

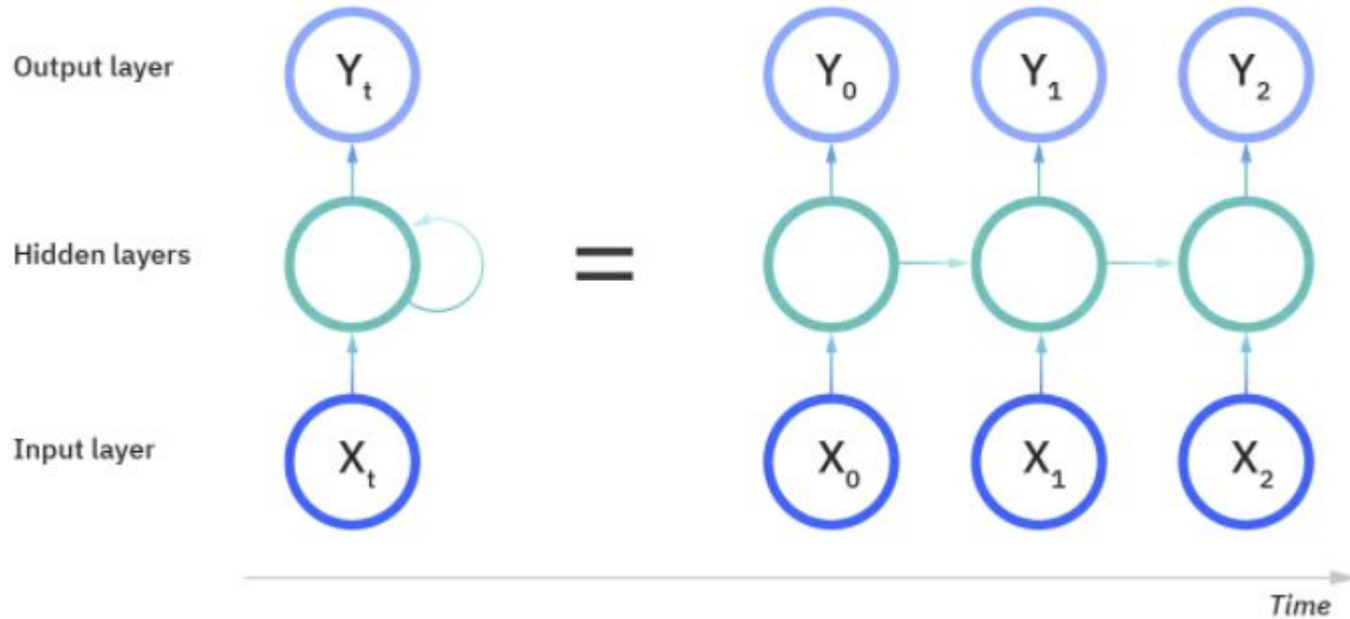
Attention and transformers



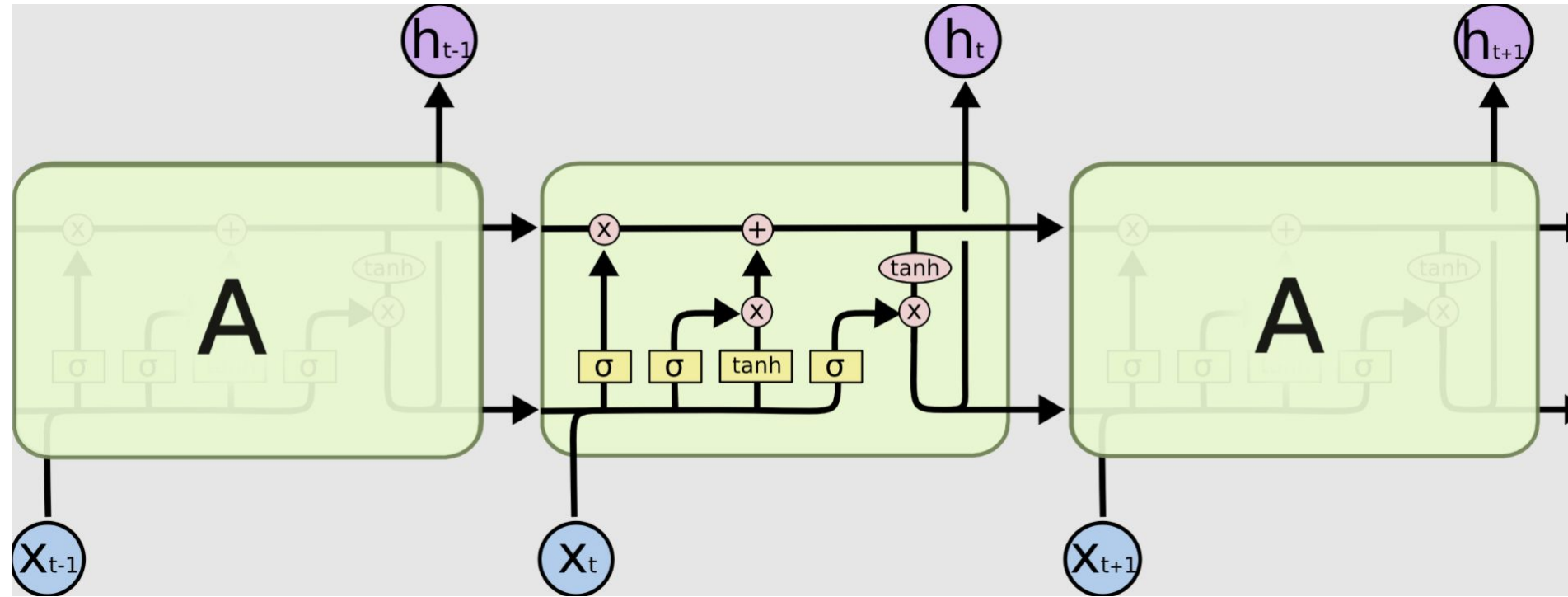
Сьогодні на лекції

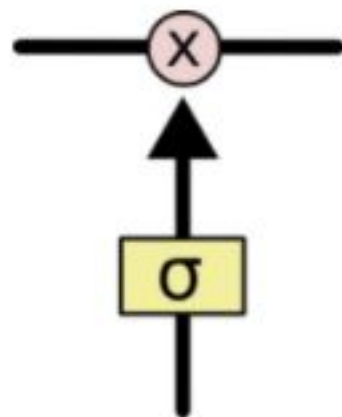
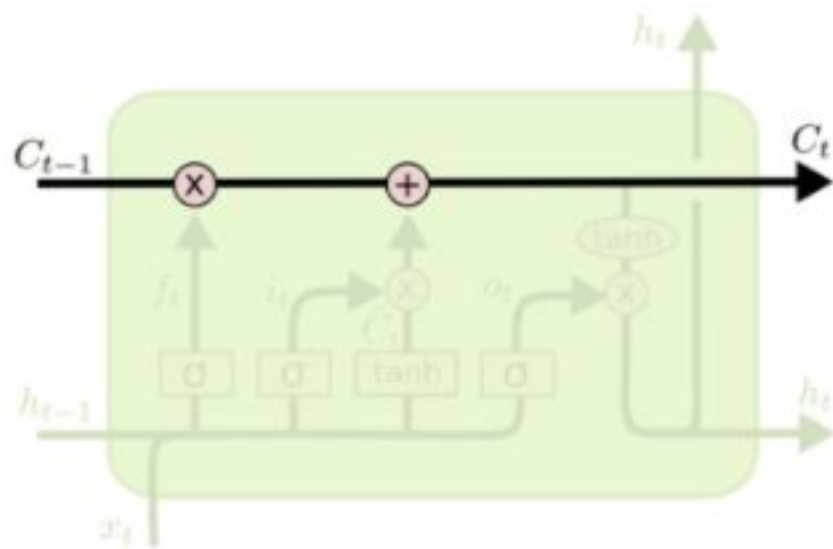


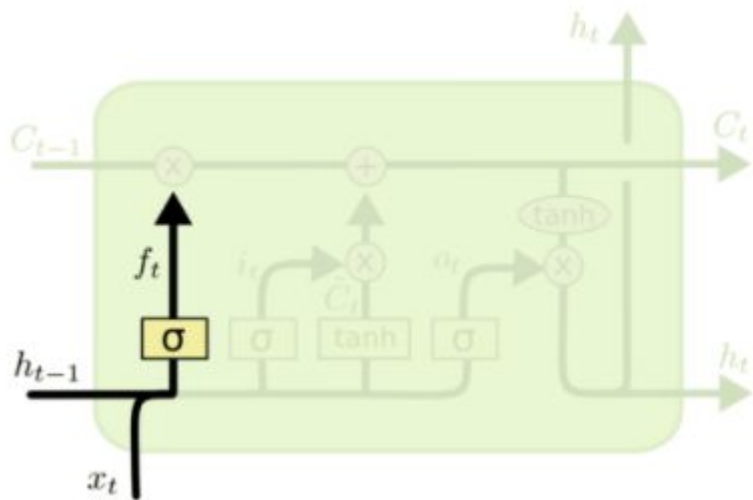
Recurrent Neural Networks



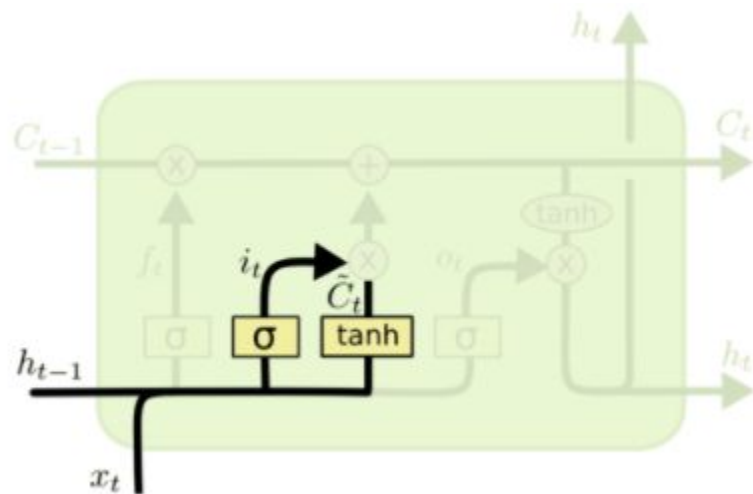
Long short-term memory (LSTM)





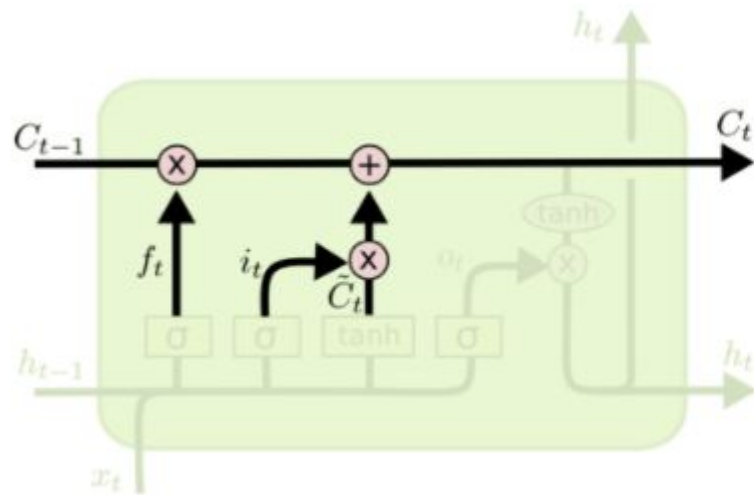


$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

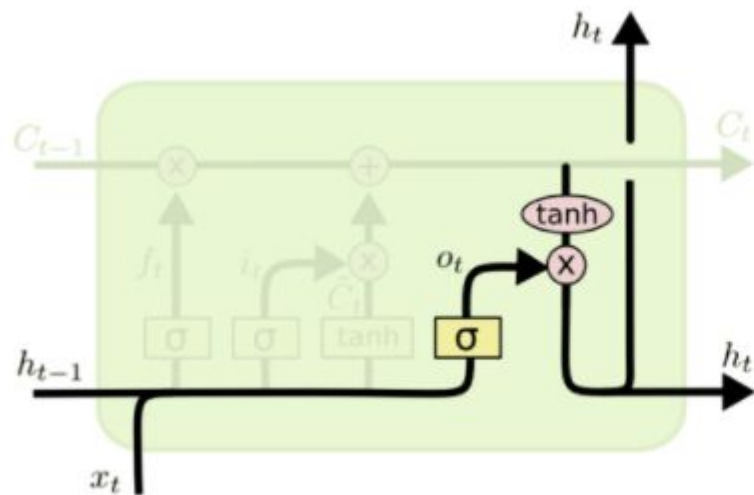


$$i_t = \sigma (W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$



$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$



В чому перевага LSTM?

Може запам'ятати попередні стани

Може забути попередні стани

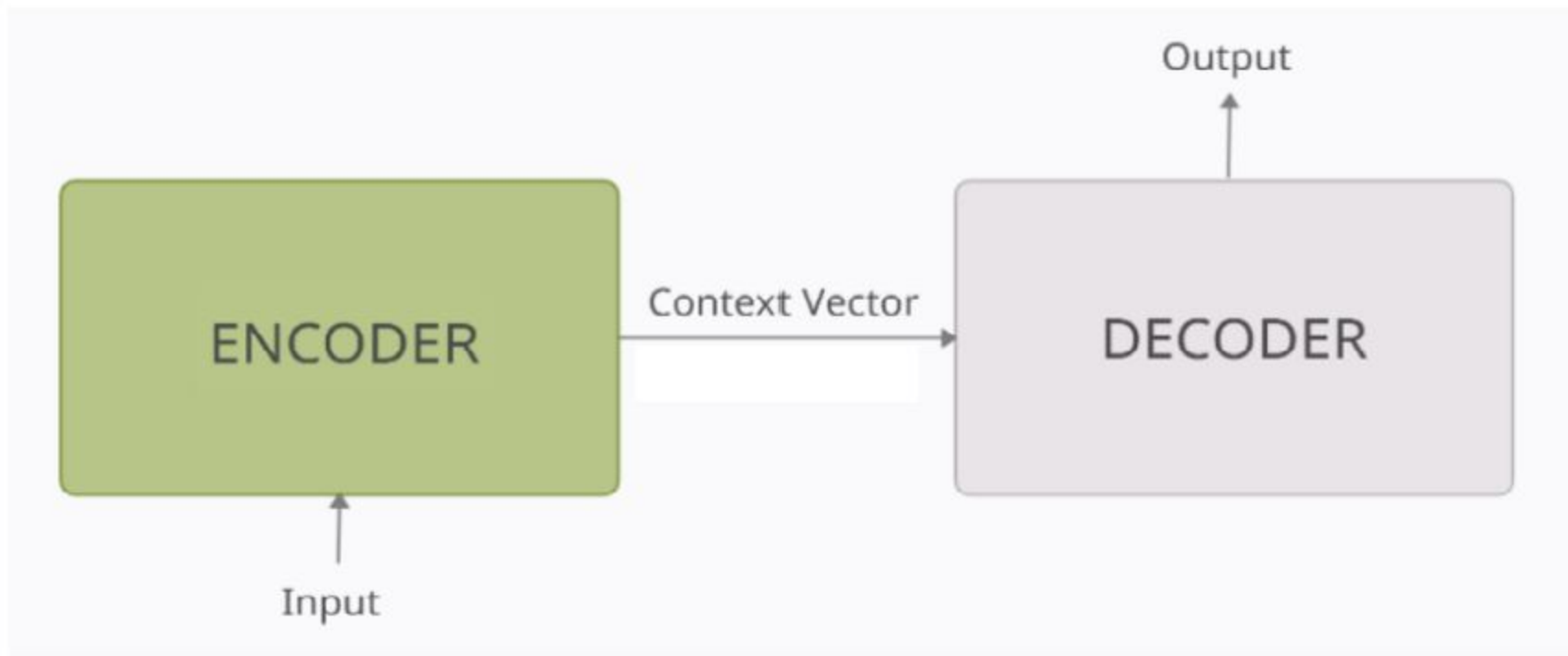
Може визначити наскільки треба
запам'ятати попередні стани

Визначає на скільки
запам'ятати/забути попередні стани

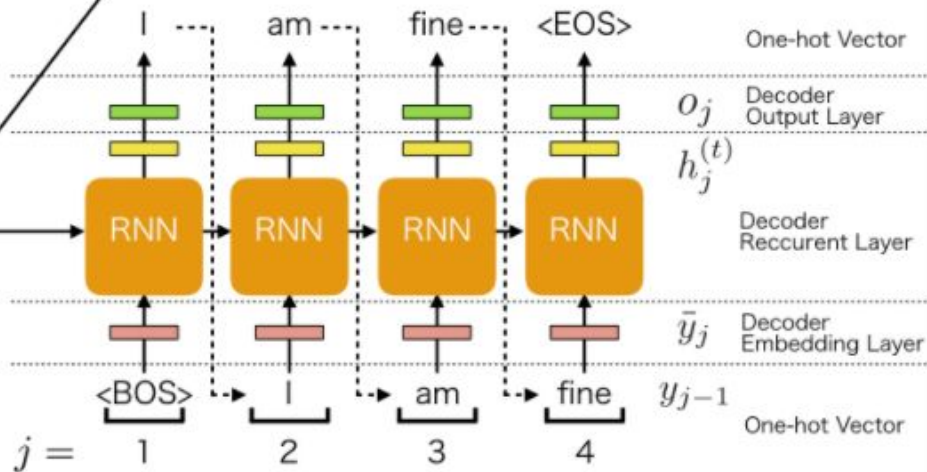
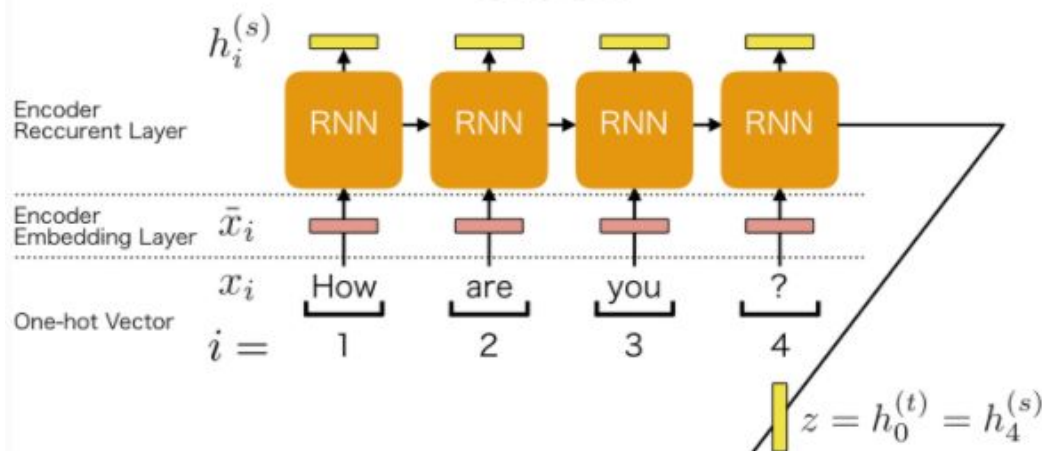
Все разом



Seq2Seq models



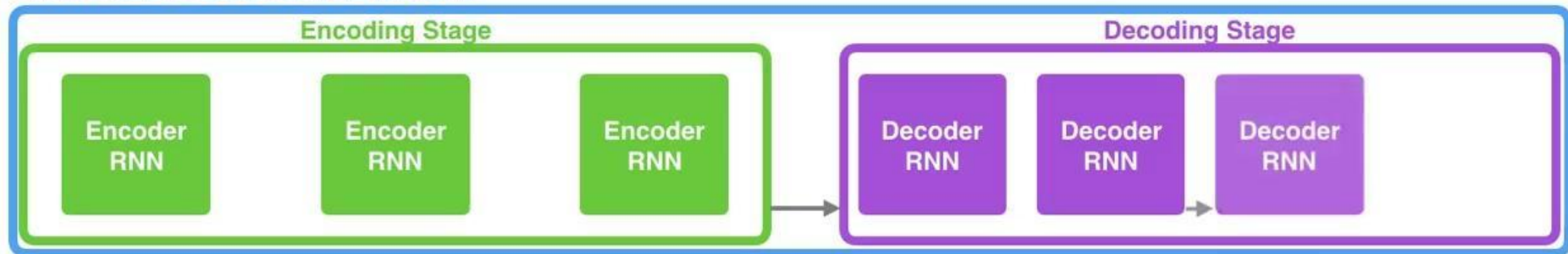
Encoder



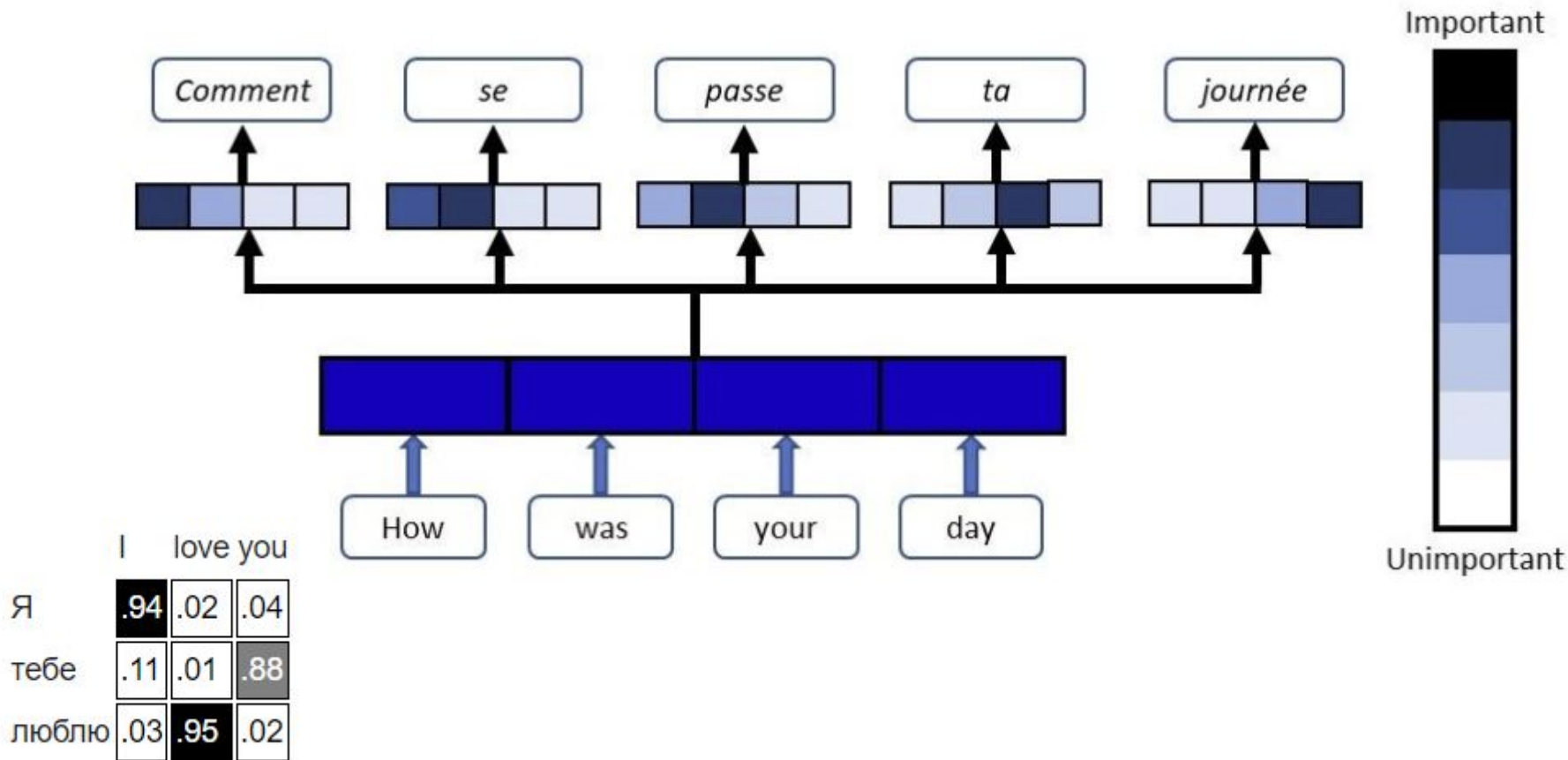
Decoder

Neural Machine Translation

SEQUENCE TO SEQUENCE MODEL



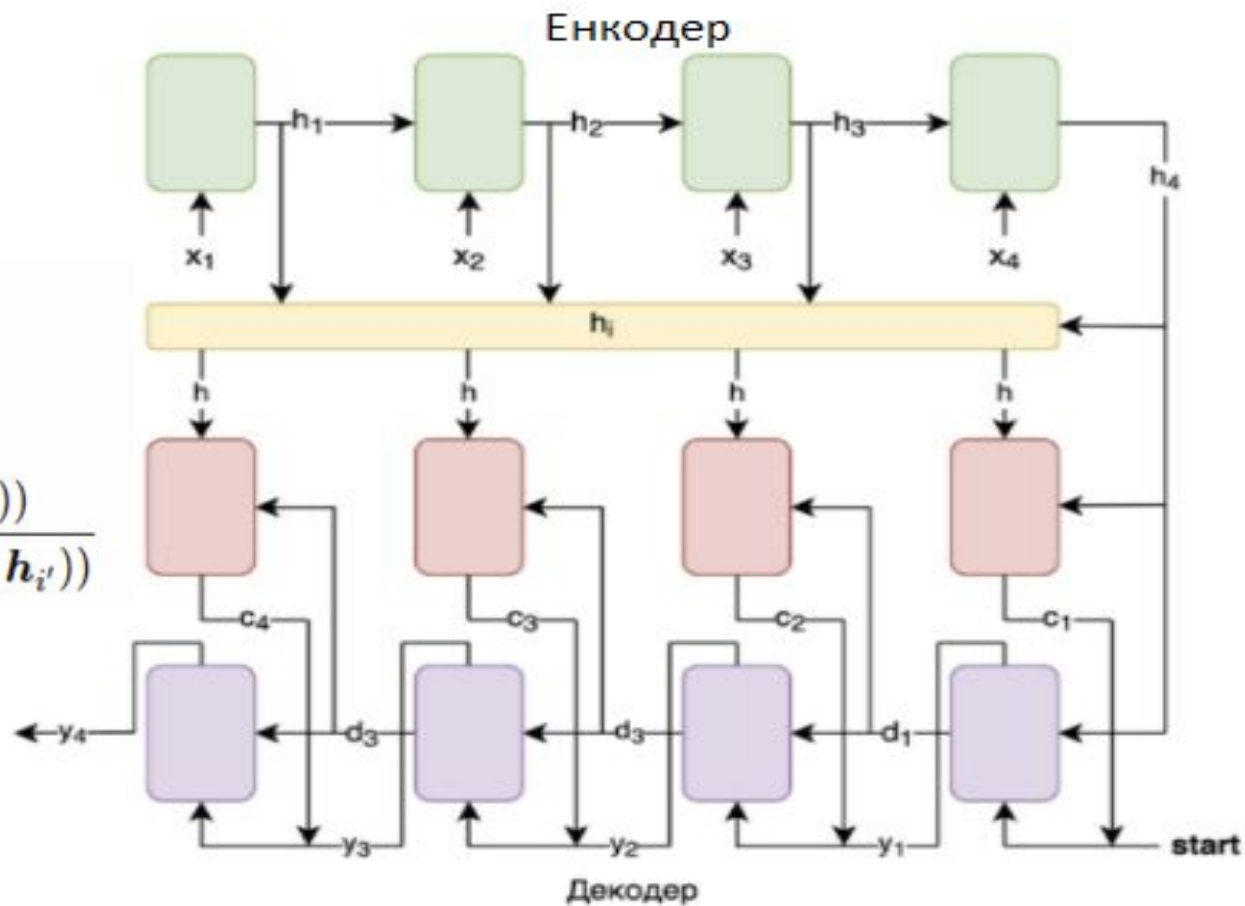
Attention mechanism



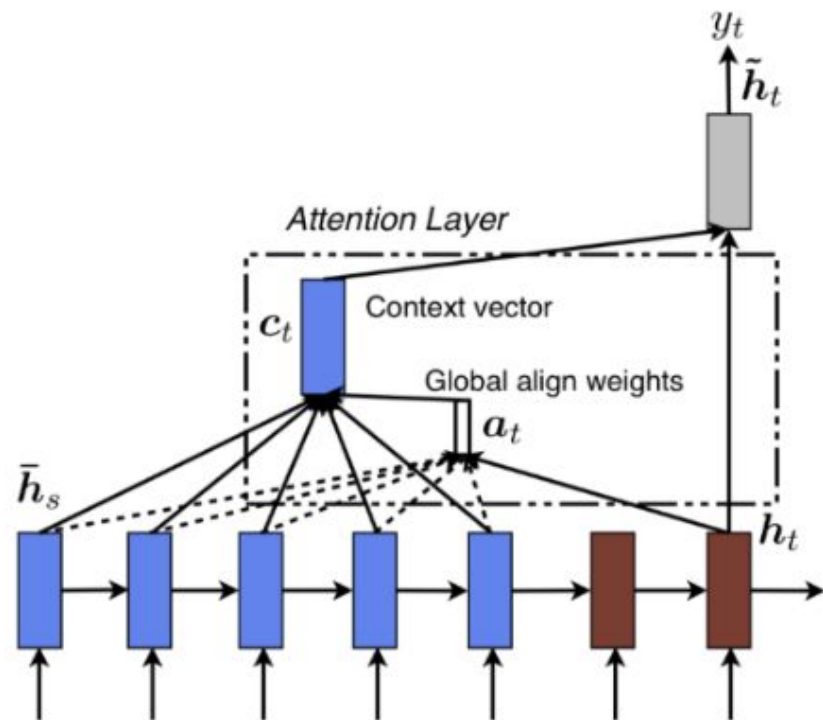
$$c_t = \sum_{i=1}^n \alpha_{t,i} h_i$$

$$\alpha_{t,i} = \text{align}(y_t, x_i)$$

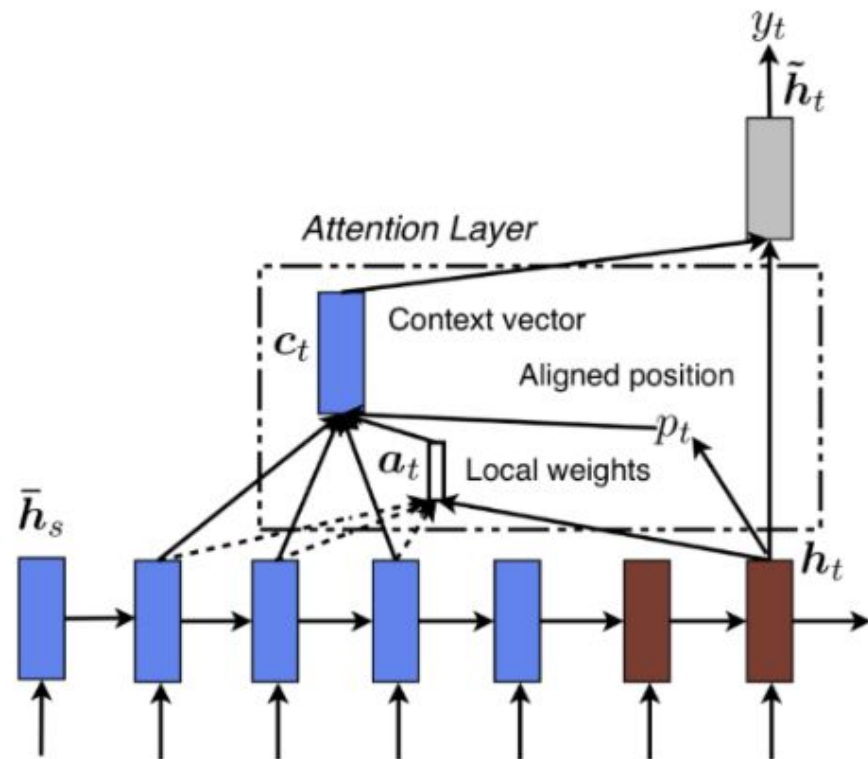
$$= \frac{\exp(\text{score}(s_{t-1}, h_i))}{\sum_{i'=1}^n \exp(\text{score}(s_{t-1}, h_{i'}))}$$



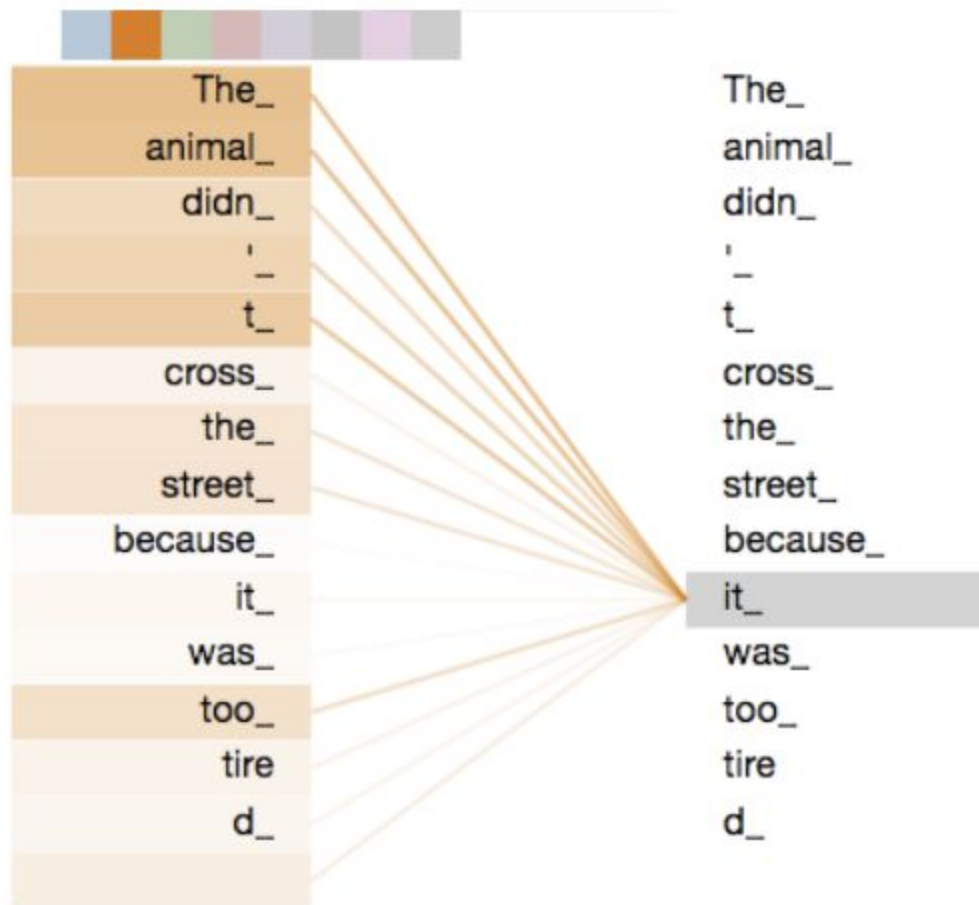
Content-base attention	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \text{cosine}[\mathbf{s}_t, \mathbf{h}_i]$
Additive(*)	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \mathbf{v}_a^\top \tanh(\mathbf{W}_a[\mathbf{s}_t; \mathbf{h}_i])$
Location-Base	$\alpha_{t,i} = \text{softmax}(\mathbf{W}_a \mathbf{s}_t)$ Note: This simplifies the softmax alignment to only depend on the target position.
Self- Attention(&)	Relating different positions of the same input sequence. Theoretically the self-attention can adopt any score functions above, but just replace the target sequence with the same input sequence.
Global/Soft	Attending to the entire input state space.
Local/Hard	Attending to the part of input state space; i.e. a patch of the input image.



Global Attention Model



Local Attention Model



Що дає механізм уваги?



Powered by  **Poll Everywhere**

Start the presentation to see live content. For screen share software, share the entire screen. Get help at pollev.com/app

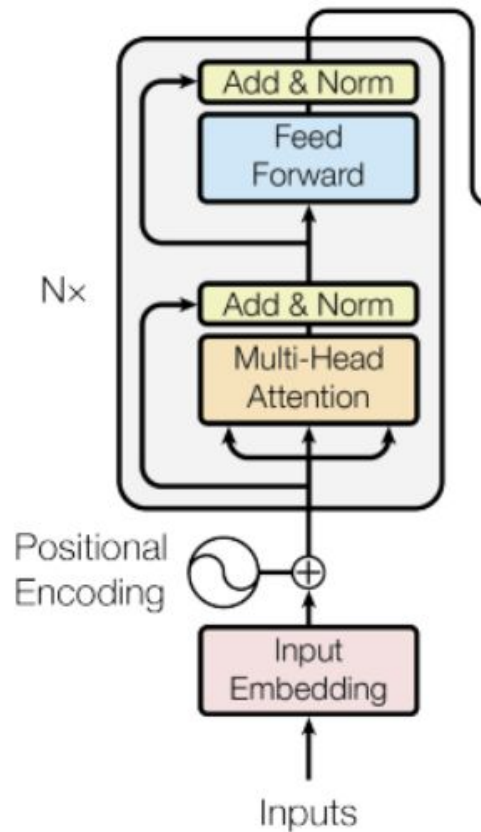
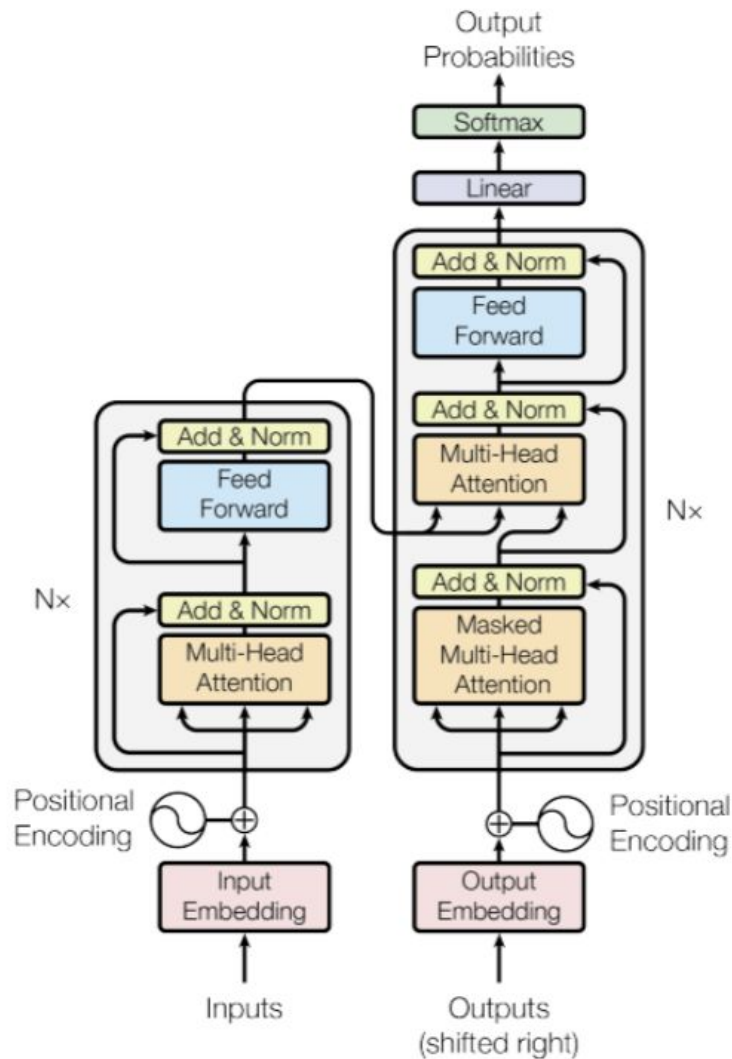
Transformers – Attention is All You Need

$$c = \sum_j \alpha_j h_j$$

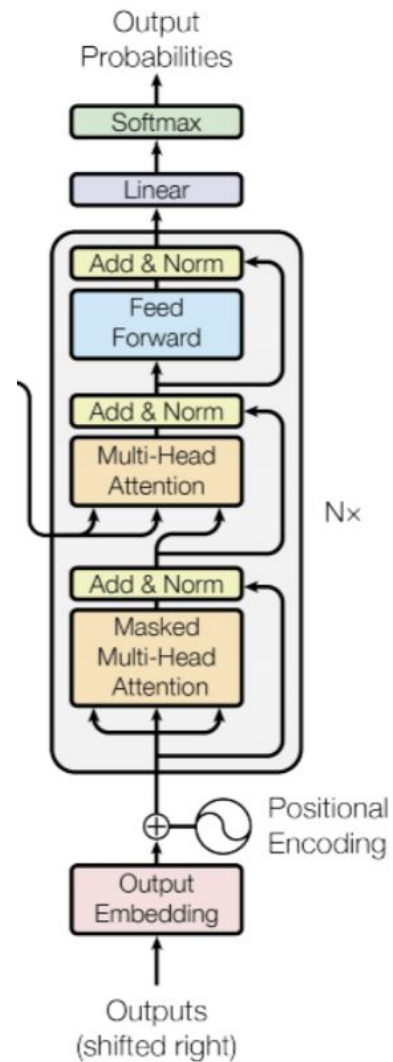
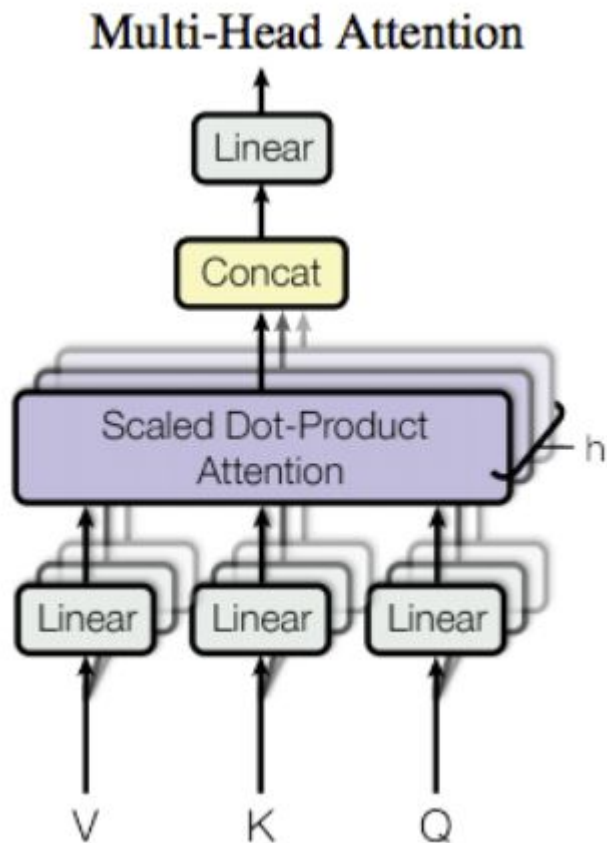
$$e_{ij} = a(s_i, h_j), \quad \alpha_{i,j} = \frac{\exp(e_{ij})}{\sum_k \exp(e_{ik})}$$

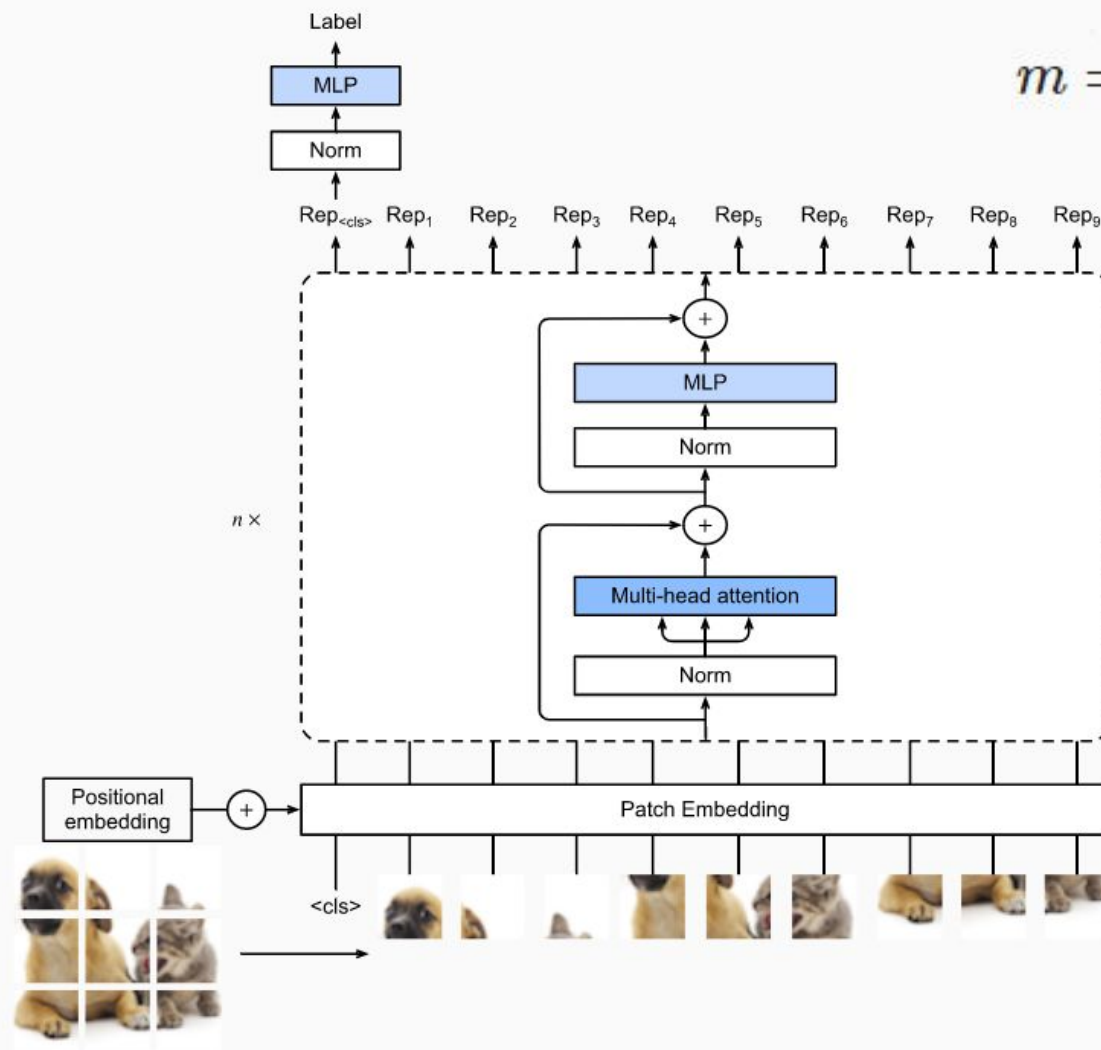
$$e_{ij} = f(s_i)g(h_j)^T$$

Attention Is All You Need



Multi-head attention





$$m = hw/p^2$$

