



Figure 1: Multi-head Attention with 2 heads.

$$\underline{A}^{[D],[\ell]} = [\underline{A}_0^{[d],[\ell]} | \underline{A}_1^{[d],[\ell]}] \quad (1a)$$

$$\underline{A}_0^{[d],[\ell]} = \text{scaled_dot_prod_att}(\underline{Q}_0^{[d],[\ell]}, \underline{K}_0^{[d],[\ell]}, \underline{V}_0^{[d],[\ell]}) \quad (1b)$$

$$\underline{A}_1^{[d],[\ell]} = \text{scaled_dot_prod_att}(\underline{Q}_1^{[d],[\ell]}, \underline{K}_1^{[d],[\ell]}, \underline{V}_1^{[d],[\ell]}) \quad (1c)$$

$$\underline{K}^{[D],[\ell]} = \underline{W}_k^{[D],[d]} \underline{e}^{[d],[\ell]} \quad (1d)$$

$$K_0^{[d],[\ell]} = \text{linear}(K^{[D],[\ell]}) \quad (\text{split, then project a component}) \quad (1e)$$

$$K_1^{[d],[\ell]} = \text{linear}(K^{[D],[\ell]}) \quad (\text{split, then project a component}) \quad (1f)$$

$$O^{[D],[\ell]} = W_{\underline{0}}^{[D],[D]} A^{[D],[\ell]} \quad (1g)$$

$$Q^{[D],[\ell]} = W_{\underline{q}}^{[D],[d]} e^{[d],[\ell]} \quad (1h)$$

$$Q_0^{[d],[\ell]} = \text{linear}(Q^{[D],[\ell]}) \quad (\text{split, then project a component}) \quad (1i)$$

$$Q_1^{[d],[\ell]} = \text{linear}(Q^{[D],[\ell]}) \quad (\text{split, then project a component}) \quad (1j)$$

$$V^{[D],[\ell]} = W_{\underline{v}}^{[D],[d]} e^{[d],[\ell]} \quad (1k)$$

$$V_0^{[d],[\ell]} = \text{linear}(V^{[D],[\ell]}) \quad (\text{split, then project a component}) \quad (1l)$$

$$V_1^{[d],[\ell]} = \text{linear}(V^{[D],[\ell]}) \quad (\text{split, then project a component}) \quad (1m)$$

$$e^{[d],[\ell]} = \text{prior} \quad (1n)$$