

$$x = [1:l]$$

$$= x[1] \cdot x[2] \cdot \dots \cdot x[l]$$

$$x[1] = \text{length} - de$$

$$x[1] \in \mathbb{R}^{d_{in} \times 1} = \mathbb{R}^{d_{ex} \times 1}$$

$$\underline{d_{in} = de}$$

$$\left\{ \begin{array}{l} x \in \mathbb{R}^{d_{ex} \times 1} \end{array} \right\}$$

$$q[i] = \underbrace{w_q}_{\mathbb{R}^{d_{atten} \times d_{in}}} \cdot x[i] \in \mathbb{R}^{d_{atten} \times 1}$$

$$\left\{ \begin{array}{l} x[i] \in \mathbb{R}^{d_{in} \times 1} \\ w_q \in \mathbb{R}^{d_{atten} \times d_{in}} \\ = w_k \end{array} \right\}$$

$$k[i] = w_k \cdot x[i] = \mathbb{R}^{d_{atten} \times d_{in}} \cdot \mathbb{R}^{d_{in} \times 1} = \mathbb{R}^{d_{atten} \times 1}$$

$$v[i] \in \mathbb{R}^{d_{out}} = \mathbb{R}^{d_{out} \times d_{in}} \cdot \mathbb{R}^{d_{in} \times 1}$$

$$= w_v \cdot x[i]$$

$$\left\{ \begin{array}{l} w_v \in \mathbb{R}^{d_{out} \times d_{in}} \end{array} \right\}$$

$$\in \mathbb{R}^{1 \times d_{hidden}}$$

$$w_1 \tilde{v}[i] = \mathbb{R}^{d_{hidden} \times d_{out}} \times \mathbb{R}^{d_{out} \times l} = \tilde{x}[i]$$

$$\begin{array}{l} q[i] \in \mathbb{R}^{1 \times d_{atten}} \\ k[i] \in \mathbb{R}^{d_{atten} \times 1} \\ \in \mathbb{R}^{1 \times 1} \end{array}$$

$$w_2 \in \mathbb{R}^{d_{hidden} \times 1}$$

$$w_1 \in \mathbb{R}^{d_{hidden} \times d_{out}}$$

$$\frac{q[i, j] \cdot v[j]}{\sum_{j=1}^l q[i, j] \cdot v[j]} \in \mathbb{R}^{d_{out} \times 1}$$

$$\frac{q[i]^T \cdot k[i]}{\sum_{i=1}^l q[i]^T \cdot k[i]} \in \mathbb{R}^{1 \times l}$$

$$\boxed{d_{out} = l} \quad \langle 1 \rangle$$

$$\tilde{v}[i] = \sum_{j=1}^l \frac{q[i, j] \cdot v[j]}{\sum_{j=1}^l q[i, j] \cdot v[j]} \in \mathbb{R}^{1 \times l} \cdot \mathbb{R}^{d_{out} \times 1}$$

$$\tilde{v}[i] \in \mathbb{R}^{1 \times l}$$

$$w_1 \cdot \tilde{v}[i] = \mathbb{R}^{d_{hidden} \times d_{out}} \cdot \mathbb{R}^{1 \times l} = \mathbb{R}^{d_{hidden} \times l}$$

$$w_2^T \in \mathbb{R}^{1 \times d_{hidden}}$$

$$\underline{\tilde{x}[i] \in \mathbb{R}^{1 \times l}}$$

$$W_q \in \mathbb{R}^{d_{\text{atten}} \times d_{\text{in}}}, W_k \in \mathbb{R}^{d_{\text{atten}} \times d_{\text{in}}} \\ W_v \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}, W_i \in \mathbb{R}^{d_{\text{hidden}} \times d_{\text{out}}} \\ W_2 \in \mathbb{R}^{d_{\text{hidden}} \times 1}, x \in \mathbb{R}^{d_{\text{ex}} \times l} = \mathbb{R}^{d_{\text{in}} \times l}$$

$$[d_{\text{in}} = d_{\text{e}}, l = d_{\text{out}}]$$

$$W_1 \in \mathbb{R}^{2 \times d_{\text{out}} = d_{\text{hidden}} \times d_{\text{out}}}$$

$$W_v \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}} = d_{\text{out}} \times d_{\text{out}}} \\ = \mathbb{R}^{l \times l}$$

$$d_{\text{e}} = l = d_{\text{out}} = d_{\text{in}}, d_{\text{hidden}} = 2$$

$$W_q = d_{\text{atten}} \times l, W_k = d_{\text{atten}} \times l$$

$$W_v = l \times l, W_1 = 2 \times l$$

$$W_2 = 2 \times 1$$

$$\rightarrow l^2 + 2l + 2l \cdot d_{\text{atten}} + 2$$

$$= l^2 + 2l \cdot (1 + d_{\text{atten}}) + 2$$

$$\tilde{y}[i] \\ = W_2^T \cdot \text{ReLU}(W_1 \cdot \tilde{v}[i])$$

$$W_1 \in \mathbb{R}^{2 \times d_{\text{out}}} = \mathbb{R}^{2 \times l}$$

$$\tilde{v}[i] \in \mathbb{R}^{l \times l} = \mathbb{R}^{d_{\text{out}} \times l}$$

$$W_1 \tilde{v}[i] \in \mathbb{R}^{2 \times l} \cdot \mathbb{R}^{l \times l}$$

$$= \boxed{4l^2}$$

$$W_2^T \in \mathbb{R}^{1 \times 2} \cdot (W_1 \tilde{v}[i] \in \mathbb{R}^{2 \times l})$$

$$= \boxed{4l}$$

$$\therefore 4l + 4l^2$$

$$= \underline{4l(1+l)}$$

feed-forward
block

<2>

3. activation $\tilde{v}[1] \tilde{v}[2] \dots \tilde{v}[l]$

$$\tilde{v}[i] \in \mathbb{R}^{l \times l} \quad i=1, 2, \dots, l$$

$$\downarrow$$

$$\underline{l \times l^2} = \text{size} \quad \langle 3 \rangle$$

2. ①

$$\alpha[i, j] = \exp(q[i]^T \cdot k[j] / \sqrt{d_{\text{atten}}})$$

② $\sum_{u=1}^l \exp(q[i]^T \cdot k[u] / \sqrt{d_{\text{atten}}})$

$$\Rightarrow \boxed{2d_{\text{atten}} \cdot (1+l) + 2l^2}$$

$$q[i]^T \in \mathbb{R}^{1 \times d_{\text{atten}}}, k[j] \in \mathbb{R}^{d_{\text{atten}} \times 1}$$

$$q[i]^T \cdot k[j] \in \mathbb{R}^{1 \times 1}$$

$$\sum_{u=1}^l \exp(q[i]^T \cdot k[u] / \sqrt{d_{\text{atten}}}) \in \mathbb{R}^{1 \times l}$$

③ $q[i]^T \cdot k[j] \in \mathbb{R}^{1 \times 1}$

$$\boxed{2d_{\text{atten}}}$$

④ $\mathbb{R}^{d_{\text{atten}}} \cdot \mathbb{R}^{d_{\text{atten}} \times l} = \boxed{2d_{\text{atten}} \cdot l}$

$$\tilde{v}[i] = \sum_{j=1}^l \alpha[i, j] \cdot v[j]$$

$$\mathbb{R}^{l \times l} \cdot \mathbb{R}^{l \times l} \in \mathbb{R}^{d_{\text{out}} \times l}$$

$$\mathbb{R}^{l \times l} \cdot \mathbb{R}^{d_{\text{out}} \times l} = \mathbb{R}^{l \times l} \cdot \mathbb{R}^{l \times l}$$

$$\Rightarrow \boxed{2l^2}$$

$$\downarrow$$

$$\text{Flop} = 2d_{\text{atten}} + 2d_{\text{atten}} \cdot l + 2l^2$$

6.

$$l(y - \hat{y}) = \sum_{i=1}^l (y[i] - \hat{y}[i])^2$$

$$\hat{y}[i] \in \mathbb{R}^{1 \times l}$$

$$y[i]^2 - 2y[i] \hat{y}[i] + \hat{y}[i]^2$$

$\hat{y}[i]$ 의 FLOPS

$$= 4l^2$$

$$w_2^T \in \mathbb{R}^{1 \times 2} \cdot w_1 \in \mathbb{R}^{2 \times d_{\text{out}}} \cdot v[i] \in \mathbb{R}^{d_{\text{out}} \times 1}$$

$$= w_2^T \in \mathbb{R}^{1 \times 2} \cdot \mathbb{R}^{2 \times l}$$

$$\boxed{4l^2 + 4l} \times 3$$

$$\Rightarrow \boxed{12l(l+1)} : y^2[i], -2y[i] \cdot \hat{y}[i], \hat{y}[i]^2 \text{ FLOPs 동일}$$

4. back ward propagation

$$w_2^T \in \mathbb{R}^{1 \times d_{\text{hidden}}} = \mathbb{R}^{1 \times 2}$$

$$w_1 \in \mathbb{R}^{d_{\text{hidden}} \times d_{\text{out}}}$$

$$= \mathbb{R}^{2 \times l}$$

<4>

$$w_2^T w_1 \in \mathbb{R}^{1 \times l} \quad [4l]$$

$$w_2^T w_1 \cdot \tilde{v}[i] = \mathbb{R}^{1 \times l} \cdot \mathbb{R}^{l \times l}$$

$$= \mathbb{R}^{1 \times l} \quad [2l^2]$$

$$\therefore 4l + 2l^2 = 2l(2+l)$$

5. $w_q \in \mathbb{R}^{d_{\text{att}} \times d_{\text{in}}} = \mathbb{R}^{d_{\text{att}} \times l}$

$$= \begin{bmatrix} 0 & \dots & 0 & 1 & 0 & \dots & 0 \\ 0 & 2M^2 & 0 & \dots & 0 & \dots & 0 \\ \vdots & & & & & & \\ 0 & \dots & \dots & 0 & 2M^2 & 0 & \dots & 0 \end{bmatrix}$$

$$w_q(t, t) = 2M^2, \quad t = 2, \dots, d_{\text{att}}$$

$$w_k \in \mathbb{R}^{d_{\text{att}} \times l}$$

$$= \begin{bmatrix} 0 & \dots & 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 & \dots & \vdots \\ 0 & 0 & 1 & \dots & 0 & \vdots \\ \vdots & & & & & \vdots \\ 0 & \dots & \dots & 0 & 1 & 0 & \dots & 0 \end{bmatrix}$$

$$w_k(t, t) = 1, \quad t = 2, 3, \dots, d_{\text{att}}$$

$$w_{V, (i, j)} \begin{cases} 1, & i = n+4, j = l. \\ \in \mathbb{R}^{1 \times l}, & 0, \text{ else where.} \end{cases}$$