

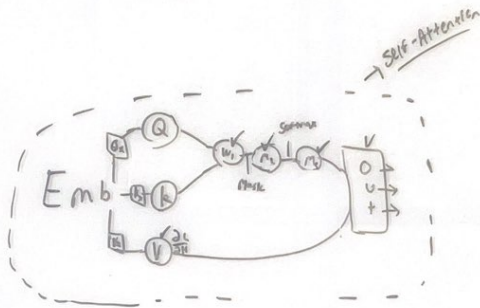
Deriving Self-Attention Head (w_1, m_1, m_2, v)

• Receive: $\frac{\partial L}{\partial \text{out}} \rightarrow \text{shape } \theta, T, C$

• Want: $\frac{\partial L}{\partial K}, \frac{\partial L}{\partial Q}, \frac{\partial L}{\partial V}$, $\rightarrow \text{output } \frac{\partial L}{\partial \text{Emb (Embed)}}$

$$\frac{\partial L}{\partial K} \rightarrow \begin{matrix} \text{key} \\ \uparrow \\ \theta, T, C \end{matrix} \otimes \begin{matrix} \text{value} \\ \uparrow \\ C, H\text{-size} \end{matrix}$$

• 3rd... affects each batch-matrix



$$B \begin{bmatrix} T \\ \vdots \end{bmatrix} \otimes T \begin{bmatrix} H \\ \vdots \end{bmatrix} \rightarrow \text{out} \begin{bmatrix} T \\ \vdots \end{bmatrix} \dots$$

$\frac{\partial L}{\partial H} \rightarrow$ • Each E_{IJ} affects all batches (sum)
• Each E_{IJ} affects I th col of M - J th col of out !

$$\frac{\partial L}{\partial E_{IJ}} \rightarrow \left(\sum_x \frac{\partial L}{\partial \text{out}_{xj}} \cdot M_{xi} \right)$$

$$\downarrow$$

$$M^T \otimes \frac{\partial L}{\partial \text{out}}, \text{sum}(0)$$

$\frac{\partial L}{\partial M} \rightarrow$ • Each E_{IJ} affects col j ! \rightarrow row i of out by row of V $\rightarrow \sum_x \frac{\partial L}{\partial \text{out}_{ix}} \cdot H_{jx}$

$$\frac{\partial L}{\partial r_2} = \left(\frac{\partial L}{\partial \text{out}} \otimes V^T \right)$$

$$\frac{\partial L}{\partial m_i} \rightarrow y = \frac{e_i^x}{\sum_j e_j^x} \rightarrow -\frac{e_i^x}{(\sum_j e_j^x)^2} + \frac{e_i^x}{\sum_j e_j^x}$$

$\frac{\partial L}{\partial w} \rightarrow$ "Soft Grad" gradients to w are not 0!

$-\frac{e_i^x \cdot e_i^x}{(\sum_j e_j^x)^2}$
 $-\frac{e_i^x}{\sum_j e_j^x} \cdot \frac{e_i^x}{\sum_j e_j^x}$

• Sum columns, multi by softmax
• for i th, add softmax!