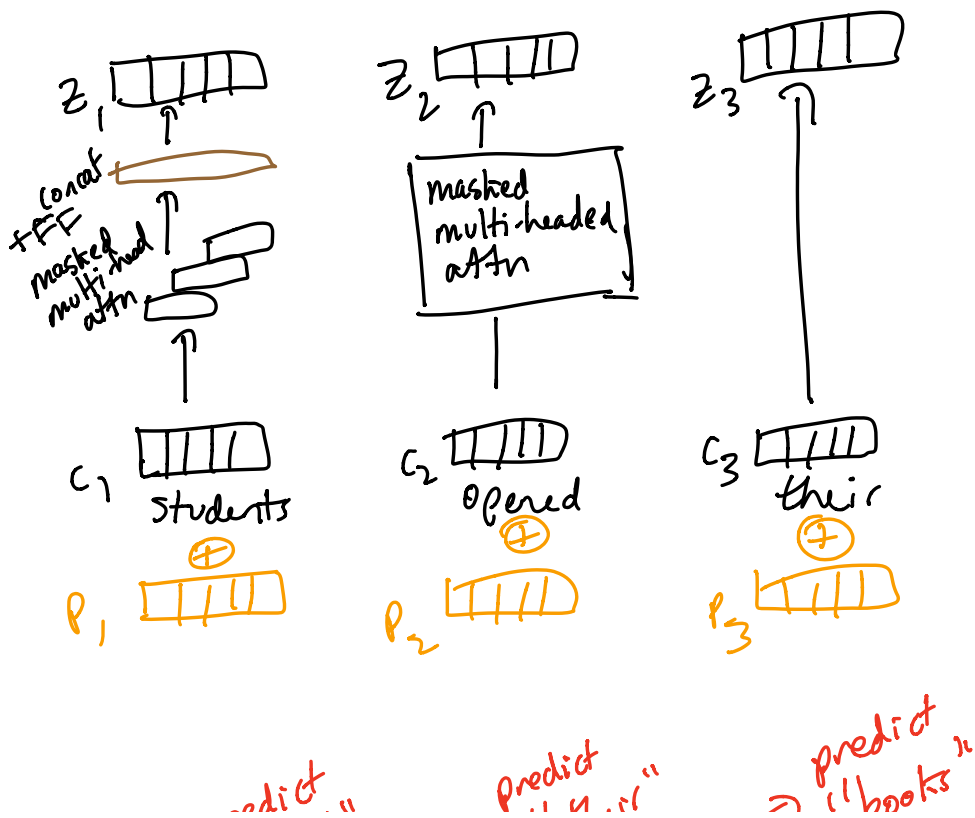
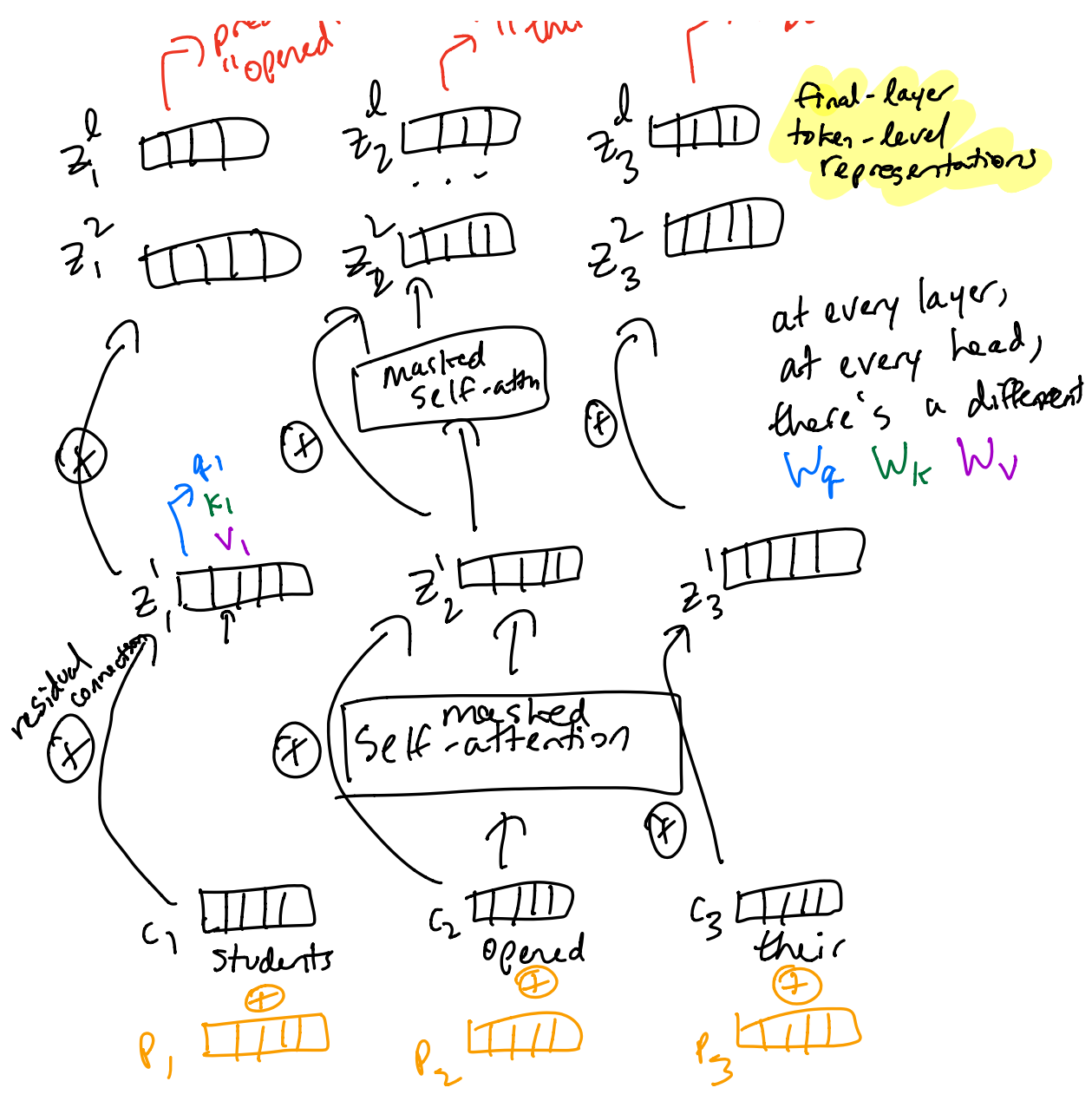


$$\left. \begin{aligned}
 q_1^1 &= f(W_{q1} c_1) \\
 q_1^2 &= f(W_{q2} c_1) \\
 q_1^n &= f(W_{qn} c_1)
 \end{aligned} \right\} \text{multi-head attn}$$

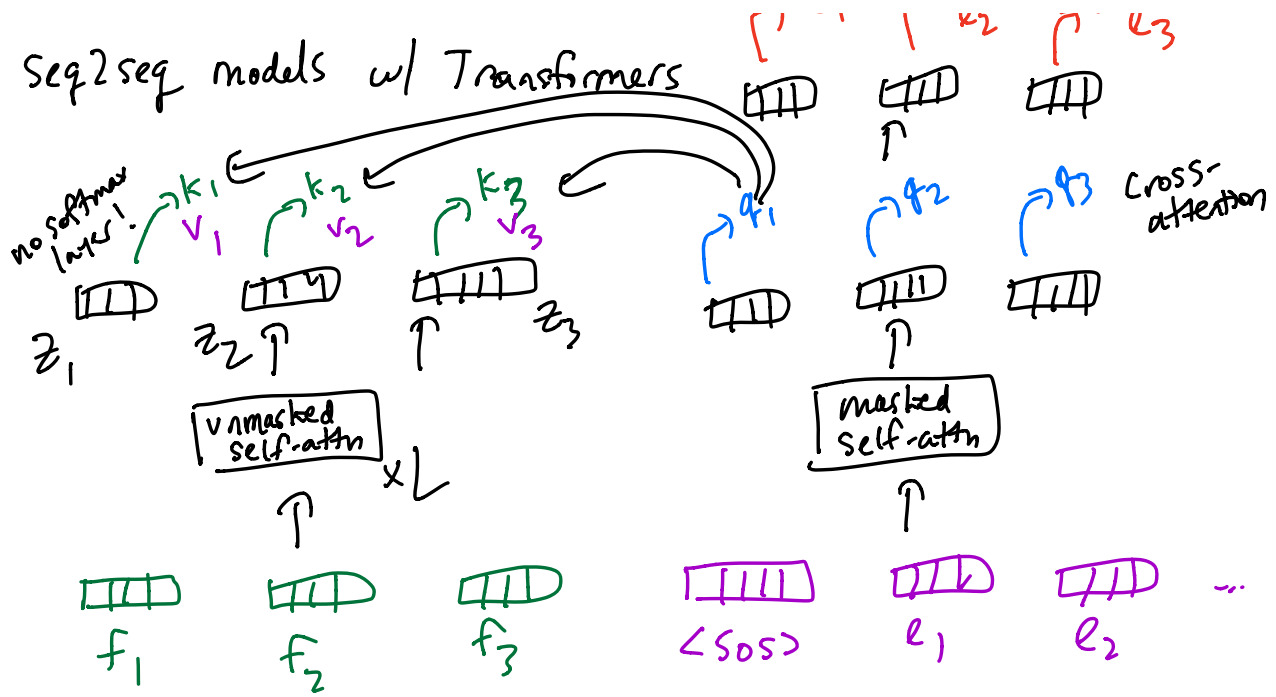
adding depth to Transformer





→ predict e_1 → predict n → predict o

Seq2seq models w/ Transformers



encoder
(encode french text)

decoder
(generating english sent)

in a decoder block, we perform both masked self-attn and cross attention