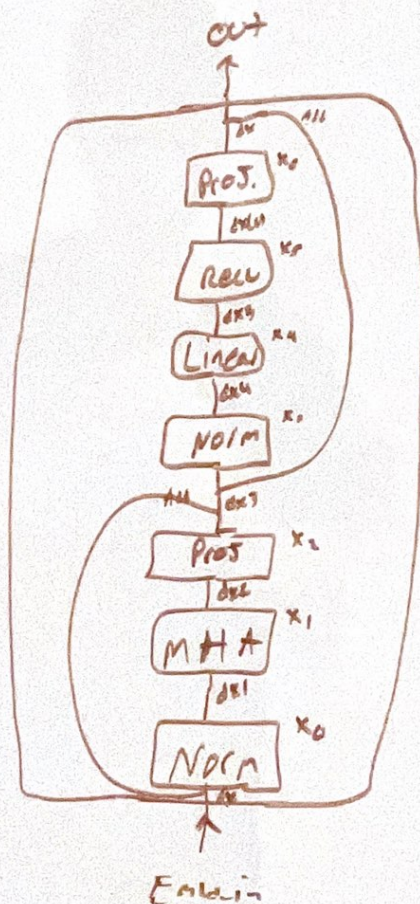


Transformer Block →  
Deriv.



$dx_n \rightarrow$  deriv  
of  $n$ th input  
(to attn block)

$$Embed = Out + dx_n$$

• Feed all  $\frac{\partial L}{\partial x_m}$  to each  $x_i$ .

- Derivatives 'skip' w/  
few connections, so add previous,  
dx thus is connects to.

↓  
This way,  
derivs pass thru  
down from net!