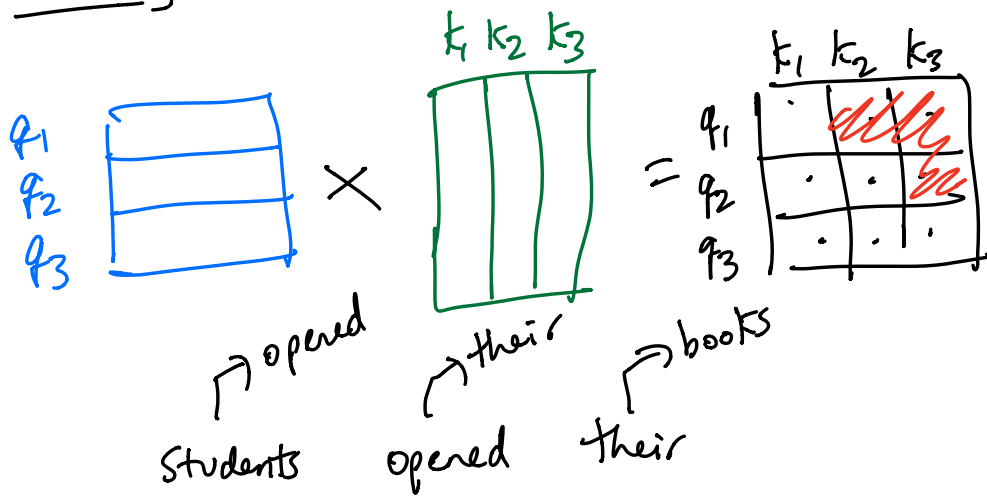


## Self-attention:

↳ training time: parallelizable

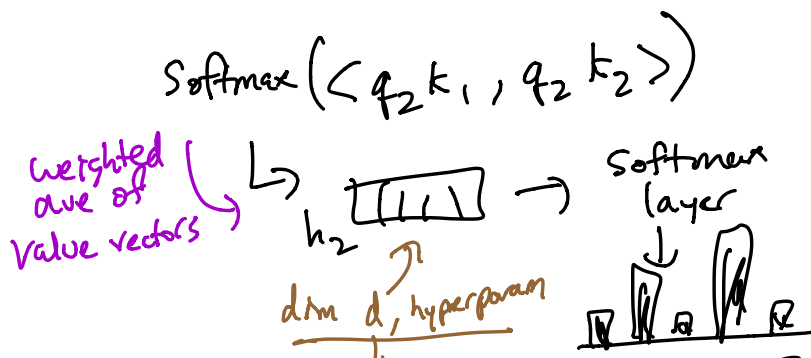
↳ test time: sequentially

### training:



- we can only parallelize this computation because we know the full sequence ahead of time

test time:  $q_1, k_1, v_1$   $q_2, k_2, v_2$  ?  
 my dog         



decoding ! <sup>dim of value vector  $a$</sup>   $\rightarrow p = \text{softmax}(w \cdot h_2)$

select a word from the pred. dist

- greedy
- sampling

$q_1 k_1 v_1$   $q_2 k_2 v_2$   $q_3 k_3 v_3$  ?  
my dog wanted ?

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

KV cache :

storing previously-computed key/value vectors so you don't have to recompute them at every timestep

$\rightarrow$  test-time

- always start w/ a  $\langle \text{SOS} \rangle$   
 $\langle \text{BOS} \rangle$  "start of seq"

---

Transformer :

$\rightarrow$  neural LM built around

multi-headed self-attention

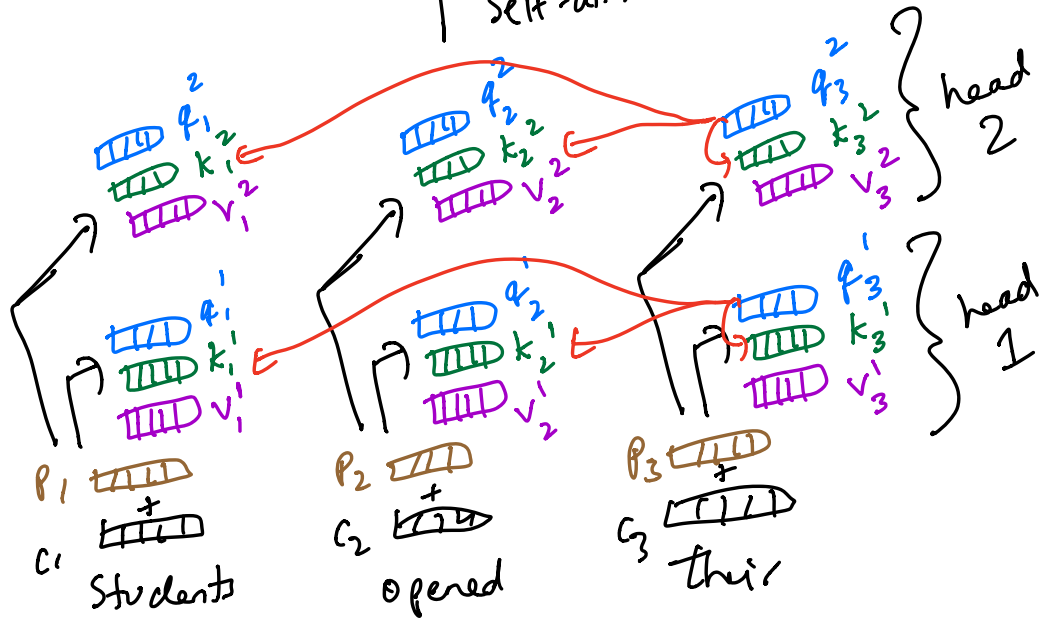
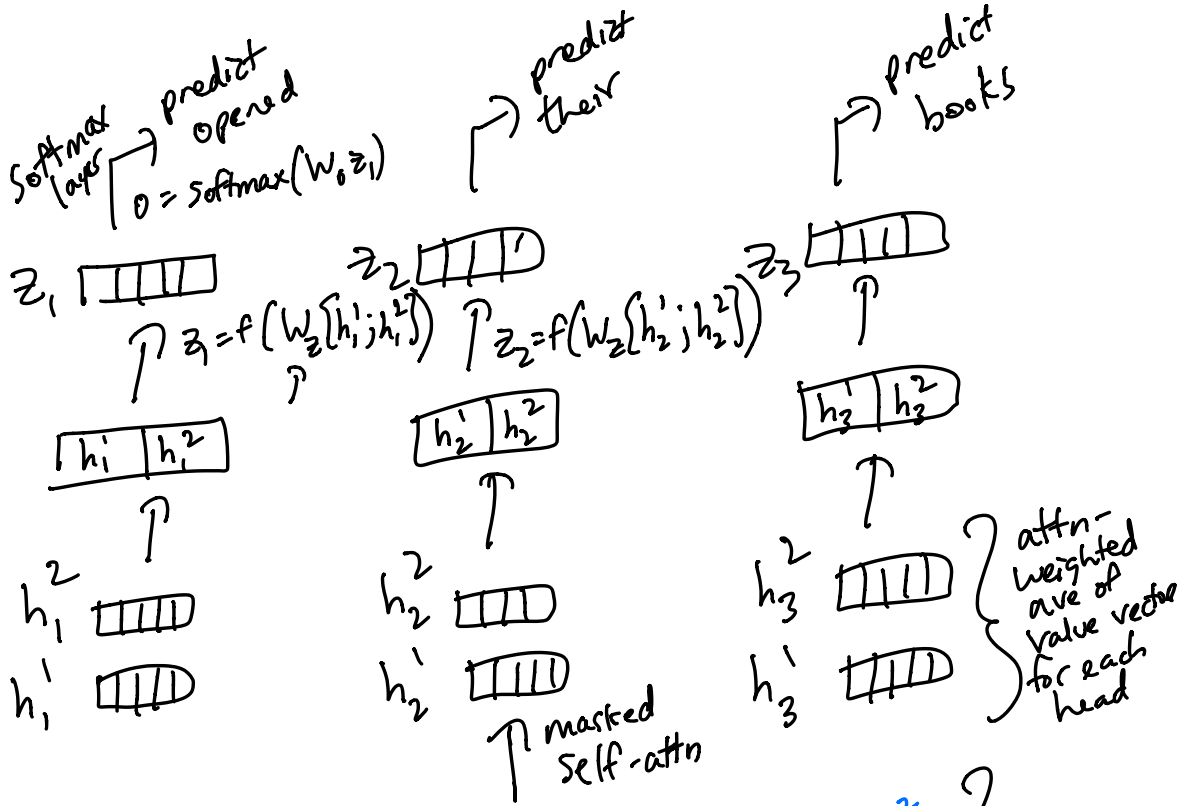
$\rightarrow$  stacked layers of attn, "deep"

multi-head self-attention :

$\rightarrow$  Vaswani et al. 2017

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

- ↳ Syntactic structures (all verbs in prefix)
- ↳ activate on certain words/phrases
- ↳ entities (dates)

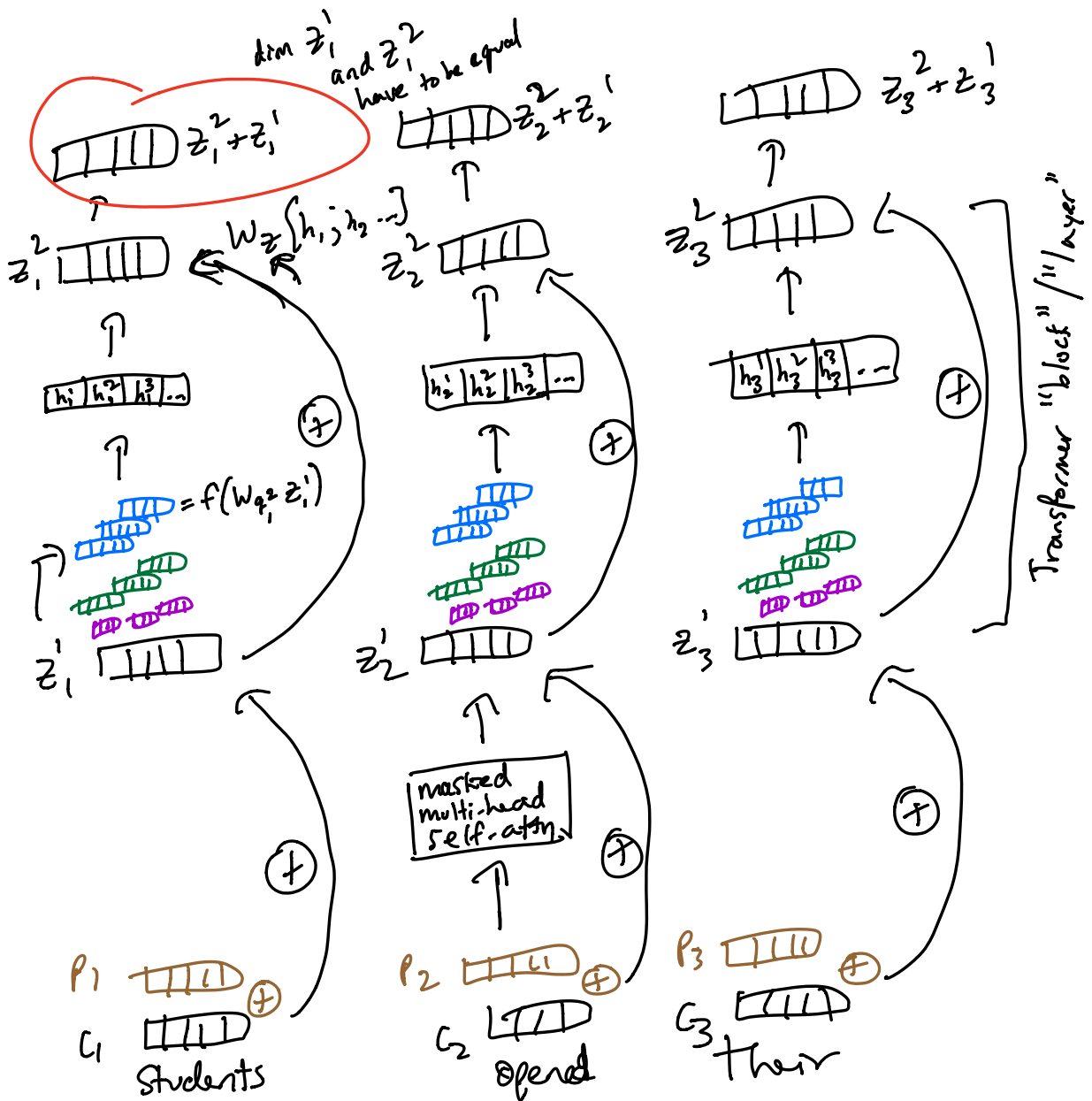


$$q_1^1 = f(W_{q^1}[P_1 + c,])$$

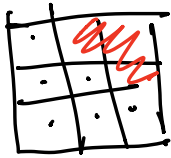
$$q_1^2 = f(W_{q^2}[P_1 + c,])$$

$$f = \text{relu}, \max(0, x)$$

Adding depth:



Decoder language model:



Sequence to -sequence models:

↳ encoder  $\Rightarrow$  source lang

↳ decoder  $\Rightarrow$  target lang

