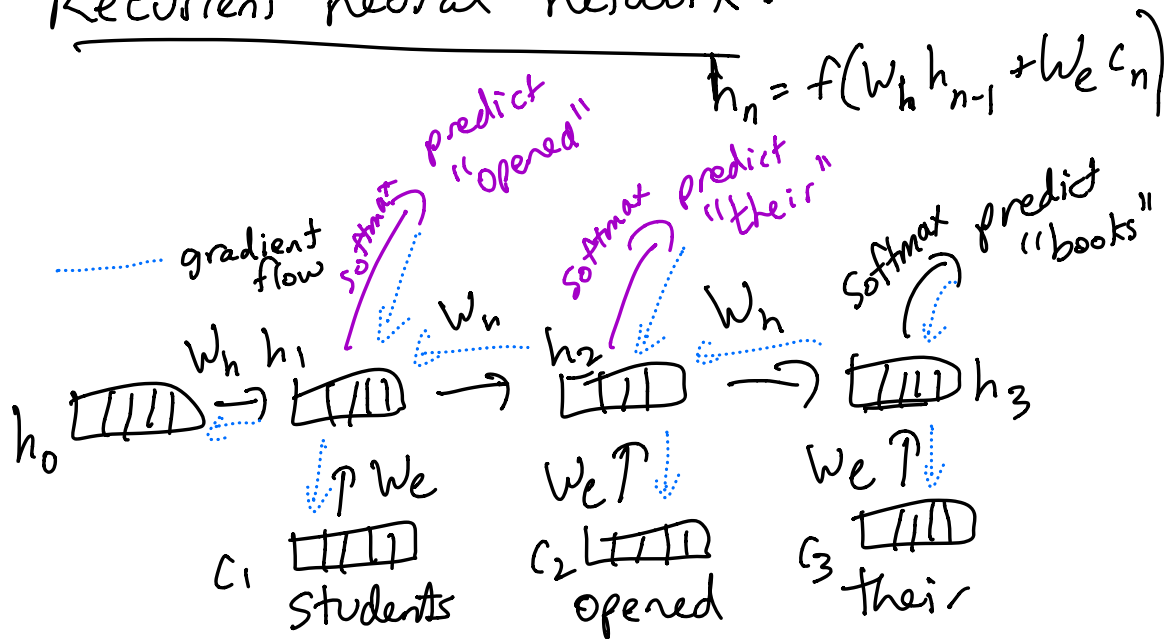


Recurrent neural network:



RNN gradient \Rightarrow "backpropagation thru time"

$$L_3 = -\log P(\text{books} \mid \text{"students opened their"})$$

$$L_2 = -\log P(\text{their} \mid \text{"students opened"})$$

$$L_1 = -\log P(\text{opened} \mid \text{students})$$

$$L = \frac{L_1 + L_2 + L_3}{3} \left. \vphantom{\frac{L_1 + L_2 + L_3}{3}} \right\} \begin{array}{l} \text{average neg. log} \\ \text{likelihood of the} \\ \text{ground-truth next word} \\ \text{over all tokens in the} \\ \text{batch} \end{array}$$

batch:

1. students opened their books
2. people walked their dog
3. the classroom fell silent

Issues w/ RNNs:

1. "bottleneck"

- Ray Mooney, 2014

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

2. lack of parallelism across timesteps

- I can only compute h_n after computing h_{n-1}

attention mechanism:

history: - developed initially for RNNs and for machine translation

- Bahdanau, Cho et al, 2014

- Vaswani et al. 2017 dropped the "recurrent" aspect and created a fully attention based architecture

- Transformer

- for machine translation

- introduced by Google

- hidden state at each timestep is independent of $h_1 \dots n-1$

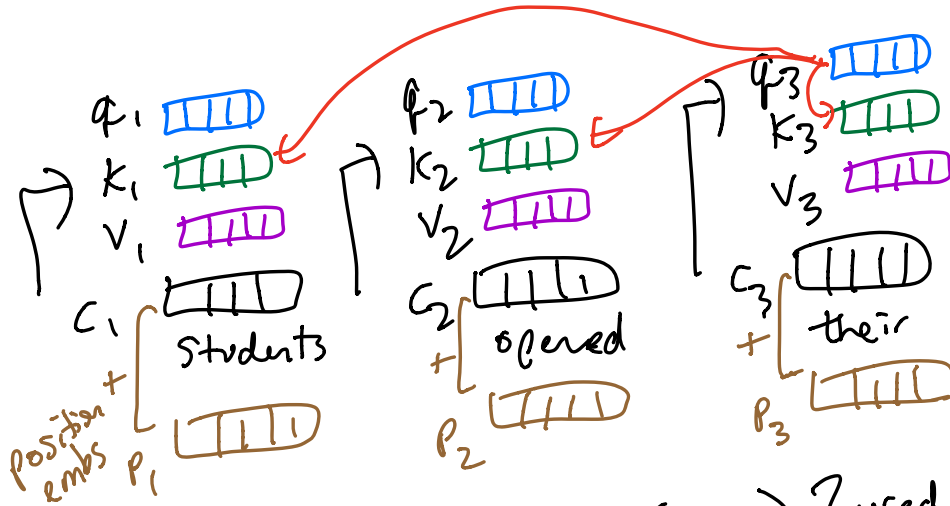
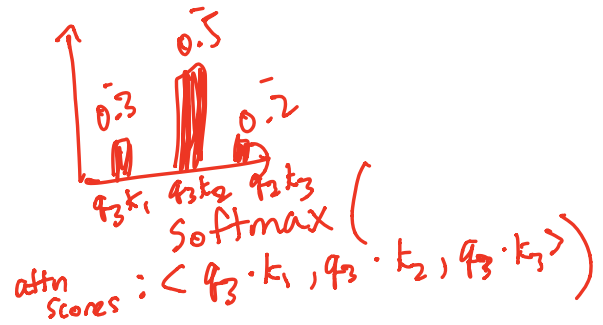
Self-attention:
parallel at training time,
sequential at test-time

Self-attention:

computation of hidden state at timestep 3:

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

↓
 [] → softmax
 ↓
 predict "books"



query: $q_1 = f(W_q c_1)$ $q_2 = f(W_q c_2)$ } used to compute "attention scores"

key: $k_1 = f(W_k c_1)$

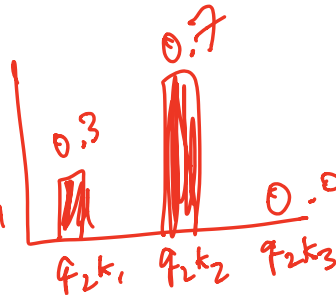
value: $v_1 = f(W_v c_1)$ } encode information used to compute hidden state

W_q, W_k, W_v are randomly initialized parameters learned during training!

computations @ second time step:

$$h_2 = 0.3 \cdot v_1 + 0.7 \cdot v_2$$

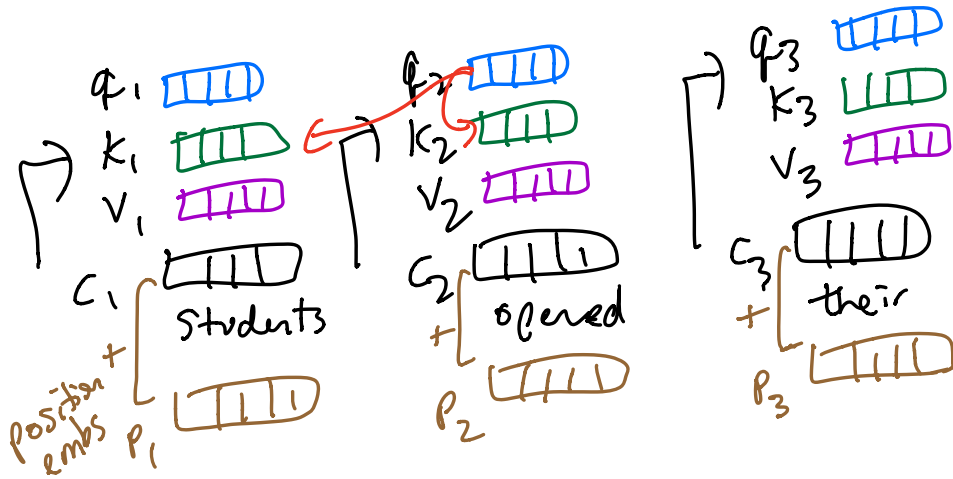
Softmax layer → predict "their"



no $q_2 k_3$
b/c of cheating!

$$\text{Softmax}(W_o h_2)$$

$$\text{attn} : \text{Softmax}(\langle q_2 \cdot k_1, q_2 \cdot k_2 \rangle)$$

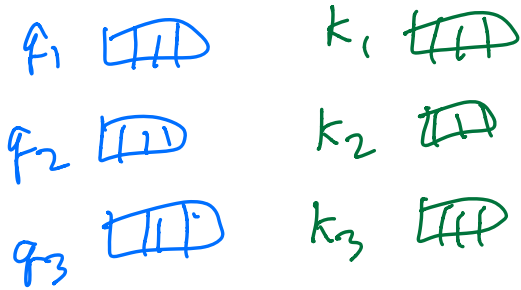


no dependencies between h_1, h_2, h_3

↳ parallelize

↳ reduce bottleneck

how to parallelize:

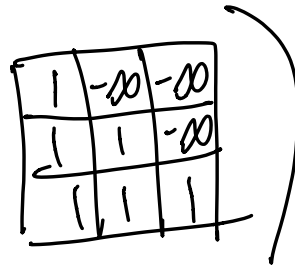
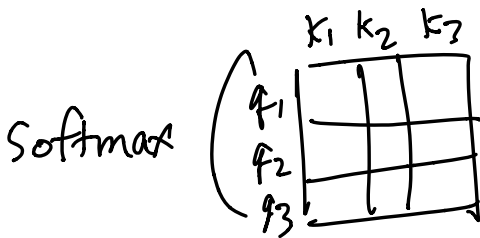
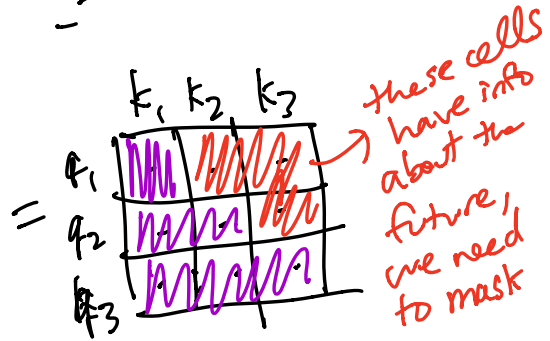
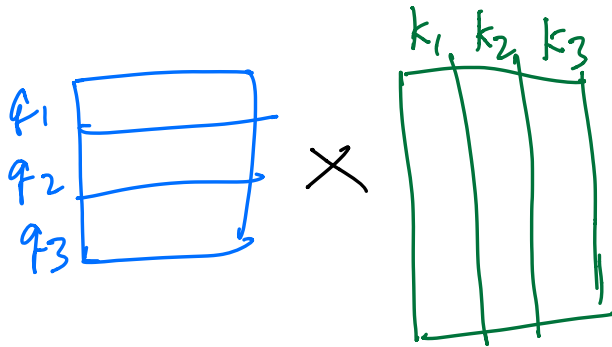


attn vectors

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

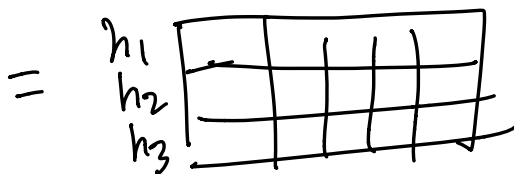
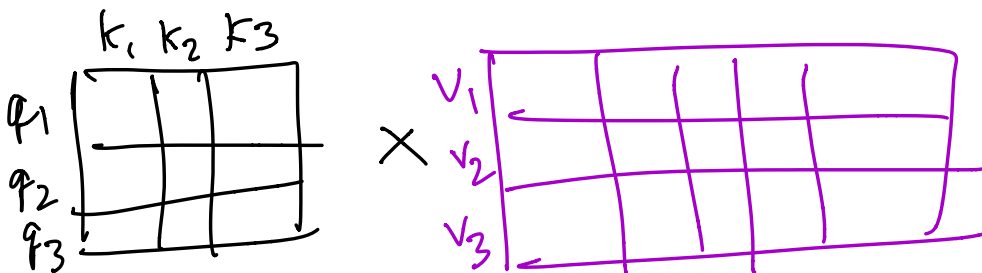
$$a_2 = \langle q_2, k_1, q_2, k_2 \rangle$$

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



\nearrow
attn scores

\curvearrowright mask matrix



-