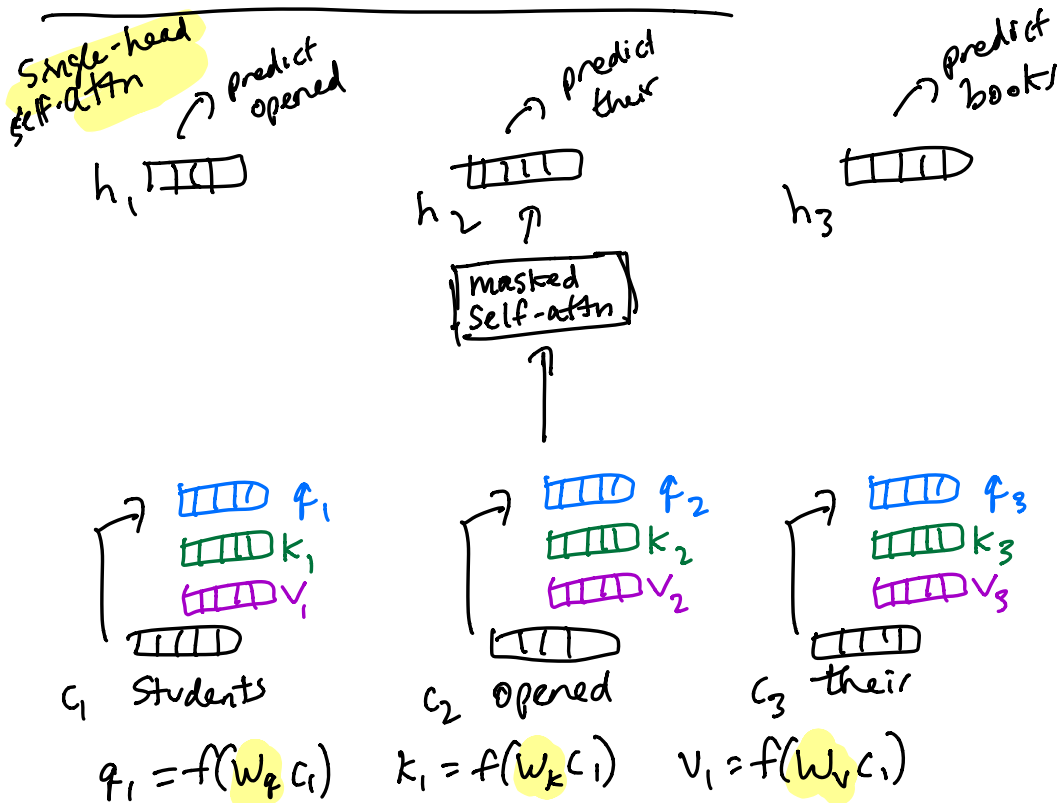# Transformers:

> a neural LM whose component is
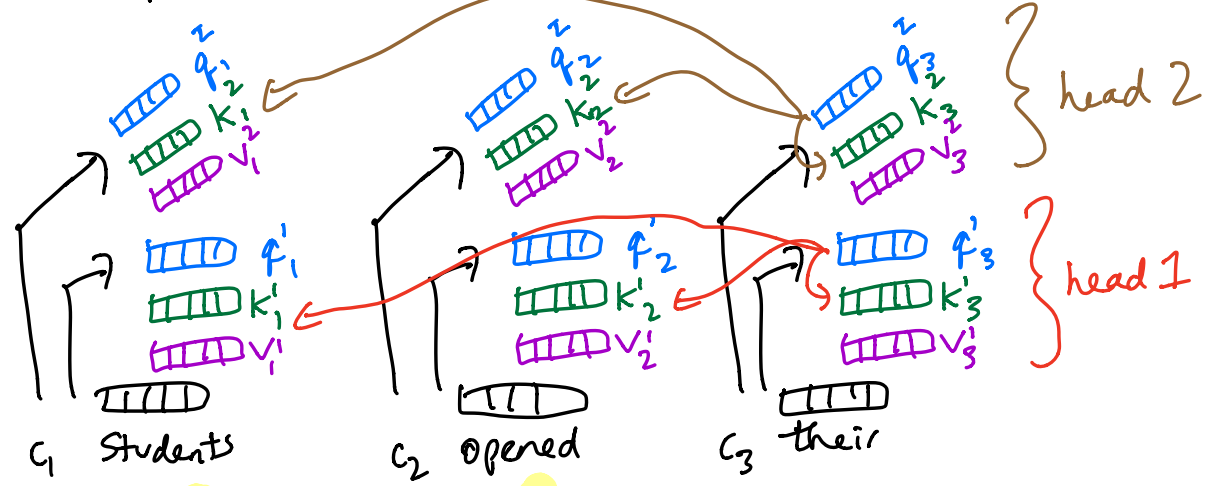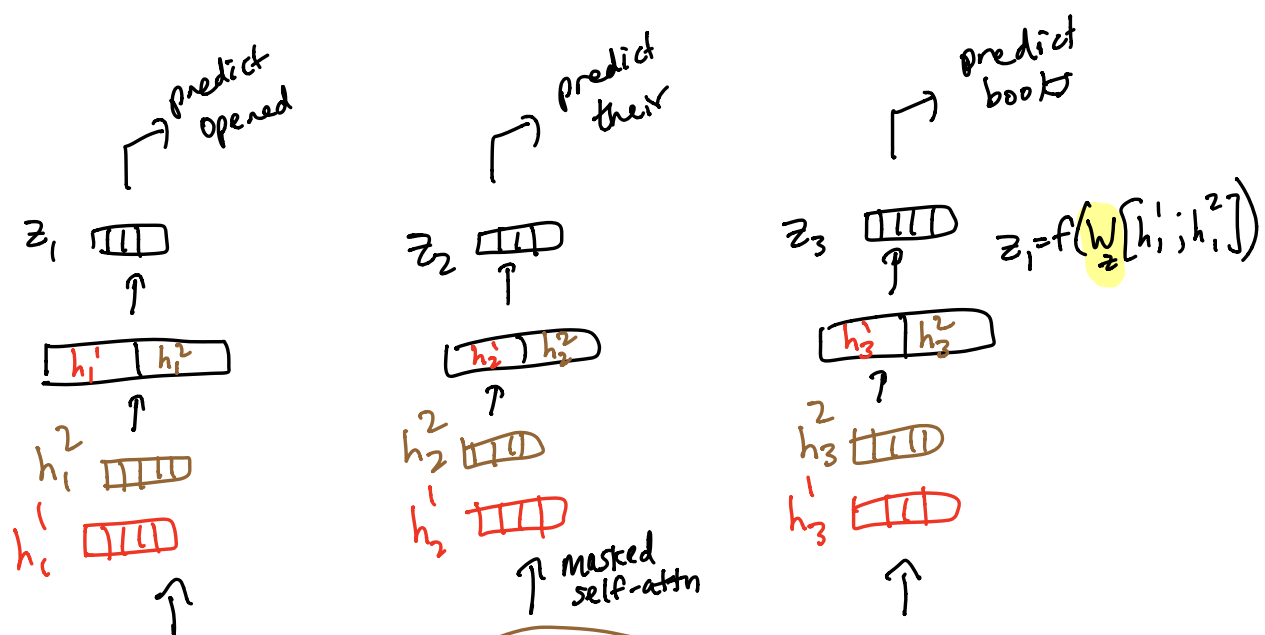> ==multi-head self-attention==

---

## multi-head self-attn:

> instead of just one set of query/key/value vectors, let's have many sets (==heads==)

> intuition: having multiple sets of $Q, K, V$ can allow each head to ==attend to== different linguistic properties of the prefix

> n-gram windows, subject/obj of sentence, discourse (global) context, entities, verbs, ...

---

==Single-head self-attn==

$h_1$ ⌗⌗⌗ ↗ predict opened

$h_2$ ↑ ⌗⌗⌗⌗ ↗ predict their

$h_3$ ⌗⌗⌗ ↗ predict books

[masked self-attn]

↑

⌗⌗⌗ $q_1$
⌗⌗⌗ $k_1$
⌗⌗⌗ $v_1$
⌗⌗⌗
$c_1$ Students

⌗⌗⌗ $q_2$
⌗⌗⌗ $k_2$
⌗⌗⌗ $v_2$
⌗⌗⌗
$c_2$ opened

⌗⌗⌗ $q_3$
⌗⌗⌗ $k_3$
⌗⌗⌗ $v_3$
⌗⌗⌗
$c_3$ their

$q_1 = f(W_q c_1)$    $k_1 = f(W_k c_1)$    $v_1 = f(W_v c_1)$

predict
opened

predict
their

predict
book

$z_1$ ⬚  $z_2$ ⬚  $z_3$ ⬚  $z_1 = f(W_z [h_1^1 ; h_1^2])$

$h_1^1$ | $h_1^2$   $h_2^1$ | $h_2^2$   $h_3^1$ | $h_3^2$

$h_1^2$ ⬚  $h_2^2$ ⬚  $h_3^2$ ⬚

$h_1^1$ ⬚  $h_2^1$ ⬚  $h_3^1$ ⬚

↑ masked
self-attn

$q_2^2$  $q_2^2$  $q_2^3$
$k_1^2$  $k_2^2$  $k_3^2$     } head 2
$v_1^2$  $v_2^2$  $v_3^2$

$q_1^1$  $q_2^1$  $q_3^1$
$k_1^1$  $k_2^1$  $k_3^1$     } head 1
$v_1^1$  $v_2^1$  $v_3^1$

$c_1$ Students   $c_2$ opened   $c_3$ their

$q_1^1 = f(W_q^1 c_1)$   $q_1^2 = f(W_q^2 c_1)$

$q_2^1 = f(W_q^1 c_2)$

Adding depth

$\rightarrow$ predict opened

$z_1^L$ 

$z_1^2$    $\{ z_1^2 = f\left(W_2^2 [h_1^1; h_{1j}^2 ...]\right)$  $\rightarrow$ ReLU   $z_2^L$   $\rightarrow$ predict their

$|h_1^1|h_1^2|h_1^3|...|$

$z_1^1$

$z_2^2$

$|h_2^1|h_2^2|h_2^3|...|$

$z_2^1$

masked multi-head self-attn

student

$P_1$

opened

$P_2$

$\rightarrow$ predict books!

$z_3^L$    } final-layer token-level representation

$z_3^2$    } second layer

$|h_3^1|h_3^2|h_3^3|...|$

residual: input to next layer is $z_3^1 + z_3^2$

$z_3^1$    } first layer

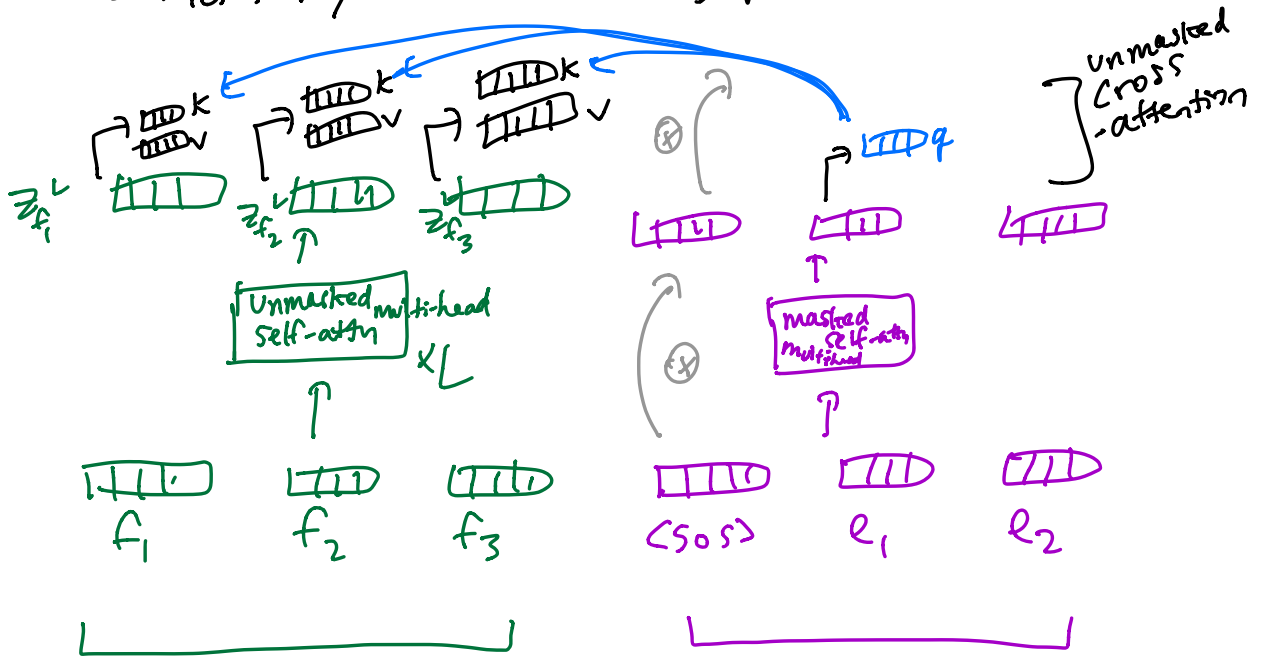residual connection: input to next layer $= z_3^1 + c_3 + P_3$

their    } embedding layer

$P_3$

what if we want to give the model
some input and have it generate a completion

↳ let's say we're translating from French to English



unmasked
cross
-attention

$z_{f_1}$  $z_{f_2}$  $z_{f_3}$

Unmasked multi-head self-attn  ×L

masked self-attn multihead

$f_1$  $f_2$  $f_3$   <SOS>  $e_1$  $e_2$

encoder
(no need to predict
the next word)

decoder:
responsible for generating text

$$P(e_n \mid e_{1 \dots n-1}, f)$$

↑
conditional LM