

Activity Sheet 05

Name: _____ ID: _____

Question 01: Multihead Self Attention

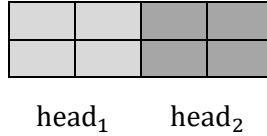
Multi-head self-attention is computed using the following formulas:

$$\mathbf{Q}^i = \mathbf{XW}_i^{\mathbf{Q}} \quad \mathbf{K}^i = \mathbf{XW}_i^{\mathbf{K}} \quad \mathbf{V}^i = \mathbf{XW}_i^{\mathbf{V}}$$

$$\text{head}_i = \text{Self-Attention}(\mathbf{A}) = \left(\text{softmax} \left(\frac{\mathbf{Q}^i \mathbf{K}^{i\mathbf{T}}}{\sqrt{d_k}} \right) \right) \mathbf{V}^i$$

$$\text{MultiHead Attention}(\mathbf{M}) = (\text{head}_1 \oplus \text{head}_2 \dots \oplus \text{head}_h) \mathbf{W}^0$$

The concatenation simply joins the results of attention heads together



As an example, consider the following results for 2 attention heads

$$\text{head}_1 = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}, \quad \text{head}_2 = \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix}$$

The concatenation will be

$$\text{head}_1 \oplus \text{head}_2 = \begin{bmatrix} 1 & 2 & 5 & 6 \\ 3 & 4 & 7 & 8 \end{bmatrix}$$

Given the input matrix:

$$\mathbf{X} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

where each row represents a word in the input sequence, and the weight matrices for two attention heads are:

Head 1:

$$\mathbf{W}_1^{\mathbf{Q}} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{W}_1^{\mathbf{K}} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \mathbf{W}_1^{\mathbf{V}} = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$$

Head 2:

$$\mathbf{W}_2^{\mathbf{Q}} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \mathbf{W}_2^{\mathbf{K}} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{W}_2^{\mathbf{V}} = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$$

CS 335: Introduction to Large Language Models

Habib University

Final Projection Matrix:

$$\mathbf{W}^O = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ -1 & 0 \\ 0 & -1 \end{bmatrix}$$

Compute the Multi-Head Self-Attention output (M).