

12/10/23

Deep Learning

## o Introduction to Attention

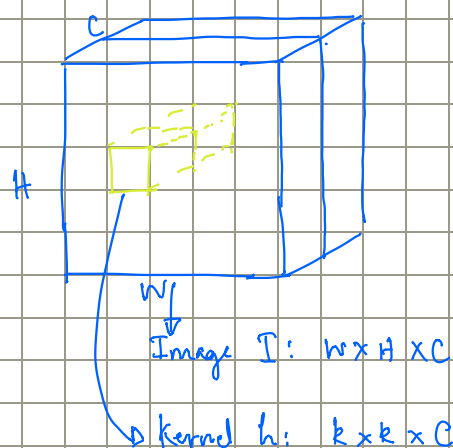
Note: References posted on classroom

- Recall Convolution
- Attention mechanism

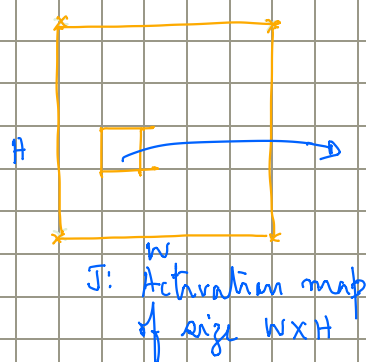
Ramachandran et al. NeurIPS 2019

## o Transformer model - Vaswani et al. NeurIPS 2017

## o Recall Convolution/Correlation



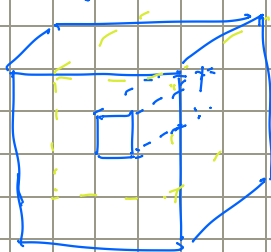
W: width  
H: height  
C: channels  
k: kernel width/height



Convolution Sum:

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

## o Attention (Self Attention):



Note: we are dealing with vectors in the definitions

In the <sup>self</sup> attention mechanism, we derive three elements - query, key and value from the same input

query:  $q(i, j) = W_q x(i, j)$

key:  $k(i, j) = W_k x(i, j)$

value:  $v(i, j) = W_v x(i, j)$

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

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

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

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

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

