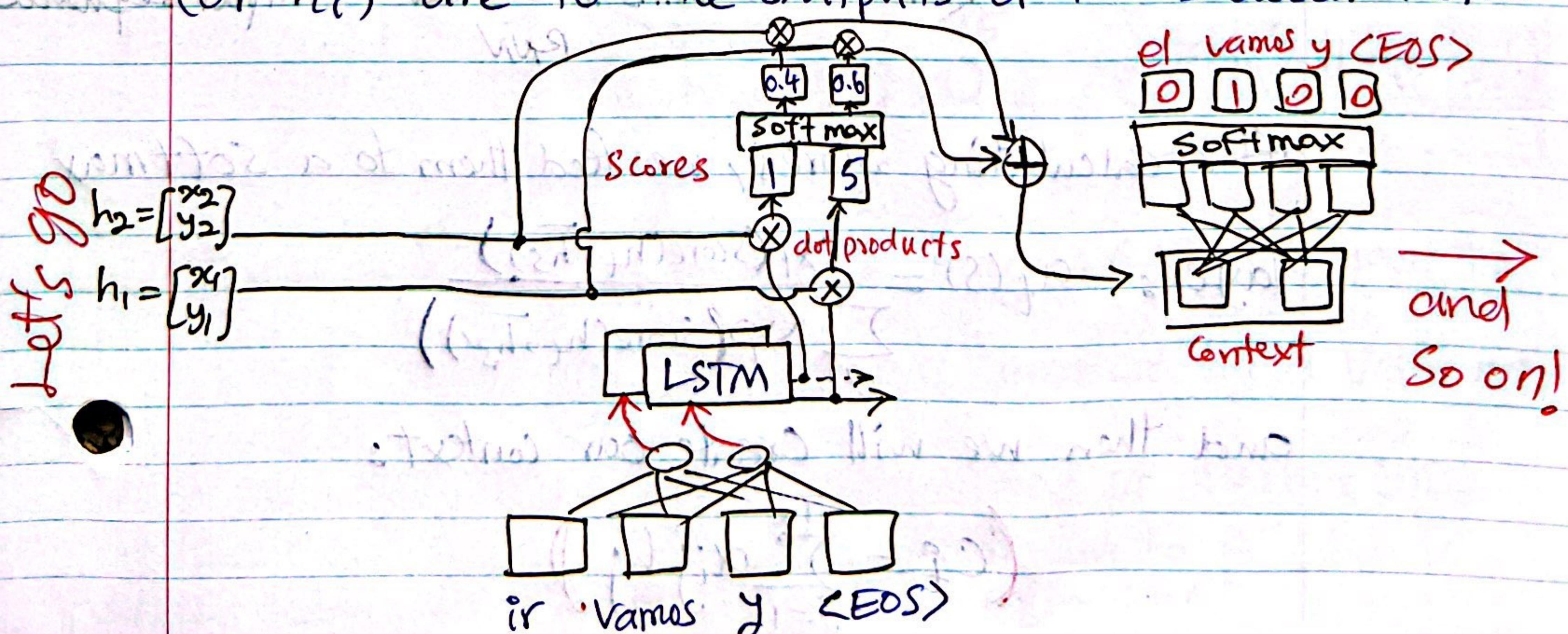
→ In this way, longer inputs have larger contexts!

another benefit is that h_i captures the essence of i^{th} element of the input sequence → we have track of all elements of input and can use them selectively!

- ① In the decoder side, we calculate a weighted sum of h_i to create context for each step → The decoder will be initialized by $\langle \text{EOS} \rangle$ input and generates the first context and then first output → The first output will be used as the input for the second step and generates new scores for h_i and new context!

- * In a lay language, the thing that attention does is to determine how similar the outputs of Encoder LSTMs (cosine loss!) ↔ dot product!!

(or h_i) are to the outputs of the Decoder LSTM



So, the whole problem is to calculate score.

functions

$$\text{Score}(h_t, \bar{h}_S) = \text{Score}(\boxed{h_t}, \boxed{\boxed{h_1 h_2 h_3}}) = \boxed{s_1 s_2 s_3}$$

target & source
(decoder) (encoder)

The length of score vector is equal to the # of h_i

In \bar{h}_S . a simple way to calculate s_i is dot product of h_t to all h_i . In other words,

$$\boxed{\text{Score}(h_t, \bar{h}_S) = h_t^T \bar{h}_S}$$

another method is the general method is

$$\boxed{\text{Score}(h_t, \bar{h}_S) = h_t^T W_a \bar{h}_S}$$

of hidden units in decoder # of hidden units in RNN # of hidden units in encoder (input sequence RNN) Ixm mxn nxr → number of elements in the input sequence

after calculating scores, we feed them to a softmax

$$\text{layer: } a_t(s) = \frac{\exp(\text{Score}(h_t, \bar{h}_S))}{\sum_{s'} \exp(\text{Score}(h_t, \bar{h}_{s'}))}$$

and then we will create our context:

$$(c_i = \sum_{j=1}^{Tx} \alpha_{ij} h_j)$$

34

and then we concatenate our context vector c_t with

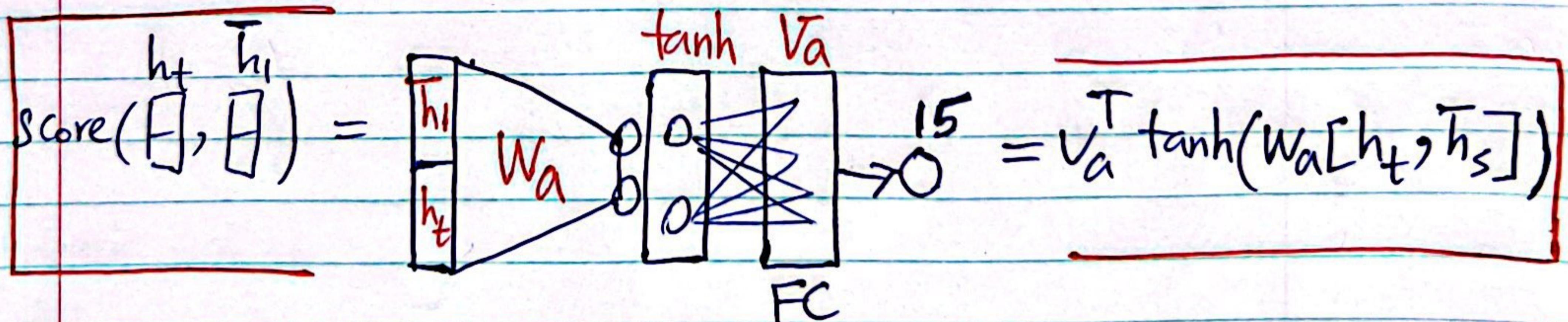
Context + hidden vector

① then we feed this to a fully connected layer

$$*\boxed{\tilde{h}_t = \tanh(W_c[c_t; h_t])}*$$

the output of this network (\tilde{h}_t) determines the output of our decoder sequence!

Another approach to calculate score between h_t and \bar{h}_i is concat method! For instance, we want to calculate the score between h_t (e.g., h_4) of decoder and \bar{h}_i in encoder.



④ Computer Vision applications: image to caption (image to context to text!) The concept is created by a VGG net.

in one layer, it produces feature map of $14 \times 14 \times 512$.

We have 512 different hidden states $\leftarrow \mathbf{h}_s = \begin{matrix} 196 \\ \text{grid} \\ 512 \end{matrix} \leftarrow 196 \times 1$ turn to vector