



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

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

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

$$F^{[d],[\ell]} = \text{feed\_forward\_nn}(N^{[d],[\ell]}) \quad (1c)$$

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

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

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

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

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

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