

Figure 1: Multi-head Attention

$$Q^{3\times 4} =) \tag{1a}$$

$$K^{3\times4} = ) \tag{1b}$$

$$V^{3\times4} = ) \tag{1c}$$

$$1^{3\times 4} = \operatorname{linear}(Q^{3\times 4}) \tag{1d}$$

$$2^{3\times 4} = \operatorname{linear}(Q^{3\times 4}) \tag{1e}$$

$$3^{3\times 4} = \operatorname{linear}(Q^{3\times 4}) \tag{1f}$$

$$4^{3\times 4} = \operatorname{linear}(K^{3\times 4}) \tag{1g}$$

$$5^{3\times4} = \operatorname{linear}(K^{3\times4}) \tag{1h}$$

$$6^{3\times 4} = \operatorname{linear}(K^{3\times 4}) \tag{1i}$$

$$7^{3\times4} = \operatorname{linear}(V^{3\times4}) \tag{1j}$$

$$8^{3\times4} = \operatorname{linear}(V^{3\times4}) \tag{1k}$$

$$9^{3\times4} = \operatorname{linear}(V^{3\times4}) \tag{11}$$

$$X^{3\times4} = \text{scaled\_dot\_prod\_att}(1^{3\times4}, 2^{3\times4}, 3^{3\times4}, 4^{3\times4}, 5^{3\times4}, 6^{3\times4}, 7^{3\times4}, 8^{3\times4}, 9^{3\times4}) \tag{1m}$$

$$Y^{3\times 4} = \text{scaled\_dot\_prod\_att}(1^{3\times 4}, 2^{3\times 4}, 3^{3\times 4}, 4^{3\times 4}, 5^{3\times 4}, 6^{3\times 4}, 7^{3\times 4}, 8^{3\times 4}, 9^{3\times 4}) \tag{1n}$$

$$Z^{3\times4} = \text{scaled\_dot\_prod\_att}(1^{3\times4}, 2^{3\times4}, 3^{3\times4}, 4^{3\times4}, 5^{3\times4}, 6^{3\times4}, 7^{3\times4}, 8^{3\times4}, 9^{3\times4}) \tag{10}$$

$$C^{3\times 4} = \text{concat}(X^{3\times 4}, Y^{3\times 4}, Z^{3\times 4})$$
 (1p)

$$L^{3\times4} = \operatorname{concat}(C^{3\times4}) \tag{1q}$$