

RNN
layer

$$y_t^T = f(x_t^T H + v^T + y_{t-1}^T G)$$

$$= \begin{bmatrix} 1 & \dots & 1 \end{bmatrix}^T + \begin{bmatrix} 1 & \dots & 1 \end{bmatrix}^T$$

mixes within current input feature vector

mixes within previous output feature vector

if G was identity matrix then current output is previous output plus a perturbation

NN
(single input)
layer

$$y^T = f(x^T H + v^T)$$

$$= \begin{bmatrix} 1 & \dots & 1 \end{bmatrix}^T + \begin{bmatrix} 1 & \dots & 1 \end{bmatrix}^T$$

mixes within input feature vector

NN
(batch input)
layer

$$Y^T = f(X^T H + 1 v^T)$$

$$\begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \begin{bmatrix} 1 & \dots & 1 \end{bmatrix}^T + \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}^T$$

mixes within input feature vectors

does not mix across input feature vectors

Self attention
(single headed)
layer

$$Y^T = A^T X^T H$$

$$\begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \begin{bmatrix} 1 & \dots & 1 \end{bmatrix}^T$$

mixes across input feature vectors

mixes within input feature vectors

$$A^T = \text{softmax}_{\text{row}} \left(\frac{X^T W_k W_v^T X}{\text{scalar}} \right)$$

self attention is function of the weighted product $X^T W X$

each row in A^T is a pmf

Encoder-attention-decoder
(linear weights)
this layer

$$y^T = a^T x^T$$

$$= \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}^T \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

used by the decoder

mixes across input feature vectors

$$a^T = \text{softmax} \left(\frac{s^T W^T X}{\text{scalar}} \right)$$

s^T comes from the decoder

many variations on this theme

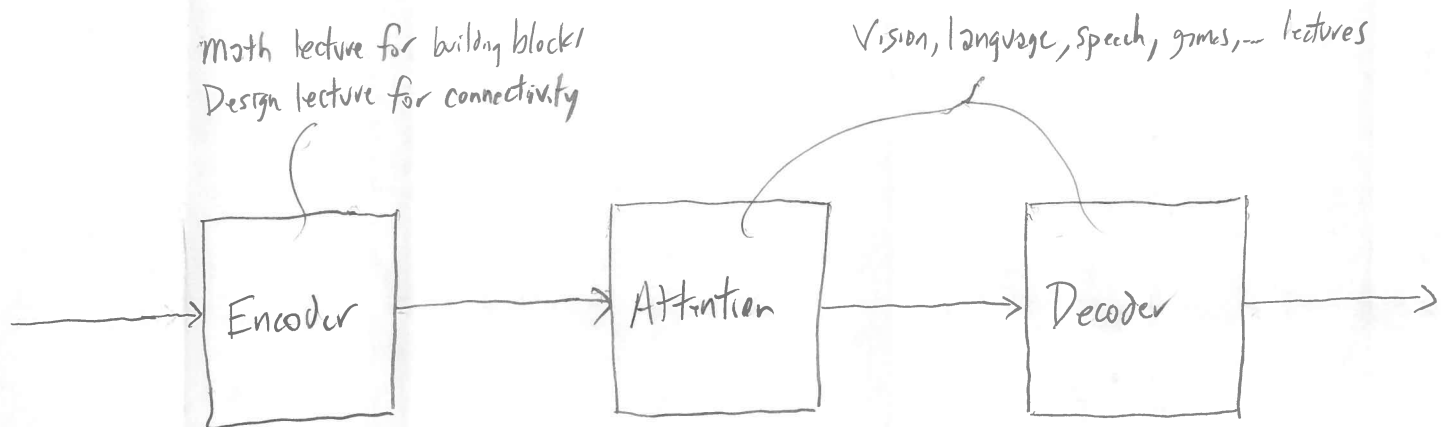
CNN
(not batched)
layer

$$Y^{3D} = f(H^{4D} \otimes X^{3D} \oplus v^{1D})$$

$$N_o \left\{ \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \right\} = \left\{ \begin{bmatrix} N_i \{ \vdots \} \\ N_i \{ \vdots \} \\ \vdots \\ N_i \{ \vdots \} \end{bmatrix} \otimes \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \right\} + \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \left\{ \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \right\} N_o$$

$H^{4D} \otimes X^{3D}$ will be implemented via lowering to matrix matrix multiplication

Encoder - attention - decoder architecture



Examples:

Input image → ResNet50 with a FPN → Features → RPN followed by RoI extraction → RoIs → Classifier and bounding box regressor → Object detection

English words → Transformer → Features → Linear weighted attention → Attended features → Transformer with masking → French words