

EE7403 2023-2024 S2

$$4. (a) \quad y_{i,j,k} = \sum_{n=1}^C \sum_{m=1}^Q \sum_{l=1}^P W_{l,m,n,k} \cdot x_{l,m,n} + b_k.$$

for each channel, there are $(PQ+1)$ learnable parameters
(including C $W_{k,n}$ and 1 b_k)
and $P \times Q$ pixels
there's D output channels, hence:

$$\text{Total number} : PQD(PQC+1) = PQDPQC + PQD$$

$$(b) \quad y_{i,j,k} = \sum_{n=1}^C \sum_{m=1}^1 \sum_{l=1}^1 W_{l,m,n,k} \cdot x_{i-l,j-m,n} + b_k$$

$$\text{each channel } (3 \times 3 \times C + 1) = 9C + 1$$

$$\text{Total} : D(9C + 1) = 9CD + D$$

$$(c) \quad y_{i,j,k} = \sum_{n=1}^C W_{n,k} x_{i,j,n} + b_k$$

$$\text{each channel } (1 \times 1 \times C + 1) = C + 1$$

$$\text{Total} : D(C + 1) = CD + D$$

$$(d) \quad y_{ik} = \sum_{n=1}^C W_{n,k} \cdot x_{i,n}$$

$$\text{each channel} : C$$

$$\text{Total} : DC.$$

$$(e) \quad \text{let } W = [W_1, W_2, \dots, W_C]^T \in \mathbb{R}^{C \times D}$$

$$Y = XW$$

In Transformer, all tokens share the same learnable parameter, the learned parameters go through all pixels, which is quite similar to CNN.