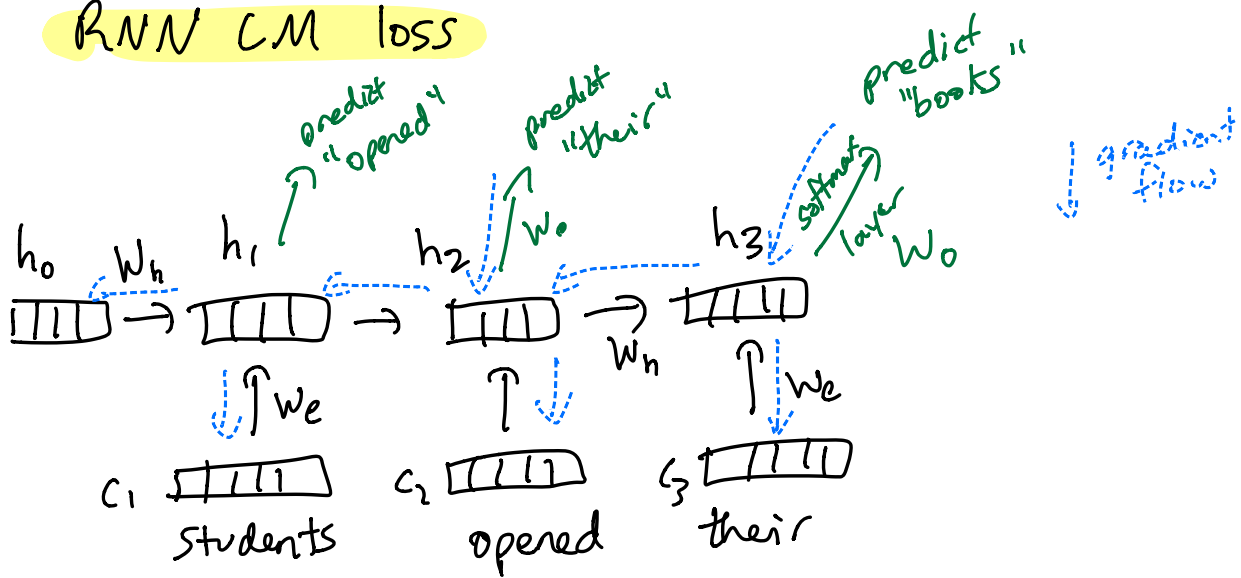


From RNNs to attention:

1. loss computation in RNN LMs
2. issues w/ RNN LM
  - ↳ bottleneck
  - ↳ cannot be parallelized across timesteps
3. self-attention

**RNN LM loss**



$$L_1 = -\log P(\text{opened} | \dots)$$

$$L_2 = -\log P(\text{their})$$

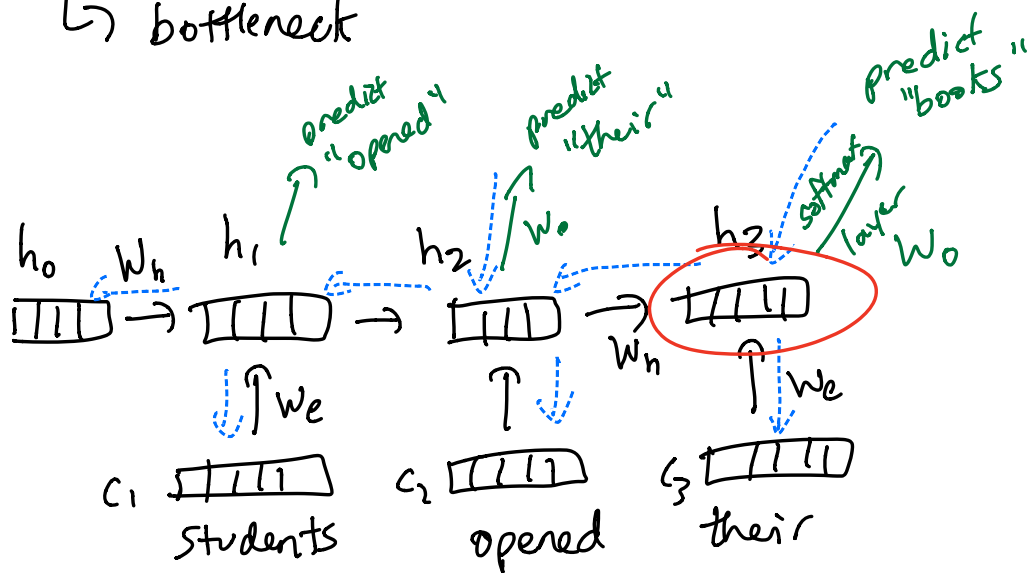
$$L_3 = -\log P(\text{books})$$

$$L = \frac{L_1 + L_2 + L_3}{3}$$

} avg. neg. log likelihood of the ground-truth next word over all tokens in the batch

## Issues of RNNs:

↳ bottleneck



$h_3$  is required to encode ALL of the useful info from the prefix to predict the next word

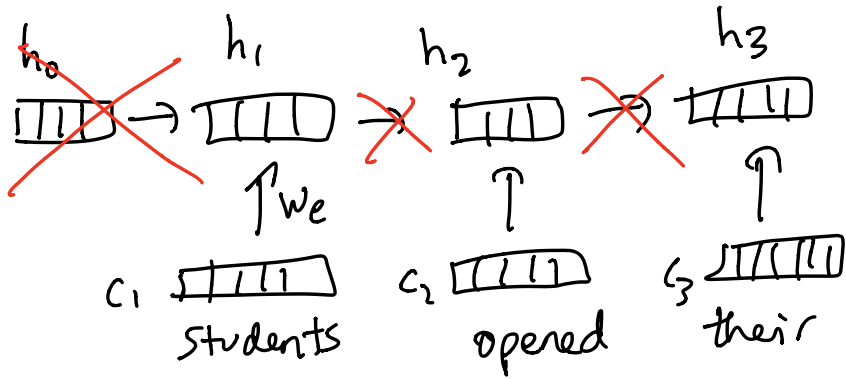
↳ with long, complex prefixes, this is an unreasonable expectation

↳ RAY MOONEY: (~2014)

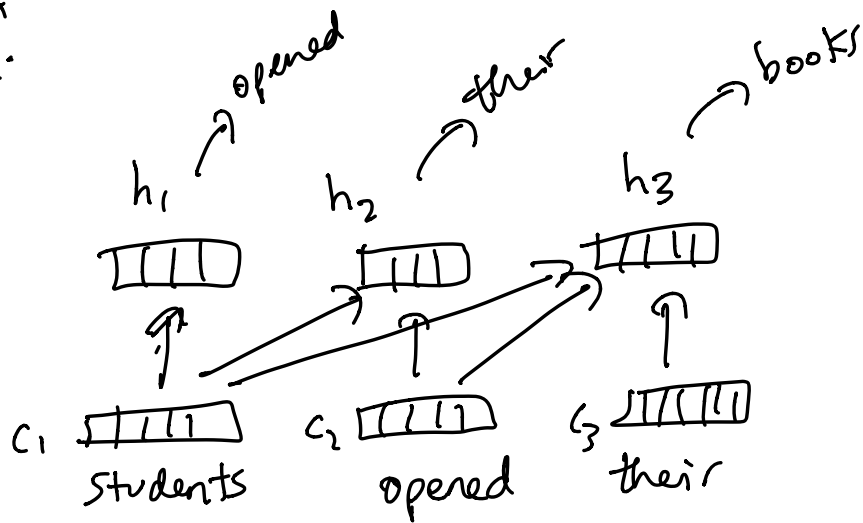
"you can't represent the meaning of a sentence in a BLEEPING vector"

## attention mechanism

- ↳ a way to get around the bottleneck in neural LMs
- ↳ intuition: attention provides a way to access hidden states that are far away
- ↳ history: developed initially for RNNs ~2014, Bahdanau, Cho
- ↳ now: "self-attention", core module of the Transformer LM
  - ↳ introduced by Google in 2017
  - ↳ reduces bottleneck effect
  - ↳ fully parallelizable across timesteps
    - ↳ hidden state at each timestep is independent of prev. hidden states



goal:

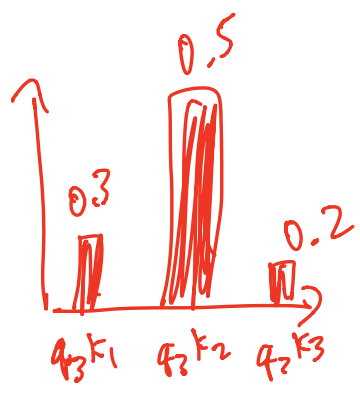


Computation of hidden state at timestep 3:

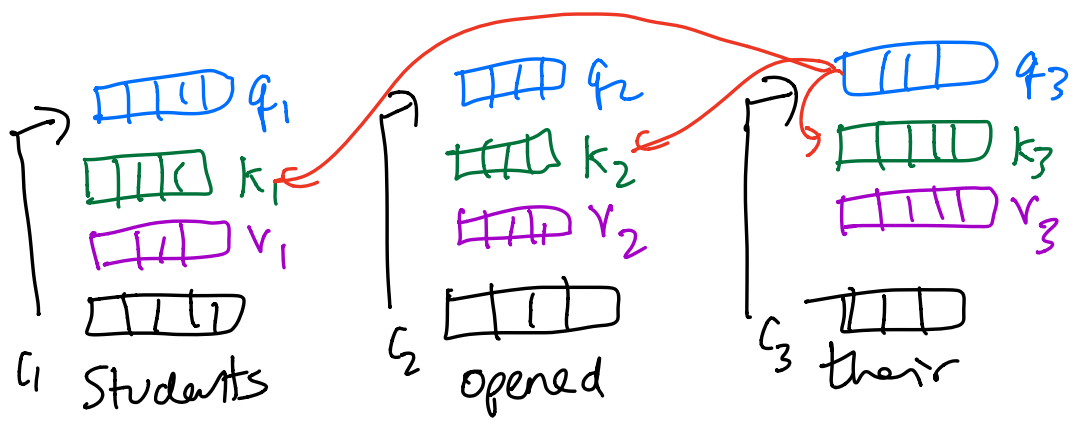
$$h_3 = 0.3v_1 + 0.5v_2 + 0.2v_3$$



→ predict  
softmax  
layer  
"books"



attn Scores:  $\text{softmax}(\langle q_3 \cdot k_1, q_3 \cdot k_2, q_3 \cdot k_3 \rangle)$



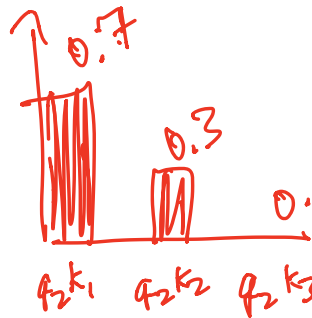
$q_i = f(W_q c_i) \Rightarrow$  query } used to compute  
 $k_i = f(W_k c_i) \Rightarrow$  key } "attention score", use dot product  $q \cdot k$   
 $v_i = f(W_v c_i) \Rightarrow$  value } encode the information used to compute the hidden state

↑  
 weight matrices  
 that are parameters  
 of the model  
 $q_2 = f(W_q c_2)$

Second timestep:

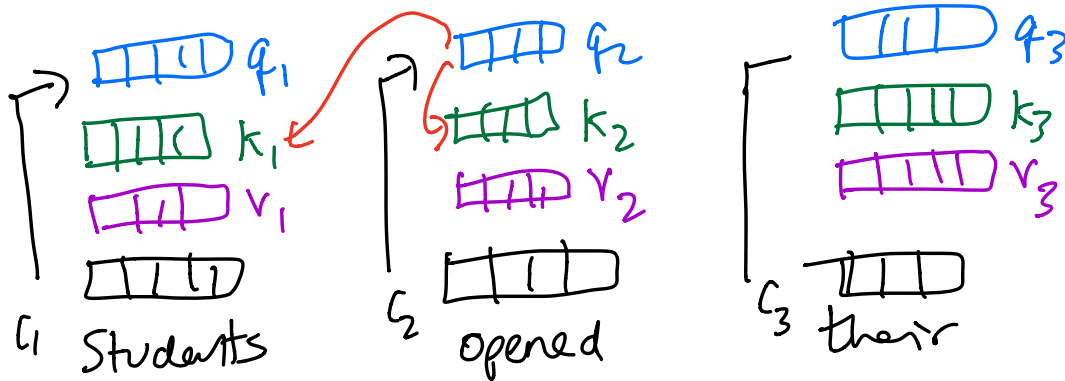
$$h_2 = 0.7v_1 + 0.3v_2$$

= [||||]   
 ↪ predict "their"



no  $q_2 k_3$   
bc of cheating!

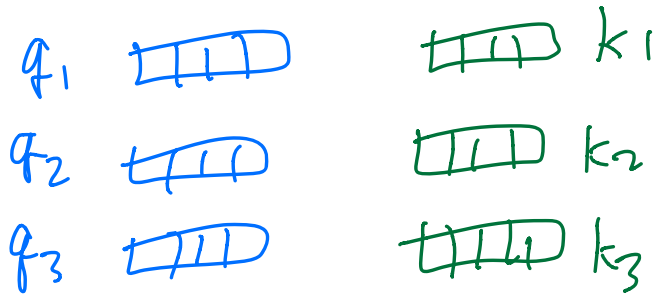
attn Scores:  $\text{softmax}(\langle q_2, k_1, q_2, k_2 \rangle)$



there are no dependencies between  $h_1, h_2, h_3!$

↳ can parallelize

↳ reduced bottleneck

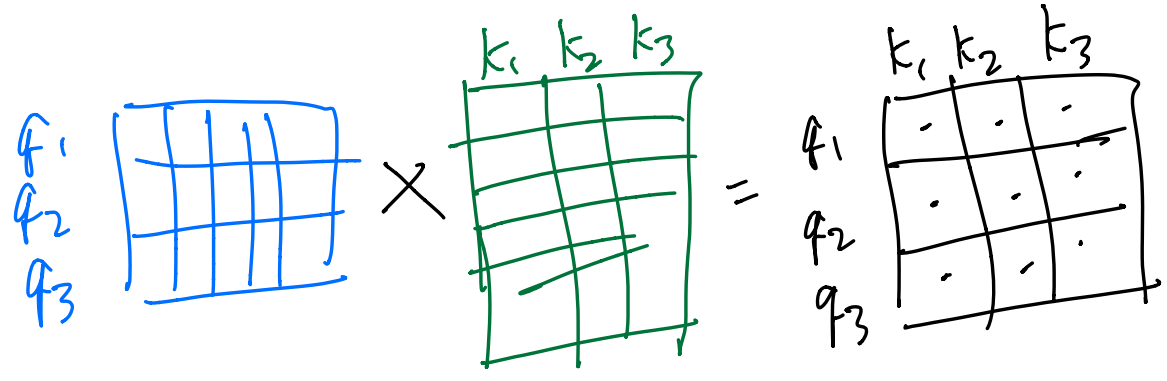


attn scores

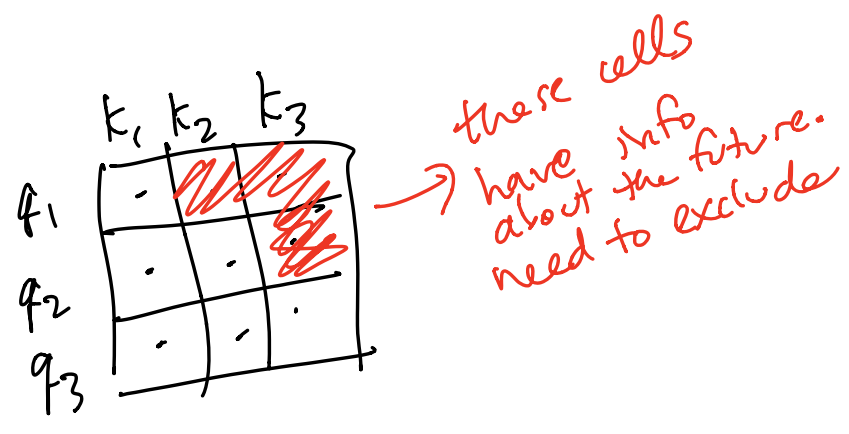
$$a_1 = \langle q_1, k_1 \rangle$$

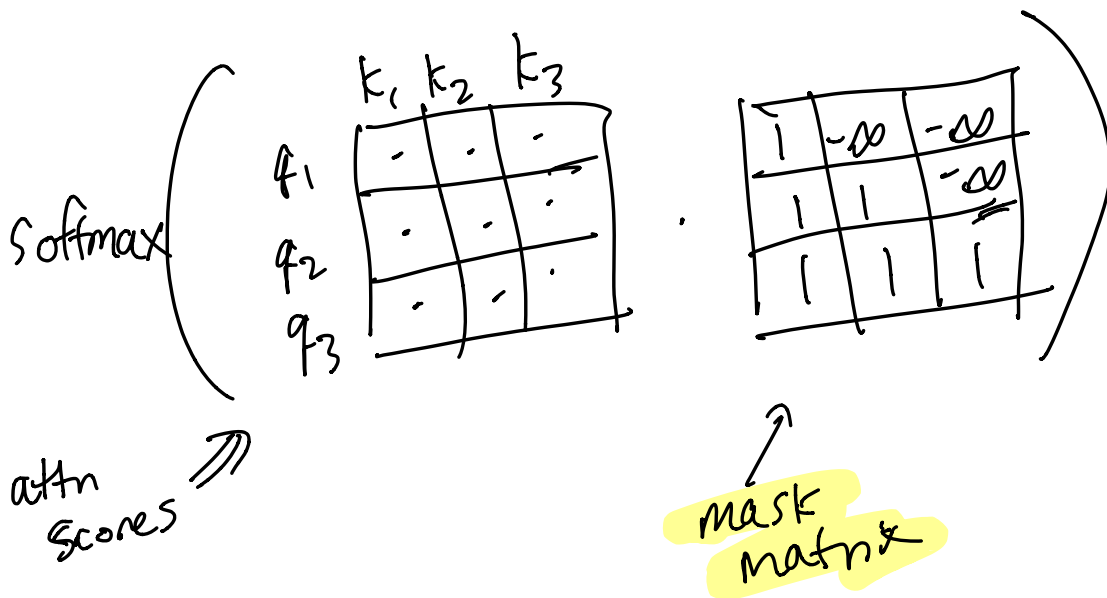
$$a_2 = \langle q_2, k_1, q_2, k_2 \rangle$$

$$a_3 = \langle q_3, k_1, q_3, k_2, q_3, k_3 \rangle$$



One matrix product to compute all  
attn scores



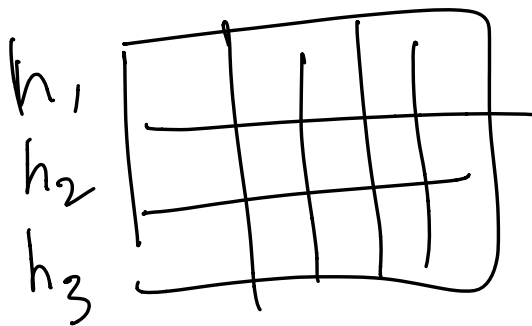


	$k_1$	$k_2$	$k_3$
$q_1$	1	0	0
$q_2$	0.3	0.7	0
$q_3$	0.1	0.5	0.4

X

$v_1$			
$v_2$			
$v_3$			

=



we just computed all three hidden states w/ just a couple matrix products

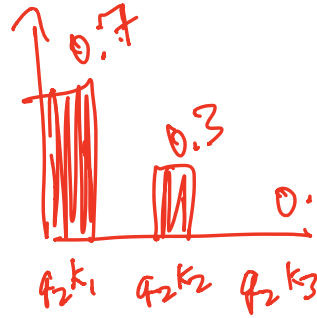


word order:

Second timestep:

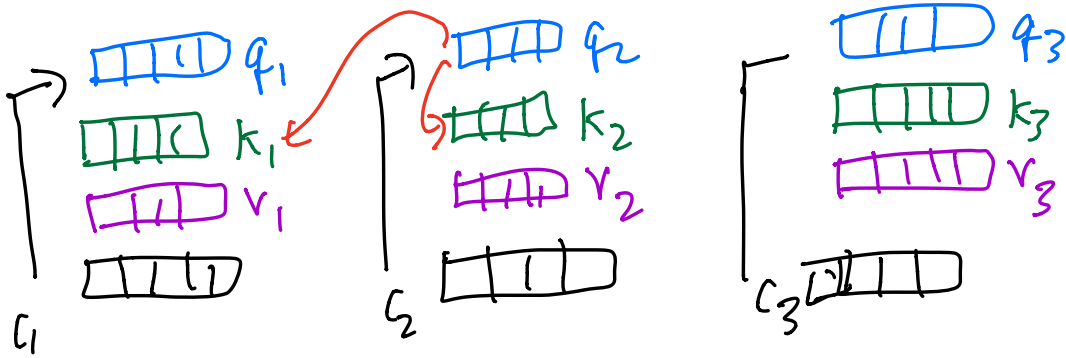
$$h_2 = 0.7v_1 + 0.3v_2$$

=    
↳ predict "their"



no  $q_2k_3$   
bc of cheating!

attn Scores:  $\text{softmax}(\langle q_2, k_1 \rangle, \langle q_2, k_2 \rangle)$



↳ positional embeddings  $\Rightarrow$  learned params

↳  $p_1$  is the same vector for the first pos. of every seq

↳  $p_2$  is the same vec for the second pos. of every seq

$$\text{emb} = c_1 + p_1$$