



Figure 1: Multi-head Attention.

$$\underline{Q}^{[L] \times [d_q]} = \text{prior} \quad (1a)$$

$$\underline{K}^{[L] \times [d_k]} = \text{prior} \quad (1b)$$

$$\underline{V}^{[L] \times [d_v]} = \text{prior} \quad (1c)$$

$$\underline{1}^{3 \times 4} = \text{linear}(\underline{Q}^{[L] \times [d_q]}) \quad (1d)$$

$$\underline{2}^{3 \times 4} = \text{linear}(\underline{Q}^{[L] \times [d_q]}) \quad (1e)$$

$$\underline{3}^{3 \times 4} = \text{linear}(\underline{Q}^{[L] \times [d_q]}) \quad (1f)$$

$$4^{3 \times 4} = \text{linear}(K^{[L] \times [d_{\mathbf{k}}]}) \quad (1\text{g})$$

$$5^{3 \times 4} = \text{linear}(K^{[L] \times [d_{\mathbf{k}}]}) \quad (1\text{h})$$

$$6^{3 \times 4} = \text{linear}(K^{[L] \times [d_{\mathbf{k}}]}) \quad (1\text{i})$$

$$7^{3 \times 4} = \text{linear}(V^{[L] \times [d_{\mathbf{v}}]}) \quad (1\text{j})$$

$$8^{3 \times 4} = \text{linear}(V^{[L] \times [d_{\mathbf{v}}]}) \quad (1\text{k})$$

$$9^{3 \times 4} = \text{linear}(V^{[L] \times [d_{\mathbf{v}}]}) \quad (1\text{l})$$

$$X^{3 \times 4} = \text{scaled_dot_prod_att}(1^{3 \times 4}, 2^{3 \times 4}, 3^{3 \times 4}, 4^{3 \times 4}, 5^{3 \times 4}, 6^{3 \times 4}, 7^{3 \times 4}, 8^{3 \times 4}, 9^{3 \times 4}) \quad (1\text{m})$$

$$Y^{3 \times 4} = \text{scaled_dot_prod_att}(1^{3 \times 4}, 2^{3 \times 4}, 3^{3 \times 4}, 4^{3 \times 4}, 5^{3 \times 4}, 6^{3 \times 4}, 7^{3 \times 4}, 8^{3 \times 4}, 9^{3 \times 4}) \quad (1\text{n})$$

$$Z^{3 \times 4} = \text{scaled_dot_prod_att}(1^{3 \times 4}, 2^{3 \times 4}, 3^{3 \times 4}, 4^{3 \times 4}, 5^{3 \times 4}, 6^{3 \times 4}, 7^{3 \times 4}, 8^{3 \times 4}, 9^{3 \times 4}) \quad (1\text{o})$$

$$C^{3 \times 4} = [X^{3 \times 4} | Y^{3 \times 4} | Z^{3 \times 4}] \quad (1\text{p})$$

$$L^{3 \times 4} = C^{3 \times 4} W_{\underline{\mathbf{c}}}^{[d_{\mathbf{v}}] \times [d]} \quad (1\text{q})$$