



Figure 1: Scaled Dot Product Attention.

$$A^{[d],[\ell]} = V^{[d],[\ell]} G^{[\ell],[\ell]} \quad (1a)$$

$$B^{[\ell],[\ell]} = (Q^{[d],[\ell]})^T K^{[d],[\ell]} \quad (1b)$$

$$G^{[\ell],[\ell]} = \text{softmax}(R^{[\ell],[\ell]}) \quad (1c)$$

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

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

$$R^{[\ell],[\ell]} = \text{mask}(Y^{[\ell],[\ell]}) \quad (1f)$$

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

$$Y^{[\ell],[\ell]} = \frac{B^{[\ell],[\ell]}}{\sqrt{d_{\underline{k}}}} \tag{1h}$$