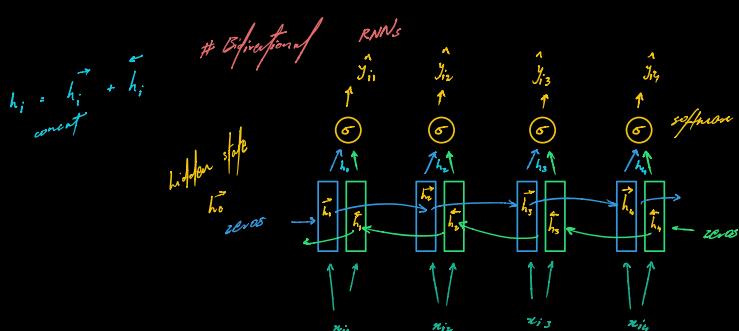
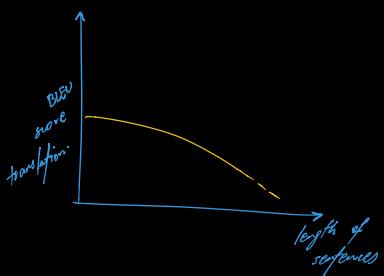


one small content vector is not enough to capture complete essence of a long sentence.

- # what humans do:
 - ① read some words -- tradeoff

- ② shift attention to next few words -- tradeoff.

This is not what seq2seq does
That's why it doesn't scale up



It doesn't matter what these boxes are, they can be RNN, LSTM, GRU, etc. whatever works

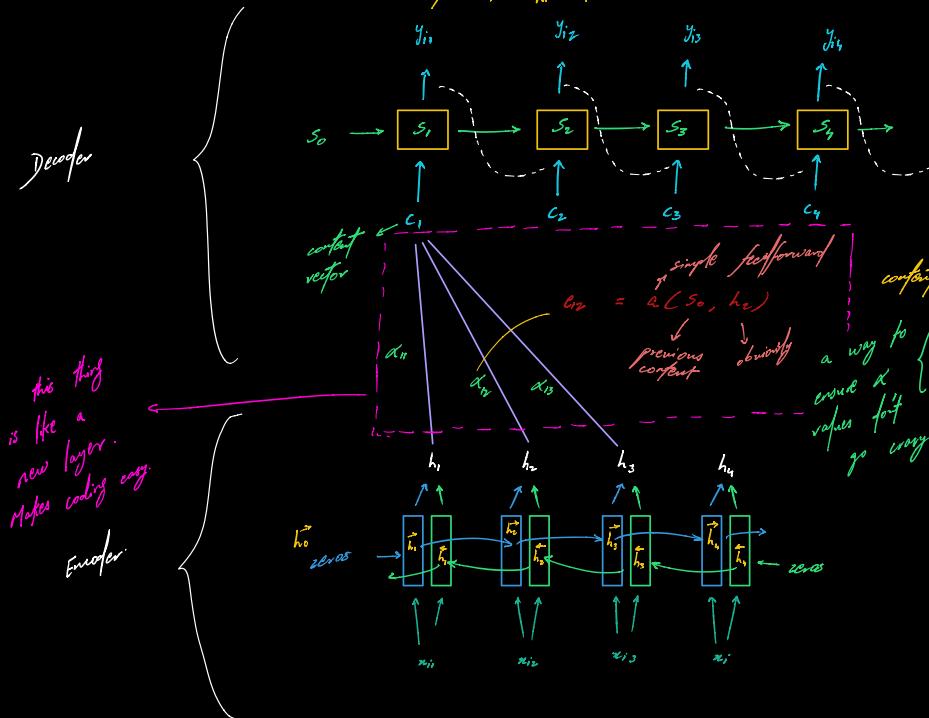
- ① forward pass as usual
- ② pass stuff backwards and collect predictions for that too
- ③ softmax over both predictions
- ④ let better than only doing forward pass.

In attention models we mostly never use vanilla pass
↳ why would we?
and everything is better if no vanilla
bi-directional, no vanilla gradient passing stuff have LSTM's & GRU's when we

Reference Paper & Motivation:

Neural Machine Translation:
- Yoshua Bengio et. al.

in different h_i values to different y_{ij} because each h_i values are taken & depends on T_m



why bidirectional encoder?
↳ output word could depend on word before & after it

why unidirectional decoder?
→ we're just generating a sequence of output + we have option. So no need for bi-directional.

① context vector is a weighted average of all encoder outputs.

② weights α_{ij} for each decoder cell.

$$c_j = \sum_{i=1}^m \alpha_{ji} h_{1i}$$

$$\sum_{i=1}^m \alpha_{ji} = 1 \quad \& \quad \alpha_{ji} \geq 0$$

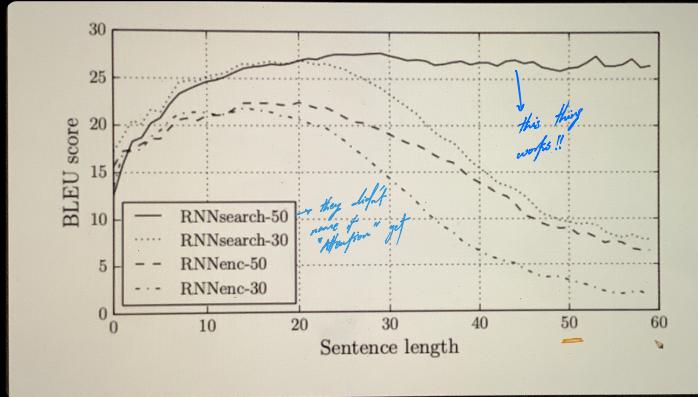
$$x_{ij} = \frac{\exp(c_{ij})}{\sum_{k=1}^m \exp(c_{ik})}$$

This will keep it so α sum = 1
This is like softmax

$c_{ij} = \alpha(s_{i-1}, h_j)$
they didn't know what f^n to use.
the models are f^n approximators. So let's put a NN there. Let the model learn it...
single N layer NN

Make architecture differentiable & let it converge

↓
given large offset,
this thing will be cool → doesn't matter how
hard it is or complex.

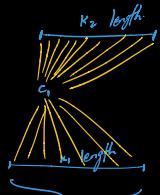


#Drawbacks

length of output

$$\text{Time Complexity} = O(k_1 \cdot k_2)$$

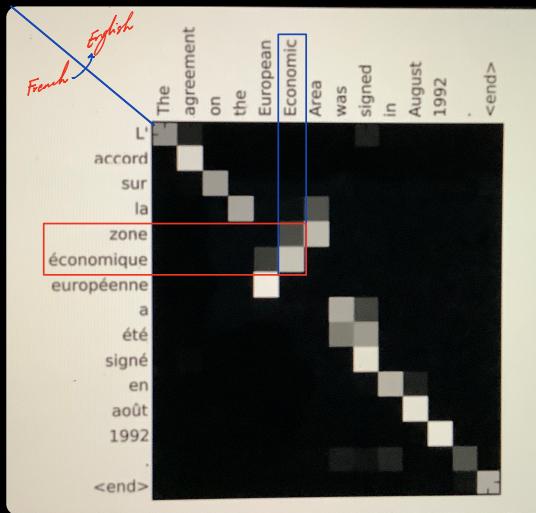
\downarrow
length of input



$T_w \rightarrow$ if T_w is large enough $\rightarrow c_*$ will depend on all inputs & T_w . If depends on c_* , then K_1, K_2

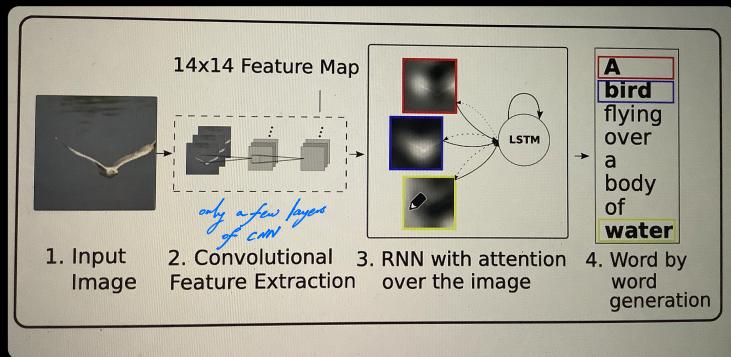
Visualizing α_{ij} → helps debug the NN

$\alpha_{ij} \rightarrow 1$
color → white



understanding what the
fuck is happening inside
the fuck box.

Show Attend & Tell . *Fisher et al. (2015)*



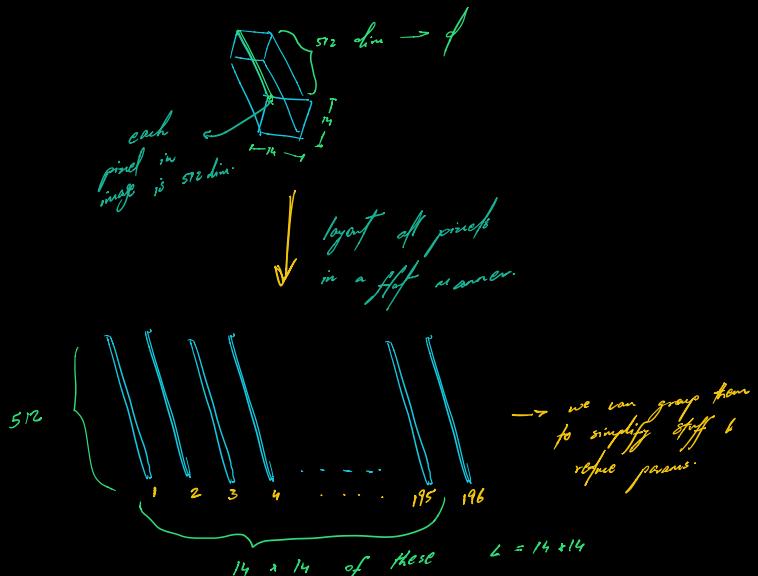
encoder → conv
decoder → same

* Divide image into sections
↳ obtain vectors for each
image region.

Try one vector for the image FC layer
is not using well the 4th layer.

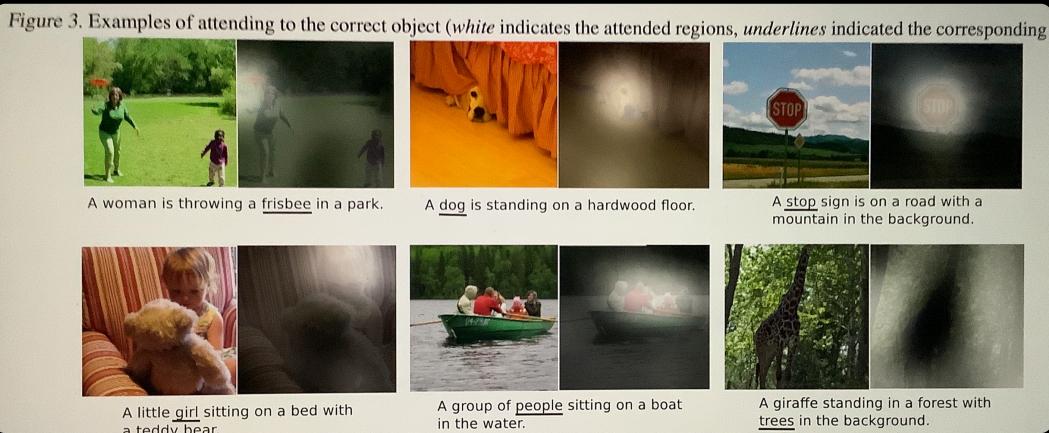
$14 \times 14 \times 512 \rightarrow$ after 4th Vocab layer

Now for each output option word
we can know which pixels got most
attention



make sense? L 14dim vectors.
196 572dim vectors

Quick note → Here a lot more make the
paper so what little is explained
here.



The model implicitly learns where to look based on the output expected

