



Geometry features will be incorporated into the attention weight matrix

$$d_{geom} = 4$$

#### Notations:

- $R$  = number of detected regions per image
  - nowhere to found (not in the paper or in the code)
- $d_{model} = 512$ , model inner dimension
- $h$  = number of heads in self-attention

One Head

1.  $X$  contains all the input vectors (e.g. region features). Let  $R = 50$ . Dimensions for all computations below:  $X: R \times 512$ ,  $W_Q, W_K, W_V: 512 \times 64$

$$Q = XW_Q, K = XW_K, V = XW_V$$

2. **Appearance attention weights** are computed as follows:

$$\Omega_A = \frac{QK^T}{\sqrt{d_k}}$$

$50 \times 50$

3. The next step is to compute **geometric attention weights**.

A displacement vector between all variations of bounding boxes  $(m, n)$  among  $R$  is computed:

$$\lambda(m, n) = \left( \log \left( \frac{|x_m - x_n|}{w_m} \right), \log \left( \frac{|y_m - y_n|}{h_m} \right), \log \left( \frac{w_n}{w_m} \right), \log \left( \frac{h_n}{h_m} \right) \right)^*$$

Afterwards, this displacement vector is passed to compute weights:

$$\omega_G^{mn} = \text{ReLU}(\text{Emb}(\lambda)W_G)**$$

4. **Combined attention weights** are computed as the next step:

$$\omega^{mn} = \frac{\omega_G^{mn} \exp(\omega_A^{mn})}{\sum_{l=1}^N \omega_G^{ml} \exp(\omega_A^{ml})}***$$

5. The output of the head is computed as the following:

$$\text{head}(X) = \text{self-attention}(Q, K, V) = \Omega V$$

Omega size:  $R \times R$

where every element in Omega is provided by combined attention weights

\* For more information, check the slides.

\*\* Check BoxRelationalEmbedding function in model implementation

\*\*\* Check box\_attention function in model implementation