



Figure 1: Decoder.

$$\underline{F}^{[D],[\ell]} = \text{feed\_forward\_nn}(\underline{j}^{[D],[\ell]}) \quad (1a)$$

$$\underline{G}^{[L],[\ell]} = \underline{I}^{[L],[\ell]} \quad (1b)$$

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

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

$$O^{[D],[\ell]} = \text{multi\_head\_attention}(Q^{[D],[\ell]}, K^{[D],[\ell]}, V^{[D],[\ell]}) \quad (1e)$$

$$Q^{[D],[\ell]} = p^{[D],[\ell]} \quad (1f)$$

$$R^{[L],[\ell]} = \quad (1g)$$

$$V^{[D],[\ell]} = p^{[D],[\ell]} \quad (1h)$$

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

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

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

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

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

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

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

$$v^{[D],[\ell]} = a^{[D],[\ell]} \tag{1p}$$