

Figure 1: Encoder.

$$F^{[D],[\ell]} = \text{feed_forward_nn}(N^{[D],[\ell]})$$
(1a)

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

$$N^{[D],[\ell]} = \text{normalize}(e^{[d],[\ell]} + O^{[D],[\ell]})$$
 (1c)

$$O^{[D],[\ell]} = \text{multi_headed_attention}(Q^{[D],[\ell]}, K^{[D],[\ell]}, V^{[D],[\ell]}) \tag{1d}$$

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

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

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

$$k^{[D],[\ell]} = n^{[D],[\ell]}$$
 (1h)

$$n^{[D],[\ell]} = \text{normalize}(N^{[D],[\ell]} + F^{[D],[\ell]})$$

$$\tag{1i}$$

$$q^{[D],[\ell]} = n^{[D],[\ell]}$$
 (1j)

$$x^{[L],[\ell]} = \text{prior} \tag{1k}$$