



Figure 1: Multi-head Attention.

$$\underline{Q}^{[D],[L]} =; prior \quad (1a)$$

$$\underline{K}^{[D],[L]} =; prior \quad (1b)$$

$$\underline{V}^{[D],[L]} =; prior \quad (1c)$$

$$\underline{Q}_0^{[d],[L]} = \text{linear}(\underline{Q}^{[D],[L]}) \quad (1d)$$

$$Q_1^{[d],[L]} = \text{linear}(Q^{[D],[L]}) \quad (1e)$$

$$Q_2^{[d],[L]} = \text{linear}(Q^{[D],[L]}) \quad (1f)$$

$$K_0^{[d],[L]} = \text{linear}(K^{[D],[L]}) \quad (1g)$$

$$K_1^{[d],[L]} = \text{linear}(K^{[D],[L]}) \quad (1h)$$

$$K_2^{[d],[L]} = \text{linear}(K^{[D],[L]}) \quad (1i)$$

$$V_0^{[d],[L]} = \text{linear}(V^{[D],[L]}) \quad (1j)$$

$$V_1^{[d],[L]} = \text{linear}(V^{[D],[L]}) \quad (1k)$$

$$V_2^{[d],[L]} = \text{linear}(V^{[D],[L]}) \quad (1l)$$

$$A_0^{[d],[L]} = \text{scaled\_dot\_prod\_att}(Q_0^{[d],[L]}, K_0^{[d],[L]}, V_0^{[d],[L]}) \quad (1m)$$

$$A_1^{[d],[L]} = \text{scaled\_dot\_prod\_att}(Q_1^{[d],[L]}, K_1^{[d],[L]}, V_1^{[d],[L]}) \quad (1n)$$

$$A_2^{[d],[L]} = \text{scaled\_dot\_prod\_att}(Q_2^{[d],[L]}, K_2^{[d],[L]}, V_2^{[d],[L]}) \quad (1o)$$

$$A^{[D],[L]} = [A_0^{[d],[L]} | A_1^{[d],[L]} | A_2^{[d],[L]}] \quad (1p)$$

$$O^{[L]} = W_{\underline{e}}^{[1][D]} A^{[D],[L]} \quad (1q)$$