



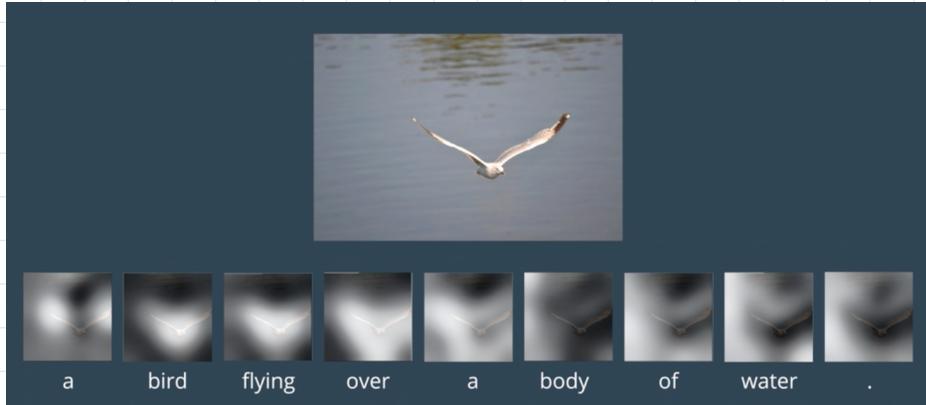
Attention

Introduction to Attention

"One important property of human perception is that one does not tend to process a whole scene in its entirety at once. Instead humans focus attention selectively on parts of the visual space to acquire information when and where it is needed, and combine information from different fixations over time to build up an internal representation of the scene, guiding future eye movements and decision making."

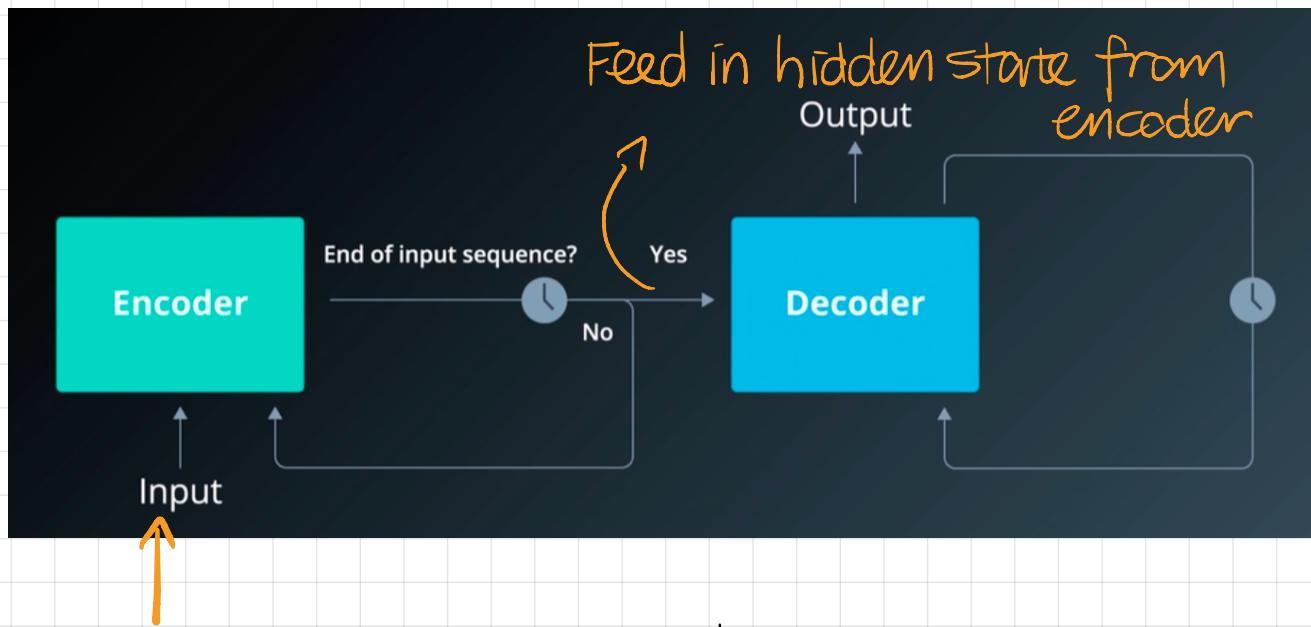
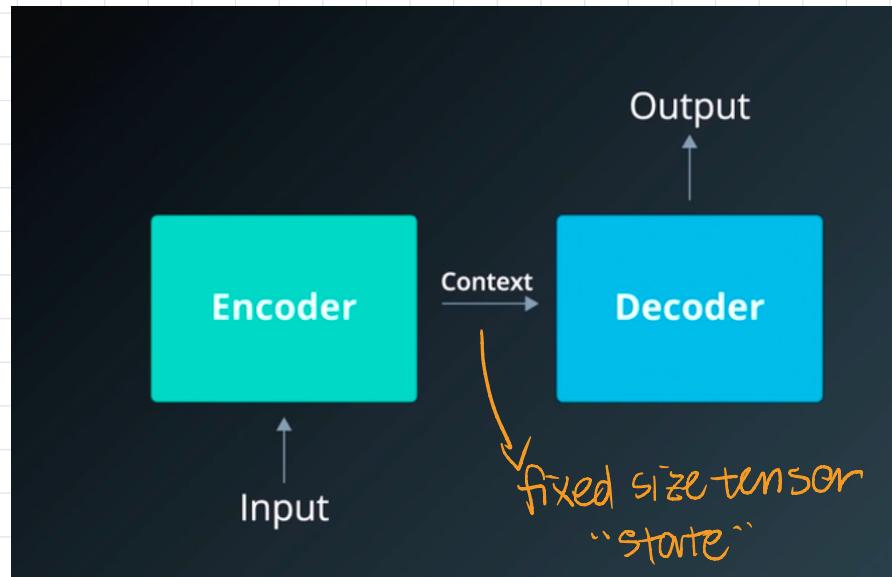
— Recurrent Models of Visual Attention, 2014

- Eyetracking device
- Classification of images
 - CNN looks at entire image to be able to class -ify.

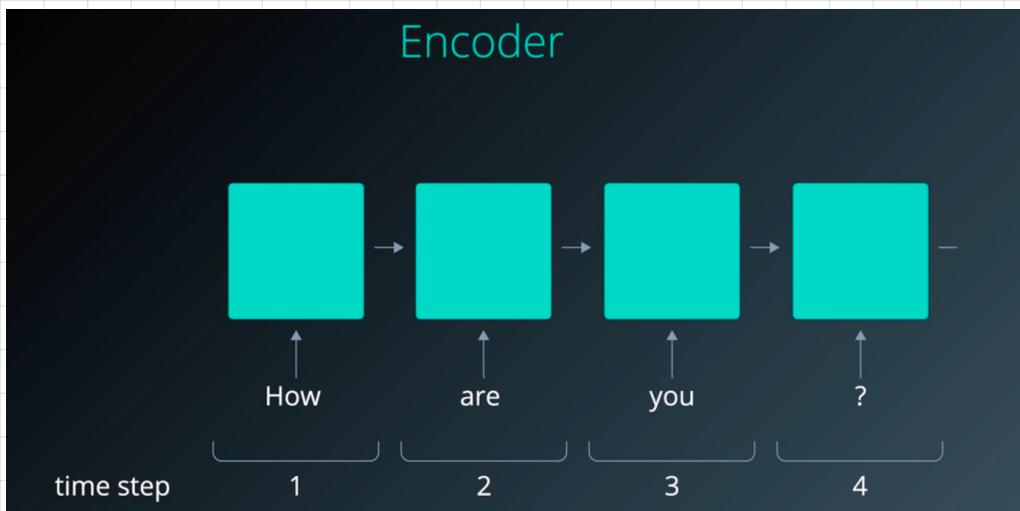


- Machine neural translation
 - Classic seq to seq models uses the entire sentence input to produce every single small output
 - Attention can use small parts of input to start translate.

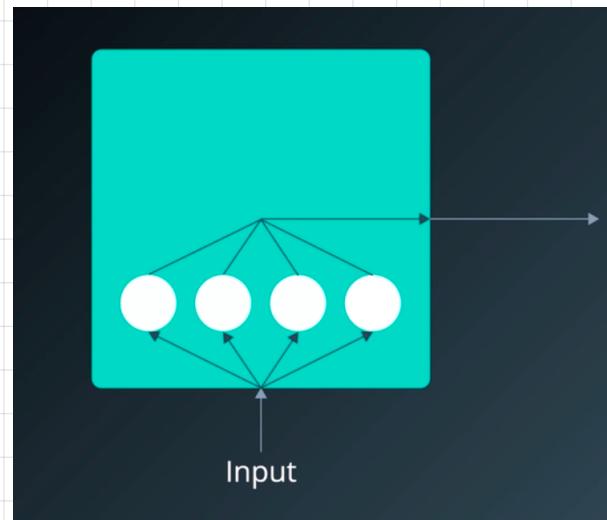
Encoders & Decoders



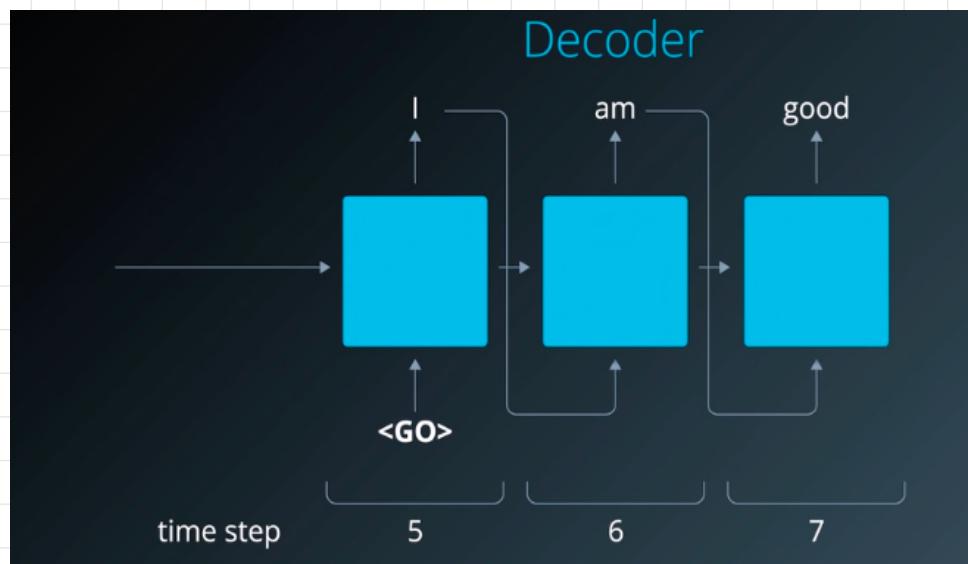
How are you? → 4 tokens, will take the RNN
4 timesteps to read



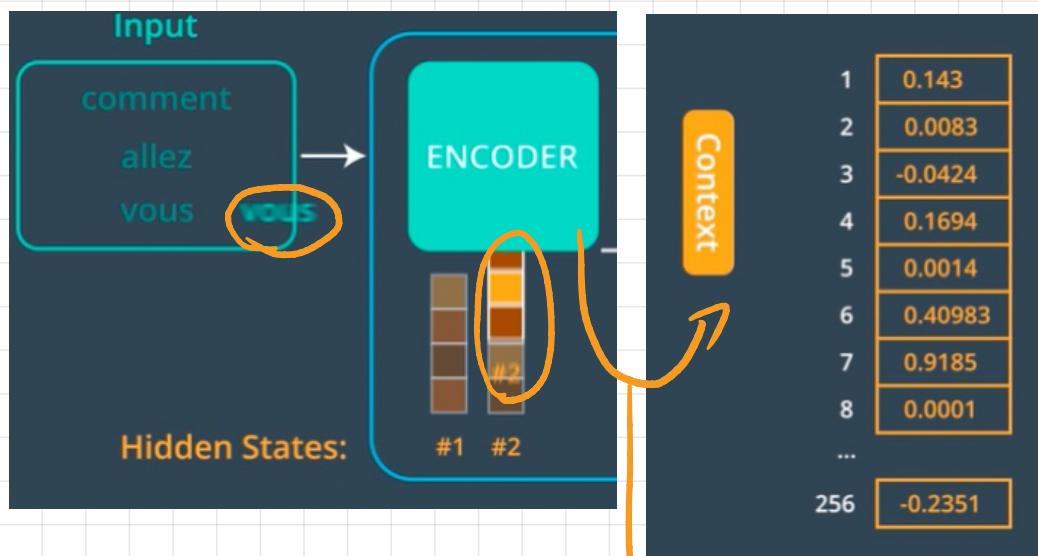
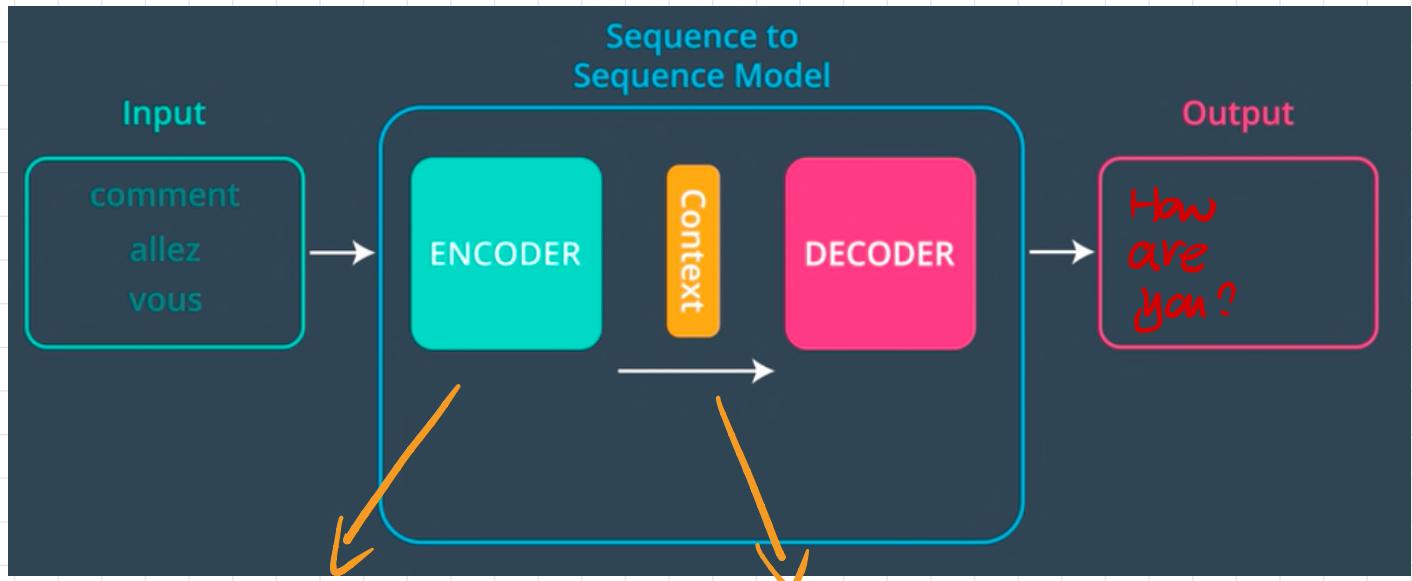
Hidden
State



Hidden-dim?



Seq to seq Model



$$h_t = f\{h_{t-1}, x_t\}$$

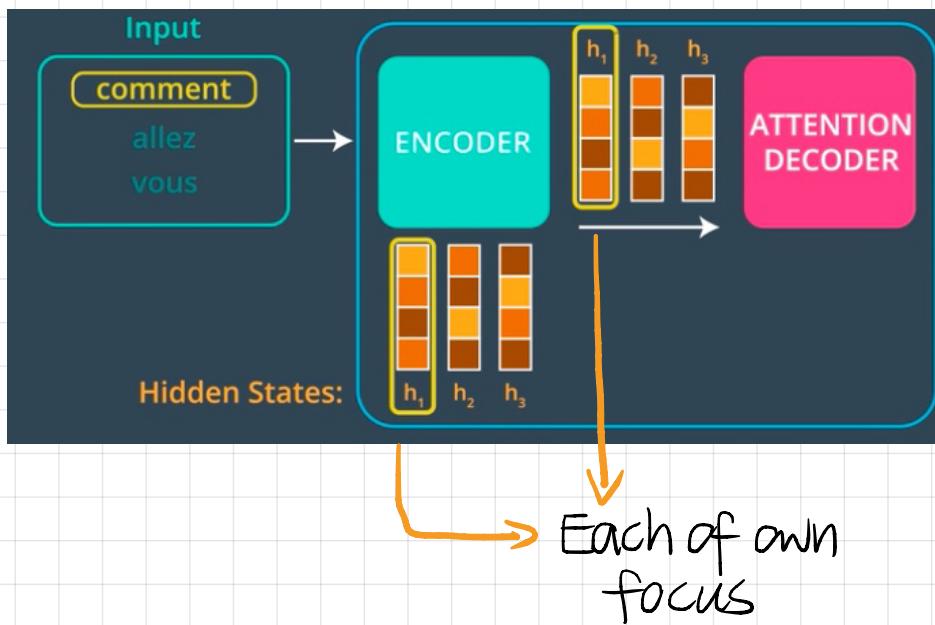
The last h generated from the input = the context vector

limitation:

- The encoder can only send a vector to decoder
- Overfitting if dim of context vector too large when dealing w/ short sequences.

Sequence to Sequence, with Attention

Encoder:



Context vector

III

All the hidden states

Longer sequence

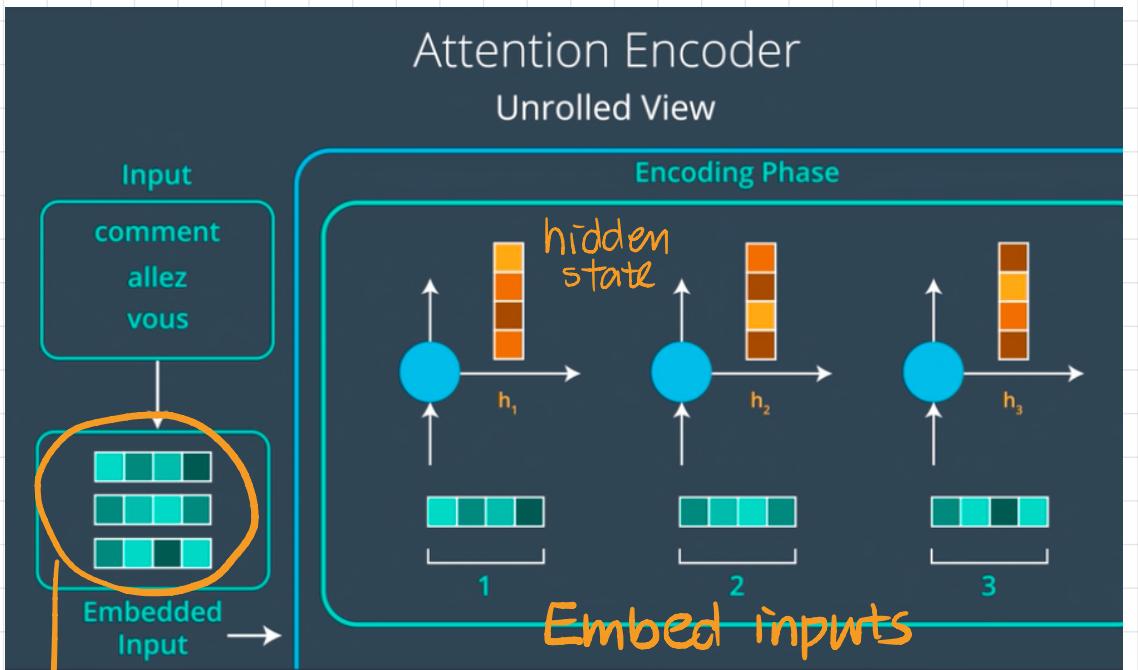
= longer context vec

Decoder:

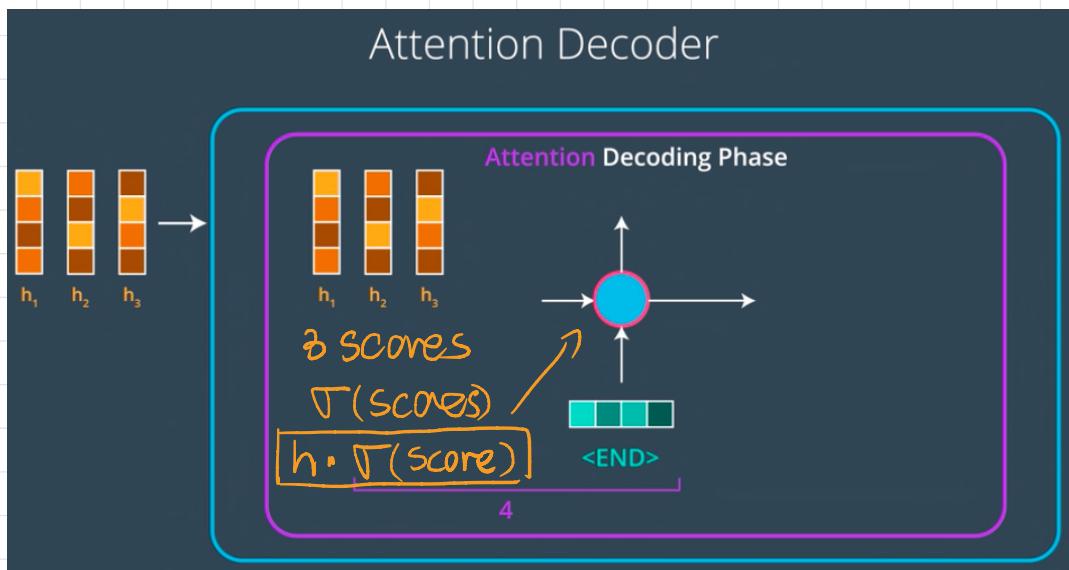


After training,

Adapt to the
grammar of
English



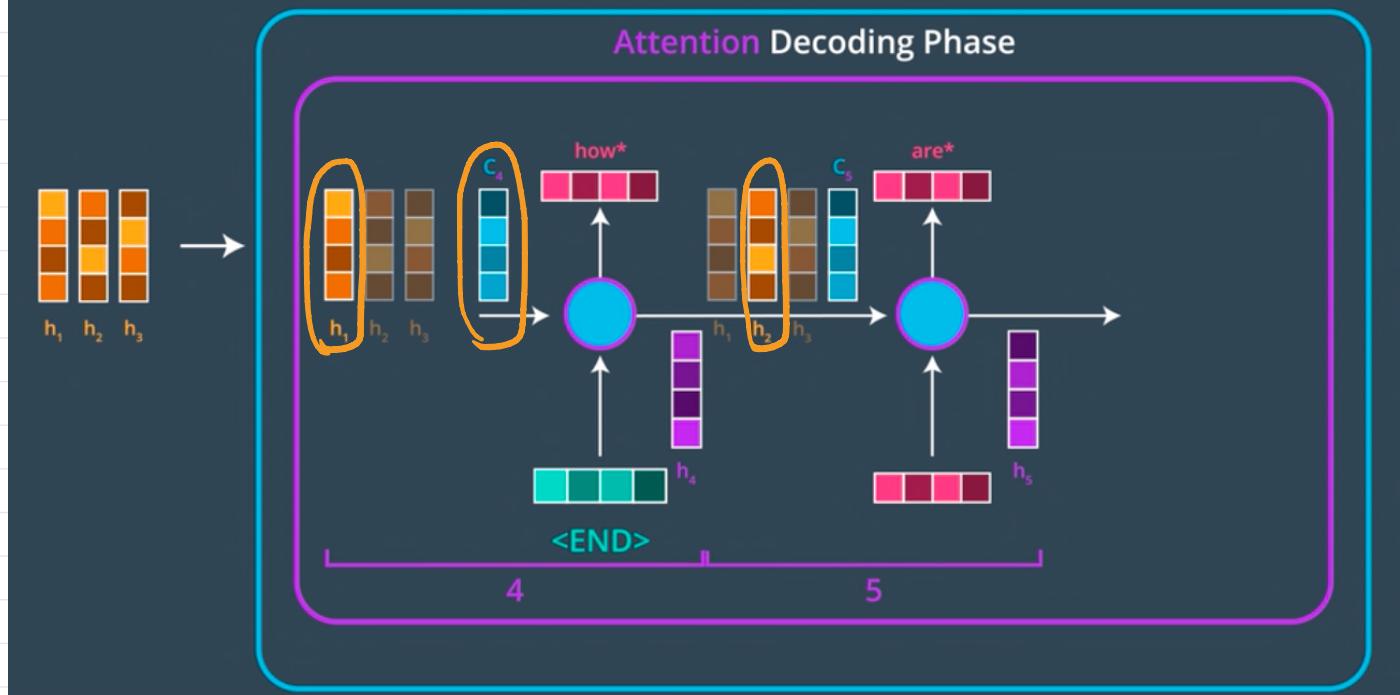
columns: embedding dim
(In real world, this will be much higher)



Before feeding hidden states into decoder:

1. Each hidden state applied with a scoring function.
2. Obtain a score for each h . n scores total
3. Apply softmax to scores
4. $J(\text{Score of each}) \times \text{each } h =$ Context vector to enter decoder

Attention Decoder



To translate :

- hidden state of each word from encoder + hidden state of decoder / output from prev. decoder cell + context word \Rightarrow output

Additive & Multiplicative Attention

Bahdanau Attention

$$e_{ij} = v_a^\top \tanh(W_a \underbrace{s_{i-1}}_{\substack{\text{hidden} \\ \text{state} \\ \text{from} \\ \text{decoder} \\ \text{in prev} \\ \text{timestep}}} + U_a \underbrace{h_j}_{\substack{\text{hidden} \\ \text{state} \\ \text{from} \\ \text{encoder}}})$$

A scoring function

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

Softmax

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$$

Weighted Sum

Luong Attention

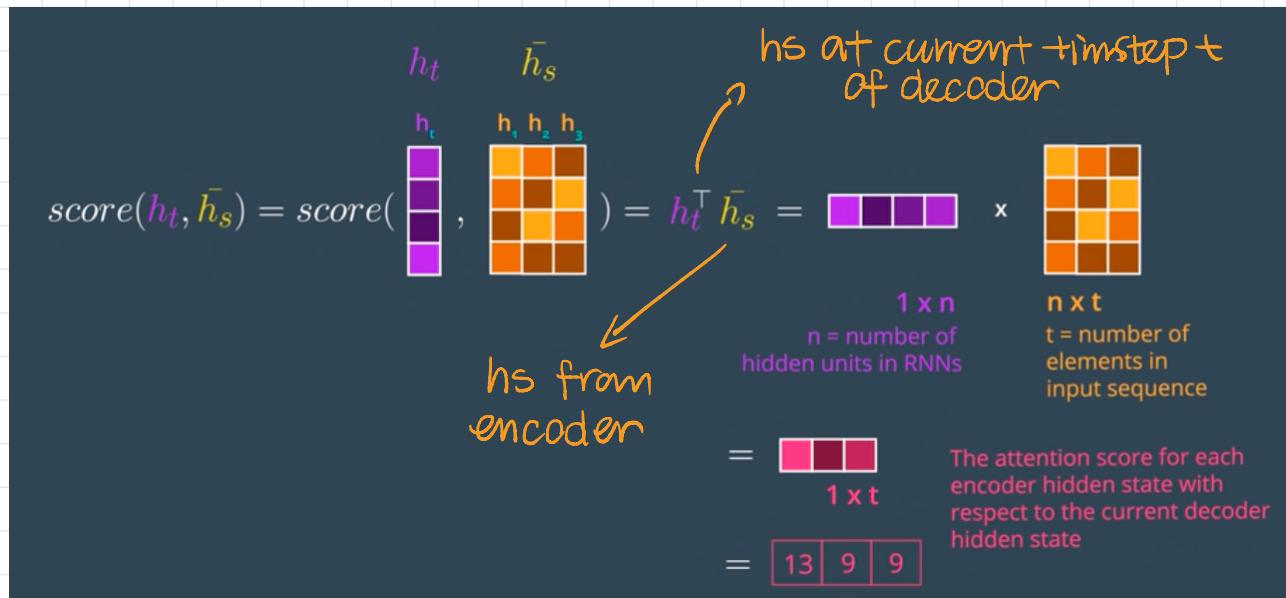
$$\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s) = \begin{cases} \mathbf{h}_t^\top \bar{\mathbf{h}}_s & \text{dot} \\ \mathbf{h}_t^\top \mathbf{W}_a \bar{\mathbf{h}}_s & \text{general} \\ \mathbf{v}_a^\top \tanh(\mathbf{W}_a[\mathbf{h}_t; \bar{\mathbf{h}}_s]) & \text{concat} \end{cases}$$

$$\begin{aligned} \mathbf{a}_t(s) &= \text{align}(\mathbf{h}_t, \bar{\mathbf{h}}_s) \\ &= \frac{\exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s))}{\sum_{s'} \exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_{s'}))} \end{aligned}$$

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_c[\mathbf{c}_t; \mathbf{h}_t])$$

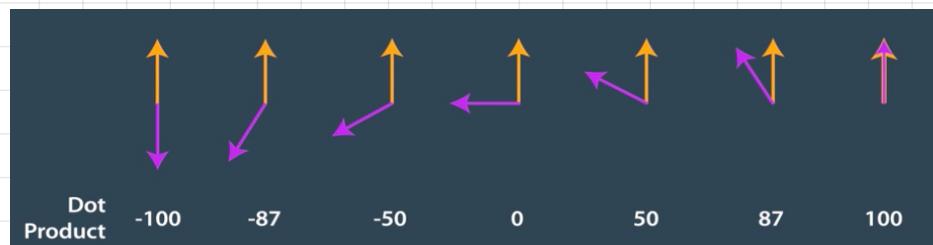
Multiplicative

Scoring method #1. Simply dot product



Intuition:

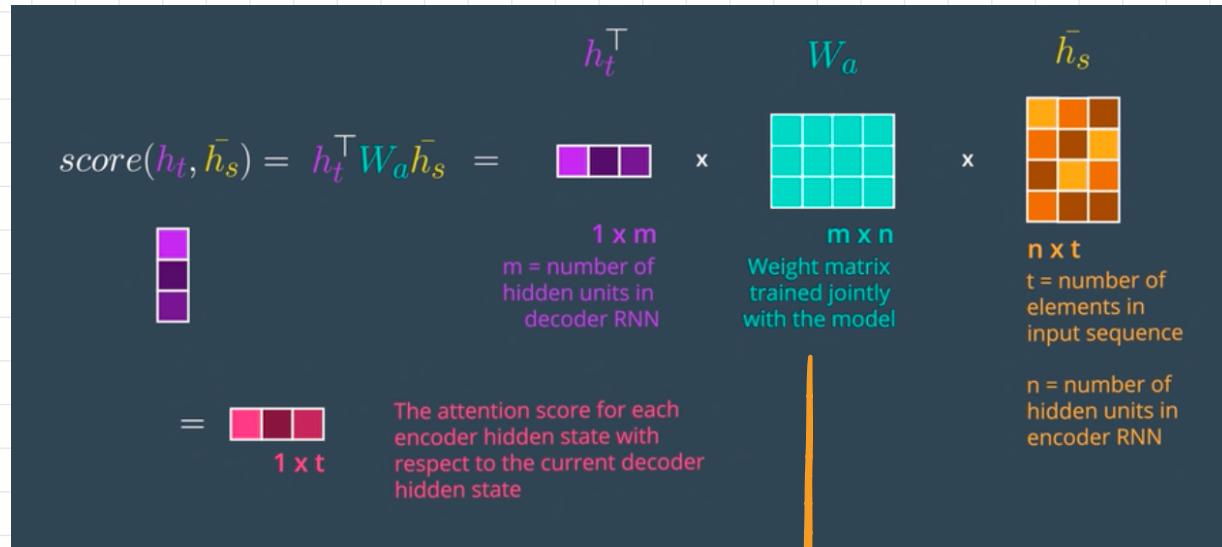
- One scoring method would be the dot product
 - DP returns a single number
 - DP \equiv similarity measurement: \downarrow angle \uparrow value / sim.



Drawback:

- Assume that encoder & decoder have same embed space
- Diff languages have diff embed space.

Scoring Method #2 : General



↓
Unite the embed dim