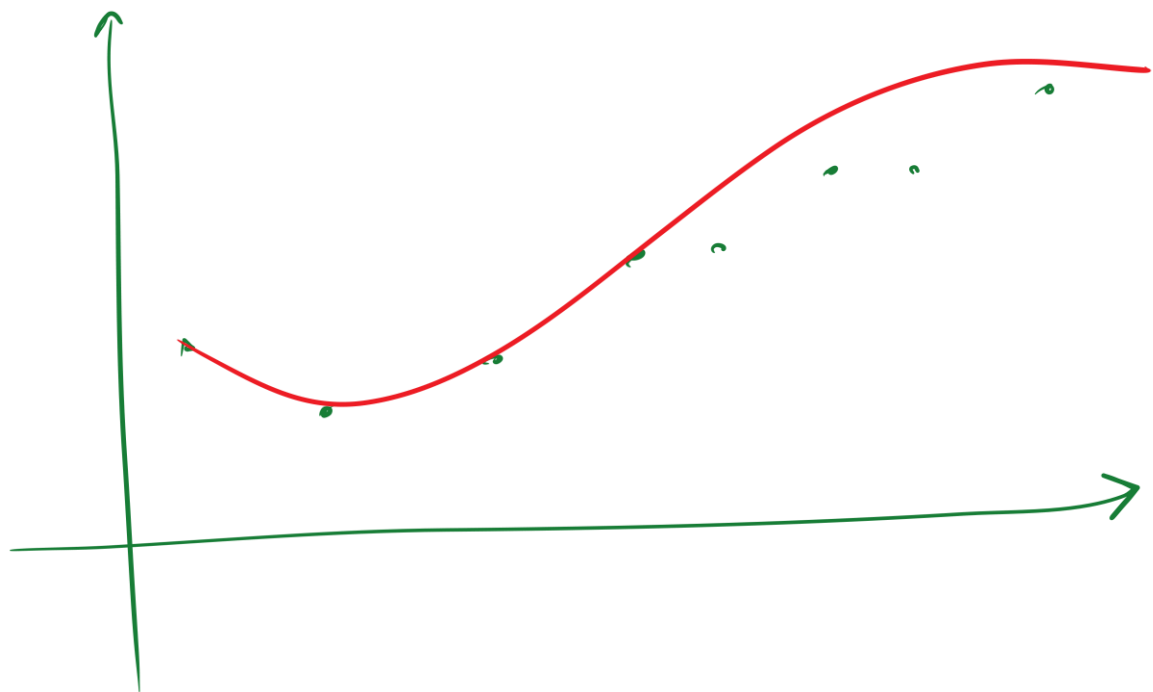


Time Series



Regression

Time Series

RNN

Recurrent

NN

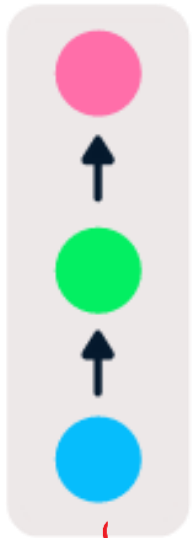
1993

image Captioning

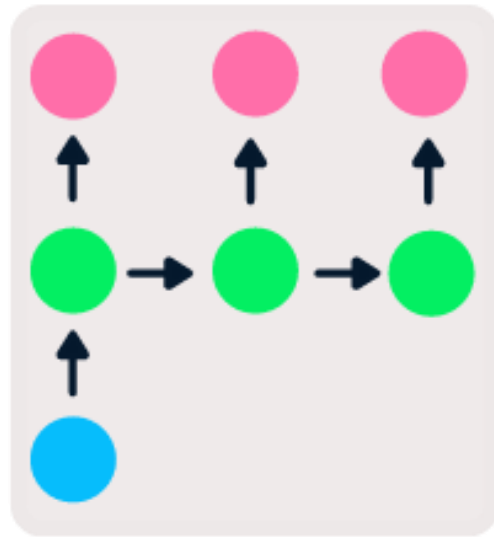
text Classification

Translate

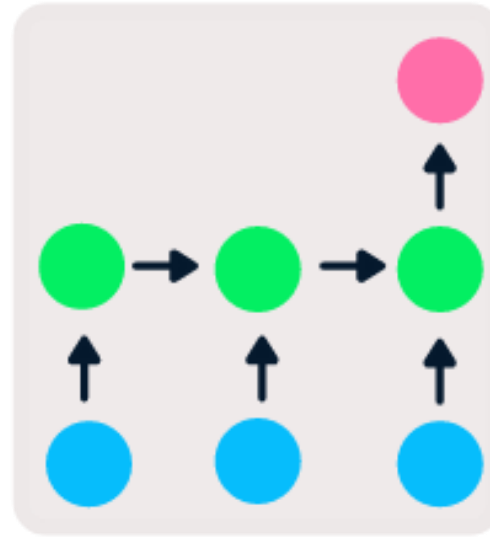
One to One



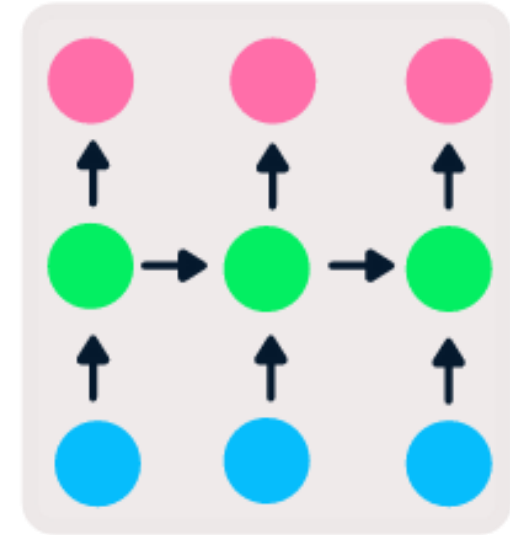
One to Many



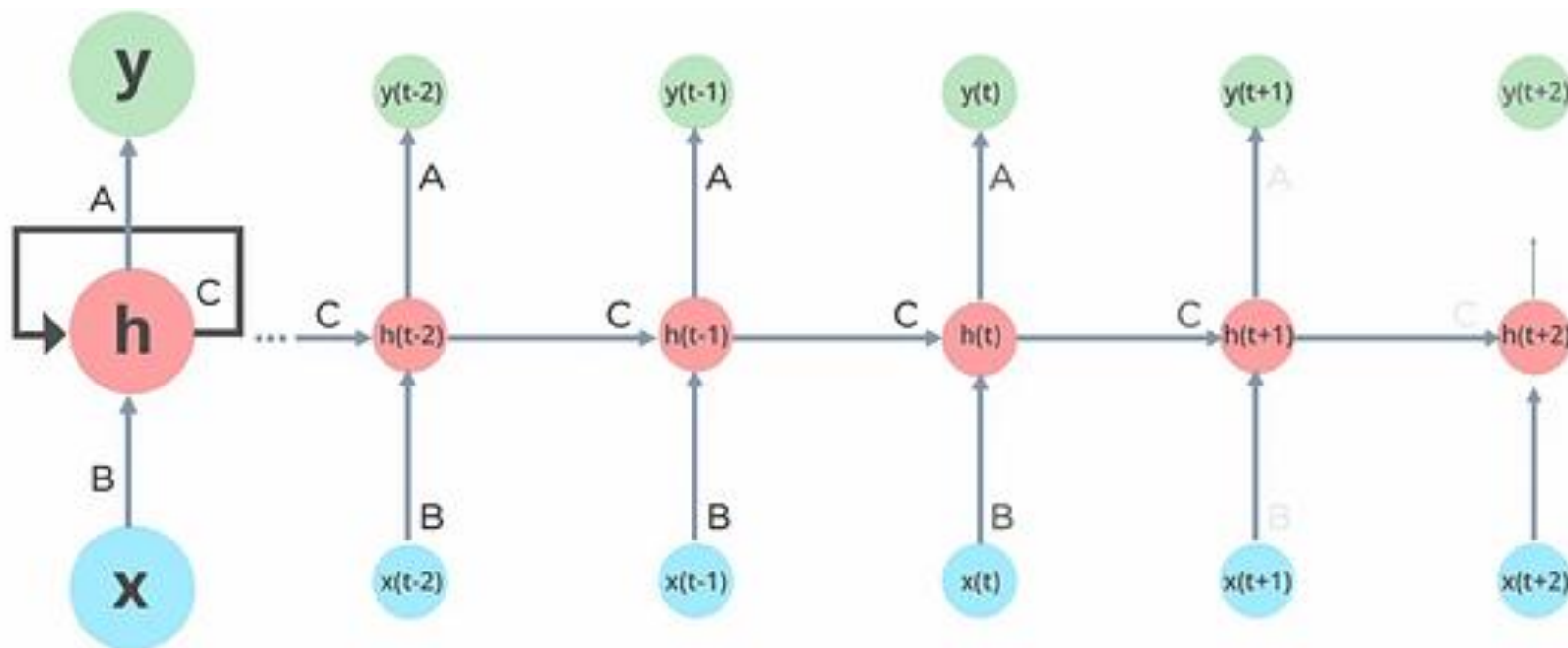
Many to One

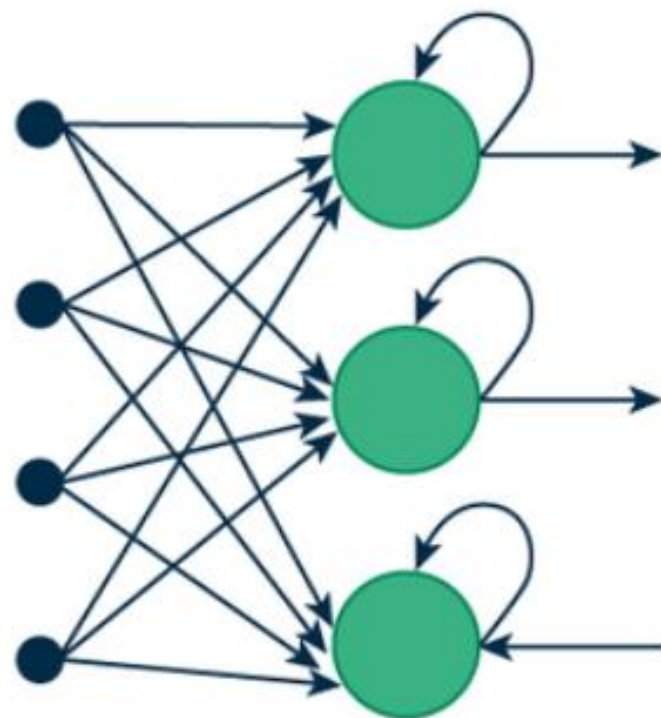


Many to Many

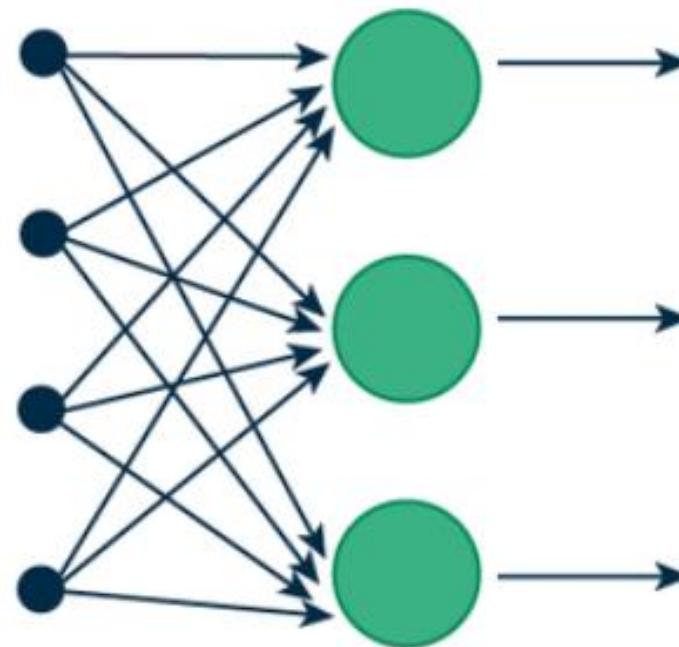


RNN

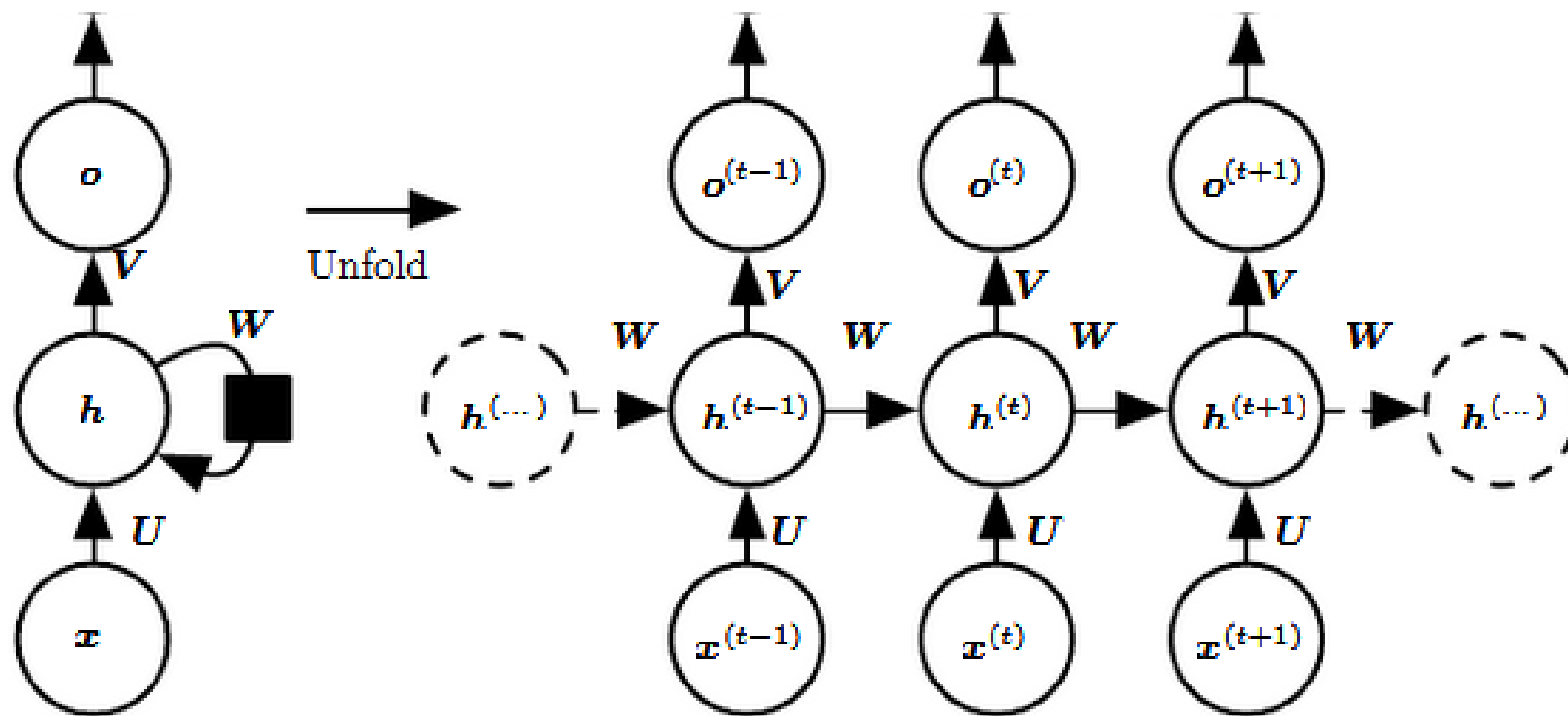




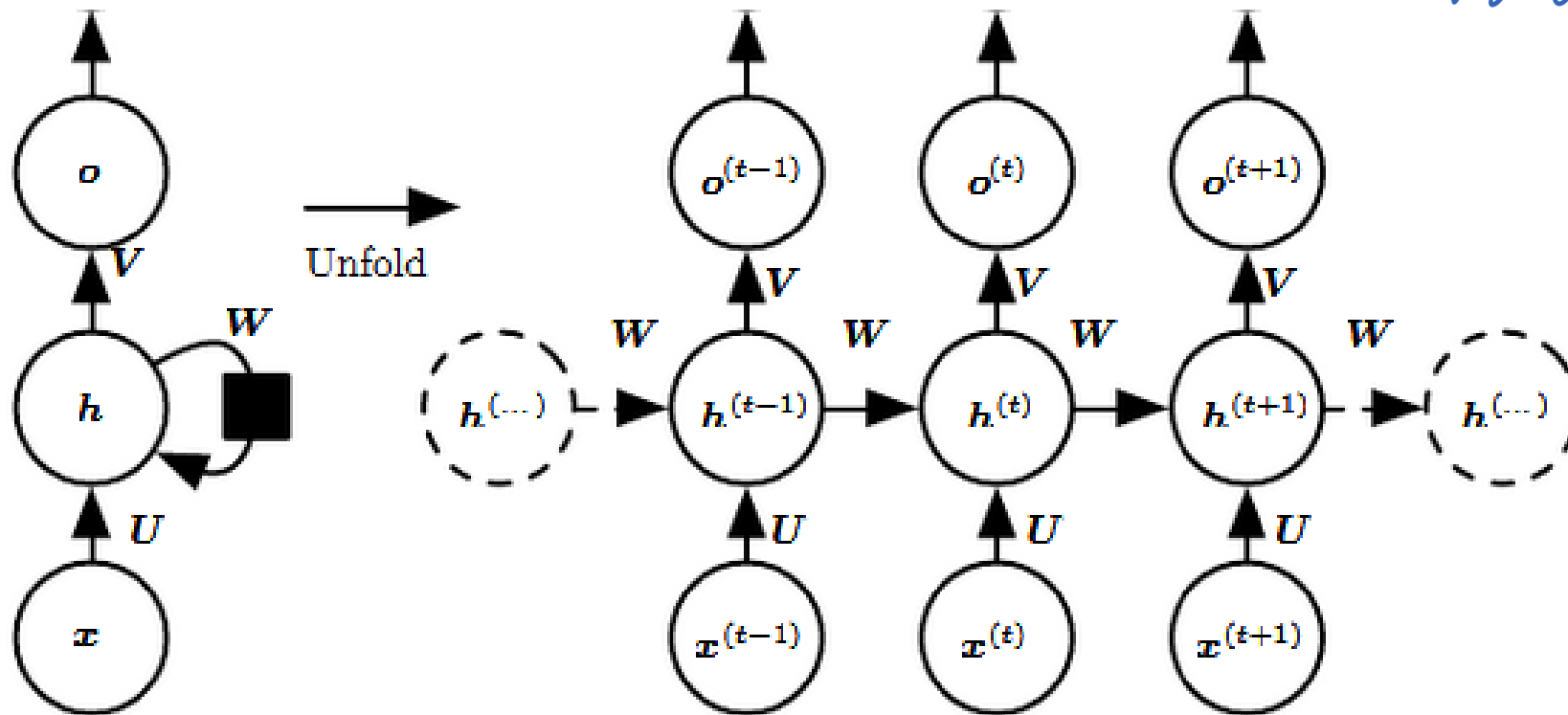
(a) Recurrent Neural Network



(b) Feed-Forward Neural Network



forward



$$a^{(t)} = b + Wh^{(t-1)} + Ux^{(t)}$$

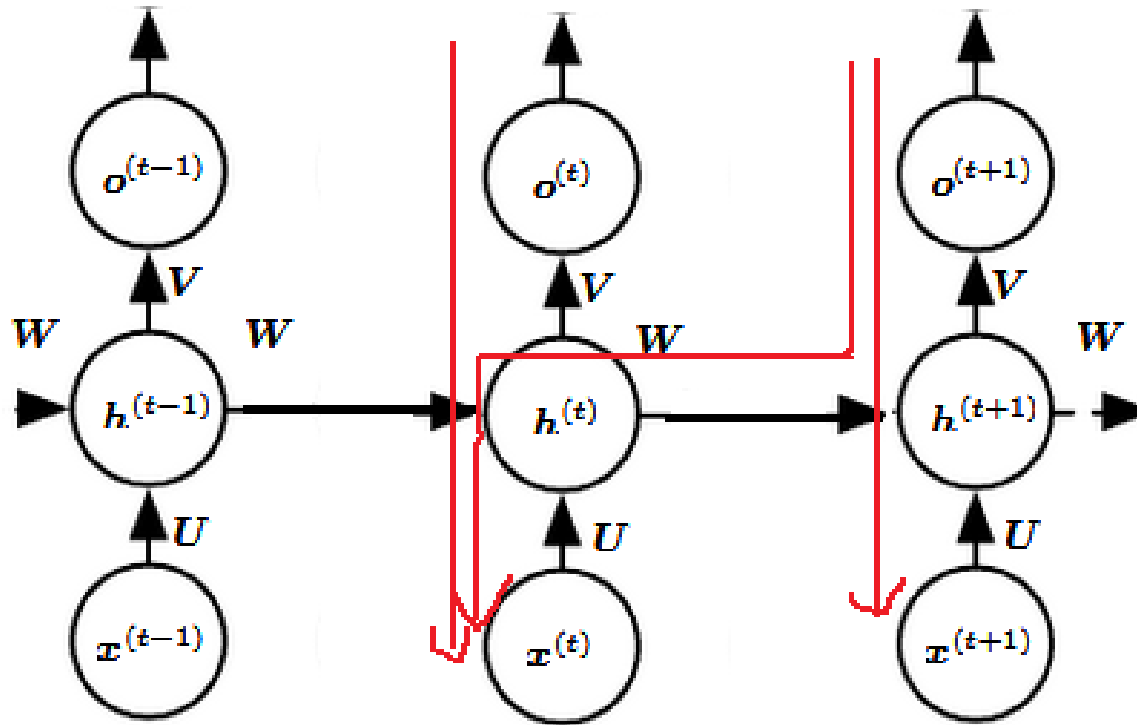
$$h^{(t)} = \tanh(a^{(t)})$$

$$o^{(t)} = c + Vh^{(t)}$$

$$\hat{y}^{(t)} = \text{softmax}(o^{(t)})$$

$$\text{units}^n + (\text{units}^n \times \text{units}^n) + (\text{units}^n \times \text{input}^n)$$

Train
(Back Propagation)



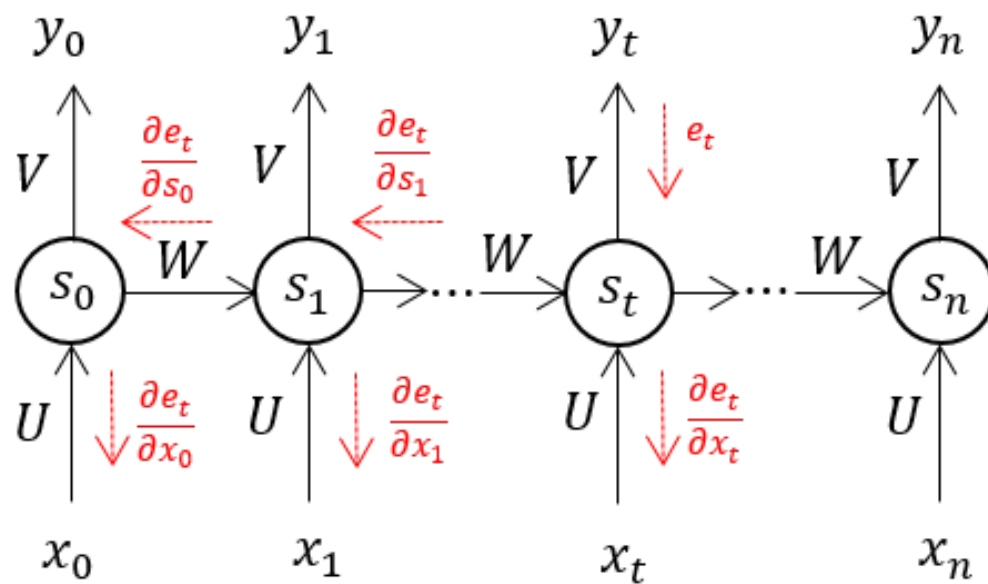
$$\nabla_{\mathbf{c}} L = \sum_t \left(\frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{c}} \right)^\top \nabla_{\mathbf{o}^{(t)}} L = \sum_t \nabla_{\mathbf{o}^{(t)}} L$$

$$\nabla_{\mathbf{b}} L = \sum_t \left(\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{b}^{(t)}} \right)^\top \nabla_{\mathbf{h}^{(t)}} L = \sum_t \text{diag} \left(1 - \left(\mathbf{h}^{(t)} \right)^2 \right) \nabla_{\mathbf{h}^{(t)}} L$$

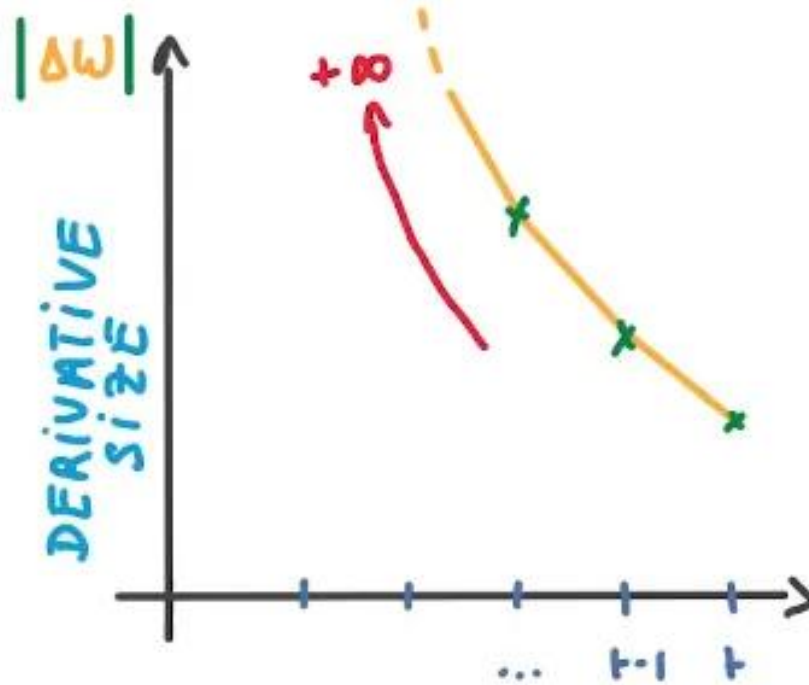
$$\nabla_{\mathbf{V}} L = \sum_t \sum_i \left(\frac{\partial L}{\partial \mathbf{o}_i^{(t)}} \right) \nabla_{\mathbf{V}^{(t)}} \mathbf{o}_i^{(t)} = \sum_t (\nabla_{\mathbf{o}^{(t)}} L) \mathbf{h}^{(t)\top}$$

$$\nabla_{\mathbf{W}} L = \sum_t \sum_i \left(\frac{\partial L}{\partial \mathbf{h}_i^{(t)}} \right) \nabla_{\mathbf{W}^{(t)}} \mathbf{h}_i^{(t)} = \sum_t \text{diag} \left(1 - \left(\mathbf{h}^{(t)} \right)^2 \right) (\nabla_{\mathbf{h}^{(t)}} L) \mathbf{h}^{(t-1)\top}$$

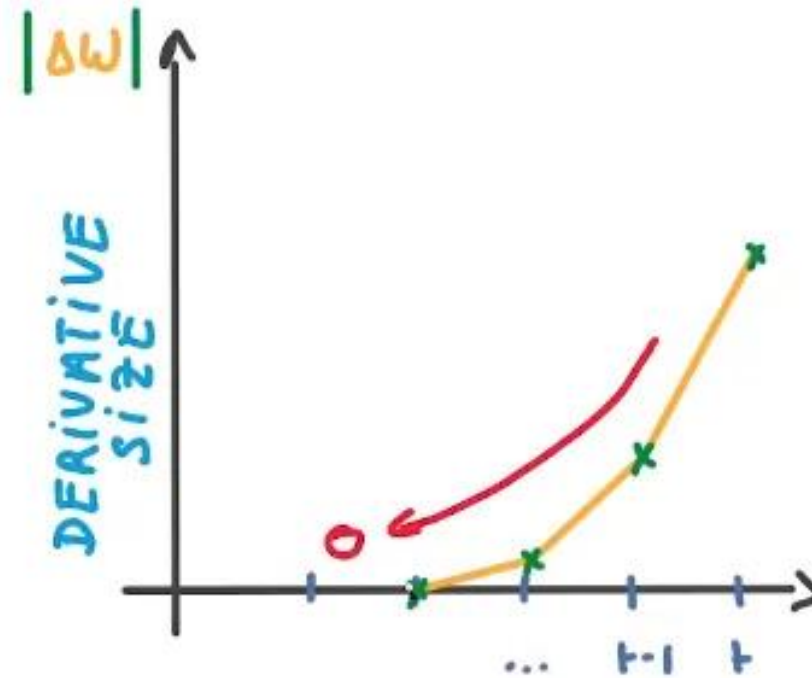
$$\nabla_{\mathbf{U}} L = \sum_t \sum_i \left(\frac{\partial L}{\partial \mathbf{h}_i^{(t)}} \right) \nabla_{\mathbf{U}^{(t)}} \mathbf{h}_i^{(t)} = \sum_t \text{diag} \left(1 - \left(\mathbf{h}^{(t)} \right)^2 \right) (\nabla_{\mathbf{h}^{(t)}} L) \mathbf{x}^{(t)\top}$$



EXPLODING GRADIENT



VANISHING GRADIENT



This is a classroom

$$Th_i \rightarrow S$$

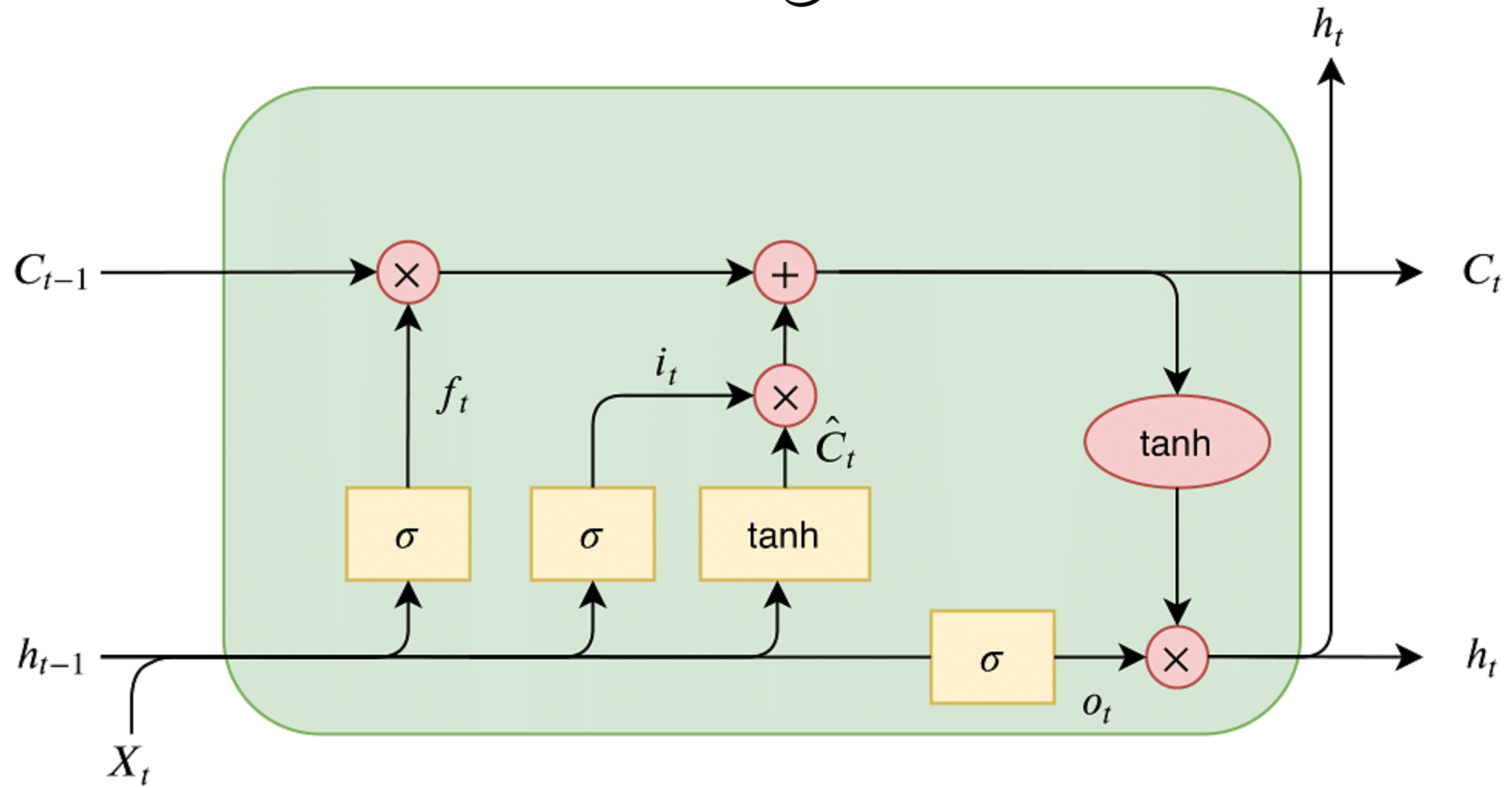
$\text{his} \rightarrow \text{'}$

is - i

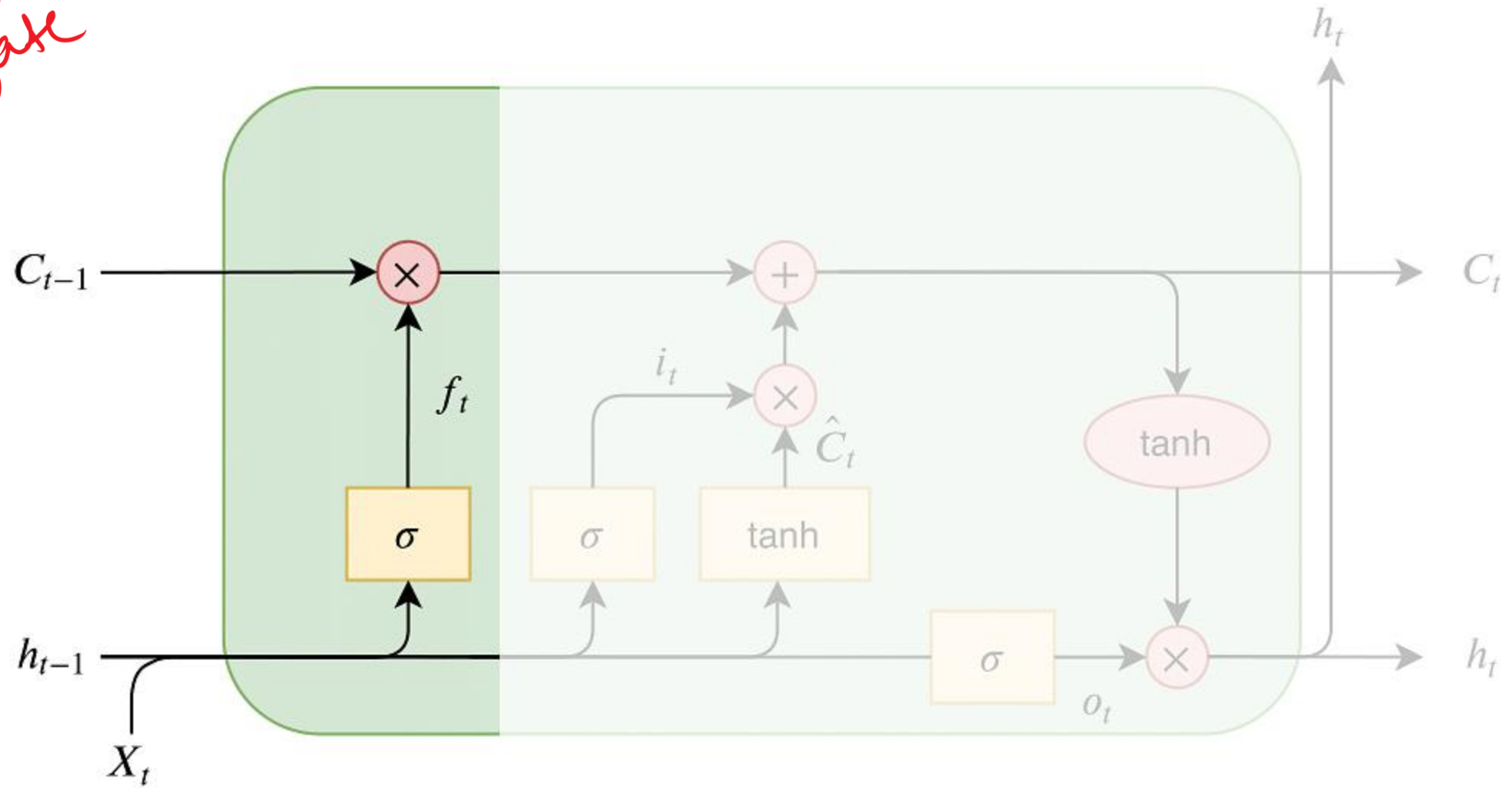
11

1997

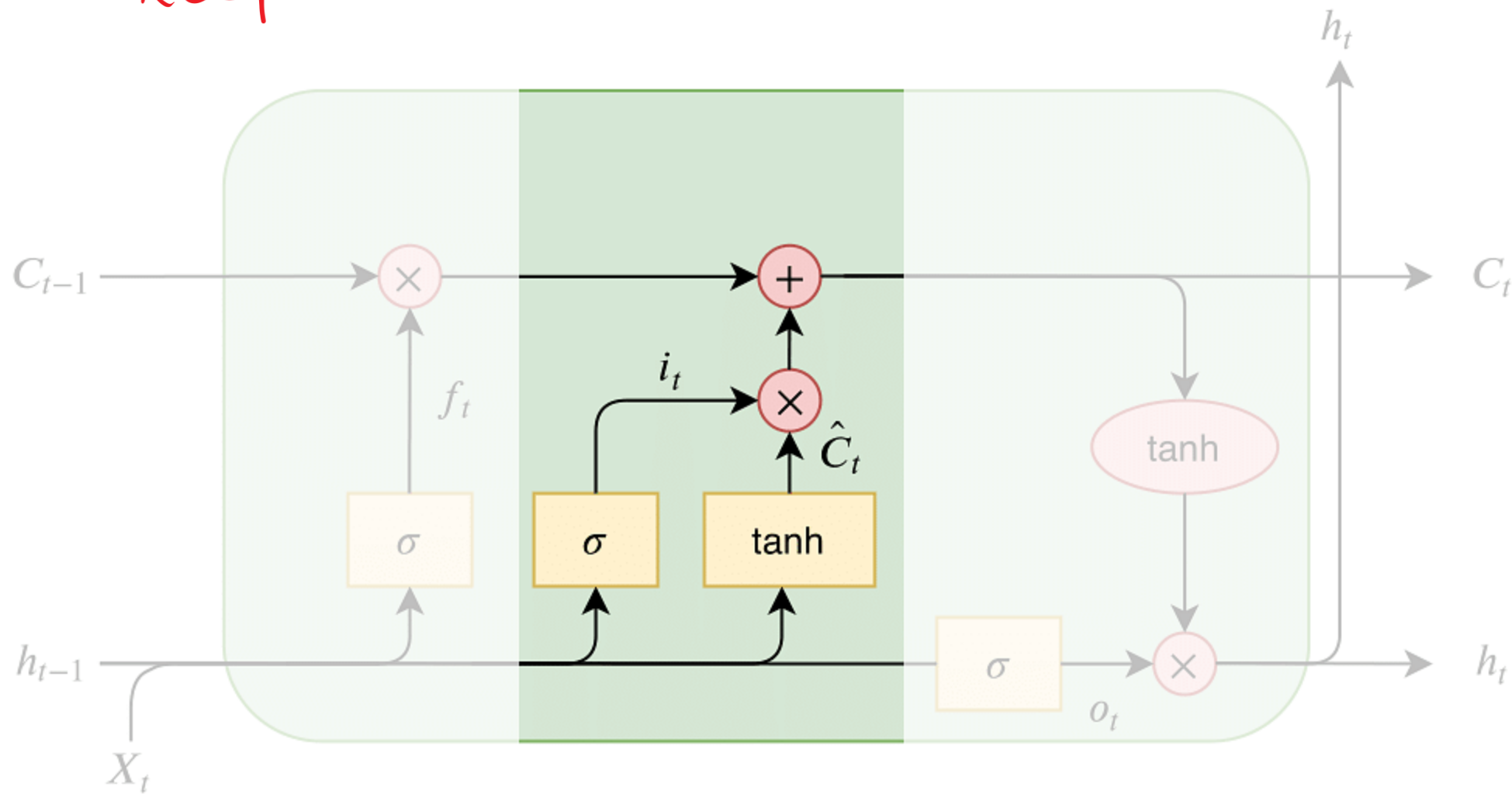
LSTM long short term memory



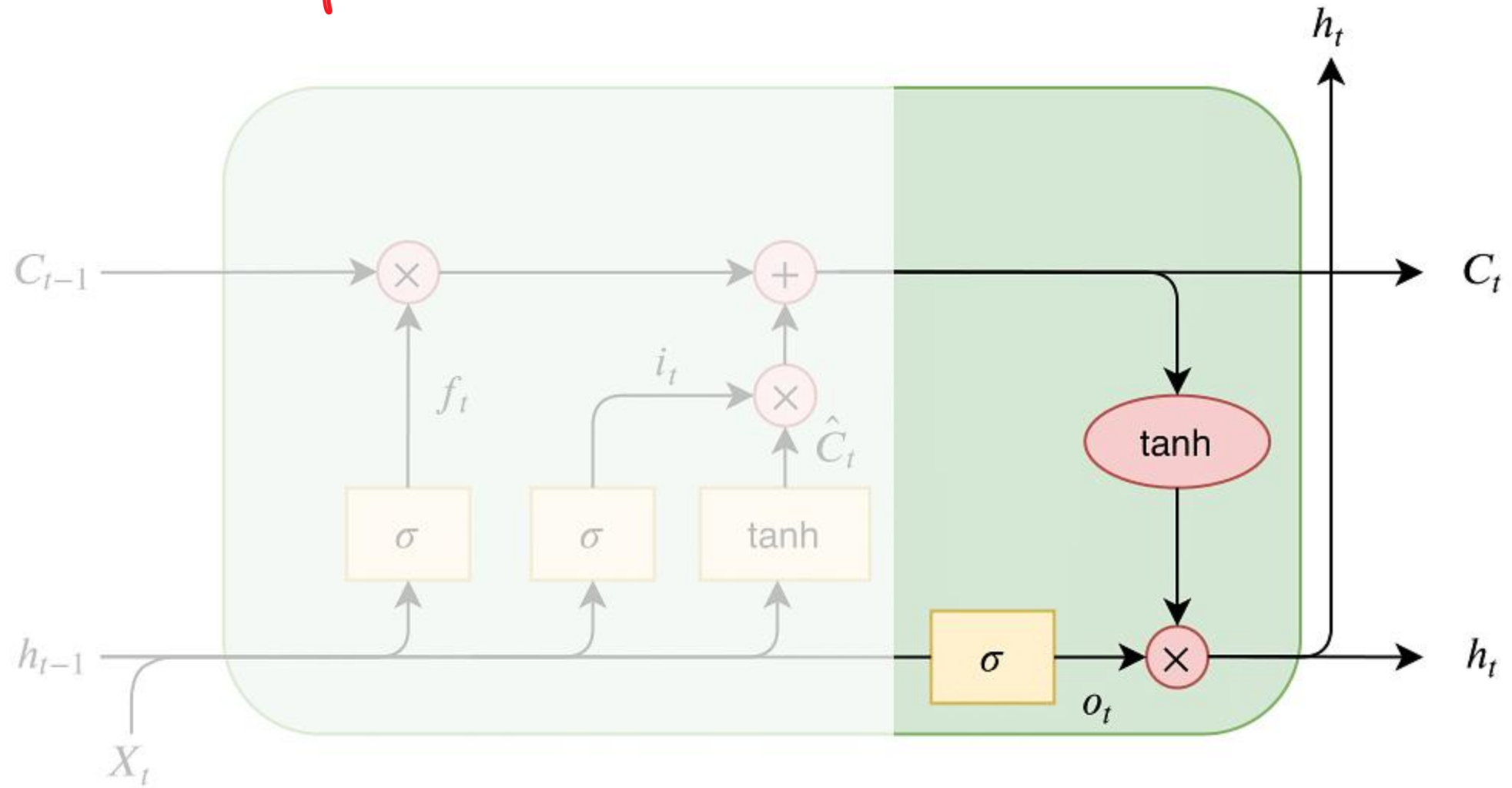
forget gate



update
keep

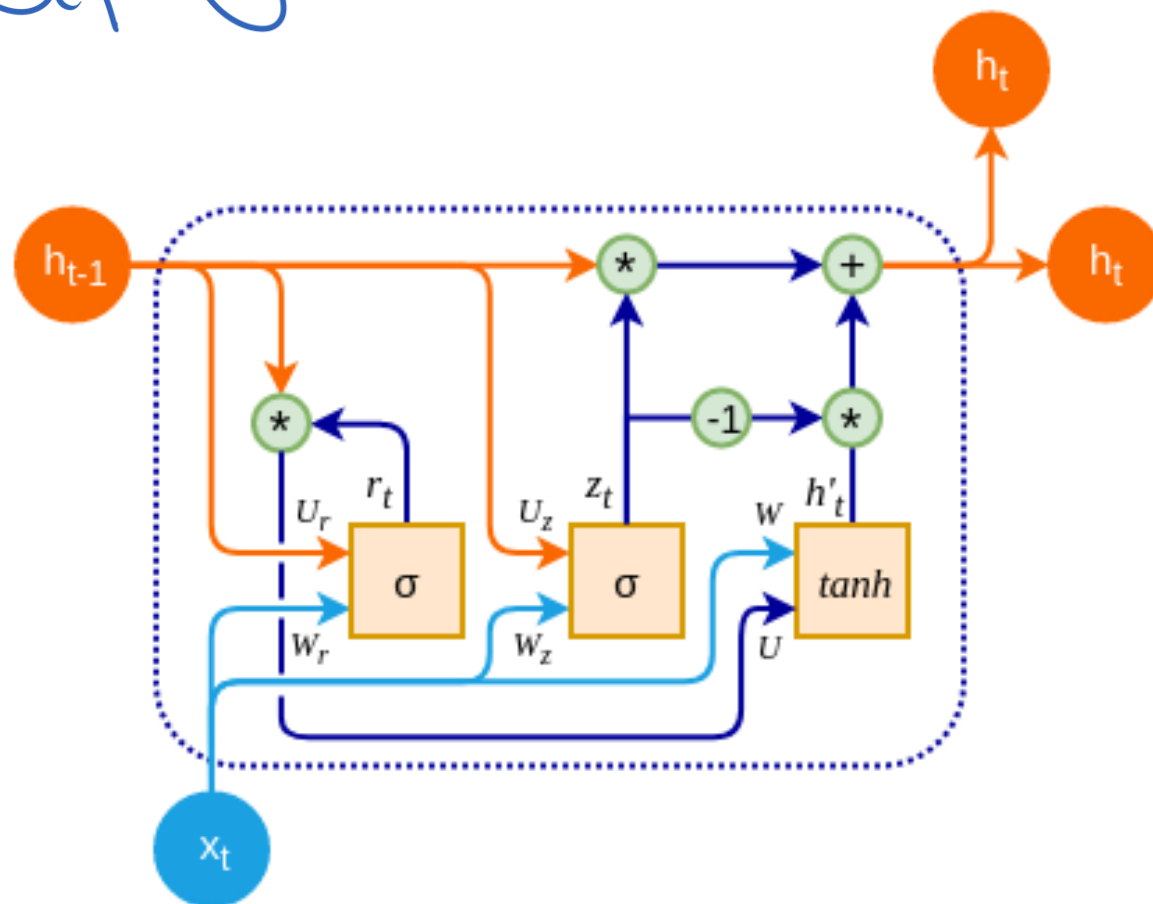


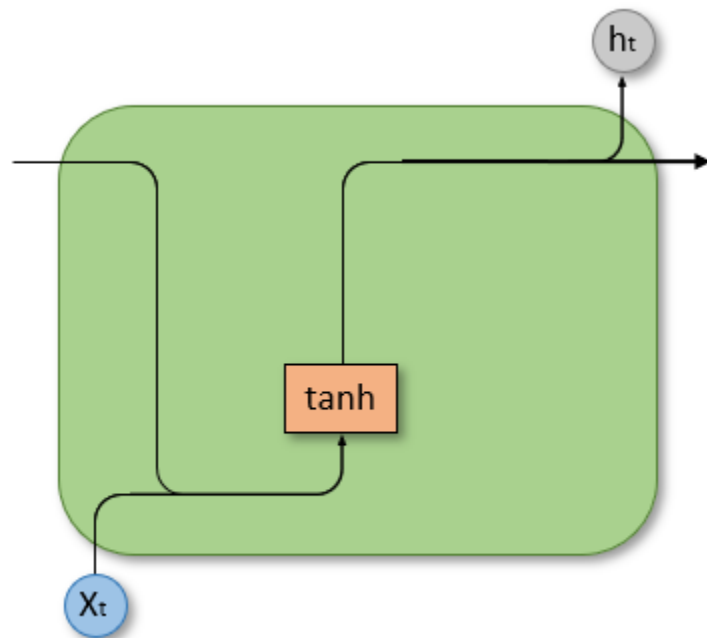
Output



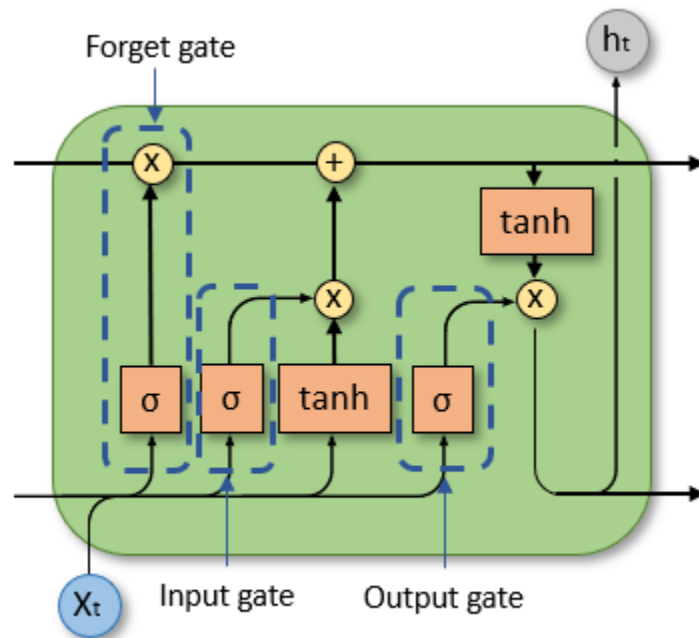
2014

GRU

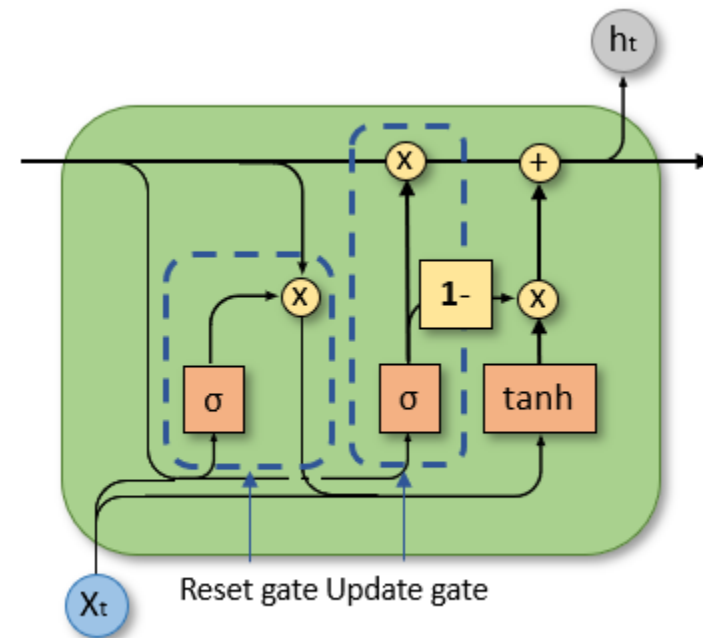




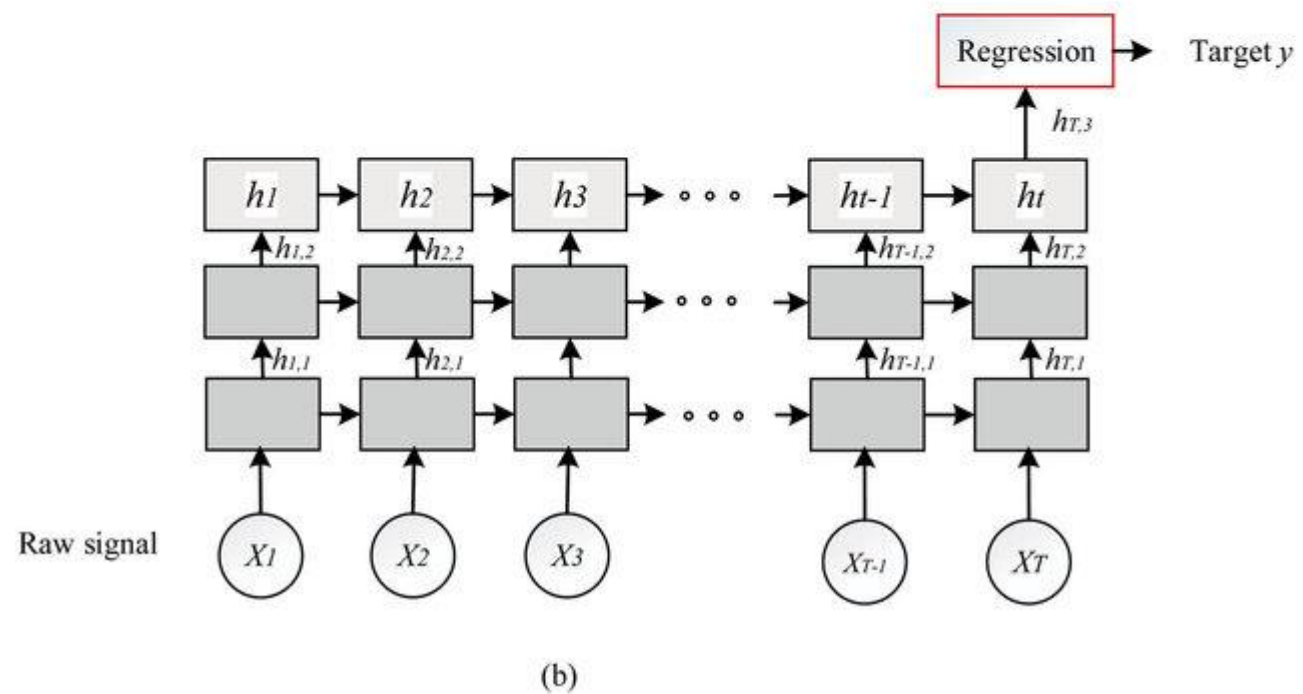
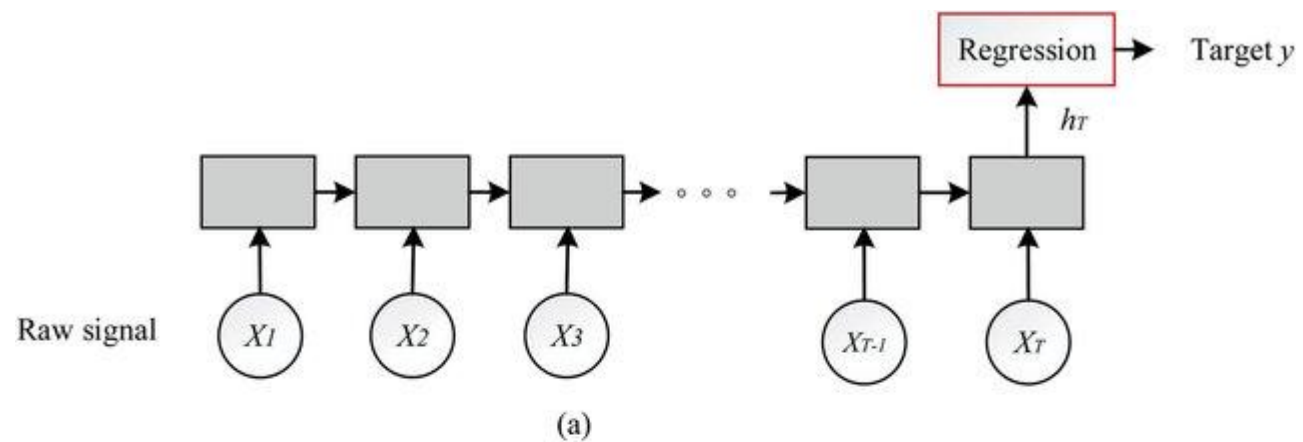
RNN

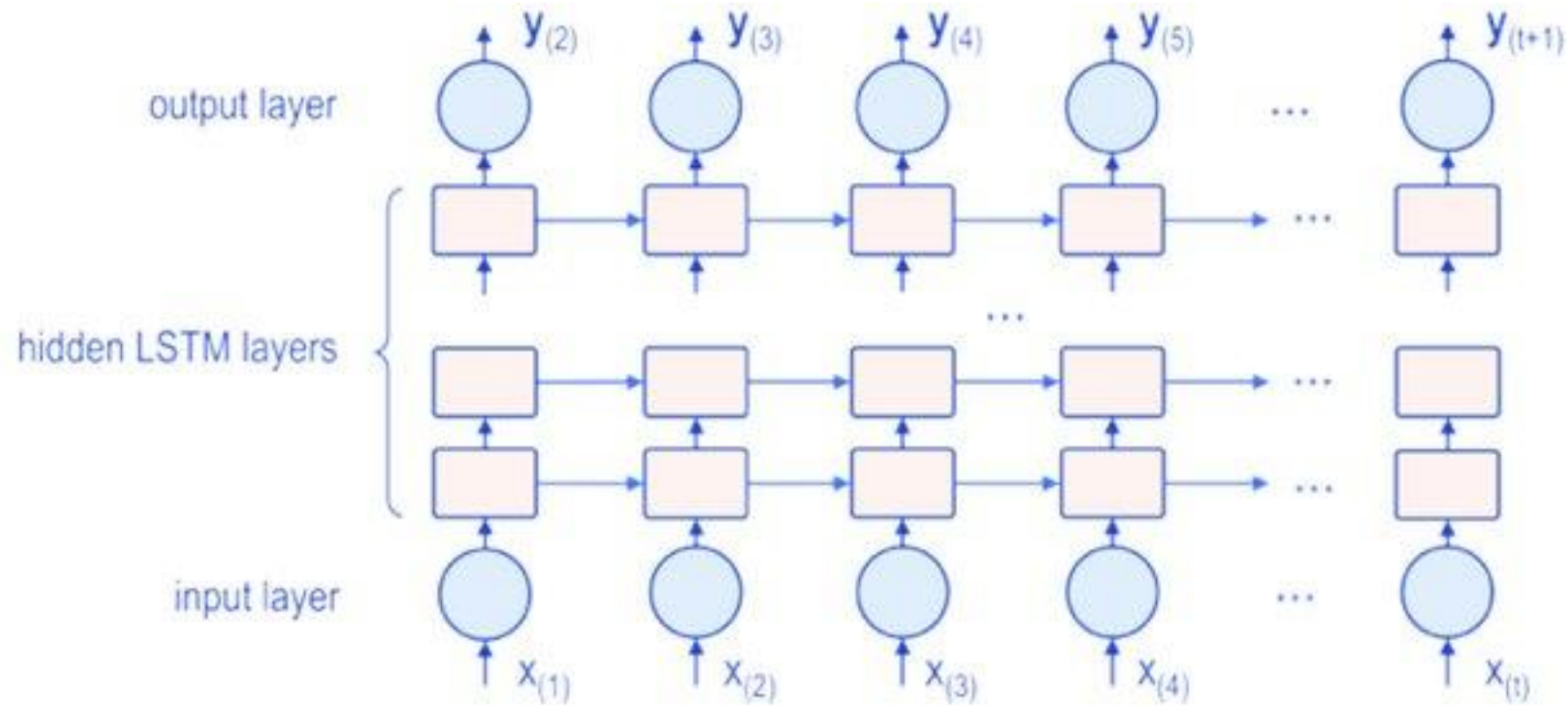


LSTM

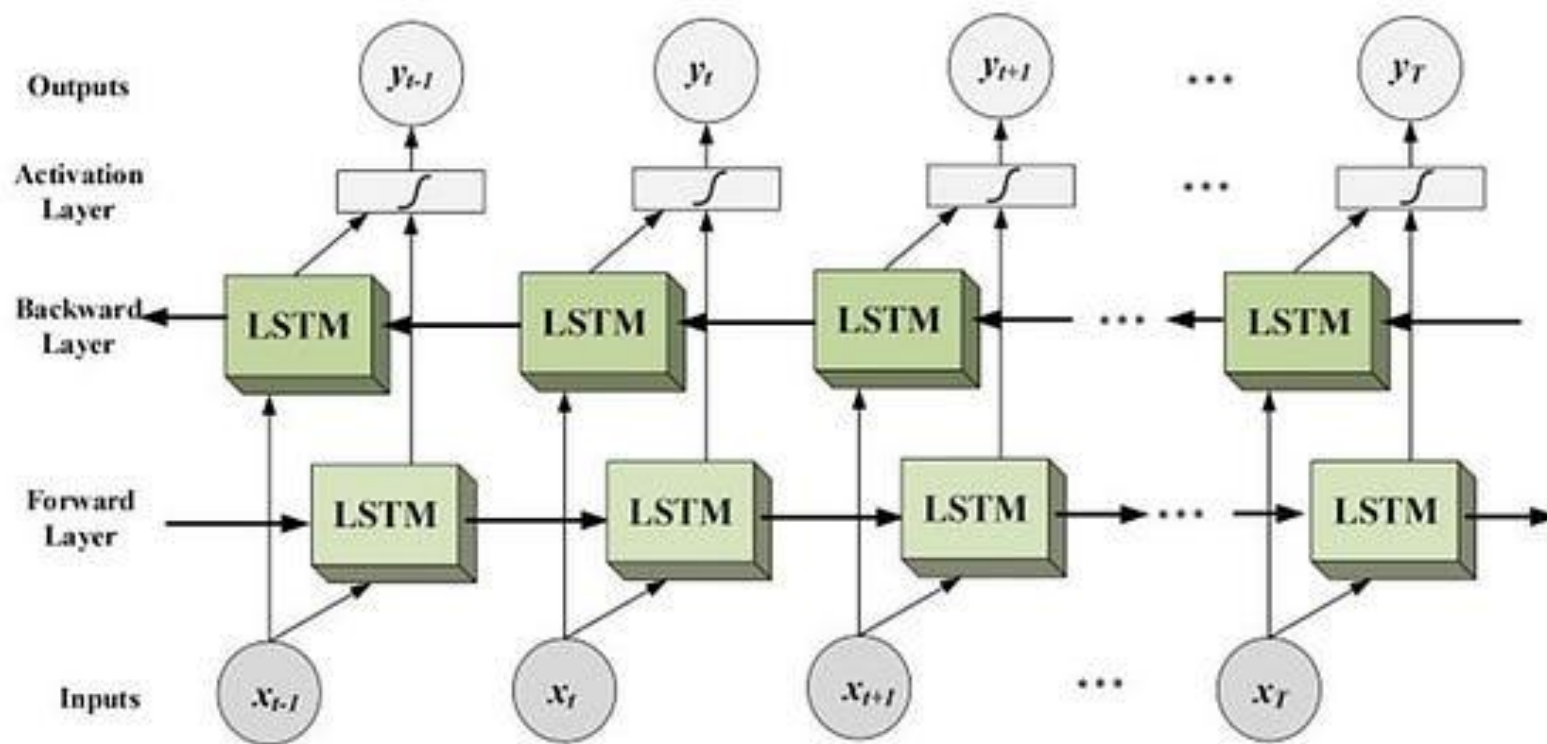


GRU



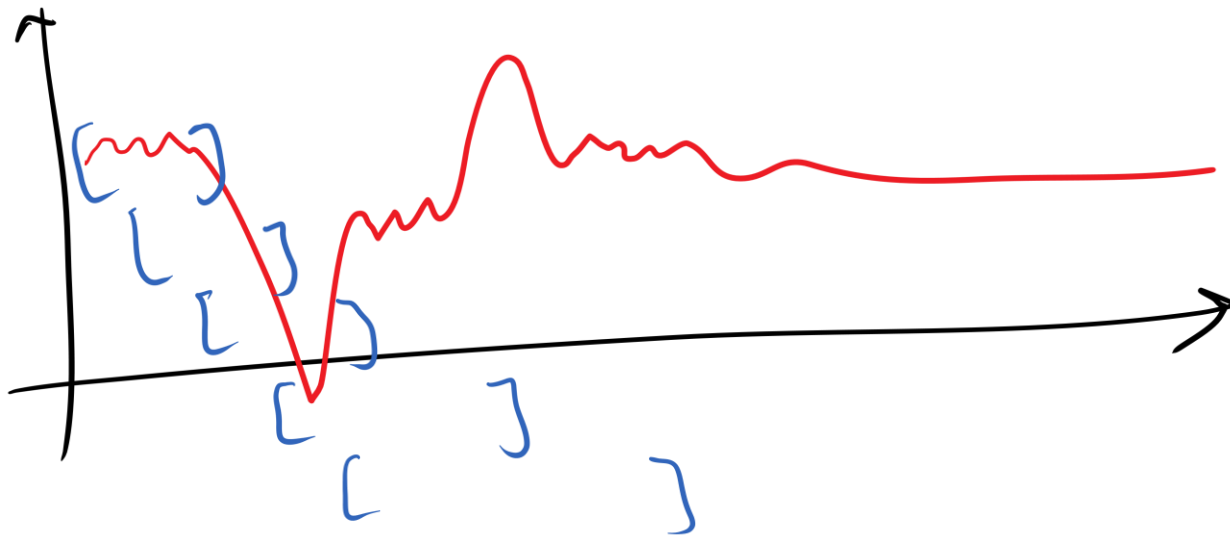


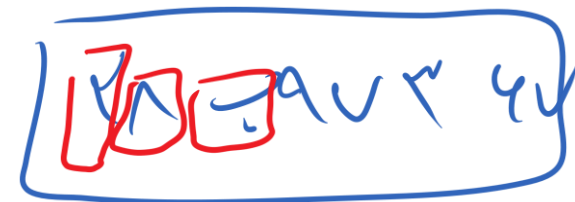
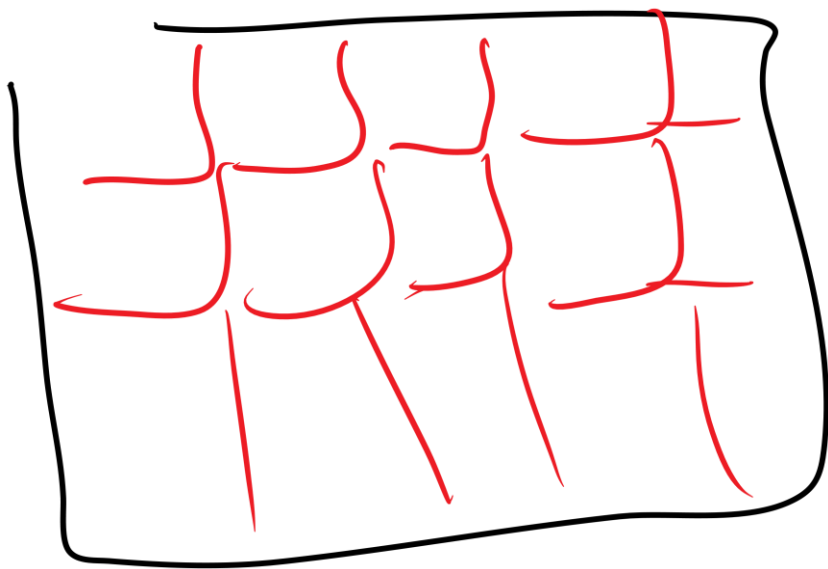
Bidirectional



CNN + RNN
LSTM

کس، مگر
tracking

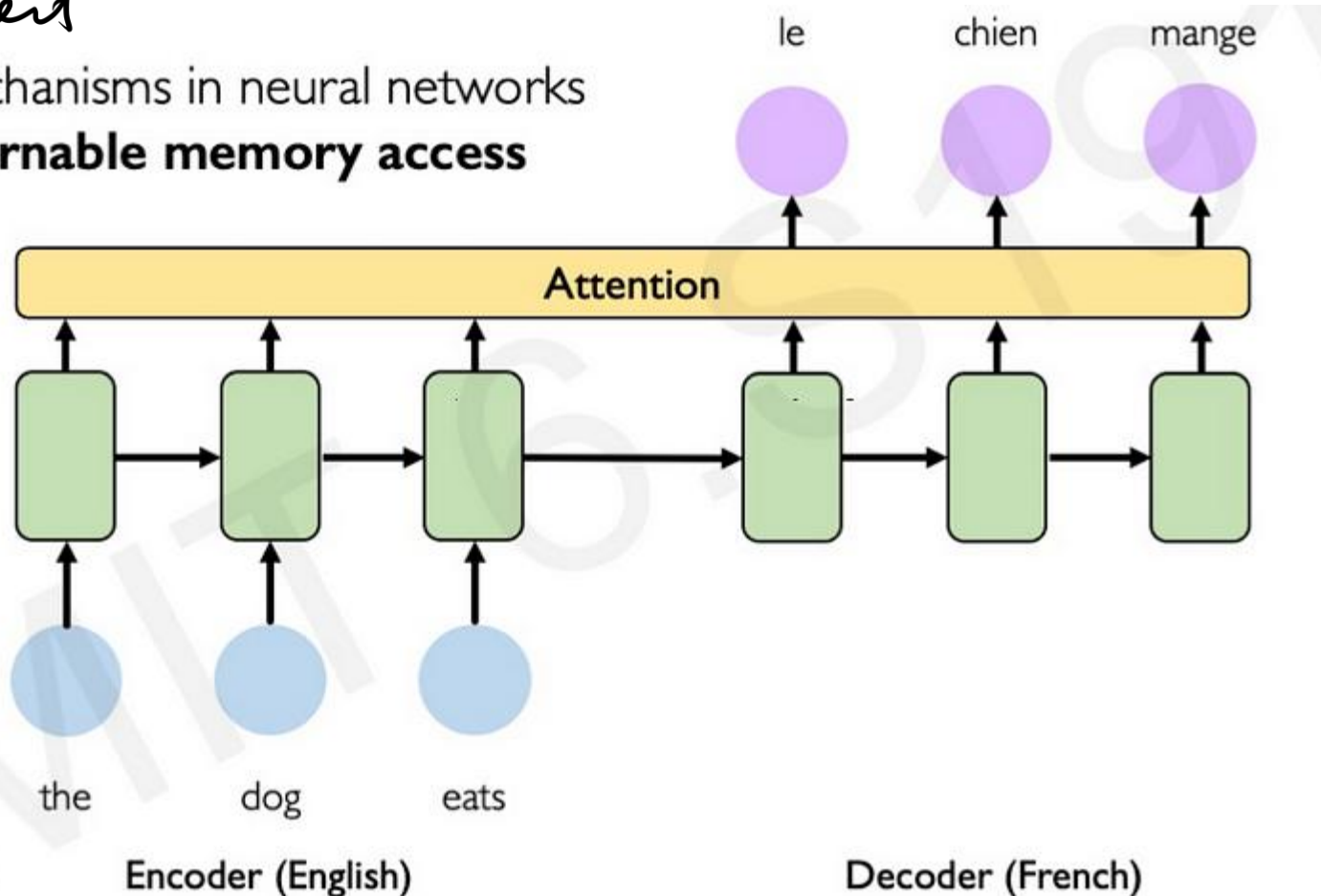


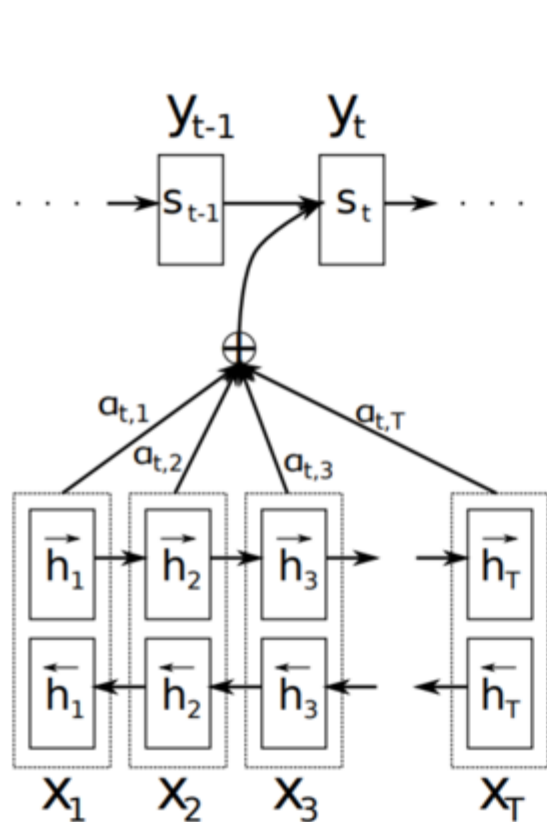


CNN + LSTM

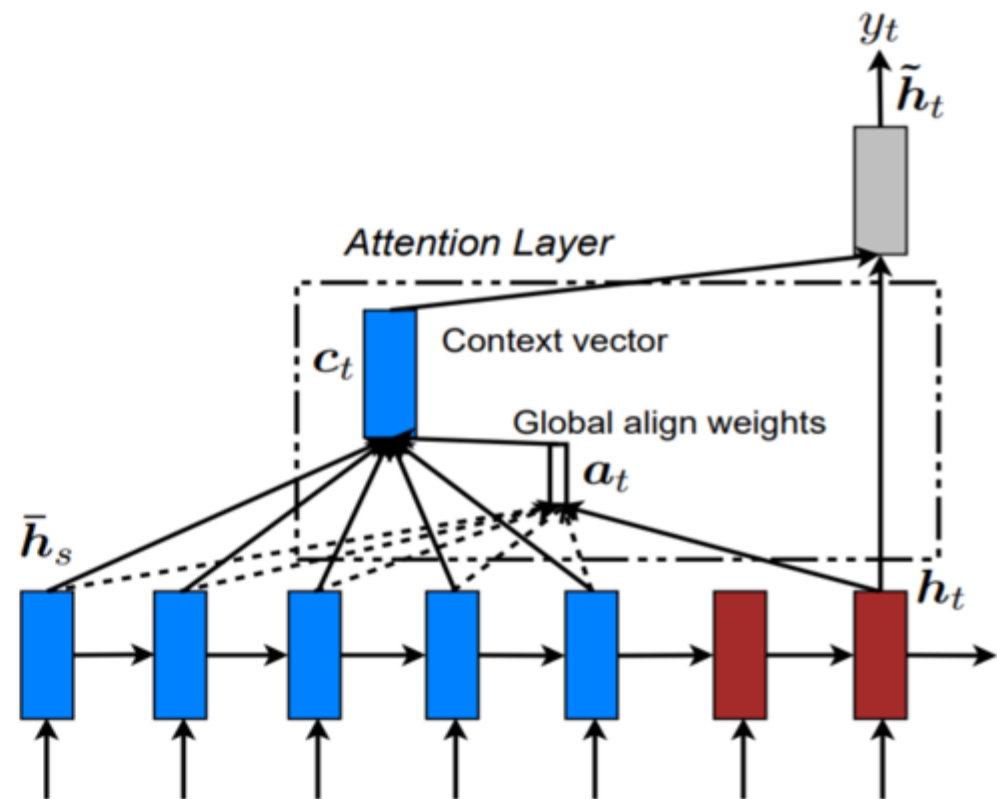
من یک دانشجو هستم
I am a student

Attention mechanisms in neural networks
provide **learnable memory access**

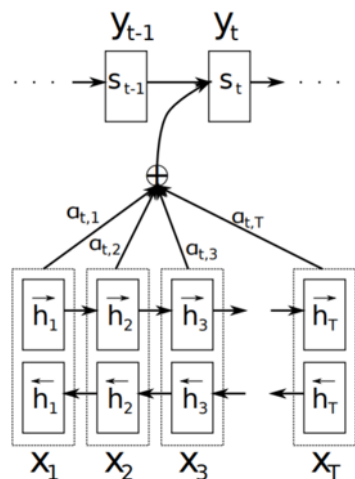




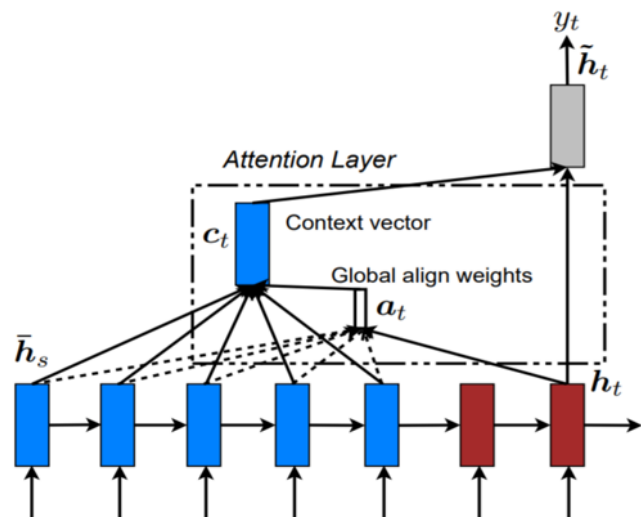
Bahdanau attention mechanism



Luong attention mechanism



Bahdanau attention mechanism



Luong attention mechanism

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

[Attention weights]

$$\mathbf{c}_t = \sum_s \alpha_{ts} \bar{\mathbf{h}}_s$$

[Context vector]

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

[Attention vector]

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

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez*[†]
University of Toronto
aidan@cs.toronto.edu

Łukasz Kaiser*
Google Brain
lukaszkaizer@google.com

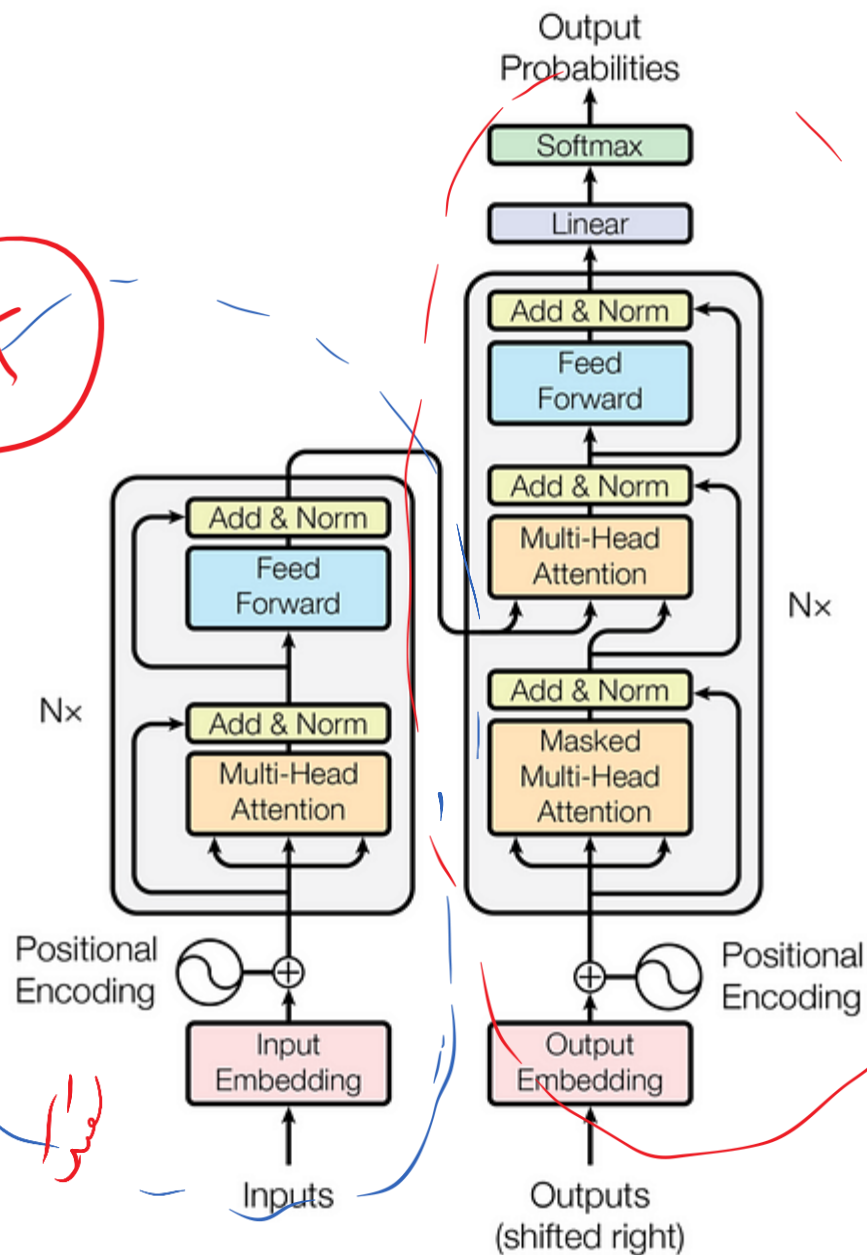
Illia Polosukhin*[‡]
illia.polosukhin@gmail.com

Abstract

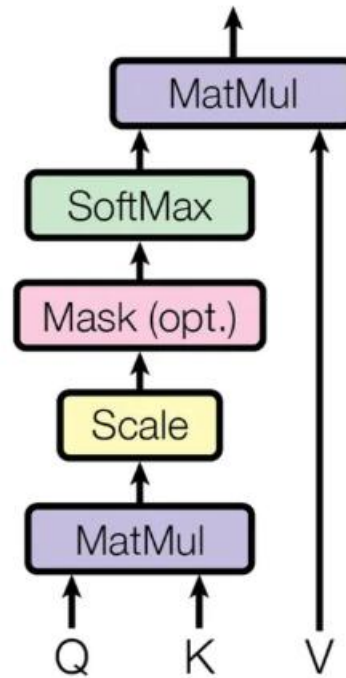
The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions

Encoder
Decoder

VIT

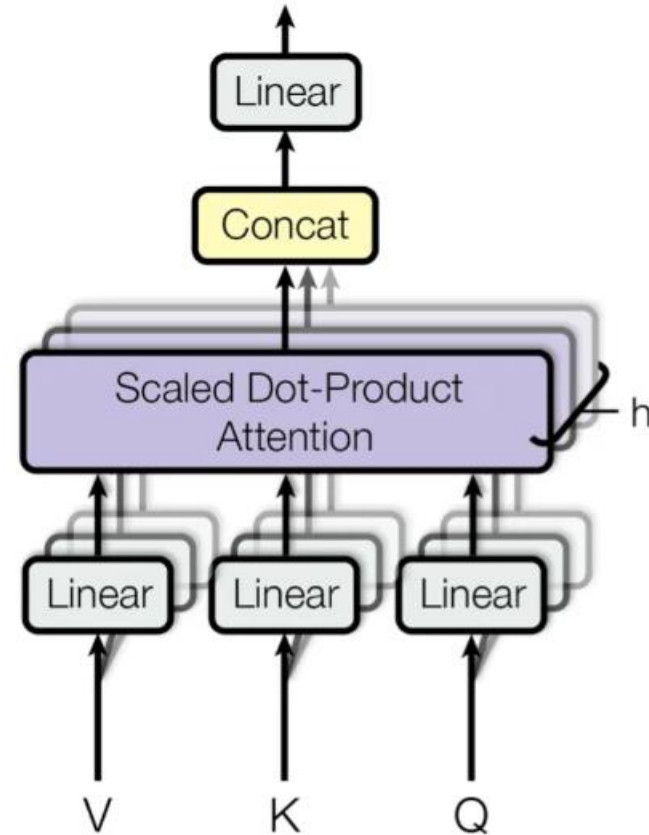


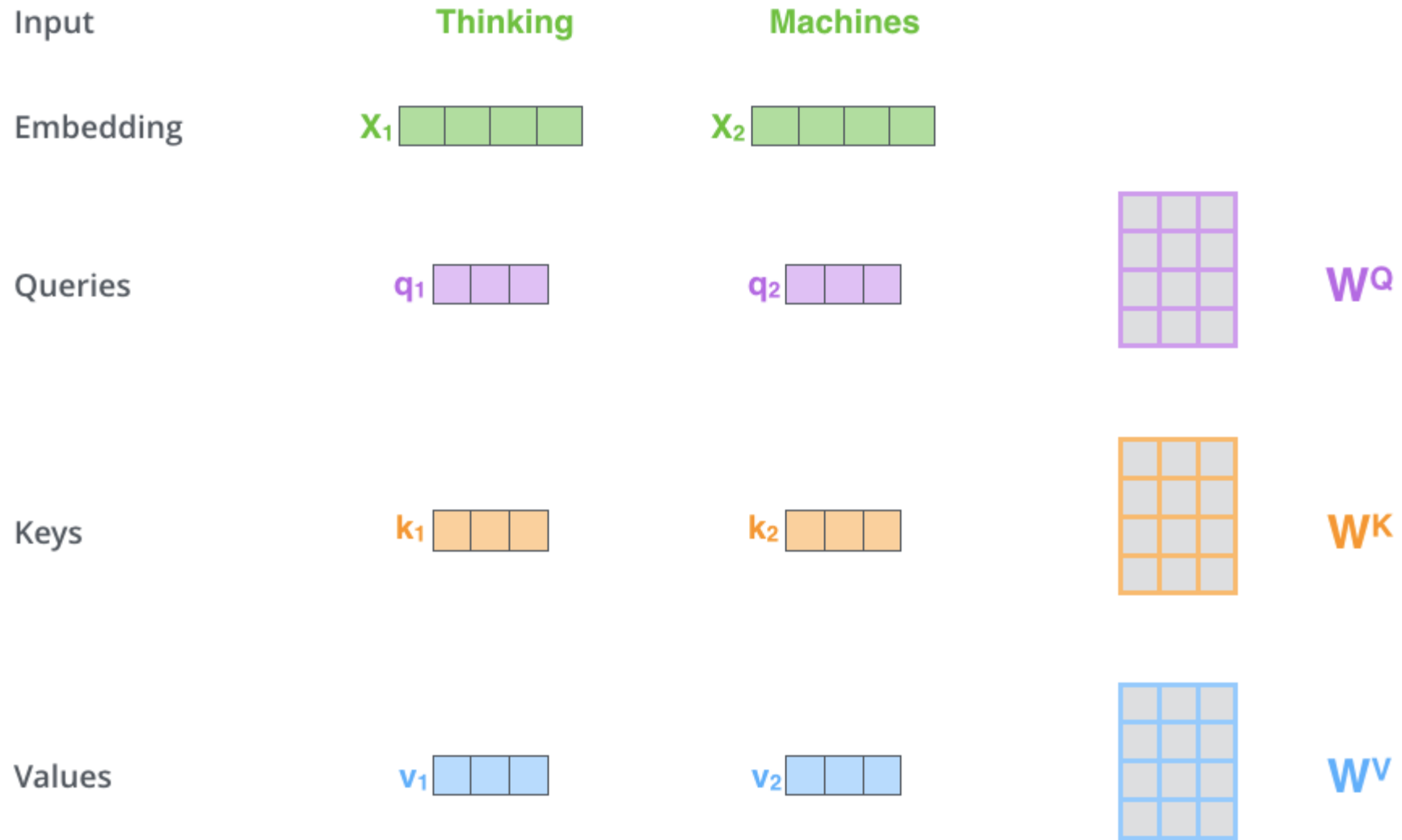
Scaled Dot-Product Attention



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

Multi-Head Attention





$$\begin{array}{c}
 \mathbf{X} \\
 \begin{array}{|c|c|c|c|} \hline \square & \square & \square & \square \\ \hline \square & \square & \square & \square \\ \hline \end{array}
 \end{array}
 \times
 \begin{array}{c}
 \mathbf{W}^Q \\
 \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array}
 \end{array}
 =
 \begin{array}{c}
 \mathbf{Q} \\
 \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array}
 \end{array}$$

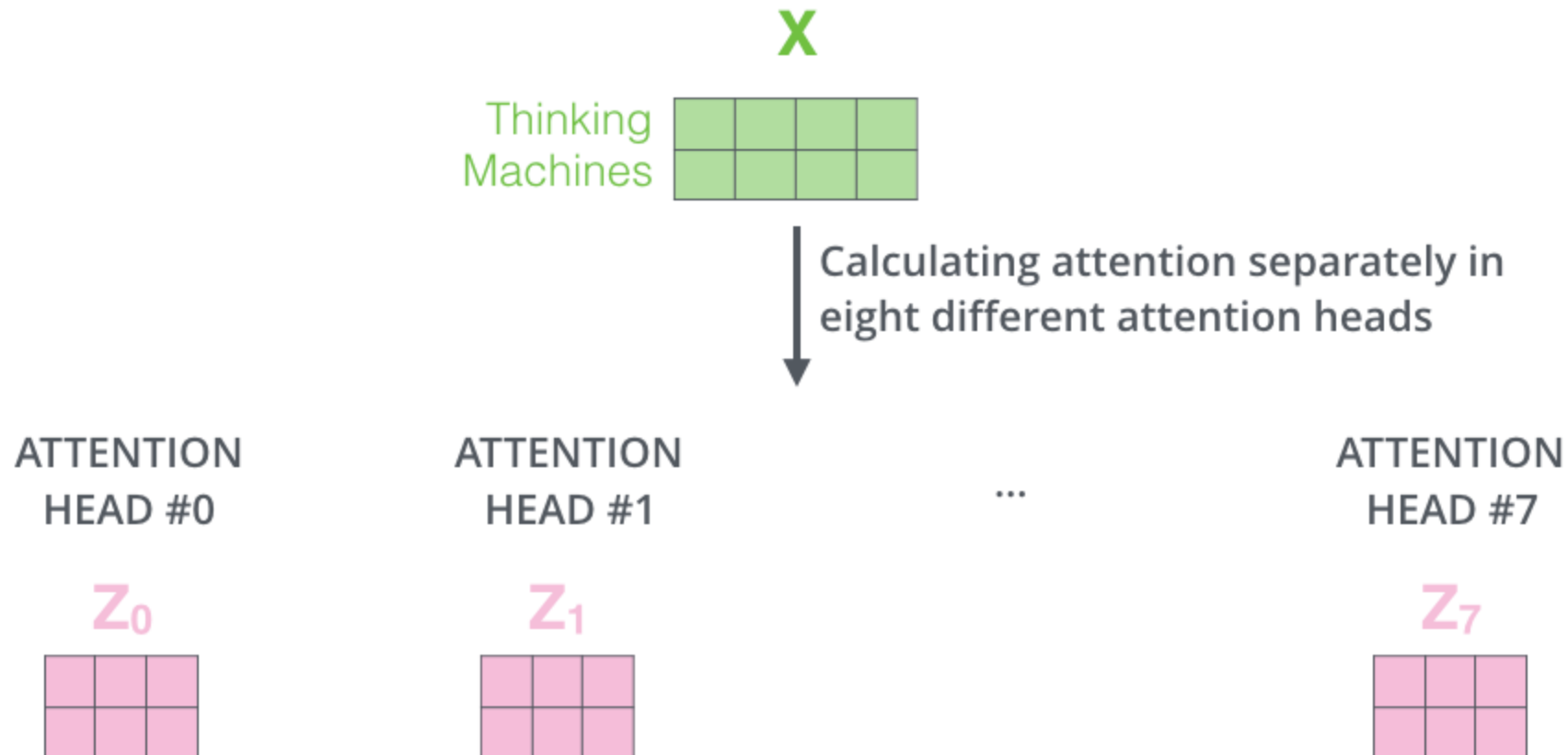
$$\begin{array}{c}
 \mathbf{X} \\
 \begin{array}{|c|c|c|c|} \hline \square & \square & \square & \square \\ \hline \square & \square & \square & \square \\ \hline \end{array}
 \end{array}
 \times
 \begin{array}{c}
 \mathbf{W}^K \\
 \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array}
 \end{array}
 =
 \begin{array}{c}
 \mathbf{K} \\
 \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array}
 \end{array}$$

$$\begin{array}{c}
 \mathbf{X} \\
 \begin{array}{|c|c|c|c|} \hline \square & \square & \square & \square \\ \hline \square & \square & \square & \square \\ \hline \end{array}
 \end{array}
 \times
 \begin{array}{c}
 \mathbf{W}^V \\
 \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array}
 \end{array}
 =
 \begin{array}{c}
 \mathbf{V} \\
 \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array}
 \end{array}$$

CSA : $\frac{A \cdot B}{(A \parallel B)}$

$$\text{softmax} \left(\frac{\underbrace{Q}_{\begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array}} \times \underbrace{K^T}_{\begin{array}{|c|c|} \hline \square & \square \\ \hline \square & \square \\ \hline \square & \square \\ \hline \end{array}}}{\sqrt{d_k}} \right) \underbrace{V}_{\begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array}}$$

$$= \underbrace{Z}_{\begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array}}$$



1) This is our input sentence*

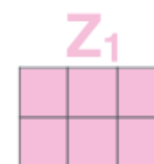
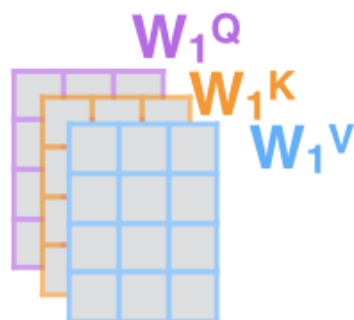
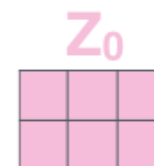
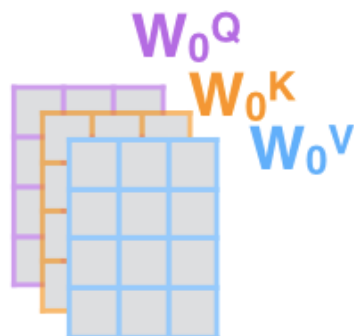
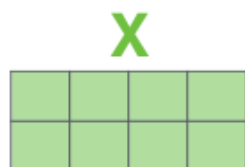
2) We embed each word*

3) Split into 8 heads. We multiply X or R with weight matrices

4) Calculate attention using the resulting $Q/K/V$ matrices

5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer

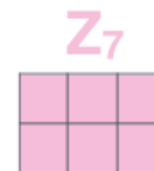
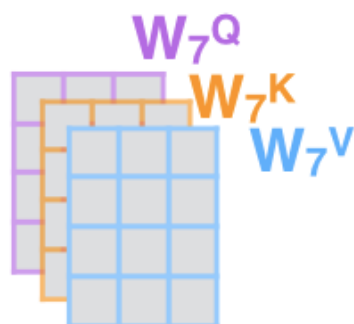
Thinking
Machines



...

...

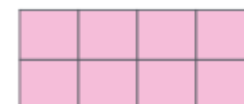
...



W^O

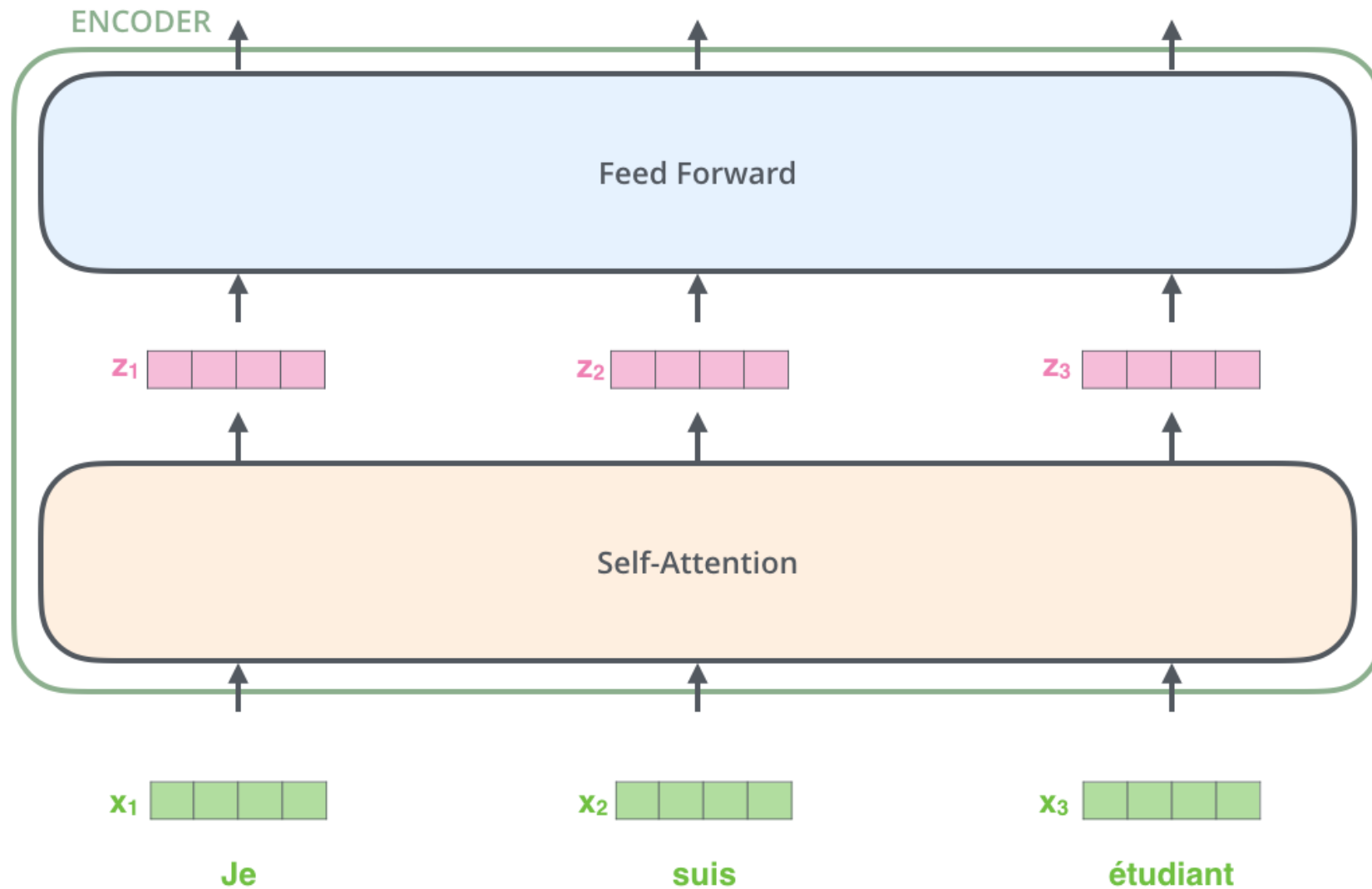


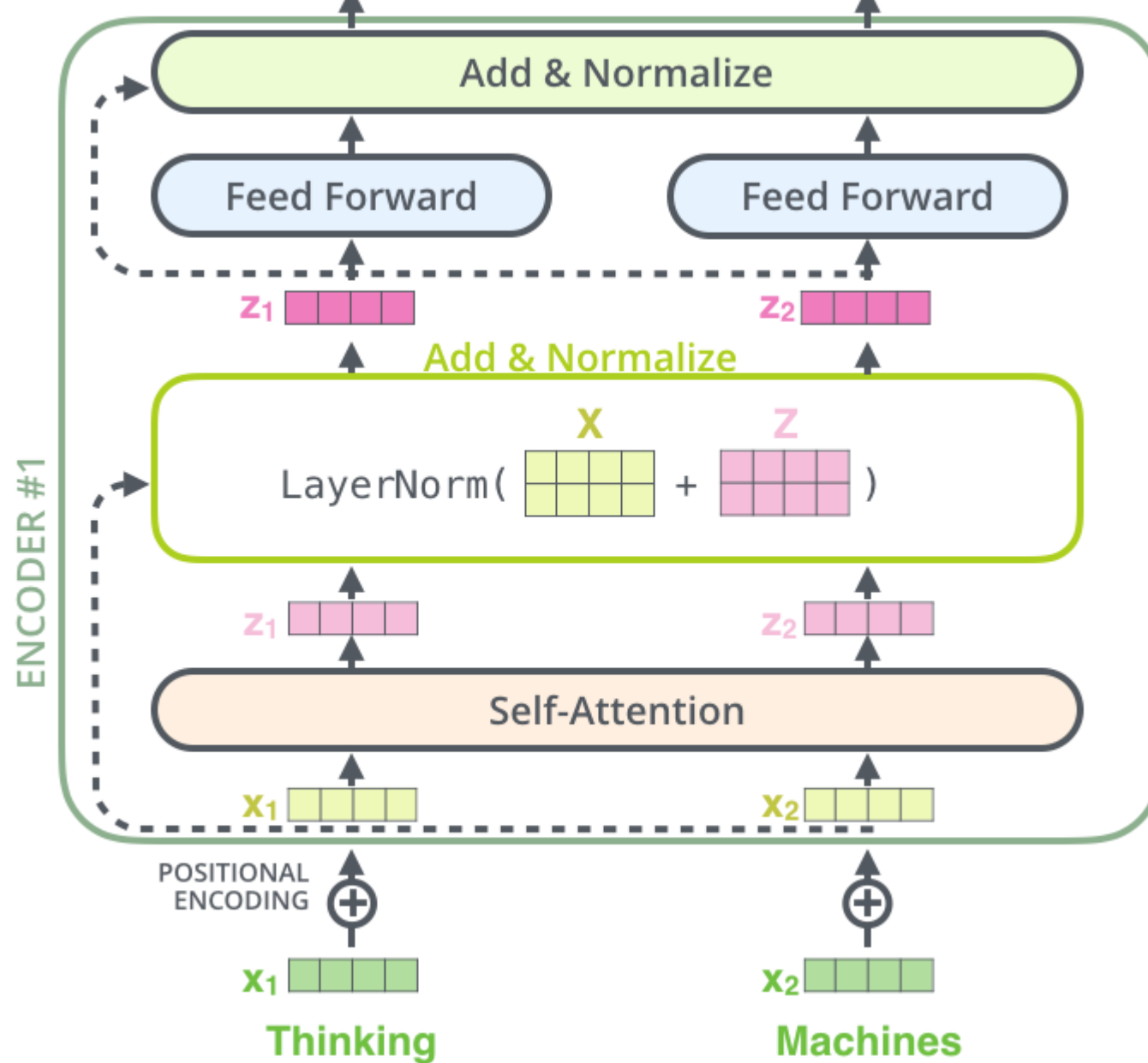
Z

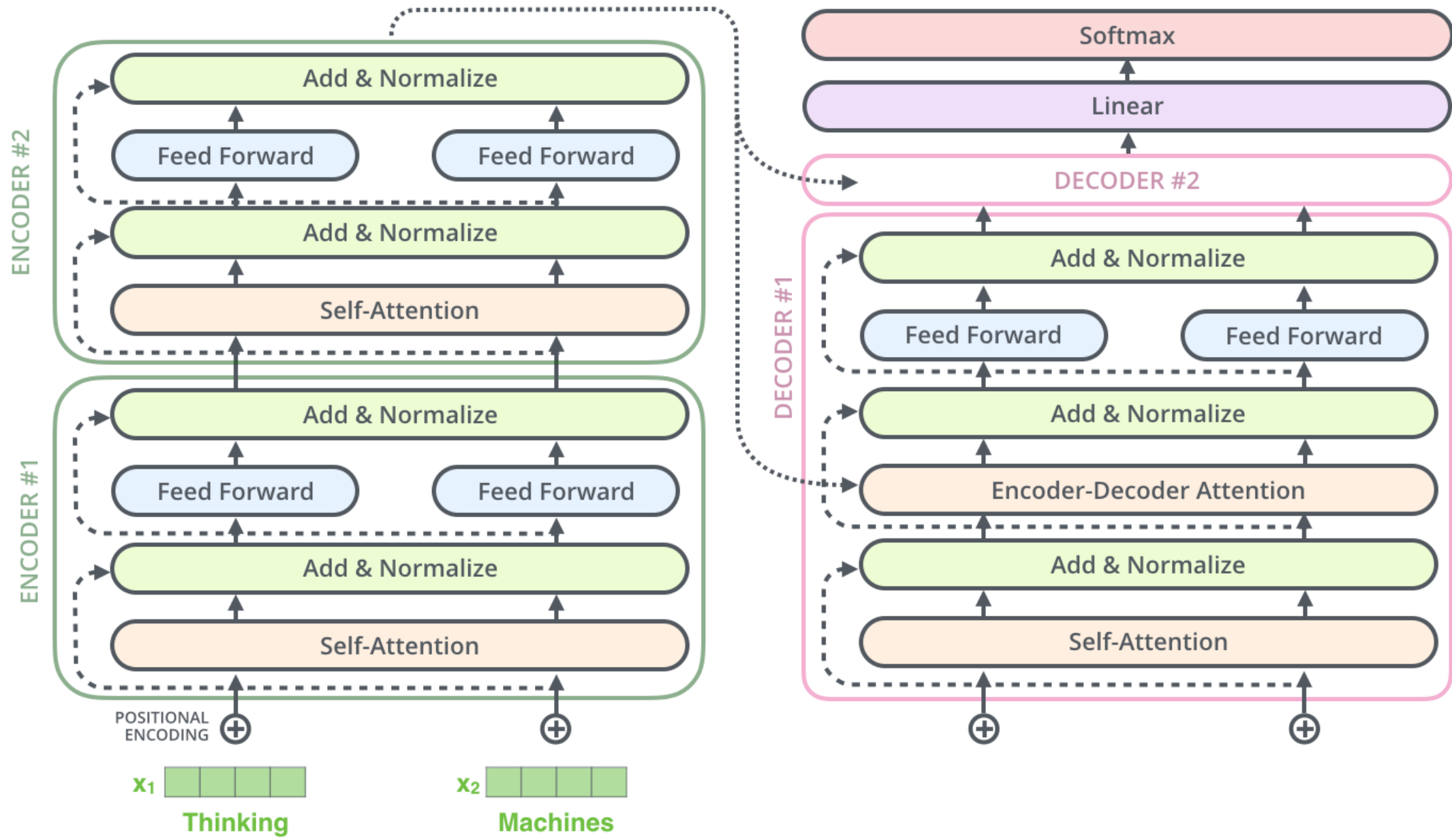


* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one





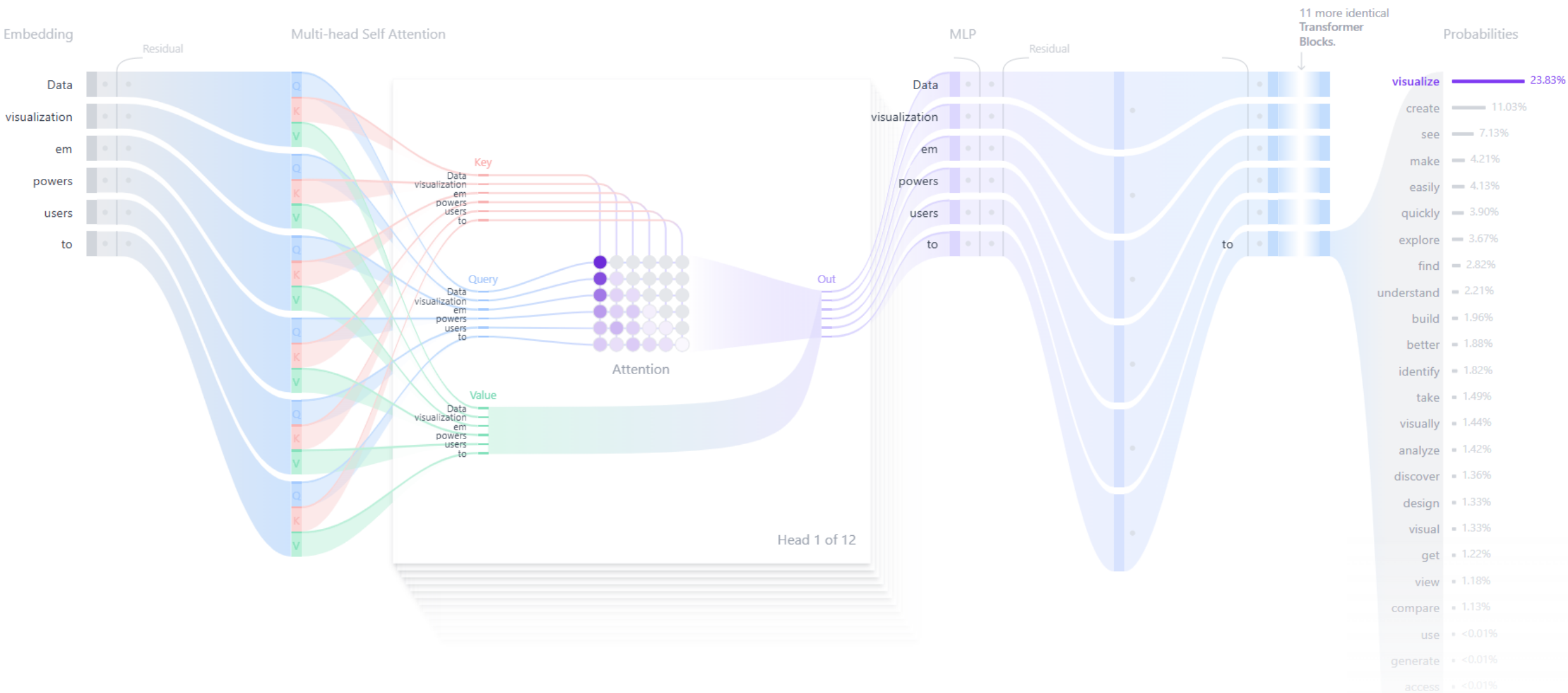




$$PE(pos, 2i) = \sin(pos/10000^{2i/d})$$

$$PE(pos, 2i + 1) = \cos(pos/10000^{2i/d})$$

Try the examples while GPT-2 model is being downloaded (600MB)



Vision Transformer ViT

