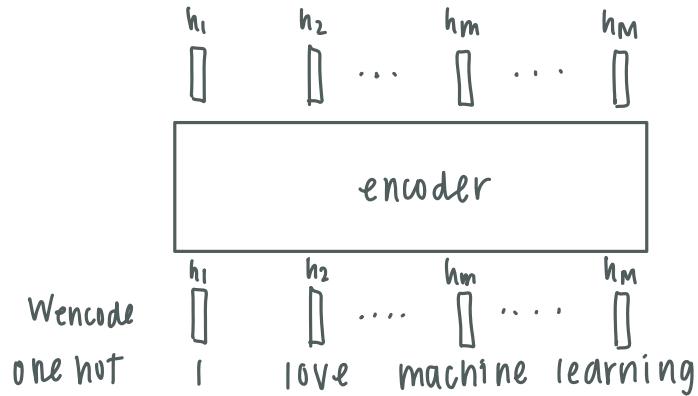
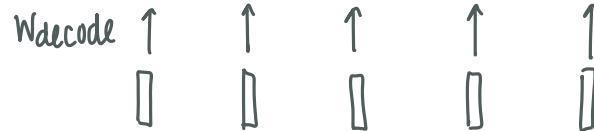


## Transformer: self attention



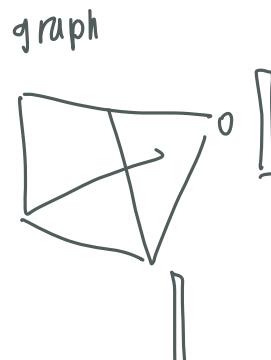
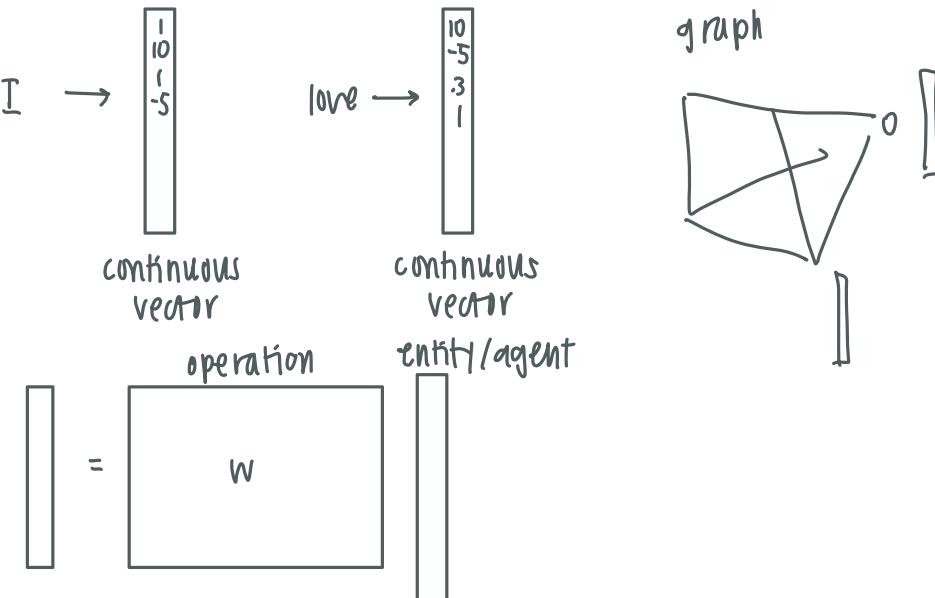
one hot Amo machinas de appendido <end>



Wencode <start> Amo machinas de appendido

problem w/ LSTM: need to do in a sufficient way, slow

## Embedding



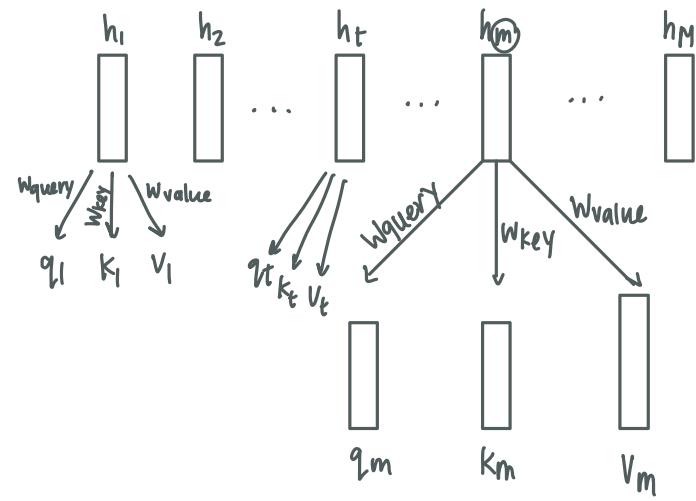
## Attention

soft max

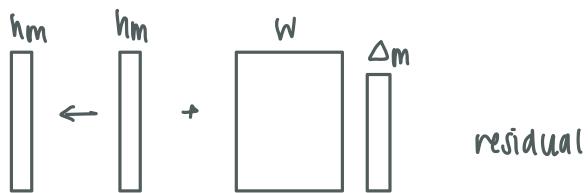
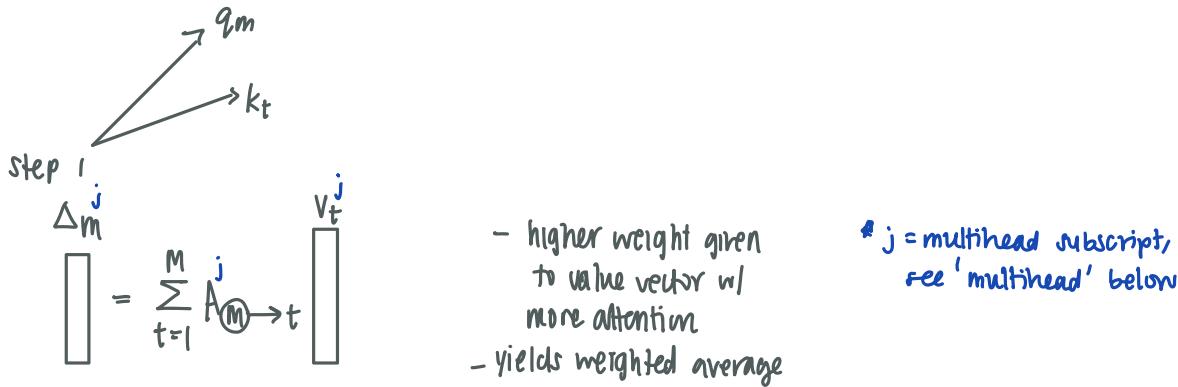
$$A_{(m) \rightarrow t} = \frac{\exp(\text{affinity}(h_m, h_t))}{\sum_{t'=1}^M \exp(\text{affinity}(h_m, h_{t'}))}$$

$$\sum_{t=1}^M A_{(m) \rightarrow t} = 1$$

Diagram below showing arrows from indices 1, 2, ..., t, ..., m, ..., M pointing to a central circle labeled "m", representing the softmax distribution over all possible targets.

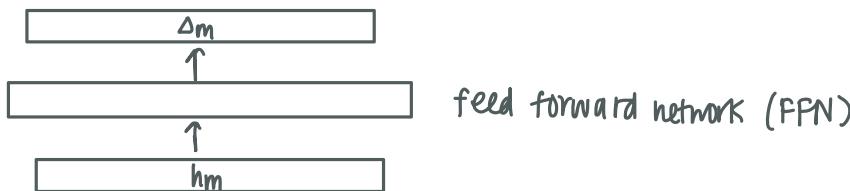


$$\text{Affinity}(h_m \otimes h_t) = \langle q_m, k_t \rangle / \sqrt{\dim}$$



contextualize  $\rightarrow$  each vector has info from all other vectors

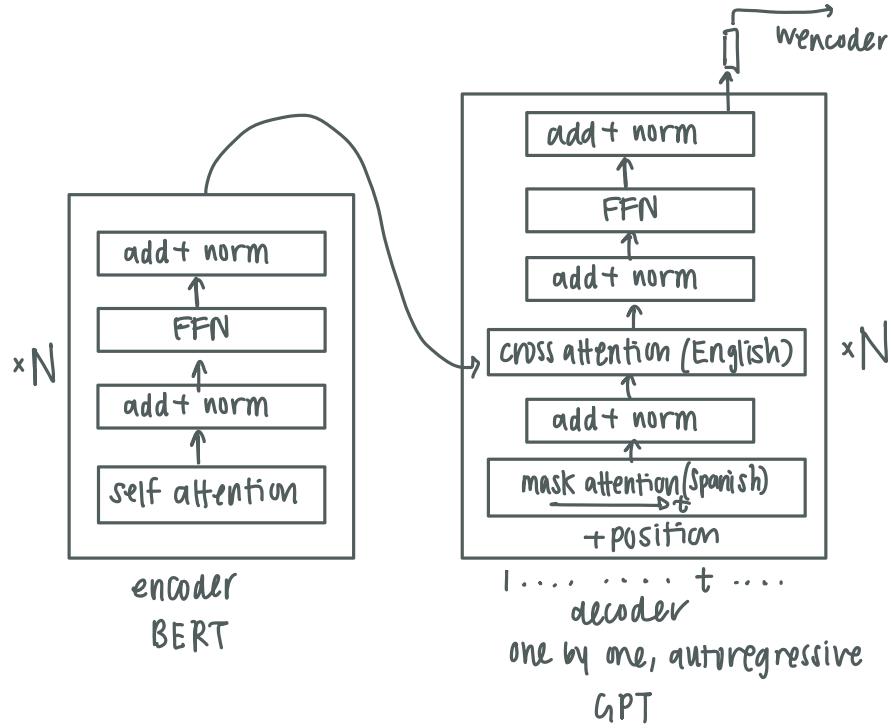
**step 2**  
Each  $h_m$  alone



$$h_m' \leftarrow h_m + \Delta_m \quad \text{residual}$$

layer normalization after each residual (not batch normalization) [so, after step 1 + step 2 residuals]

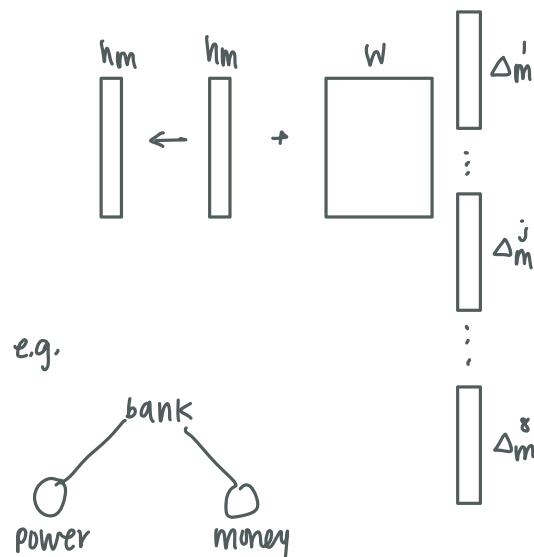
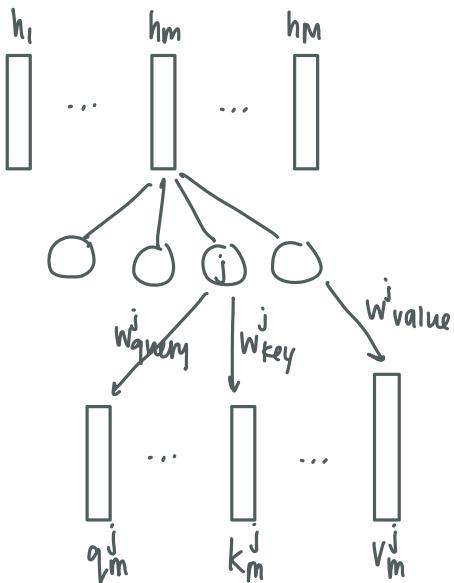




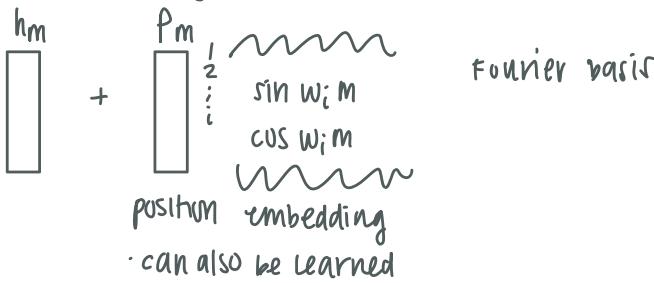
GPT-3	170 billion	scaling law:
PaLM	540 billion	$\propto$ size of model
		$\propto$ performance

**decoder**  
train AR model  
conditioned on the  
prompt.  
doesn't need to be retrained  
due to AR nature

### Multilead \*



### (initial) embedding



### Decoding

