



Figure 1: Decoder.

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

$$I^{[L],[\ell]} = W^{[L],[D]} Y^{[D],[\ell]} \quad (1b)$$

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

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

$$P^{[L],[\ell]} = \text{softmax}(I^{[L],[\ell]}) \quad (\sum_{\alpha \in [\ell]} P^{[L],\alpha} = 1) \quad (1e)$$

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

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

$$Y^{[D],[\ell]} = \text{normalize}(F^{[D],[\ell]} + a^{[D],[\ell]}) \quad (1h)$$

$$a^{[D],[\ell]} = \text{normalize}(O^{[D],[\ell]} + e^{[D],[\ell]}) \quad (1i)$$

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

$$j^{[D],[\ell]} = \text{normalize}(o^{[D],[\ell]} + a^{[D],[\ell]}) \quad (1k)$$

$$k^{[D],[\ell]} = \quad (1l)$$

$$o^{[D],[\ell]} = \text{multi_head_attention}(q^{[D],[\ell]}, k^{[D],[\ell]}, v^{[D],[\ell]}) \quad (1m)$$

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

$$v^{[D],[\ell]} = a^{[D],[\ell]} \quad (1o)$$

$$x^{[L],[\ell]} = \tag{1p}$$