

Self-Attention (Transformer)

I PITY THE FOOL



USING RECURRENT LAYERS

Self-Attention





Convolution vs Self-Attention



$$C = [W_v X] \left([W_k X]^T [W_q X] \right)_{\text{ssmax}}$$
$$c_i = \sum_j \langle q_i, k_j \rangle v_j$$

Convolution:

$$c_i = \sum_j \frac{h_{i-j} x_j}{j}$$

$$W_{v} = W_{k} = W_{q} = I$$
$$C = [X] ([X]^{T} [X])_{ssmax}$$
$$c_{i} = \sum_{j} \langle x_{i}, x_{j} \rangle x_{j}$$

With ssmax:
$$c_i = \sum_j \frac{\exp(\langle x_i, x_j \rangle)}{\sum_k \exp(\langle x_i, x_k \rangle)} x_j$$

Convolution vs Self-Attention



Convolution vs Self-Attention



Local and fixed during inference





Local and fixed during inference



Local and fixed during inference



Local and fixed during inference

Global and variable during inference



Attention in Image Processing

- Image Classification & Detection:
 - CBAM (Convolutional Block Attention Module)
 - Dual Attention (Spatial and Channel)
 - ViT (Vision Transformer)
 - CoAtNet (Convolution with Self-Attention)

Convolutional Block Attention Module





Dual Attention





Fu et al., DANet, 2019

Dual Attention



A. Position attention module



B. Channel attention module

The Transformer

- Multi-Headed Self- and Cross-Attention
- Masked Multi-Headed Self-Attention
- Layer Normalization
- Linear + ReLU
- Positional Encoding
- Residual Connection +
- Dropouts



Vision Transformer (ViT)







CoAtNet

• Marry convolution and attention

$$c_{i} = \sum_{j} \frac{\exp(\langle x_{i}, x_{j} \rangle)}{\sum_{k} \exp(\langle x_{i}, x_{k} \rangle)} x_{j}$$
$$c_{i} = \sum_{j} \frac{\exp(\langle x_{i}, x_{j} \rangle + w_{i-j})}{\sum_{k} \exp(\langle x_{i}, x_{k} \rangle + w_{i-k})} x_{j}$$



Bianco et al., IEEE Access, 2018





https://bit.ly/zoi-sem01