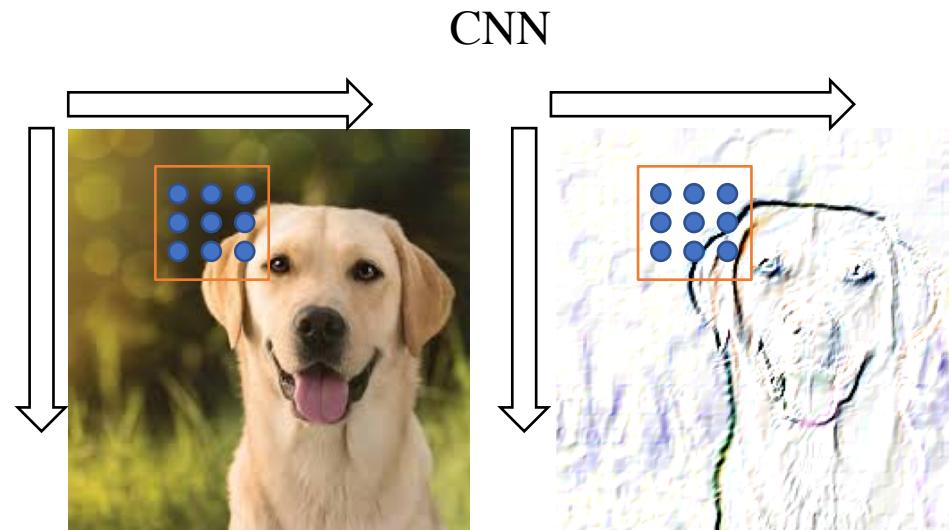


Attention in deep learning



- Color, surface, posture
 - Edge of the object against the background
 - Edge of the eyes, noise, mouth...
- a dog

RNN

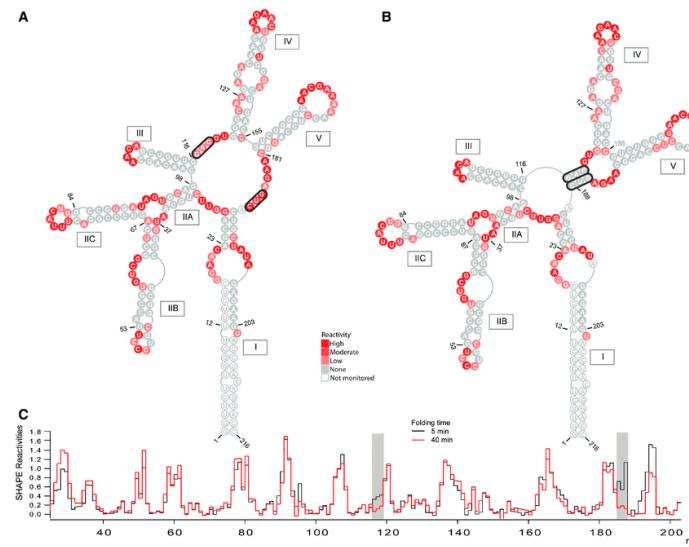
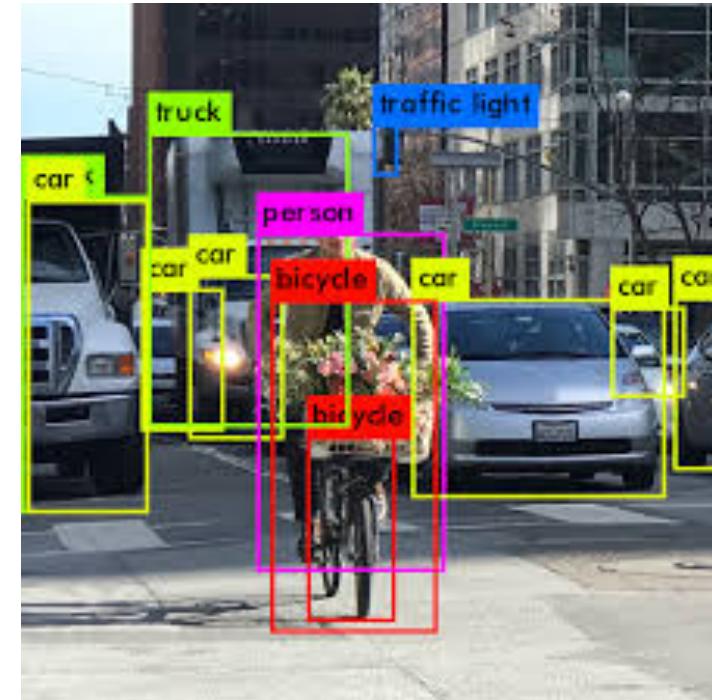


- Goes up at this state with this input
 - Goes down at that state with that input
- correlations

But these are the easy ones

Why we need attention

- Many objects:
- Convolved correlations



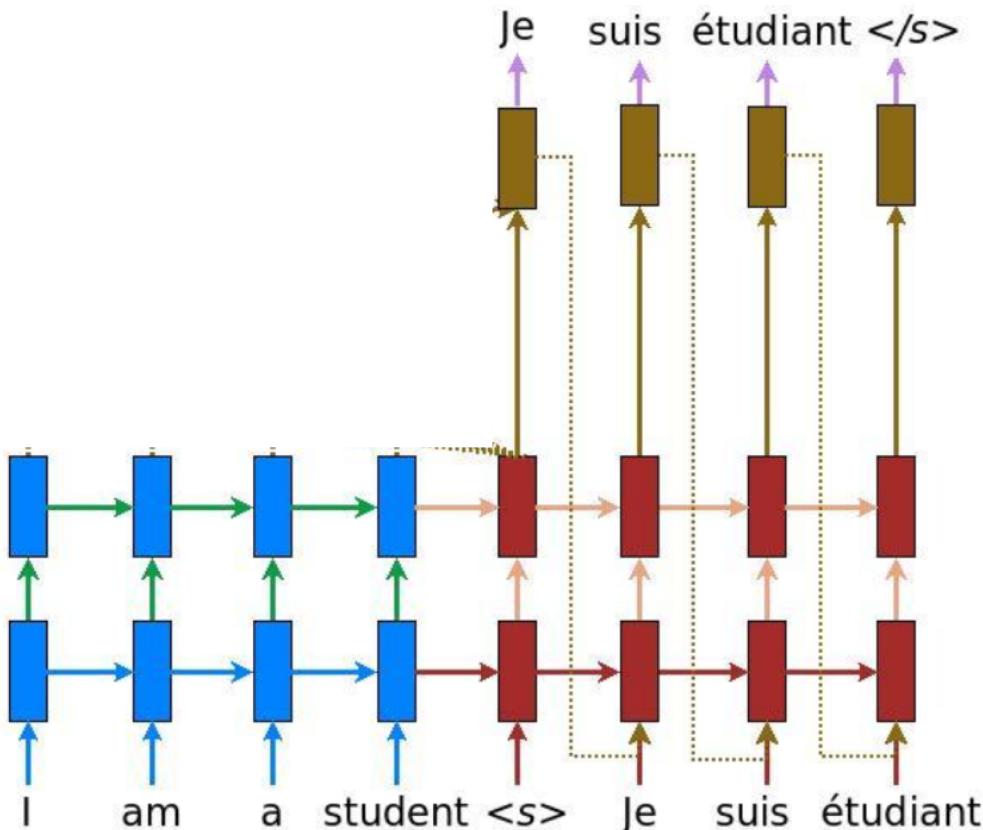
https://www.researchgate.net/figure/RRE2-structures-obtained-by-SHAPE-as-a-function-of-RNA-folding-time-RRE-2-RNA-assumes-fig1_236608422

Summarize spatio-temporal information

Pooling/down-sampling/aggregation etc.

- Unparameterized:
 - `reduce_max`, `reduce_mean`, `reduce_sum`
 - Local pooling with sliding window
 - Slicing and scaling
 - etc.
- Weighted:
 - **Attention**
 - `set2set`
 - etc. and etc.

Vanilla attention



$$\alpha_{ts} = \frac{\exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s))}{\sum_{s'=1}^S \exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_{s'}))} \quad (1)$$

$$\mathbf{c}_t = \sum_s \alpha_{ts} \bar{\mathbf{h}}_s \quad (2)$$

$$\mathbf{a}_t = f(\mathbf{c}_t, \mathbf{h}_t) = \tanh(\mathbf{W}_c[\mathbf{c}_t; \mathbf{h}_t]) \quad (3)$$

[Attention weights] (1)

[Context vector] (2)

[Attention vector] (3)

Tensorflow seq2seq, Neural Machine Translation (Luong et al., 2015, Bahdanau et al., 2015)

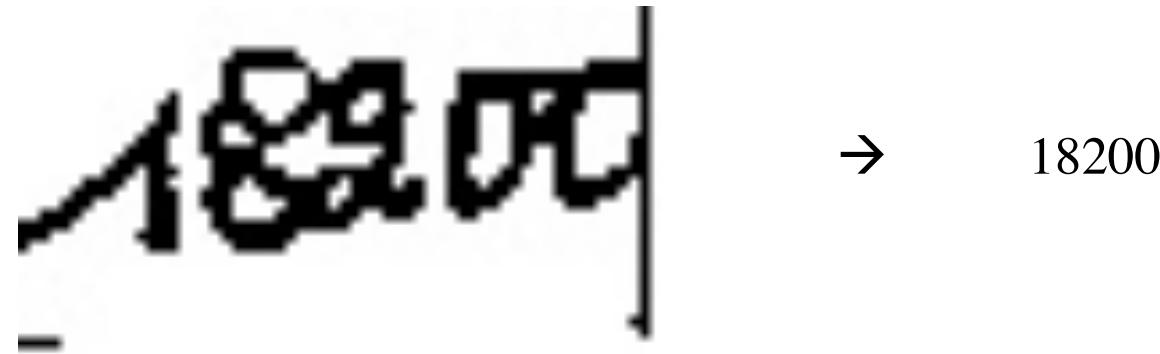
The gist of attention

$$\alpha_{ts} \propto \text{score}(h_t, h_s)$$

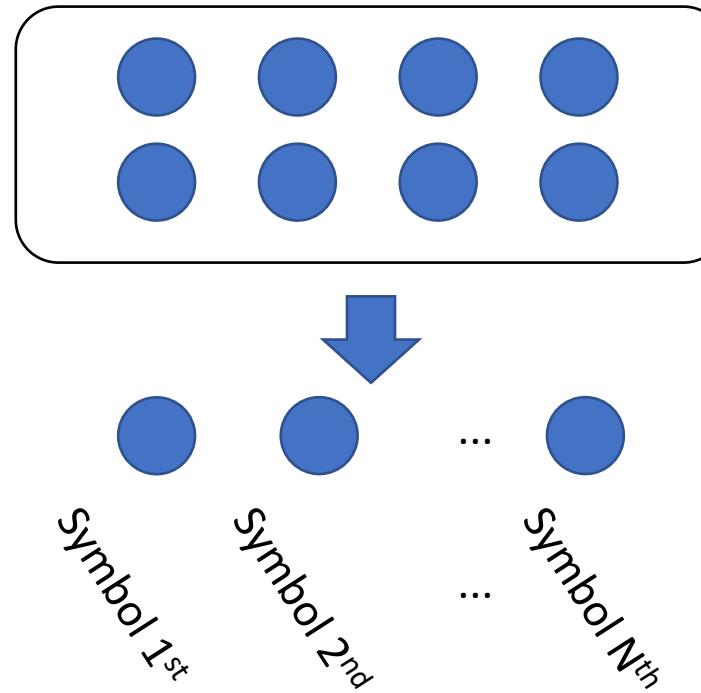
$$\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s) = \begin{cases} \mathbf{h}_t^\top \mathbf{W} \bar{\mathbf{h}}_s & [\text{Luong's multiplicative style}] \\ \mathbf{v}_a^\top \tanh(\mathbf{W}_1 \mathbf{h}_t + \mathbf{W}_2 \bar{\mathbf{h}}_s) & [\text{Bahdanau's additive style}] \end{cases}$$

Attention in image to sequence translation

- OCR2Text



- Q. What's the problem of a naive mapping from the source to the target?



- A. It may not align well.
- Therefore use attention.
- Demo

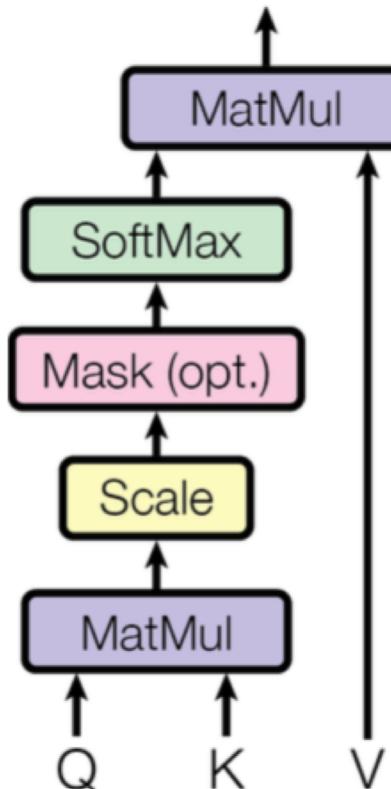
Self-attention: coping with the loss of reference

- Generalized attention function:
- Q — querys (target hidden state h_t)
 - [batch_size, d_k]
- K — keys (source hidden states $h_s, s = 1..N$)
 - [batch_size, N, d_k]
- V — memory (source hidden states $h_s, s = 1..N$)
 - [batch_size, N, d_v]
- Note: multiplicative style
- Scaling factor accounts for gradient saturation

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

A. Vaswani et al. Attention is all you need. NIPS. 2017

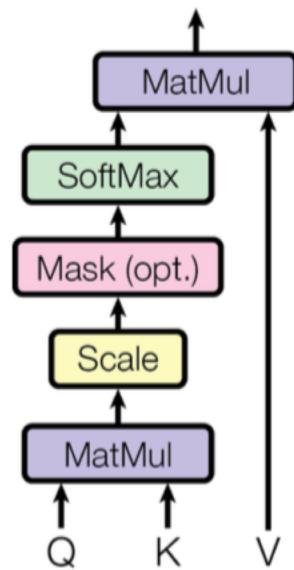
Scaled Dot-Product Attention



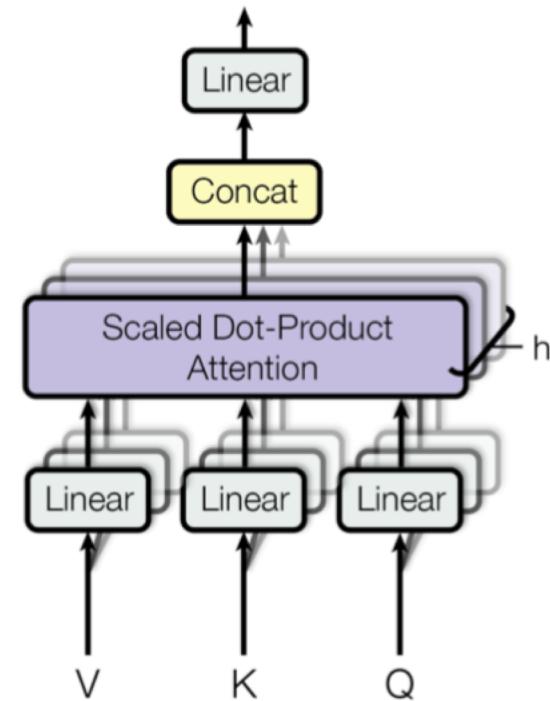
- Self-attention is when Q, K, V are the same → comparing to itself

Mutli-head attention

Scaled Dot-Product Attention

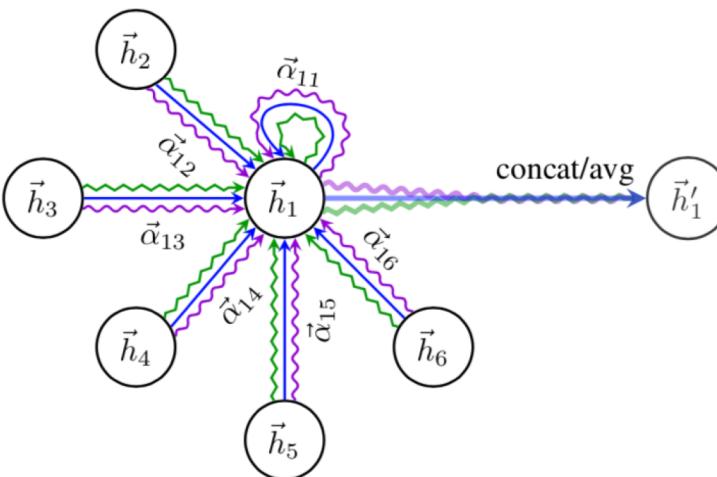
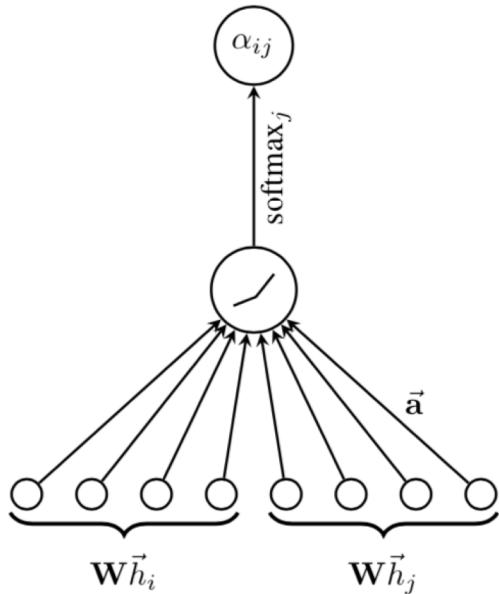


Multi-Head Attention



- * have h parallel&independet self-attention modules from the last slide
- * then concatenate/average their features

Attention is pervasive



Q: What's the style of this attention?

A: additive

$$\alpha_{ij} = \frac{\exp \left(\text{LeakyReLU} \left(\vec{a}^T [\vec{W}\vec{h}_i \| \vec{W}\vec{h}_j] \right) \right)}{\sum_{k \in \mathcal{N}_i} \exp \left(\text{LeakyReLU} \left(\vec{a}^T [\vec{W}\vec{h}_i \| \vec{W}\vec{h}_k] \right) \right)}$$

where \cdot^T represents transposition and $\|$ is the concatenation operation.