



Figure 1: Scaled Dot Product Attention.

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

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

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

$$B^{[L] \times [L]} = Q^{[L] \times [d_q]} (K^{[L] \times [d_k]})^\dagger \quad (1d)$$

$$Y^{[L] \times [L]} = \frac{B^{[L] \times [L]}}{\sqrt{d_k}} \quad (1e)$$

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

$$G^{[L] \times [L]} = \text{softmax}(R^{[L] \times [L]}) \quad (1g)$$

$$P^{[L] \times [d_v]} = G^{[L] \times [L]} V^{[L] \times [d_v]} \quad (1h)$$