



Figure 1: Encoder of Vanilla Transformer Net. N copies of the boxed part are connected in series.

$$\underline{A}^{[D],[\ell]} = \text{Attention}(\underline{Q}^{[D],[\ell]}, \underline{K}^{[D],[\ell]}, \underline{V}^{[D],[\ell]}) \quad (1a)$$

$$\underline{F}^{[d],[\ell]} = \text{feed_forward_nn}(\underline{N}^{[d],[\ell]}) \quad (1b)$$

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

$$N^{[d],[\ell]} = \text{normalize}(e^{[d],[\ell]} + W_{\underline{a}}^{[d],[D]} A^{[D],[\ell]}) \quad (1d)$$

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

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

$$e^{[d],[\ell]} = E^{[d],[L]} x^{[L],[\ell]} \quad (1g)$$

$$n^{[d],[\ell]} = \text{normalize}(N^{[d],[\ell]} + F^{[d],[\ell]}) \quad (1h)$$

$$x^{[L],[\ell]} = \text{prior} \quad (1i)$$