

- Recap: Attention and Convolution
- Transformer

Recap:  $J(i,j) = \sum_c \sum_m \sum_n I(i-m, j-n, c) \cdot h(m, n, c)$  (Convolution)

$J(i,j) = \sum_{a,b \in \mathcal{N}(i,j)} \text{softmax}_{ab}(\langle q(i,j), k(a,b) \rangle) v(a,b)$  (Attention)

Self Attention

$$\begin{cases} q(i,j) = W_q x(i,j) \\ k(i,j) = W_k x(i,j) \\ v(i,j) = W_v x(i,j) \end{cases}$$

$$x(i,j) \in \mathbb{R}^{d \times c}$$

$$W_q \in \mathbb{R}^{d \times c}$$

$$W_k \in \mathbb{R}^{d \times c}$$

$$W_v \in \mathbb{R}^{d \times c}$$

Note: If  $W_q = W_k = W$ ,

$$\langle W x(i,j), W x(a,b) \rangle$$

$$= (W x(i,j))^T \cdot W x(a,b),$$

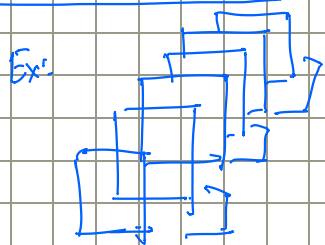
o we would have a high value if the inner product at  $(a,b) = (i,j)$

This in turn would mean that  $v(i,j)$  is always going to be assigned a high weight.

o This is not ideal.

$\therefore$  we have  $W_q$  and  $W_k$  that are not the same.

o Multi Head Attention (MHA): In MHA,  $c' = c/h$  where  $h$  is the number of heads



$$c = 6, h = 3 \Rightarrow c' = 6/3 = 2$$

o We now work with  $c'$  channels per head.

o Suppose we have  $d'$  to be the output dimension,  $\tilde{d} = h \cdot d'$

ie. the outputs of each head are concatenated. If  $\tilde{d} \neq d$

a projection matrix  $W_o \in \mathbb{R}^{d \times \tilde{d}}$  is applied

# The Transformer Architecture

