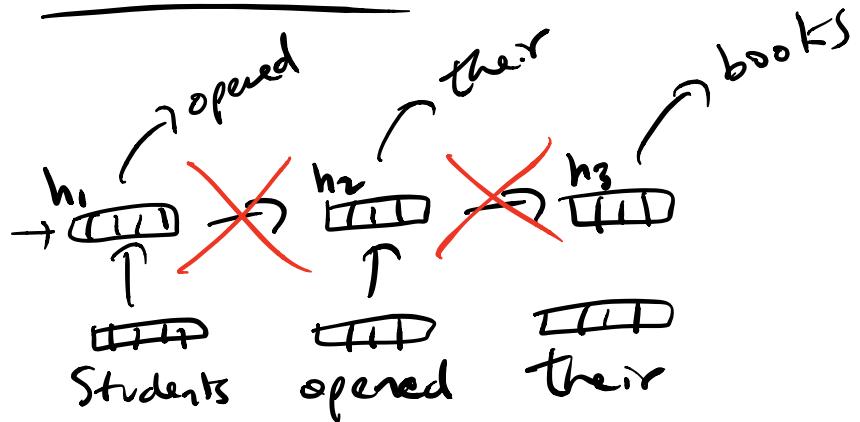
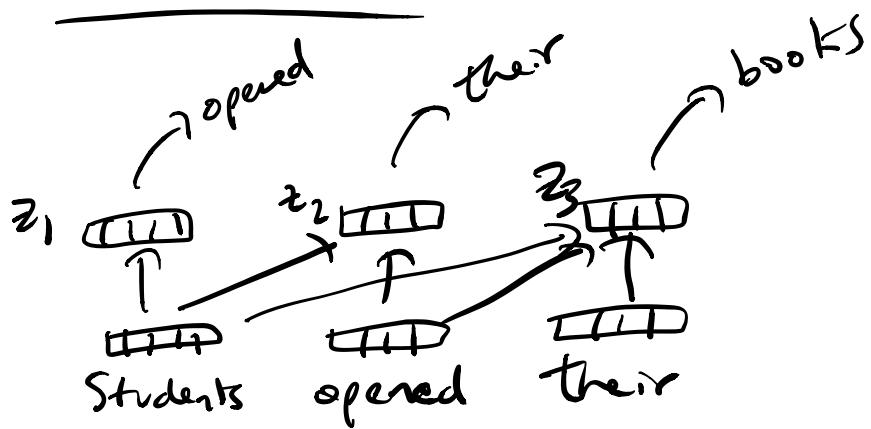


Self-attention:



if i can get rid of recurrence at
training, i can compute all h_i in parallel

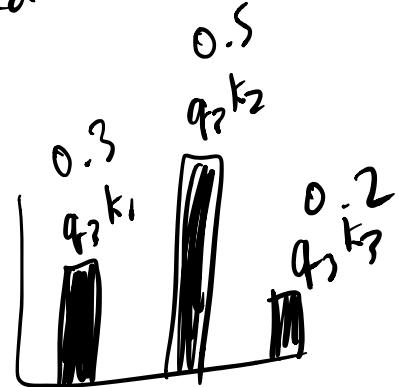
goal:



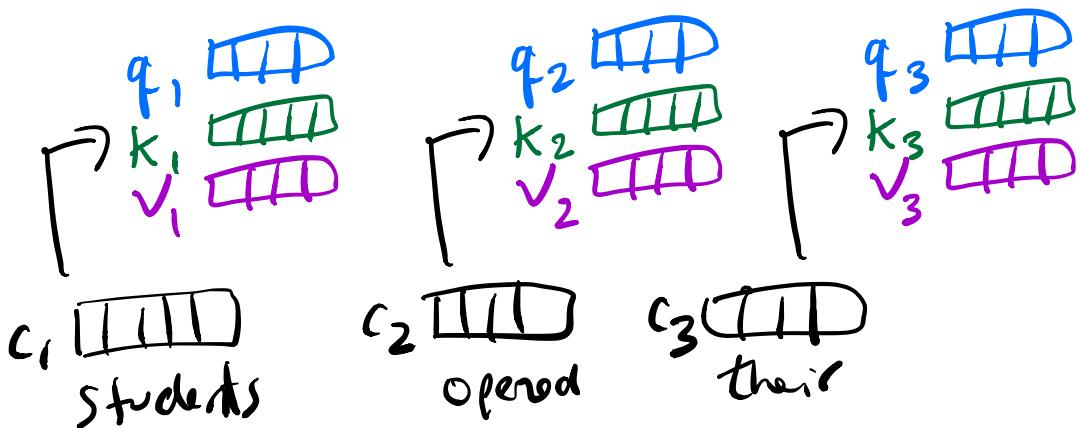
$$\hat{z}_3 = 0.3 \cdot v_1 + 0.5 \cdot v_2 + 0.2 \cdot v_3$$

$\hat{z}_3 = \boxed{\text{||||}} \rightarrow \text{predict books}$

use attn scores to compute weighted ave. of value vectors



$$\text{softmax}(q_3 \cdot k_1, q_3 \cdot k_2, q_3 \cdot k_3)$$



$$q_i = f(W_q, c_i)$$

$\xrightarrow{W_k, W_v}$

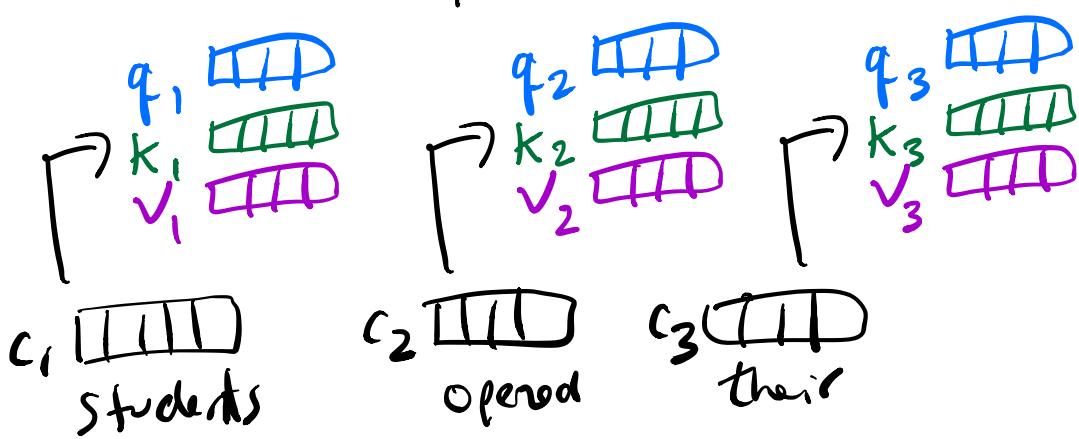
$$z_2 = 0.7 \cdot v_1 + 0.3 \cdot v_2$$

$= \boxed{\text{----}}$ \rightarrow predict their

0.7
0.3
0.0

$\langle q_2 k_1, q_2 k_2 \rangle$

\rightarrow no $q_2 k_3$
b/c of cheating



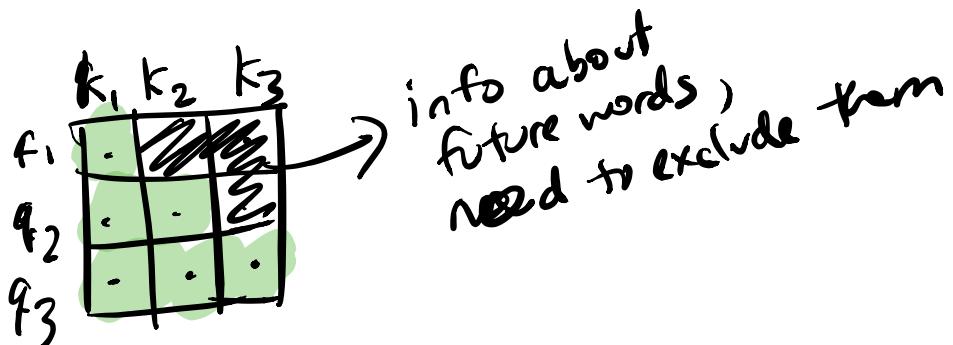
no dependency between z_1, z_2, z_3

how can we compute z_i in parallel?

attn scores

q_1		k_1		$a_1 = \langle q_1, k_1 \rangle$
q_2		k_2		$a_2 = \langle q_2, k_1, q_2, k_2 \rangle$
q_3		k_3		$a_3 = \langle q_3, k_1, q_3, k_2, q_3, k_3 \rangle$

q_1		\times		$\begin{matrix} k_1 & k_2 & k_3 \\ f_1 & \cdot & \cdot & \cdot \\ f_2 & \cdot & \cdot & \cdot \\ f_3 & \cdot & \cdot & \cdot \end{matrix}$
-------	--	----------	--	--



	k_1	k_2	k_3
f_1	-	-	.
f_2	-	-	:
f_3	-	-	.

Mask matrix

1	$-\infty$	$-\infty$
1	1	$-\infty$
1	1	1

Output goes to zero after softmax

	k_1	k_2	k_3
f_1	-	0	0
f_2	-	-	0
f_3	-	0	.