



Figure 1: Sentence Ax tranet. $D = dn_h$ where $d = 768$ is the hidden dimension per head, and $n_h = 12$ is the number of heads. 2 copies of dashed box connected in series. 5 copies of plain box connected in series.

$$\underline{A}^{[D],[105]} = \text{Attention}(\underline{Q}^{[D],[105]}, \underline{K}^{[D],[105]}, \underline{V}^{[D],[105]}) \quad (1a)$$

$$\underline{B}^{[768],[105]} = \text{BERT}() \quad (1b)$$

$$\underline{E}^{[768],[105]} = \underline{X}^{[768],[105]} \quad (1c)$$

$$K^{[D],[105]} = W_{\underline{k}}^{[D],[768]} B^{[768],[105]} \quad (1d)$$

$$L^{[6],[84]} = \text{ilabel}(M^{[768],[105]}) \quad (1e)$$

$$M^{[768],[105]} = \text{merge}(S^{[768],[105]}) \quad (1f)$$

$$n^{[768],[105]} = W_{\underline{a}}^{[768],[D]} A^{[D],[105]} \quad (1g)$$

$$Q^{[D],[105]} = W_{\underline{q}}^{[D],[768]} B^{[768],[105]} \quad (1h)$$

$$S^{[768],[105]} = X^{[768],[105]} + n^{[768],[105]} \quad (1i)$$

$$V^{[D],[105]} = W_{\underline{v}}^{[D],[768]} B^{[768],[105]} \quad (1j)$$

$$X^{[768],[105]} = 0 \quad (1k)$$