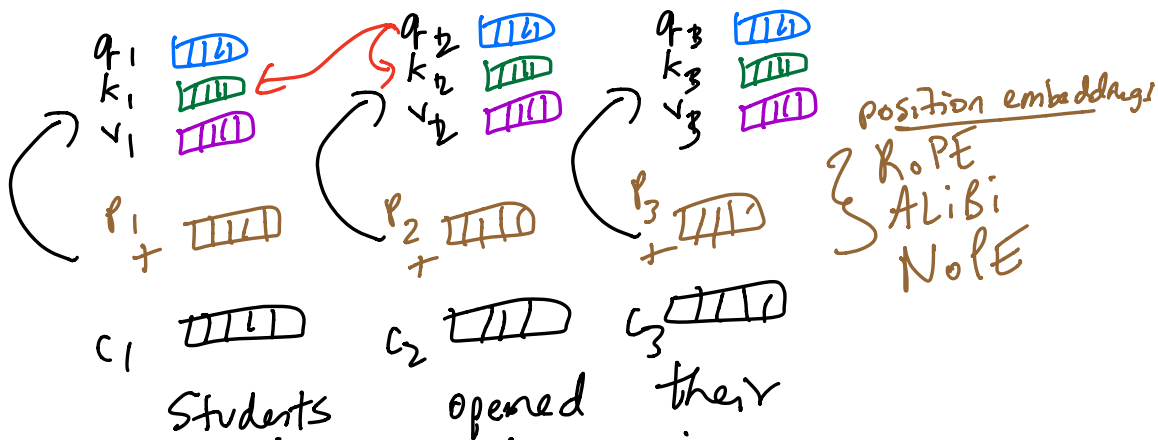




$$h_2 = 0.3v_1 + 0.7v_2$$

$$\text{attn: } \text{softmax}(c q_2 k_1, q_2 k_2)$$



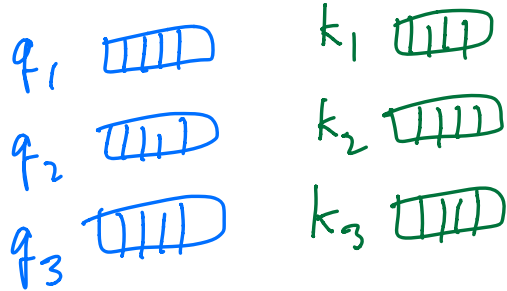
$$q_1 = f(W_q [c_1 + p_1])$$

→ no dependencies between  $h_1, h_2, h_3$

↳ parallelize

↳ reduce bottleneck

# how to parallelize a self-attn computation:

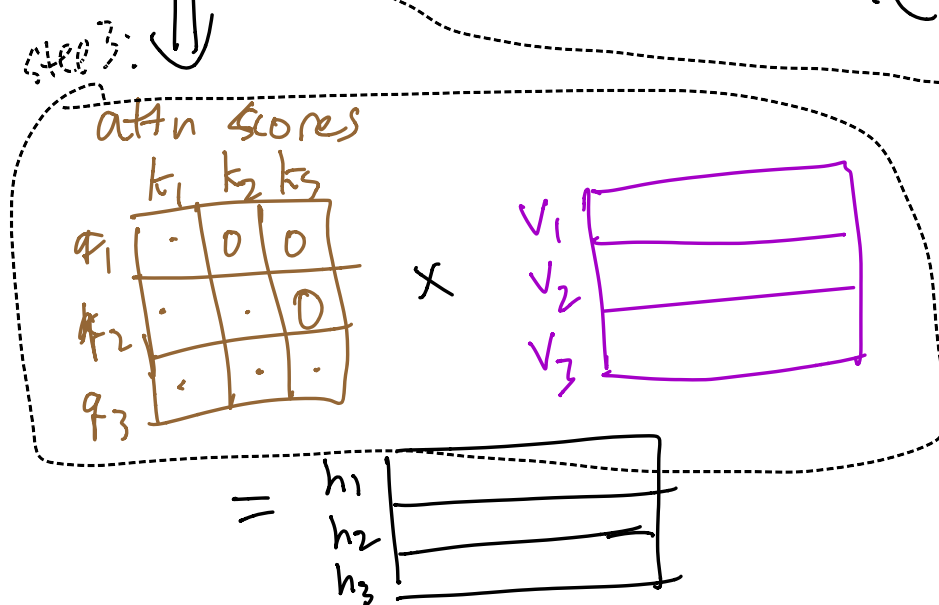
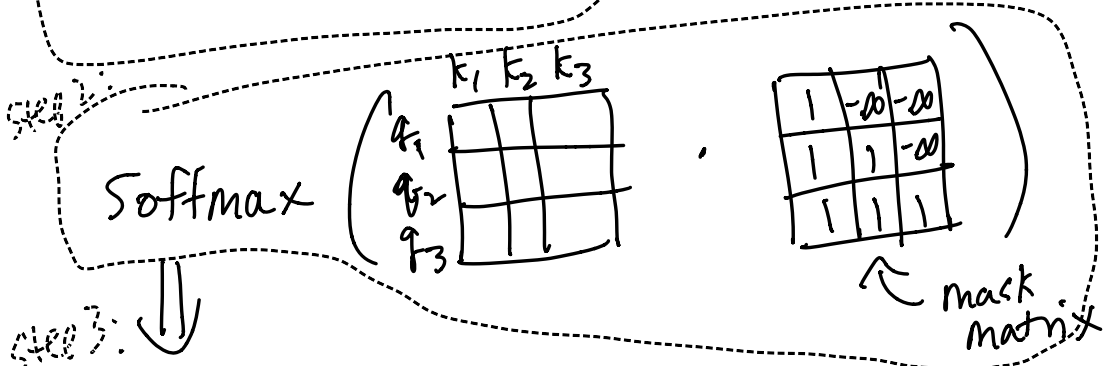
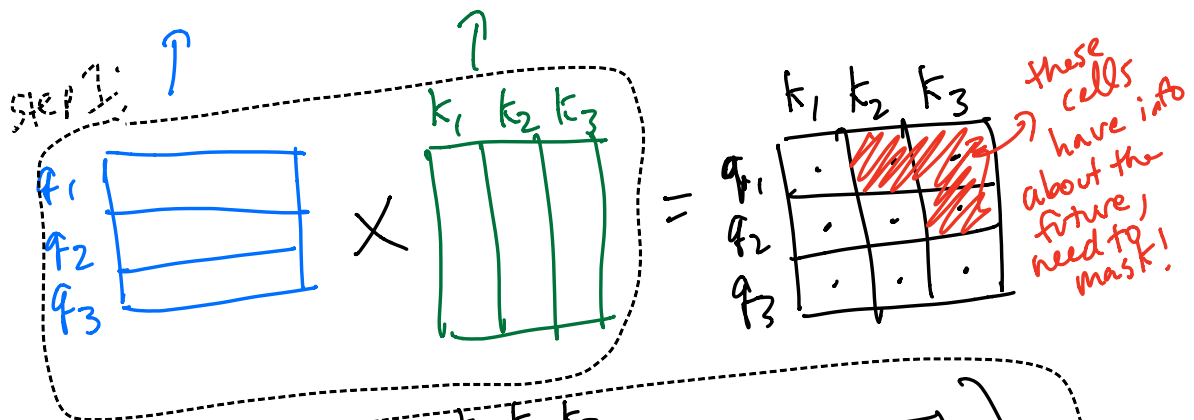


attn vectors

$$a_1 : \langle \underline{q_1}, k_1 \rangle$$

$$a_2 : \langle \underline{q_2}, k_1, q_2, k_2 \rangle$$

$$a_3 : \langle \underline{q_3}, k_1, q_3, k_2, q_3, k_3 \rangle$$



## Self-attention:

↳ training time: parallelizable

↳ test-time: sequentially

⇒ we can only parallelize computations of hidden states when we know the full sequence ahead of time

## test-time

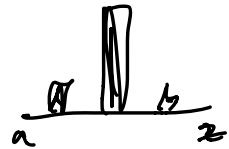
$q_1, k_1, v_1$       $q_2, k_2, v_2$   
↑                    ↑  
students     opened     ?

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

multiply  
by  
value  
vecs

↳  $h_2$  

→ softmax layer



↑  
sample "laptops"

$q_1, k_1, v_1$      . . .  
↑  
students     opened laptops     ?

## KV cache

↳ storing prev. computed keys/values so you don't have to recompute every time step

# Transformer

↳ neural LM built on  
multi-head self-attn

↳ "deep", stacked layers (or "blocks")  
of MH self-attn

↳ Vaswani et al. 2017

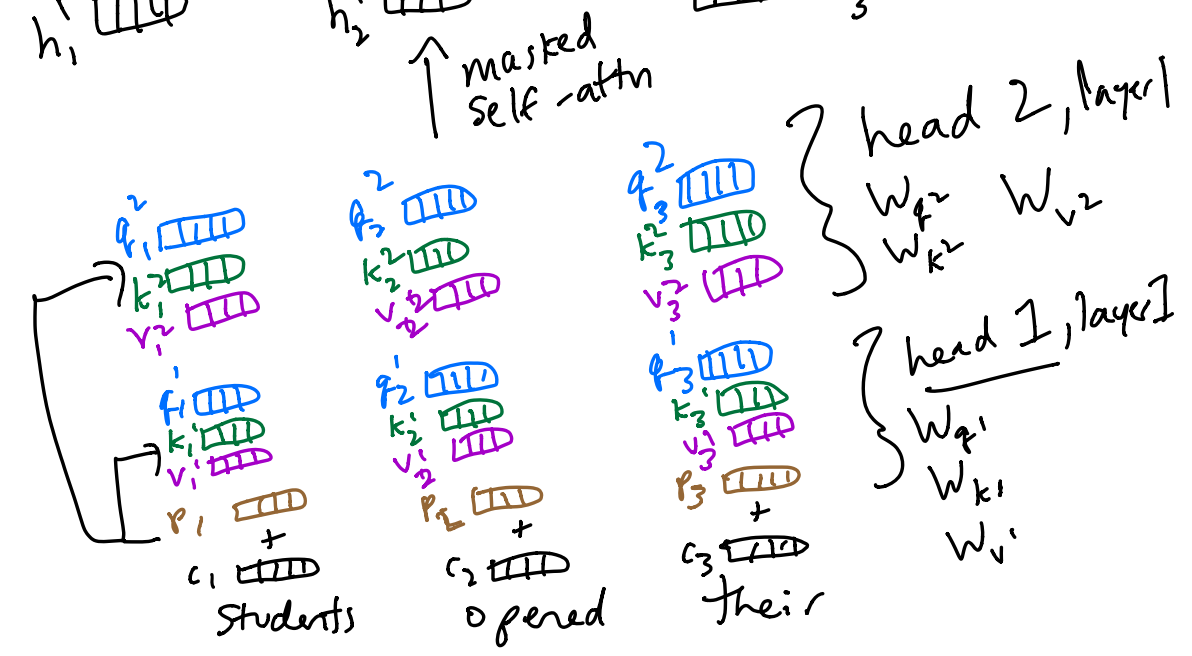
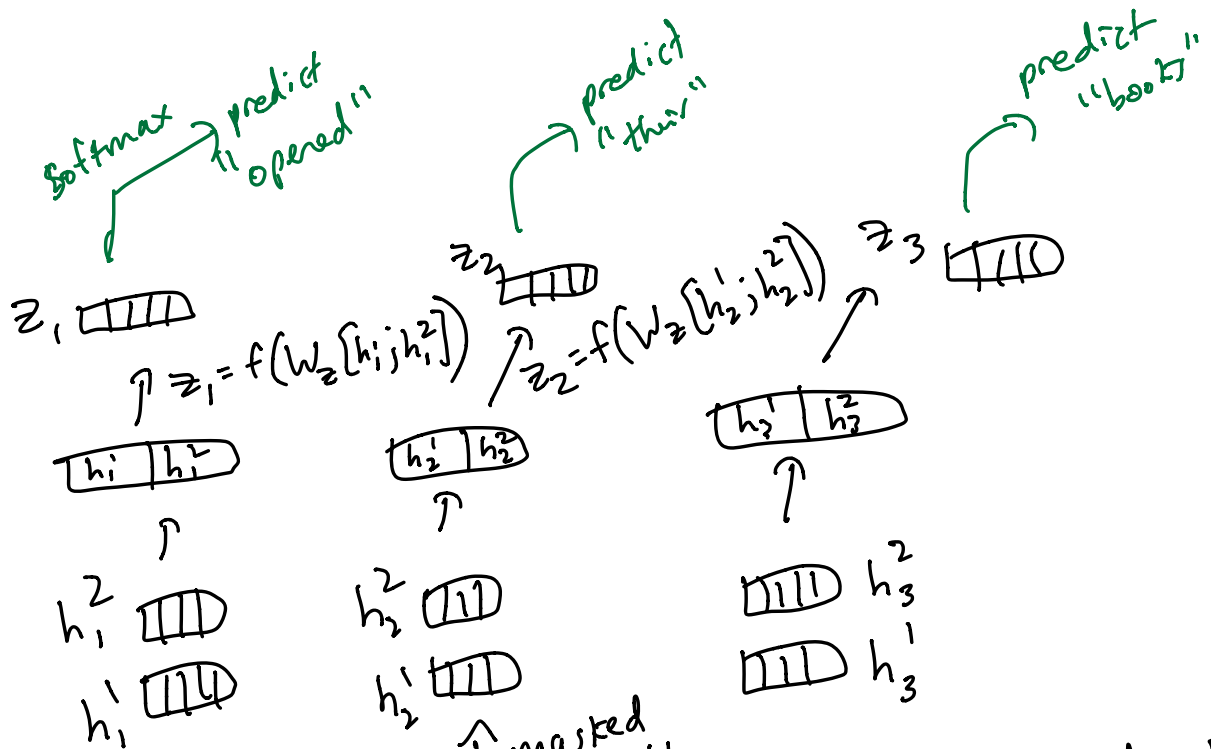
↳ intuition:

→ let's have multiple sets of  $q, k, v$   
vectors for each token,  
then each "head" can focus on  
a specific linguistic property

↳ syntax (e.g. all verbs in phrase)

↳ activate on specific words/phrases

↳ entities/dates



# Adding depth:

