

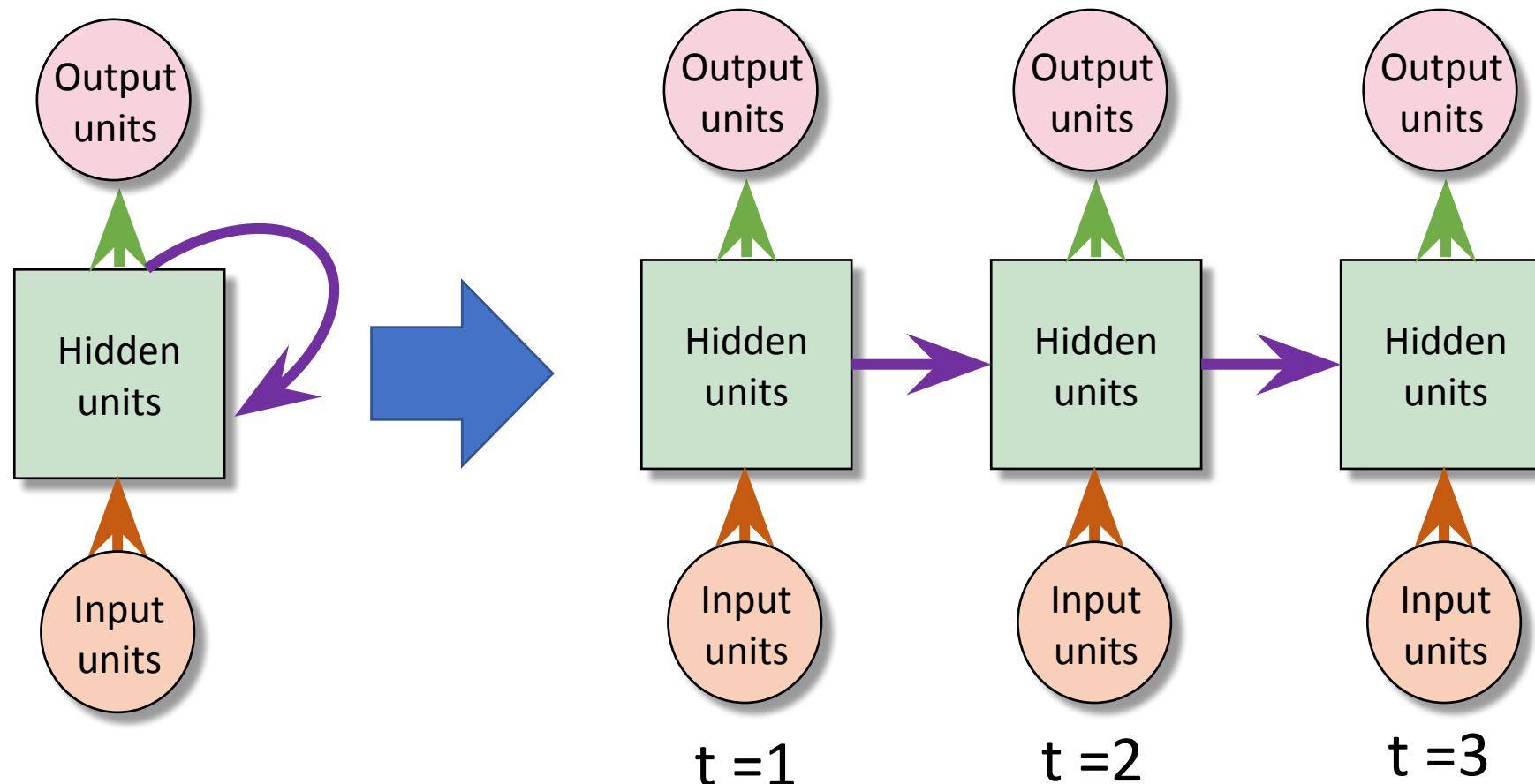
Variants of Recurrent Neural Network

Contents

- Recap! Recurrent Neural Network (RNN)
- Bidirectional RNN
- Long Short-Term Memories Cell (LSTM)
- Gate Recurrent Units (GRU)
- Attention Mechanism
- Self-Attention
- Transformer
- Summary

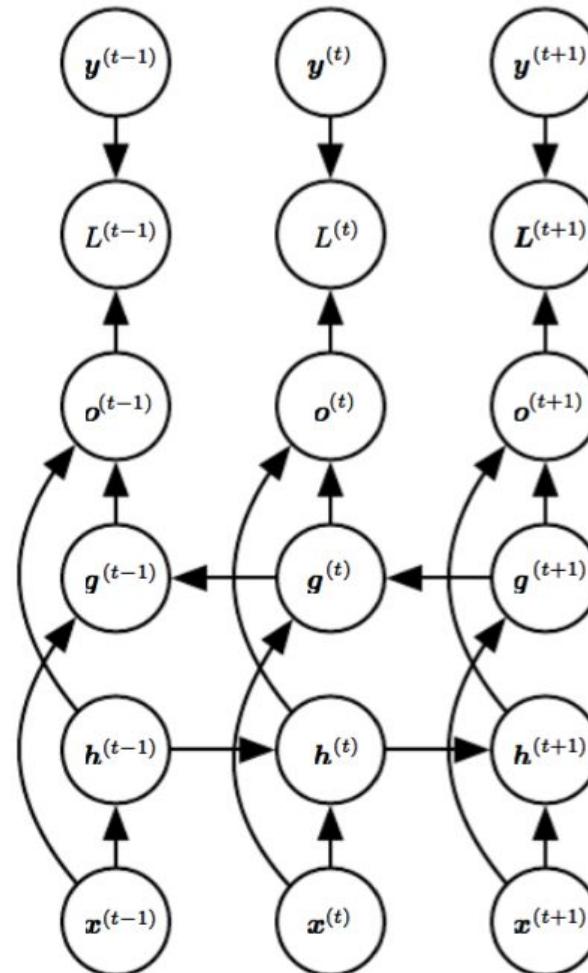
Recap! Recurrent Neural Network (RNN)

- Unfold the RNN (e.g., t=3)



Bidirectional RNNs

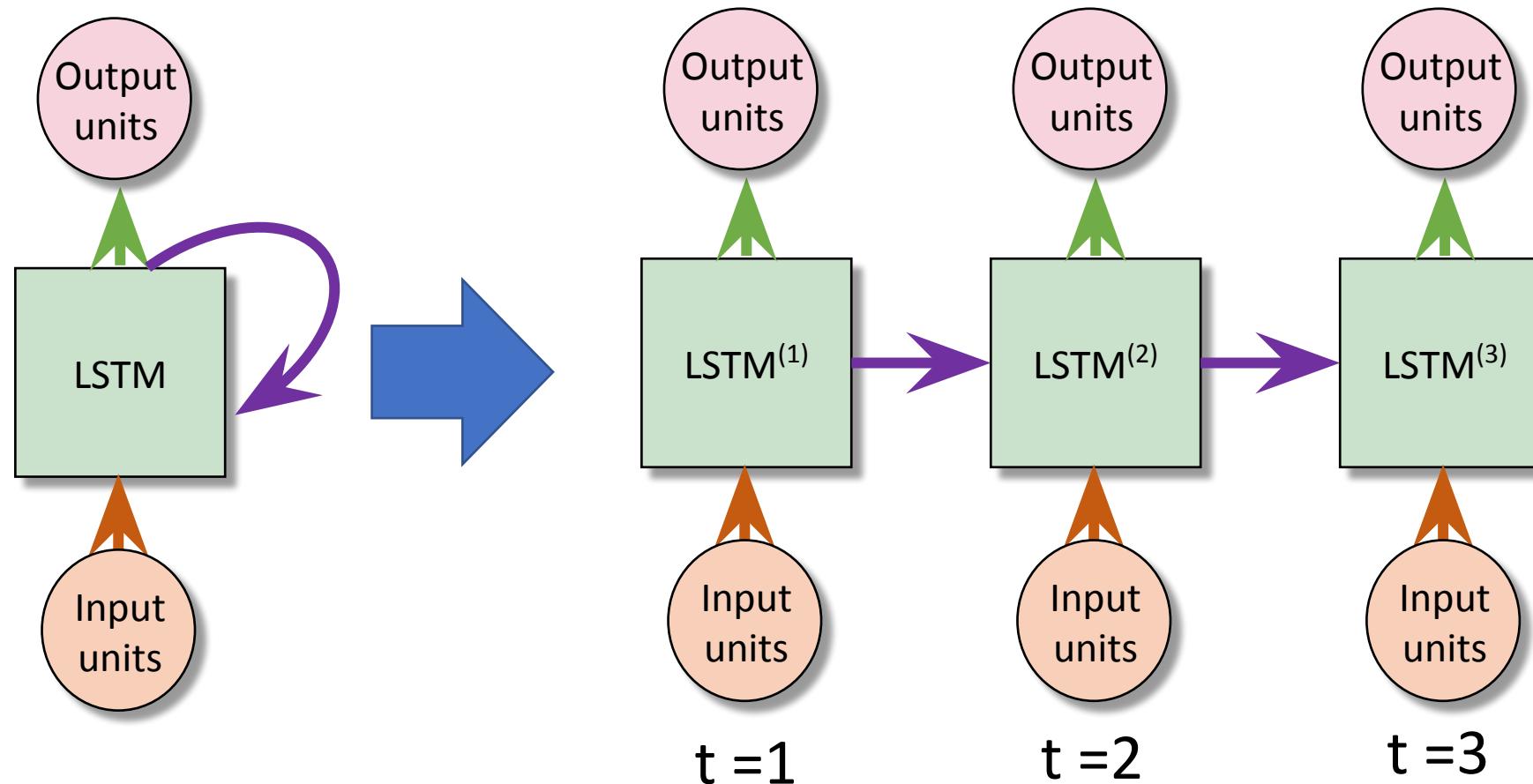
- We can think it combines **two RNN**
- One moves forward through time beginning from the start of the sequence while **another RNN** that **moves backward** through time beginning from the end of the sequence.



LSTM – Introduce a Memory Cell

- LSTM was proposed as a **solution to the vanishing gradients** problem.
- On step t, there is a **hidden state** and a **cell state**
 - The cell stores **long-term information**
 - The LSTM can **erase**, **write** and **read** information from the cell
- The selection of which information is **erased/written/read** is controlled by **three corresponding gates**
 - On each timestep, each element of the gates can be **open** (1), **closed** (0), or somewhere **in-between**.
 - The gates are **dynamic**: their value is computed based on the current context

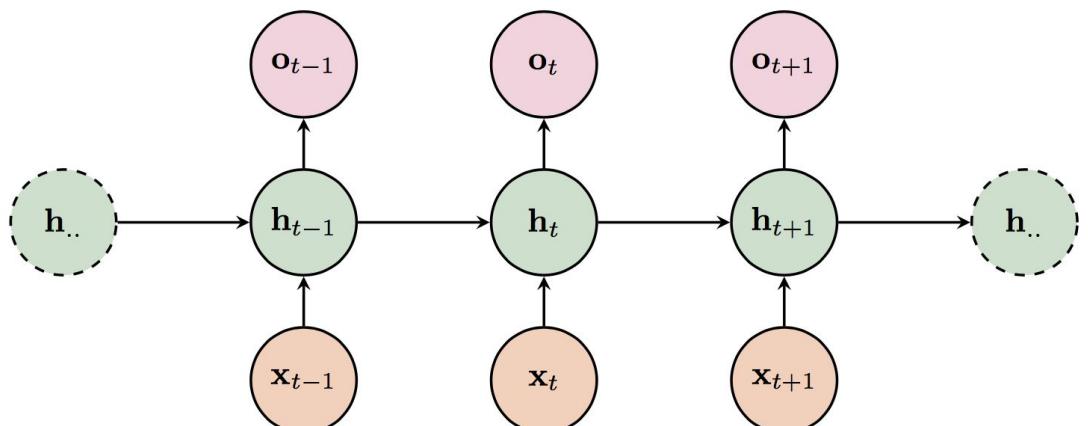
LSTM – Introduce a Memory Cell



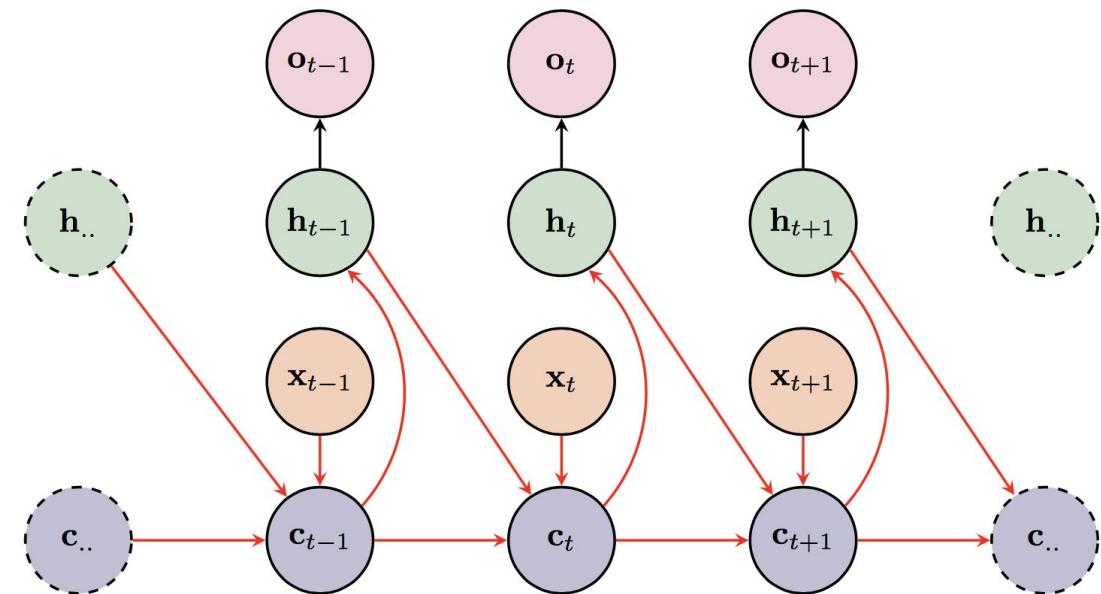
LSTM cell contains a **hidden state** and a **cell state (Memory cell)**

RNN vs LSTM

Graphical Structure

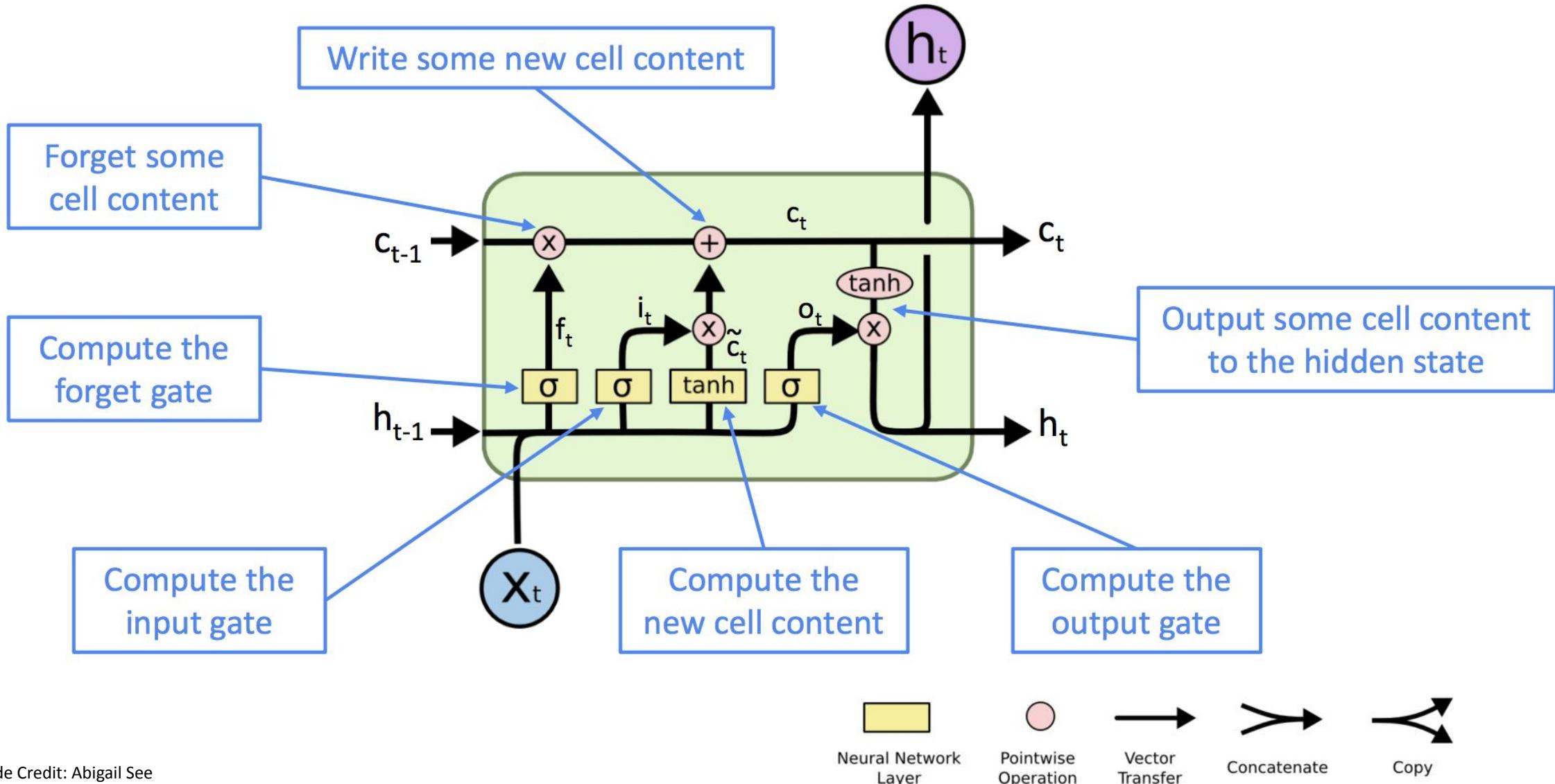


RNN



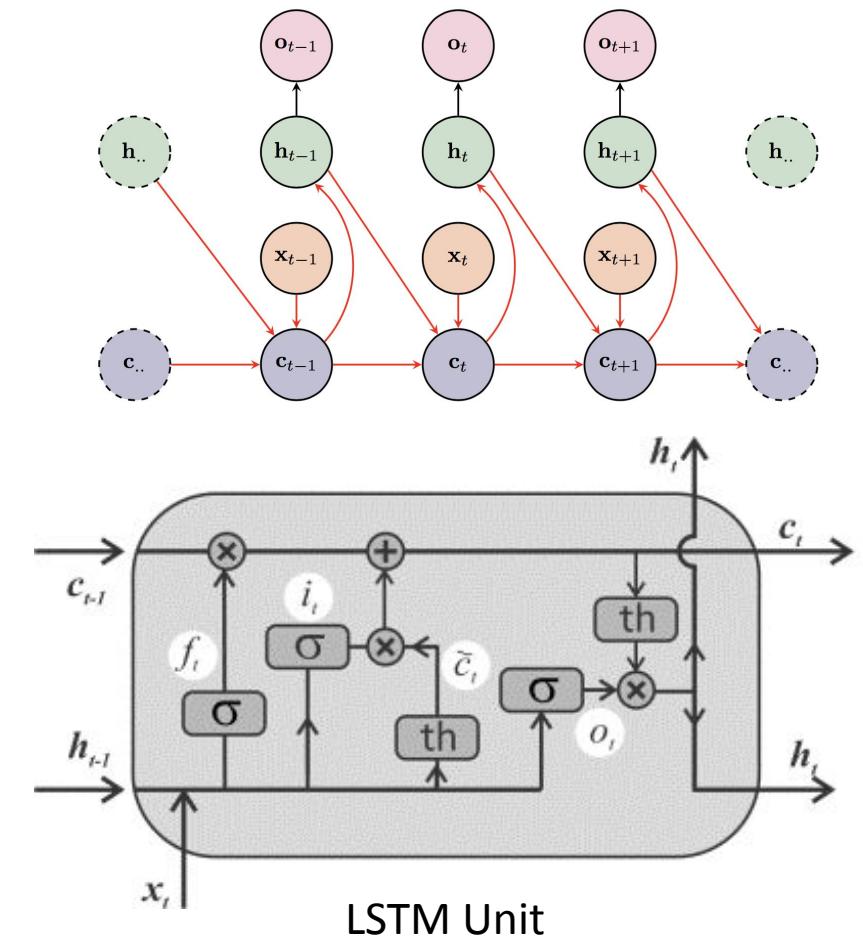
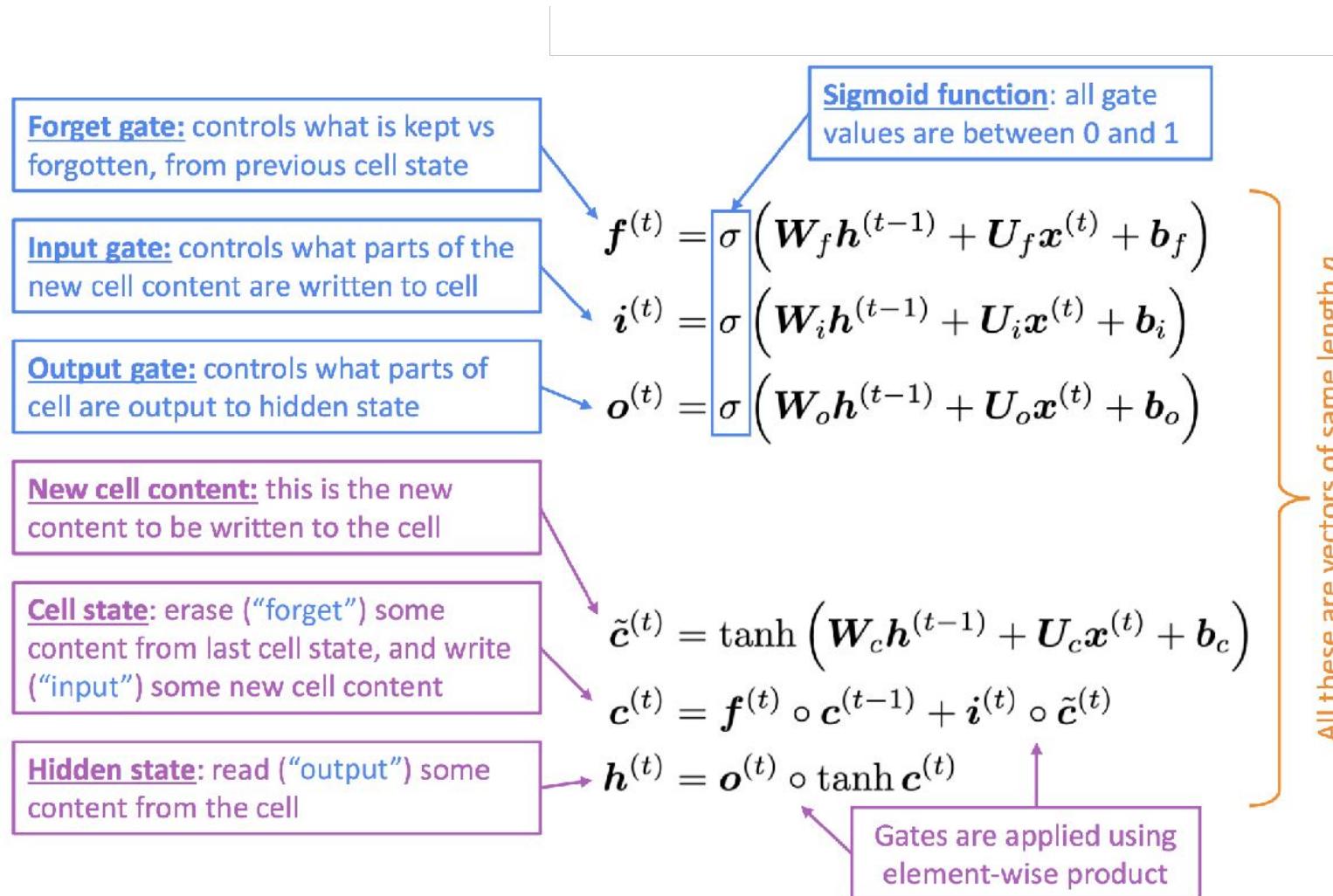
LSTM

LSTM – Memory Cell

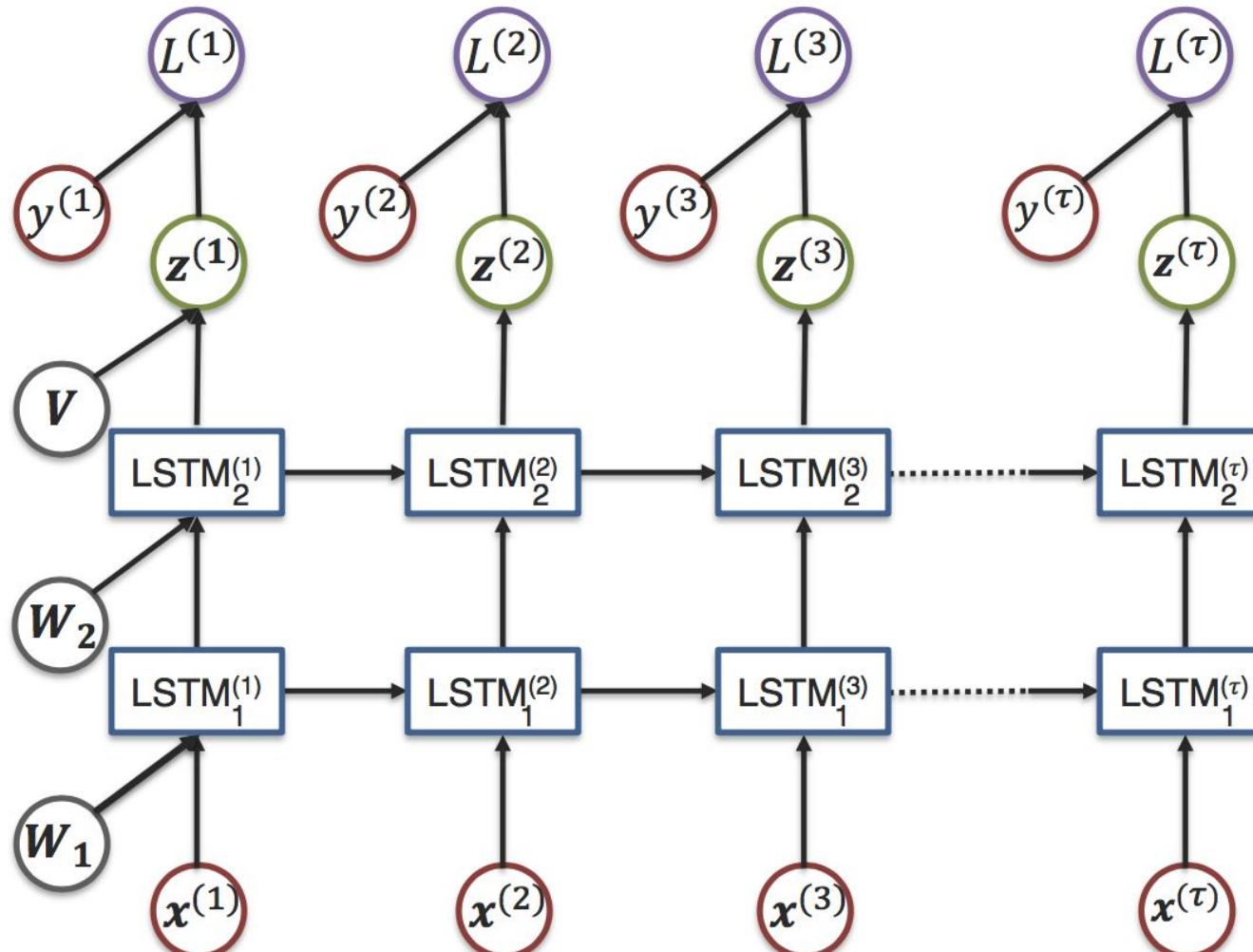


LSTM – Memory Cell

On time t : We have a sequence of inputs $x^{(t)}$, and we will compute a sequence of hidden states $h^{(t)}$ and cell states.



Deep LSTM



Gated Recurrent Units (GRUs)

- Proposed by Cho et al. in 2014 as a simpler alternative to the LSTM.
- On each timestep t : we have input and hidden state (no cell state).

Update gate: controls what parts of hidden state are updated vs preserved

Reset gate: controls what parts of previous hidden state are used to compute new content

New hidden state content: reset gate selects useful parts of prev hidden state. Use this and current input to compute new hidden content.

Hidden state: update gate simultaneously controls what is kept from previous hidden state, and what is updated to new hidden state content

$$u^{(t)} = \sigma(W_u h^{(t-1)} + U_u x^{(t)} + b_u)$$

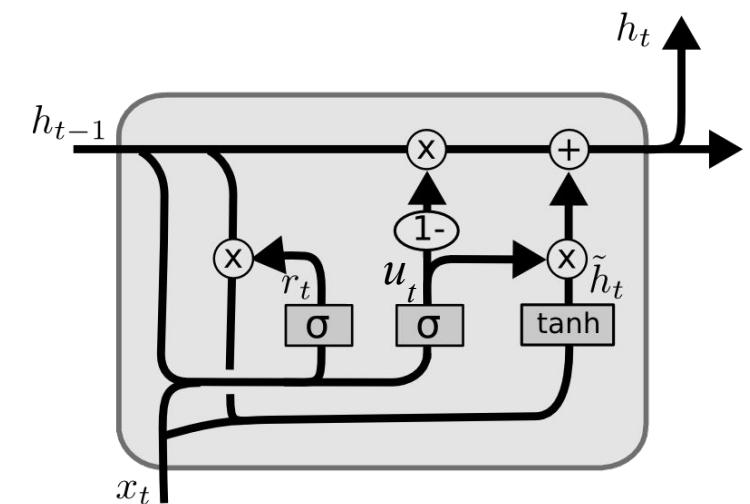
$$r^{(t)} = \sigma(W_r h^{(t-1)} + U_r x^{(t)} + b_r)$$

$$\tilde{h}^{(t)} = \tanh(W_h(r^{(t)} \circ h^{(t-1)}) + U_h x^{(t)} + b_h)$$

$$h^{(t)} = (1 - u^{(t)}) \circ h^{(t-1)} + u^{(t)} \circ \tilde{h}^{(t)}$$

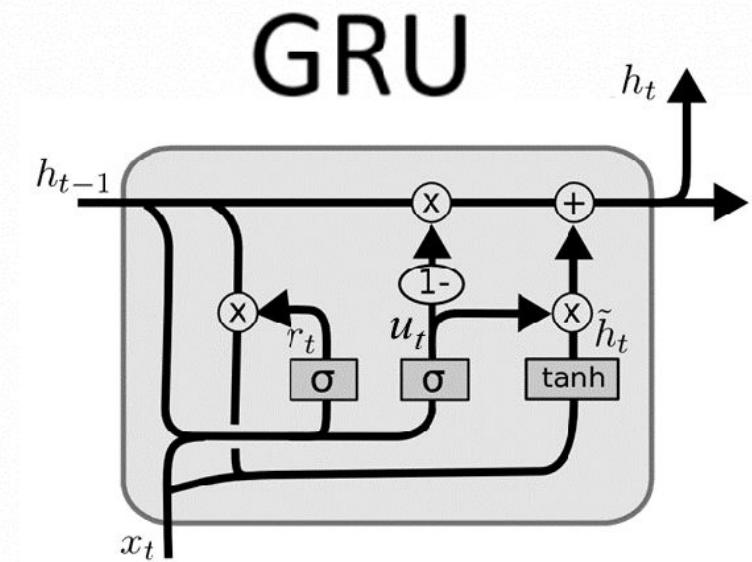
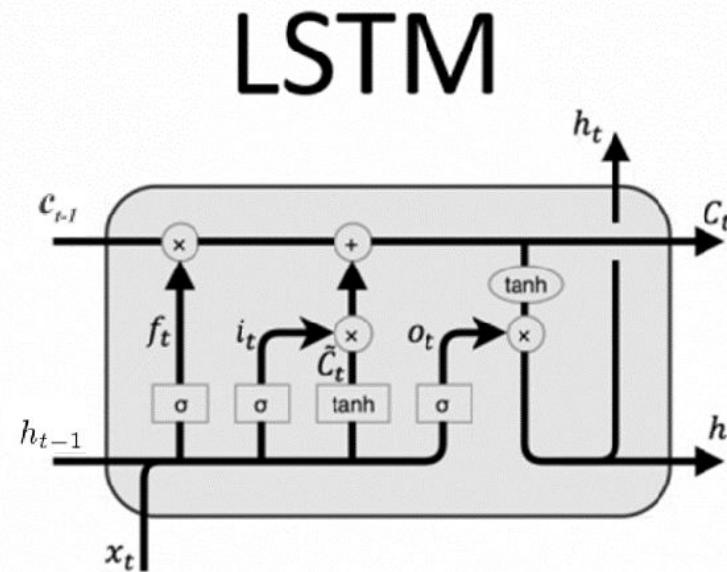
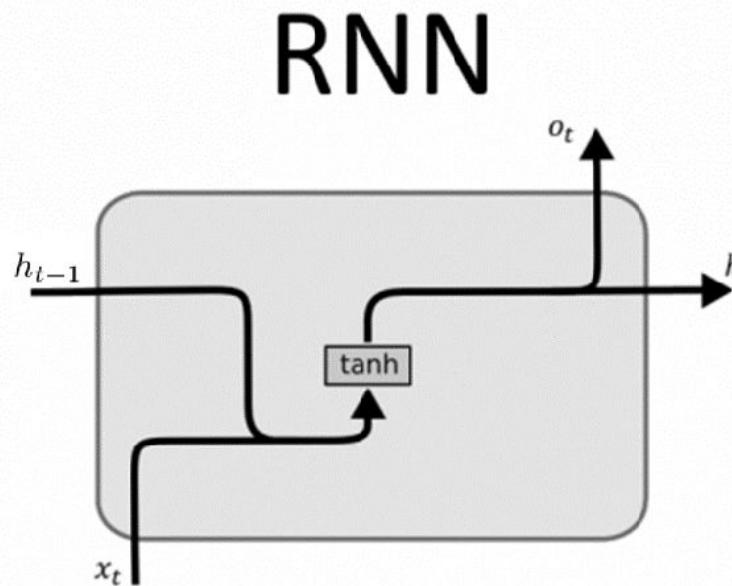
How does this solve vanishing gradient?

Like LSTM, GRU makes it easier to retain info long-term (e.g. by setting update gate to 0)



GRU Unit

RNN, LSTM and GRU



LSTM vs GRU

- Researchers have proposed many gated RNN variants, but LSTM and GRU are the most widely-used.
- The biggest difference is that GRU is quicker to compute and has fewer parameters
- There is no conclusive evidence that one consistently performs better than the other
- LSTM is a good default choice (especially if your data has particularly long dependencies, or you have lots of training data)
- Rule of thumb: start with LSTM, but switch to GRU if you want something more efficient

Transformers: Why new architecture?

- **Challenges with RNN**

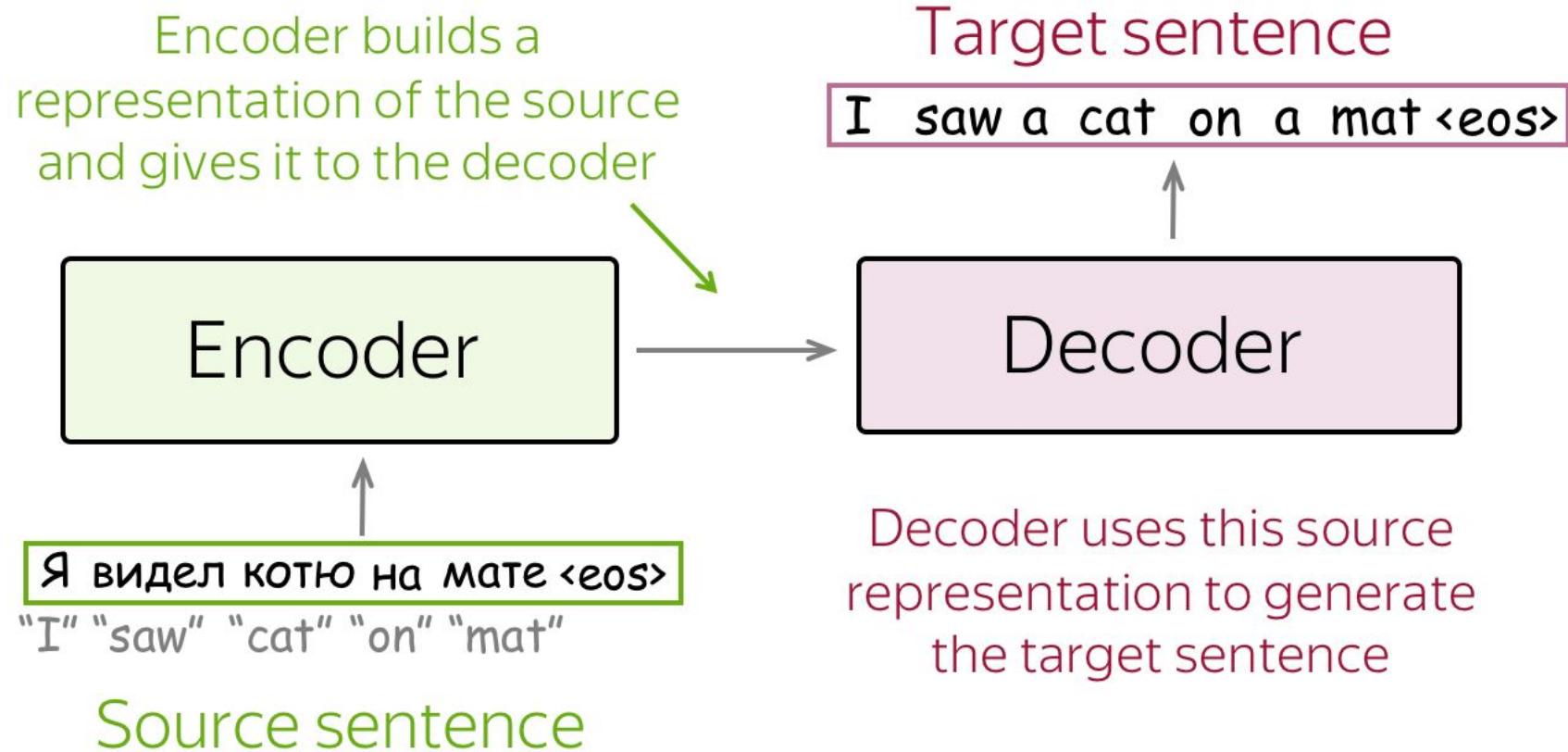
- Long range dependencies
- Gradient vanishing and explosion
- Large number of training steps
- Recurrence prevents parallel computation

- **Transformer Networks**

- Facilitate long range dependencies
- No gradient vanishing and explosion
- Fewer training steps
- No recurrence that facilitate parallel computation

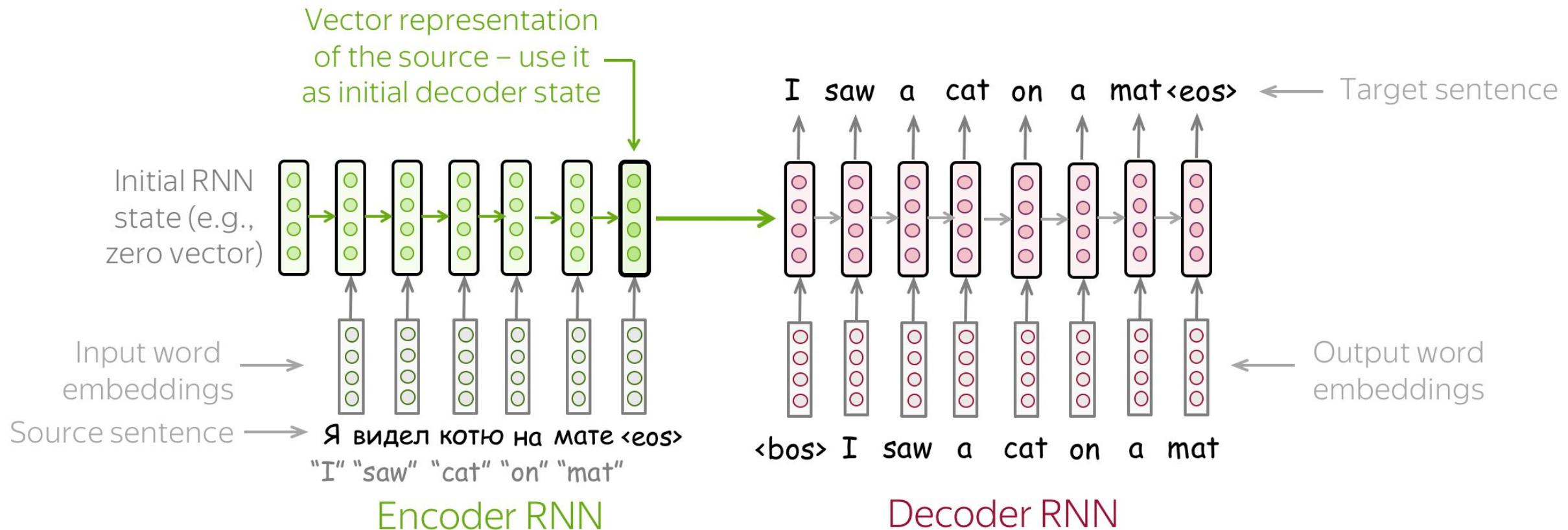
Transformer Networks

First understand Sequence to Sequence Model



The Simplest Model

Two RNNs for Encoder and Decoder



Transformer Networks

I arrived at the **bank** after crossing thestreet? ...river?

What does **bank** mean in this sentence?



I've no idea: let's wait
until I read the end

RNNs

$O(N)$ steps to process a
sentence with length N

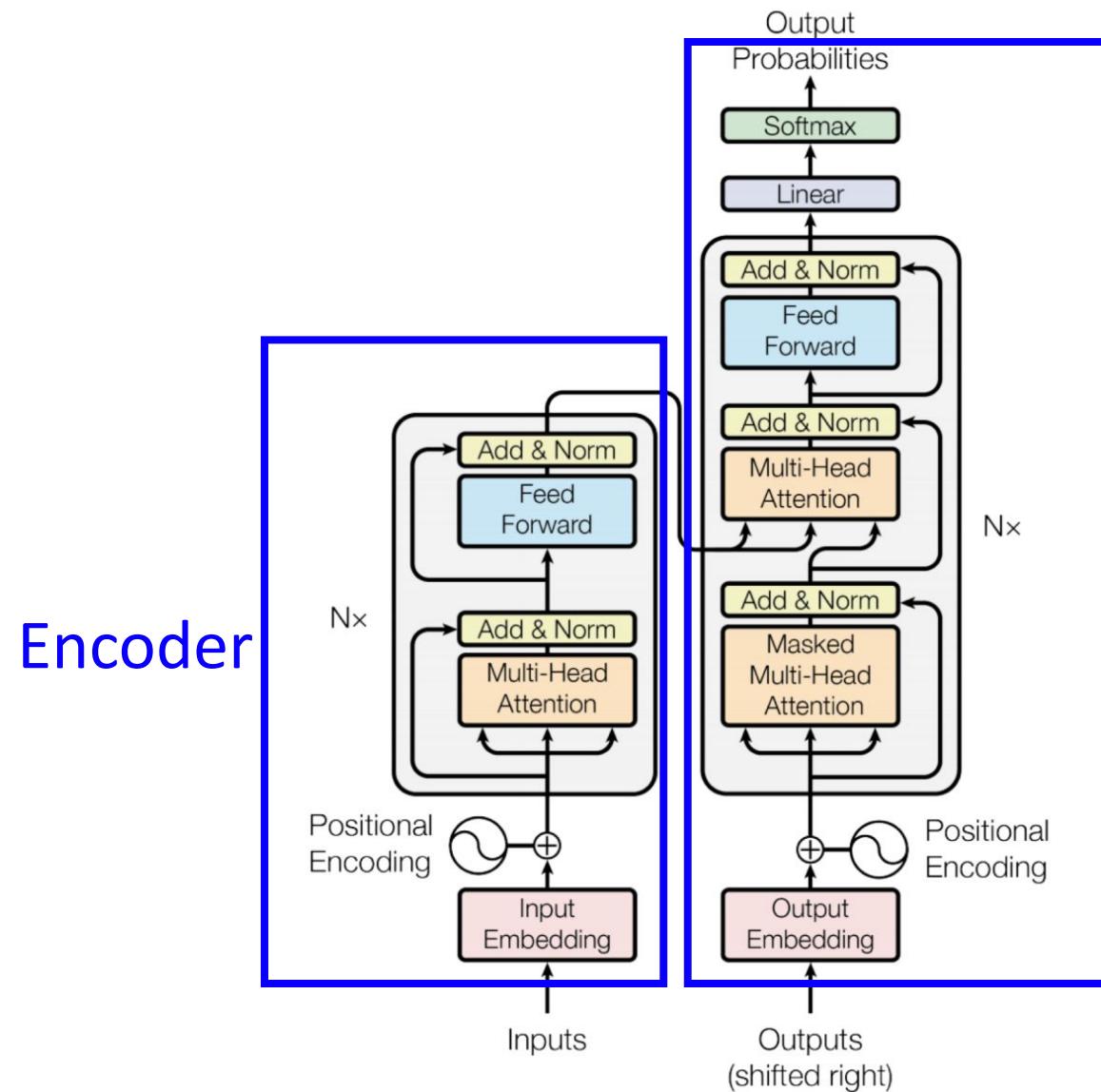


I don't need to wait - I
see all words at once!

Transformer

Constant number of steps
to process any sentence

Transformer Basics



- The basic building blocks of transformer networks

Self-Attention layers!

Decoder

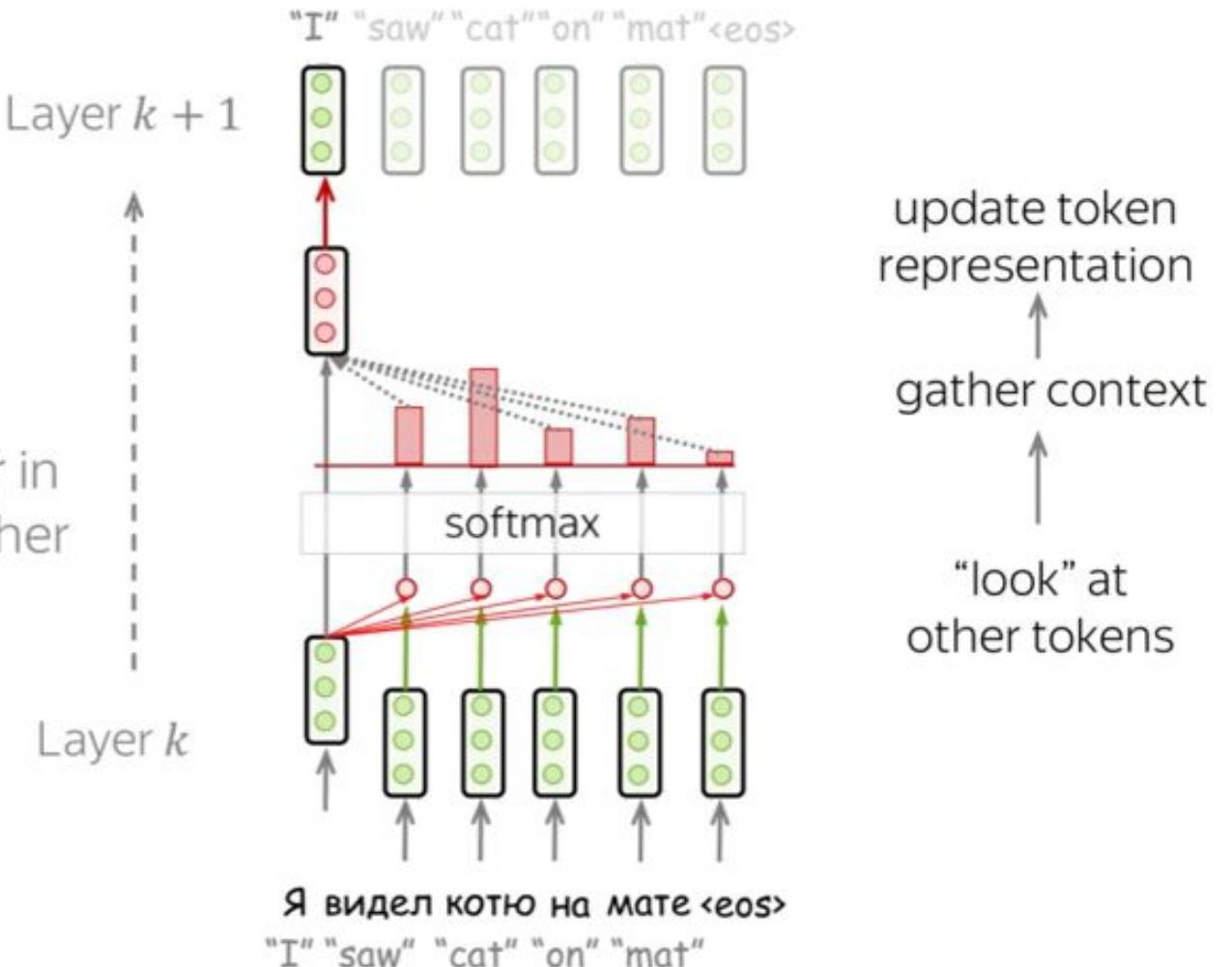
Self-Attention

$$\text{Attention}(q, k, v) = \text{softmax}\left(\frac{qk^T}{\sqrt{d_k}}\right)v$$

from \swarrow to \nearrow

Attention weights
vector dimensionality of K, V

Tokens try to understand themselves better in context of each other



Self-Attention

$$\text{Attention}(q, k, v) = \text{softmax}\left(\frac{qk^T}{\sqrt{d_k}}\right)v$$

from to
Attention weights
vector dimensionality of K, V

Each vector receives three representations (“roles”)

$$[W_Q] \times \begin{bmatrix} \text{I} \\ \text{saw} \\ \text{cat} \\ \text{on} \\ \text{mat} \\ \text{<eos>} \end{bmatrix} = \begin{bmatrix} \text{I} \\ \text{saw} \\ \text{cat} \\ \text{on} \\ \text{mat} \\ \text{<eos>} \end{bmatrix}$$

Query: vector from which the attention is looking

“Hey there, do you have this information?”

$$[W_K] \times \begin{bmatrix} \text{I} \\ \text{saw} \\ \text{cat} \\ \text{on} \\ \text{mat} \\ \text{<eos>} \end{bmatrix} = \begin{bmatrix} \text{I} \\ \text{saw} \\ \text{cat} \\ \text{on} \\ \text{mat} \\ \text{<eos>} \end{bmatrix}$$

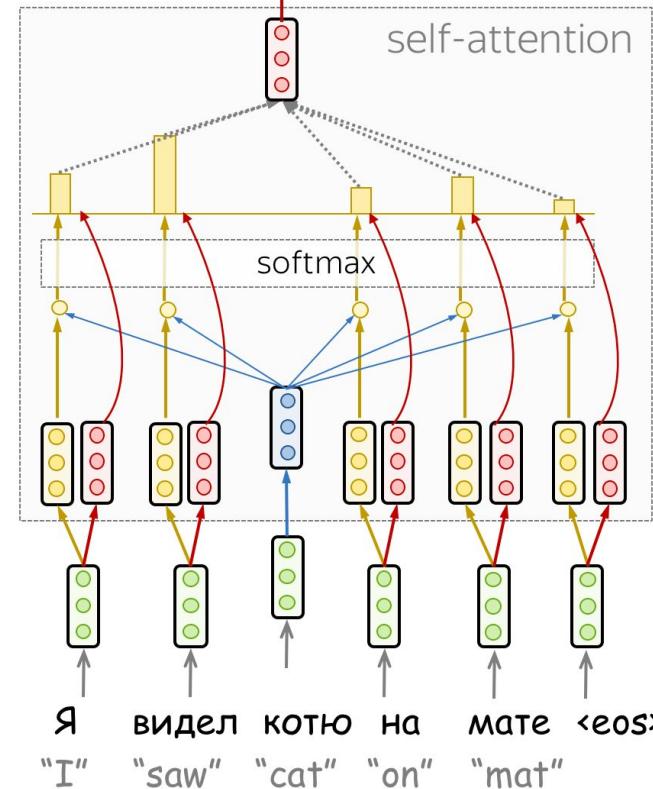
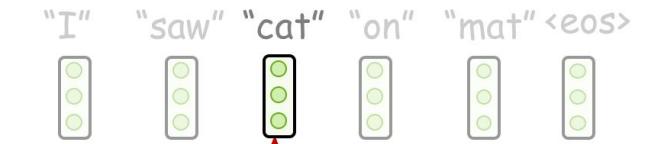
Key: vector at which the query looks to compute weights

“Hi, I have this information – give me a large weight!”

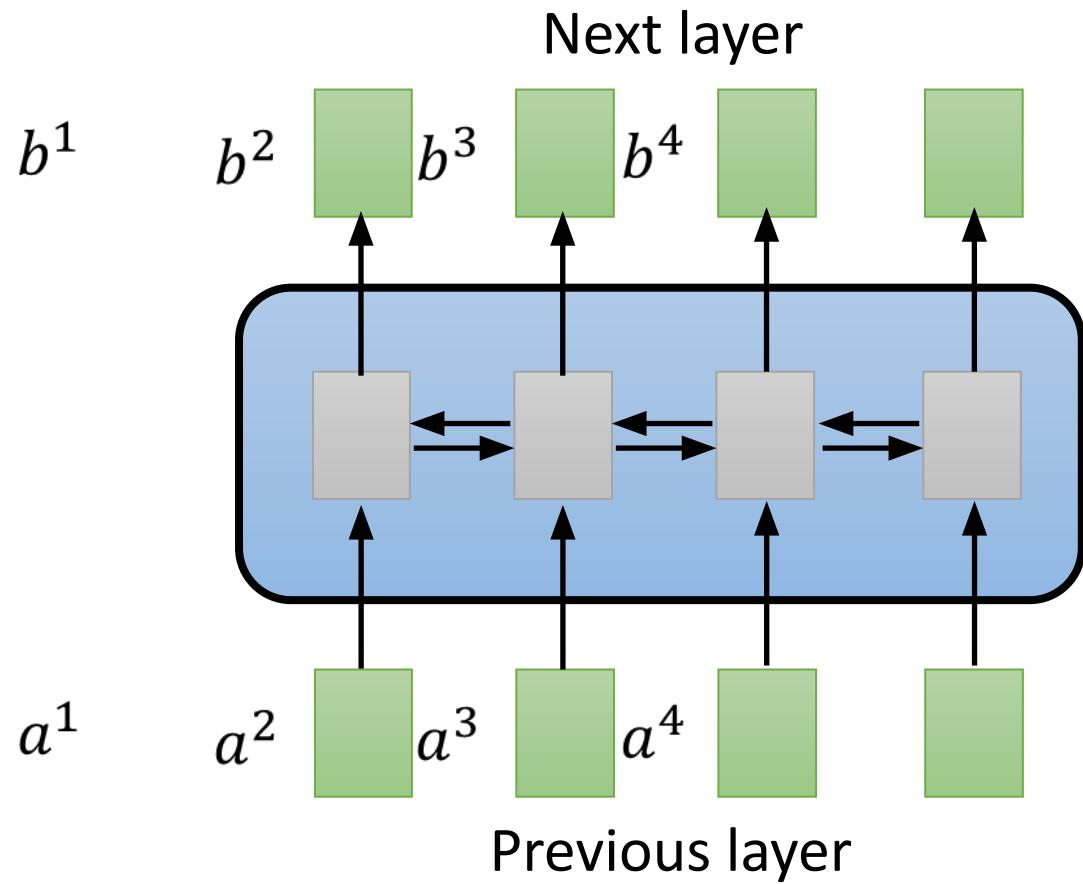
$$[W_V] \times \begin{bmatrix} \text{I} \\ \text{saw} \\ \text{cat} \\ \text{on} \\ \text{mat} \\ \text{<eos>} \end{bmatrix} = \begin{bmatrix} \text{I} \\ \text{saw} \\ \text{cat} \\ \text{on} \\ \text{mat} \\ \text{<eos>} \end{bmatrix}$$

Value: their weighted sum is attention output

“Here’s the information I have!”



Self-Attention

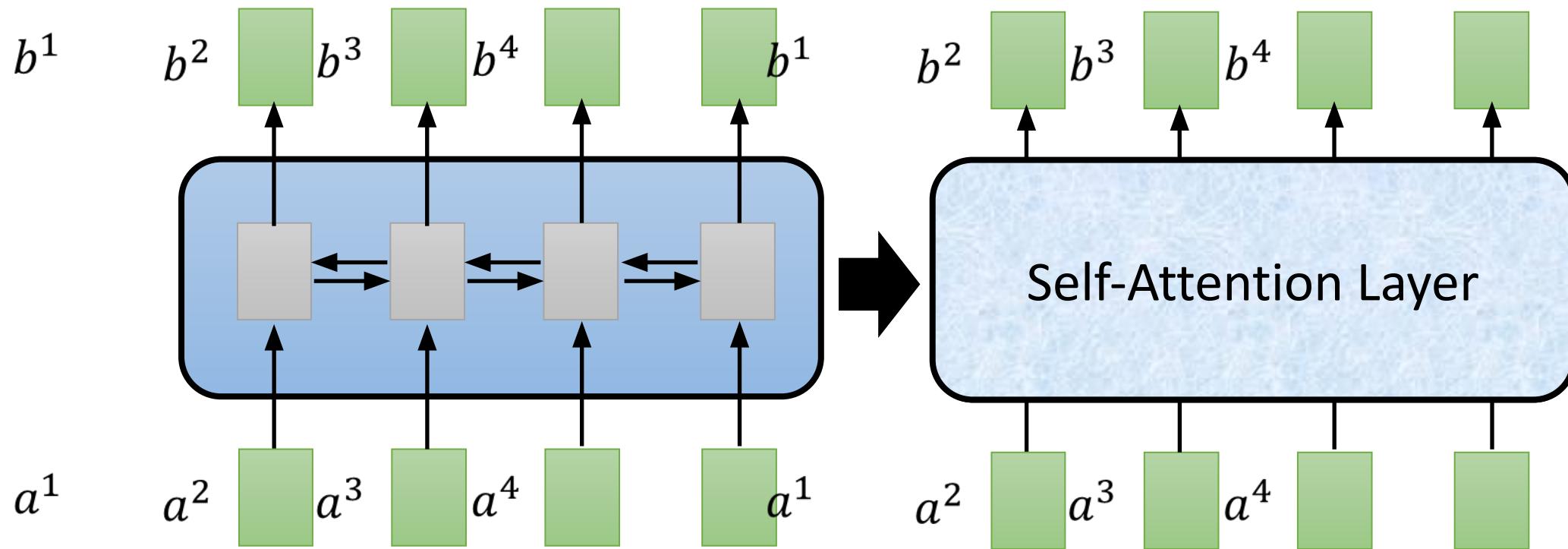


Hard to parallelize

Self-Attention

b^i is obtained based on the whole input sequence.

b^1, b^2, b^3, b^4 can be parallelly computed.



You can try to replace anything that has been done by RNN with self-attention.

Self-Attention

Each input token in **self-attention** receives three representations corresponding to the roles it can play:

- query (q) – asking for information
- key(k) – saying that it has some information
- value(v) – giving the information.

$$\text{Attention}(q, k, v) = \frac{\text{softmax}\left(\frac{qk^T}{\sqrt{d_k}}\right)v}{\text{Attention weights}}$$

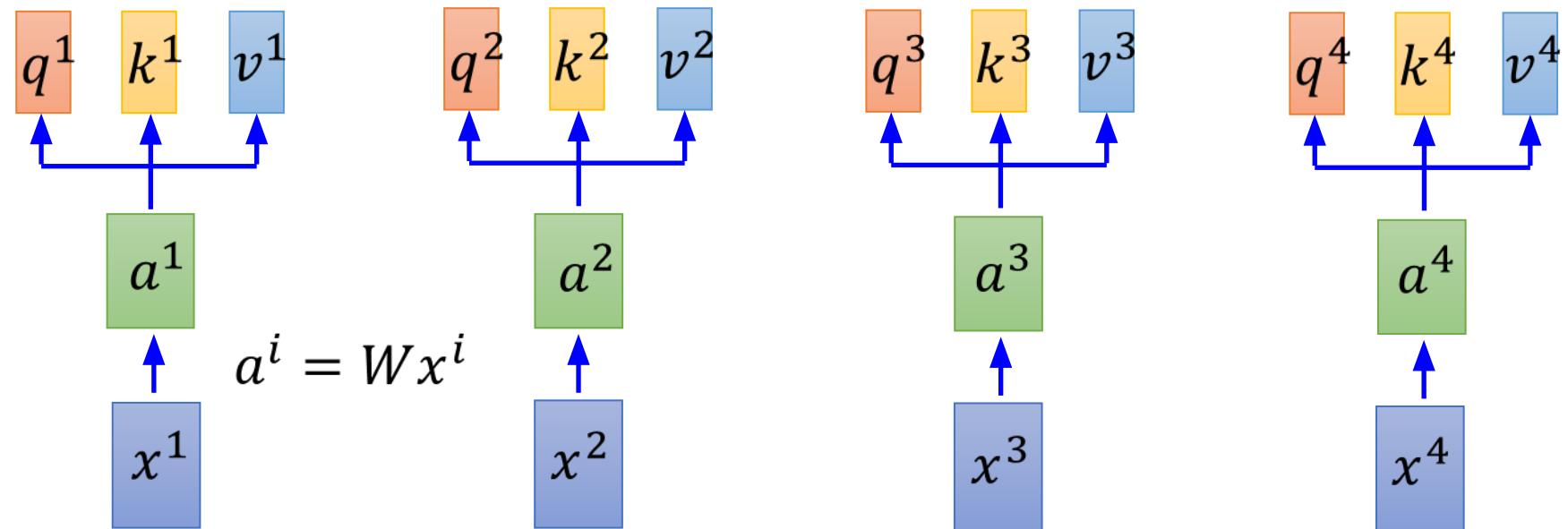
from to

vector dimensionality of K, V

$$q^i = W^q a^i$$

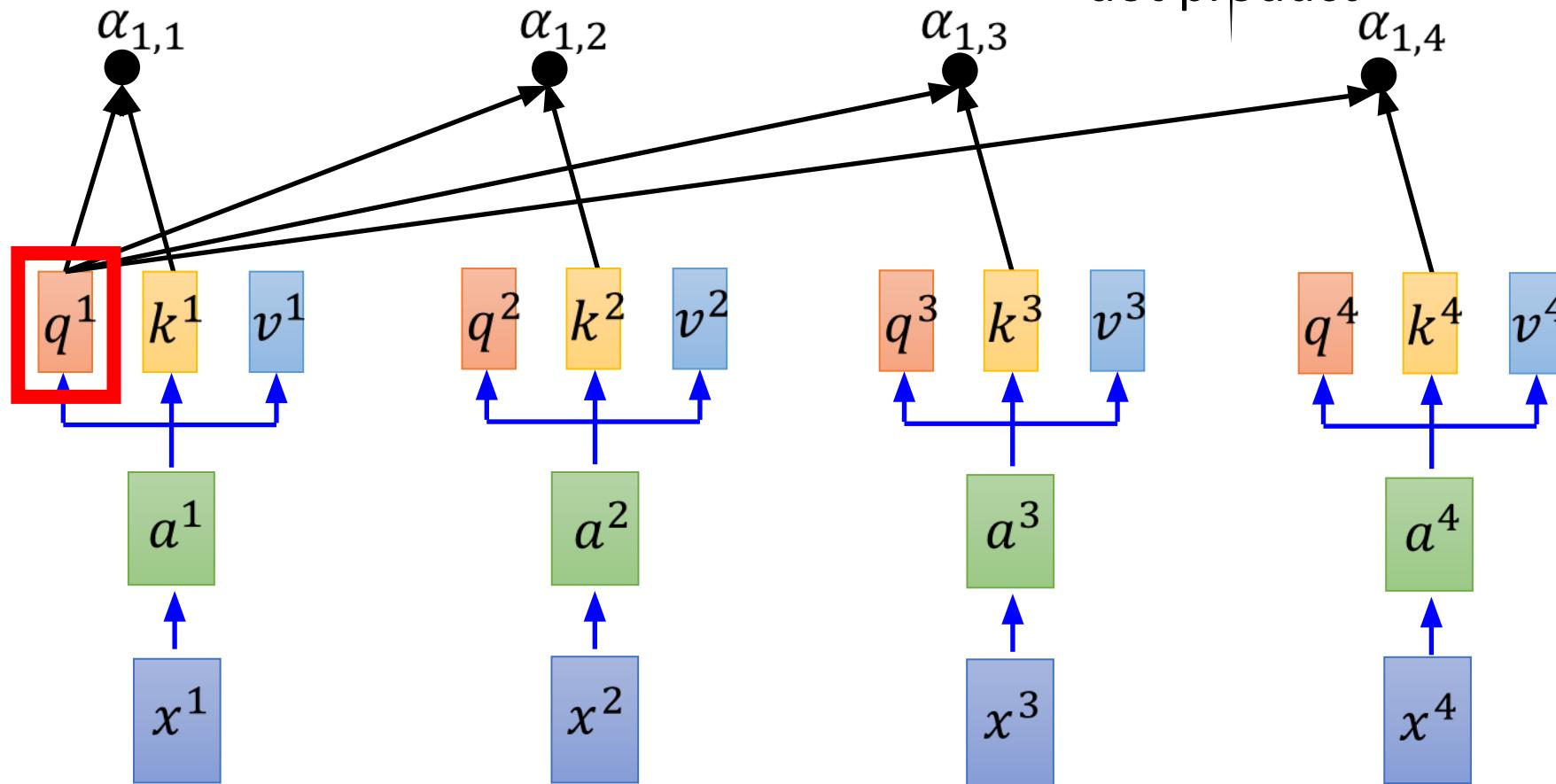
$$k^i = W^k a^i$$

$$v^i = W^v a^i$$



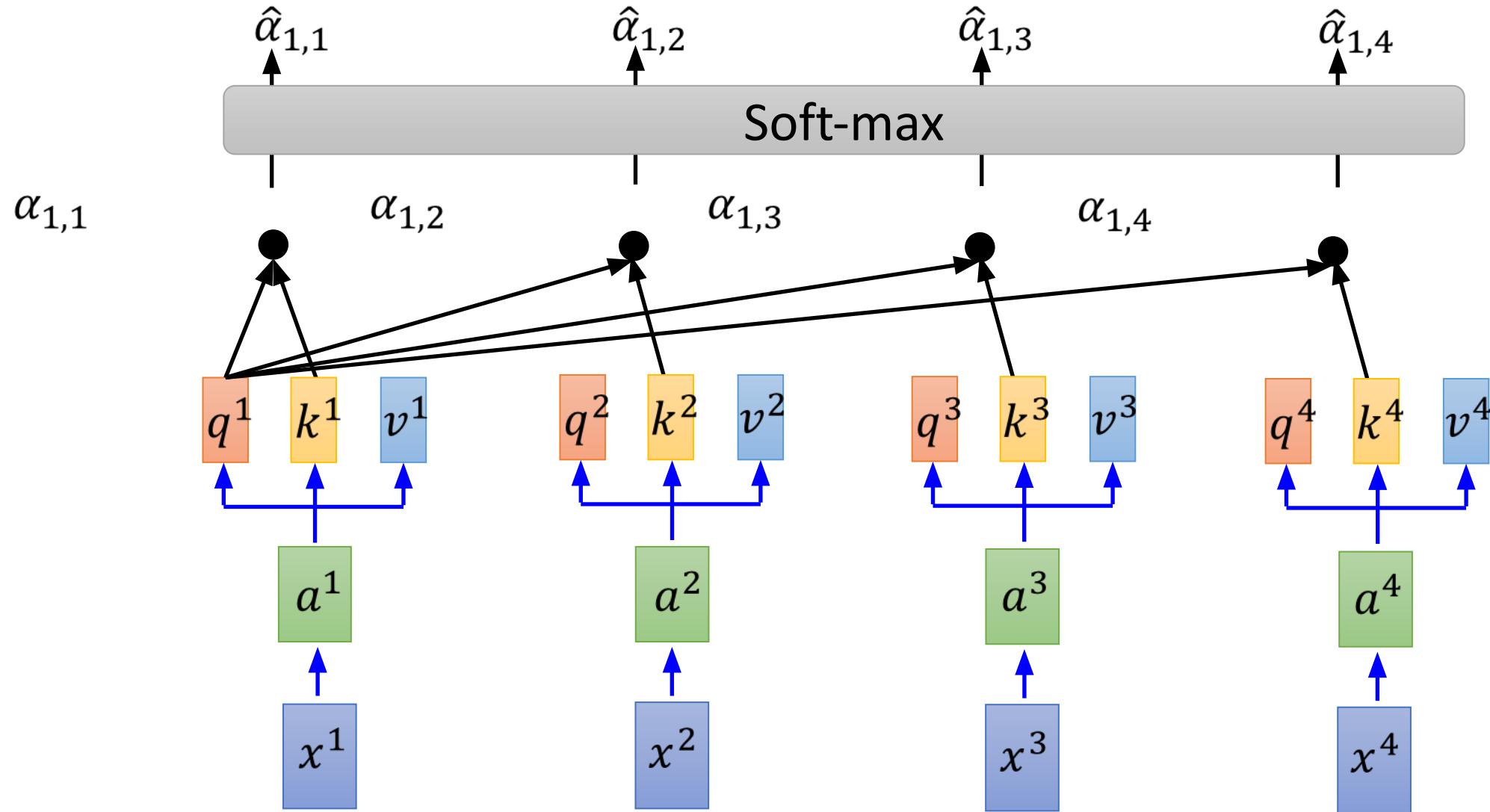
Self-Attention

Scaled Dot-Product Attention: $\alpha_{1,i} = \underbrace{q^1 \cdot k^i}_{\text{dot product}} / \sqrt{d}$



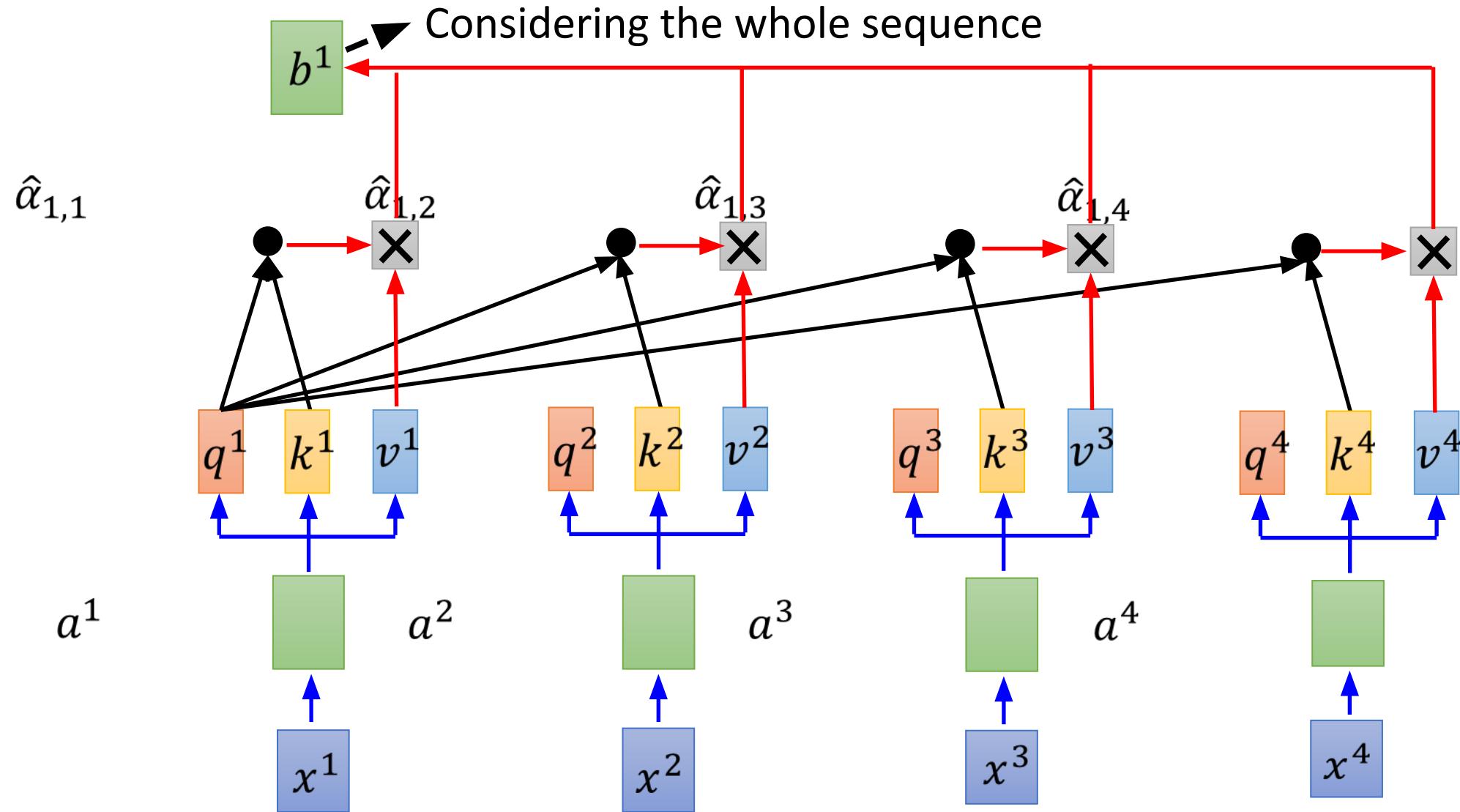
Self-Attention

$$\hat{\alpha}_{1,i} = \exp(\alpha_{1,i}) / \sum_j \exp(\alpha_{1,j})$$



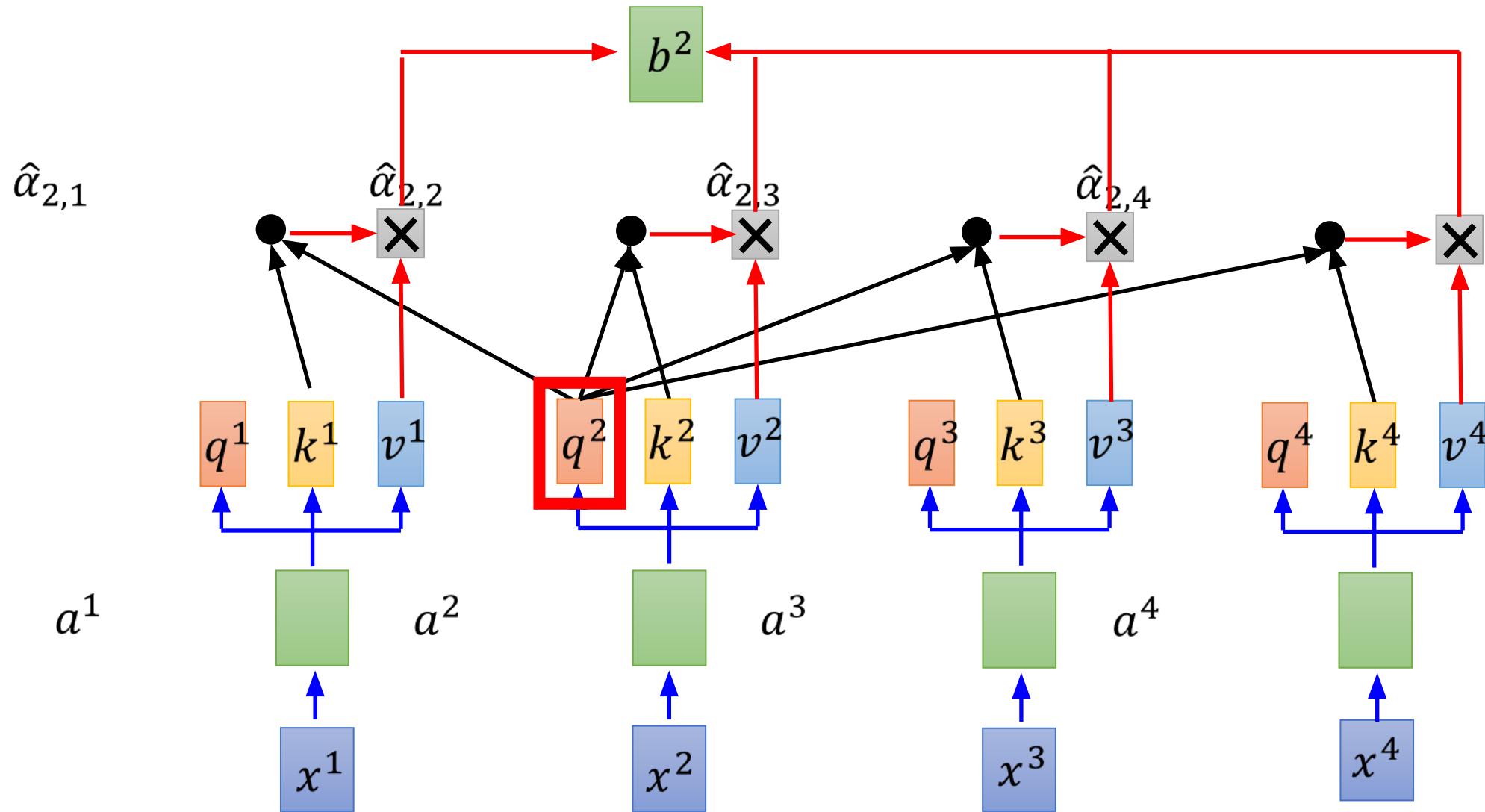
Self-Attention

$$b^1 = \sum_i \hat{\alpha}_{1,i} v^i$$



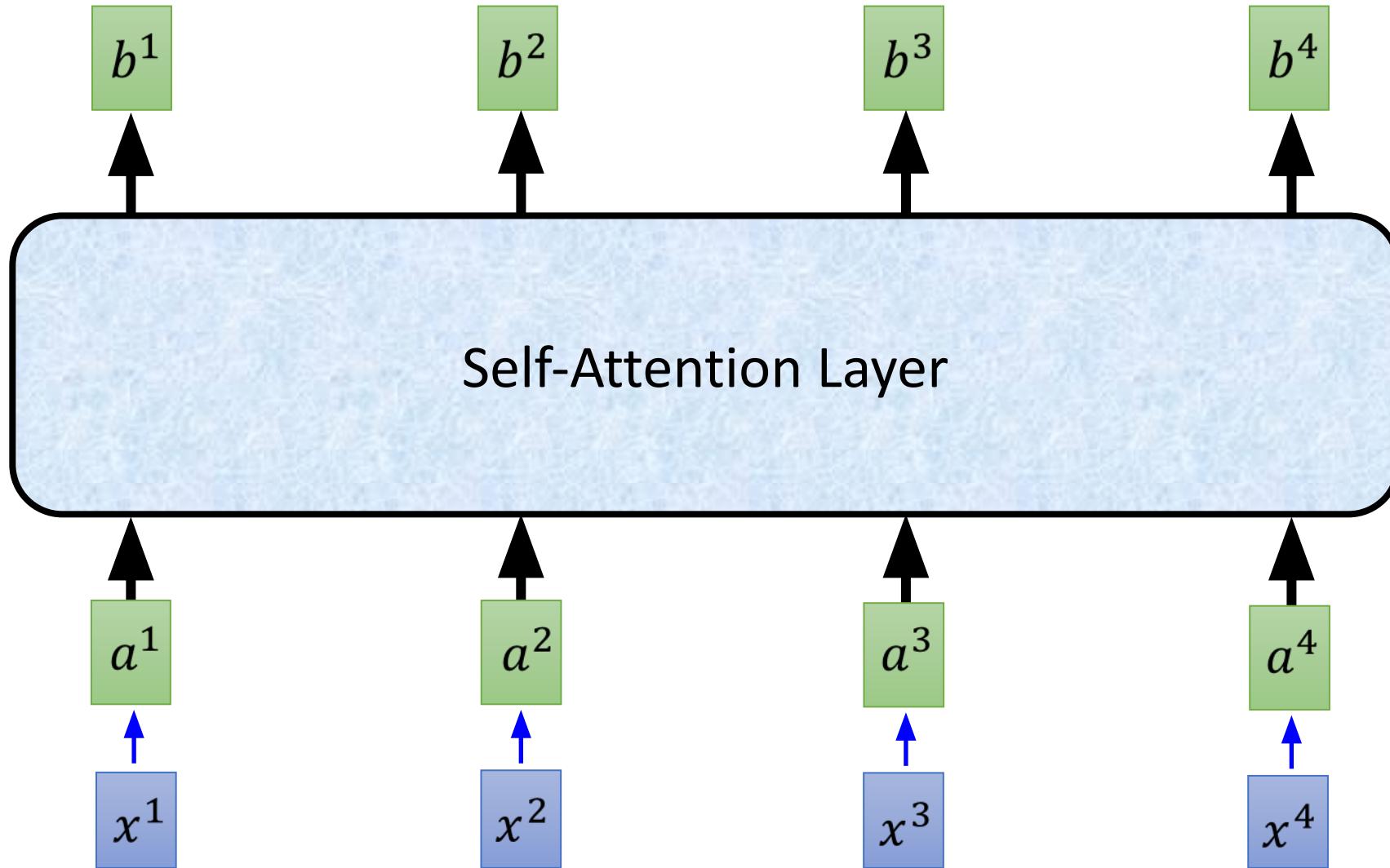
Self-Attention

$$b^2 = \sum_i \hat{\alpha}_{2,i} v^i$$



Self-Attention

b^1, b^2, b^3, b^4 can be parallelly computed.



Self-Attention

$$q^i = W^q a^i$$

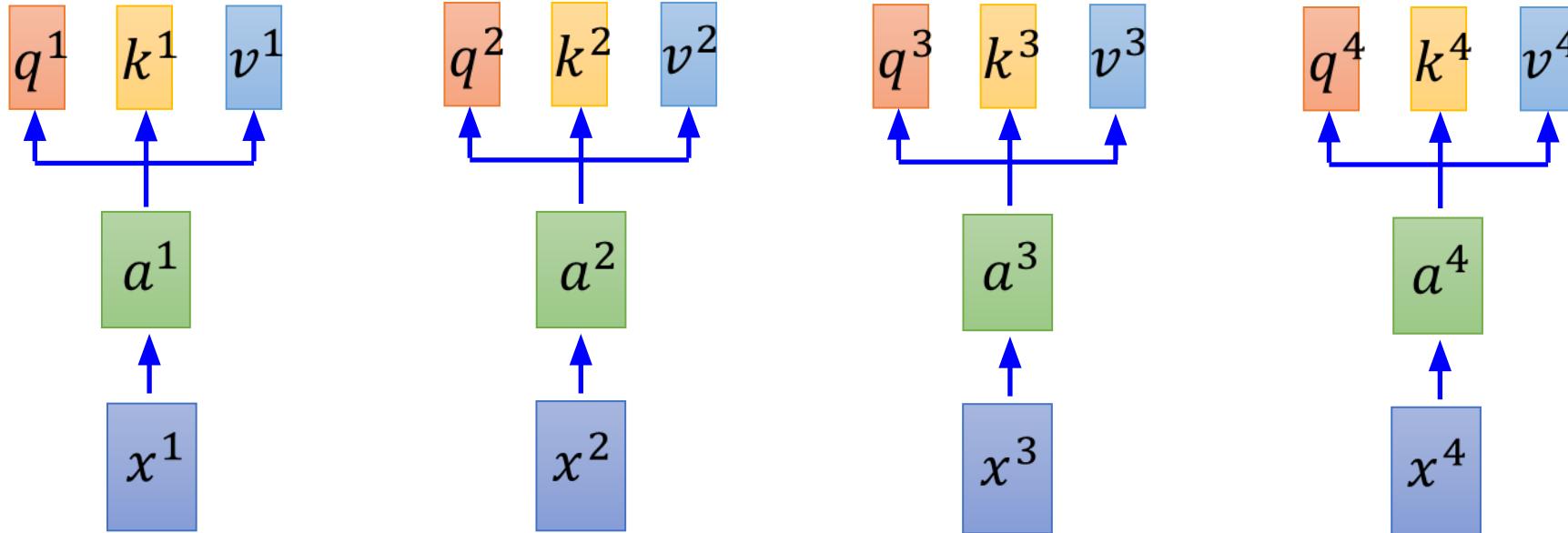
$$k^i = W^k a^i$$

$$v^i = W^v a^i$$

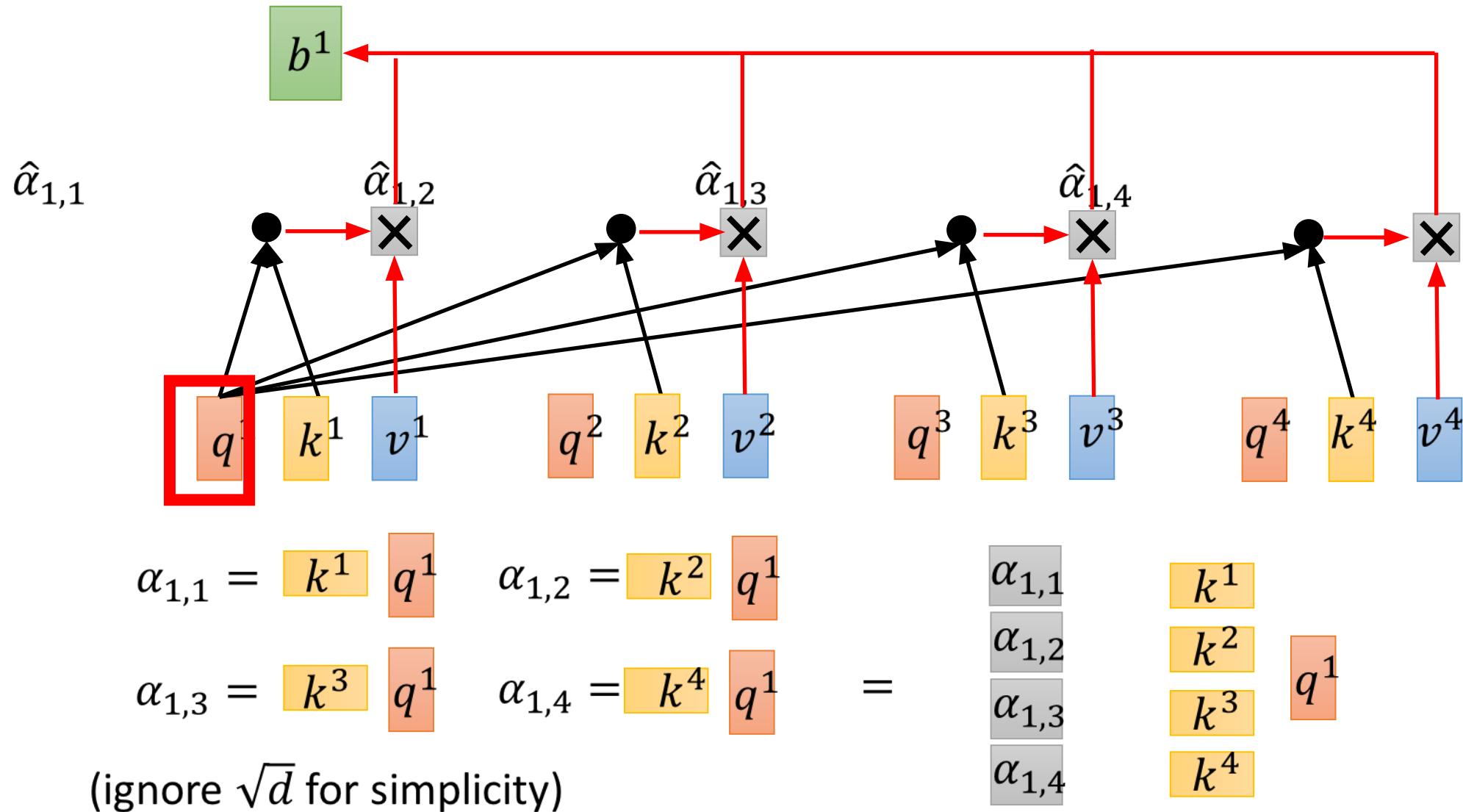
$$Q \quad q^1 \boxed{q^2} q^3 q^4 = \boxed{W^q} \quad a^1 \boxed{a^2} a^3 a^4 \\ I$$

$$K \quad k^1 \boxed{k^2} k^3 k^4 = \boxed{W^k} \quad a^1 \boxed{a^2} a^3 a^4 \\ I$$

$$V \quad v^1 v^2 \boxed{v^3} v^4 = \boxed{W^v} \quad a^1 \boxed{a^2} a^3 a^4 \\ I$$

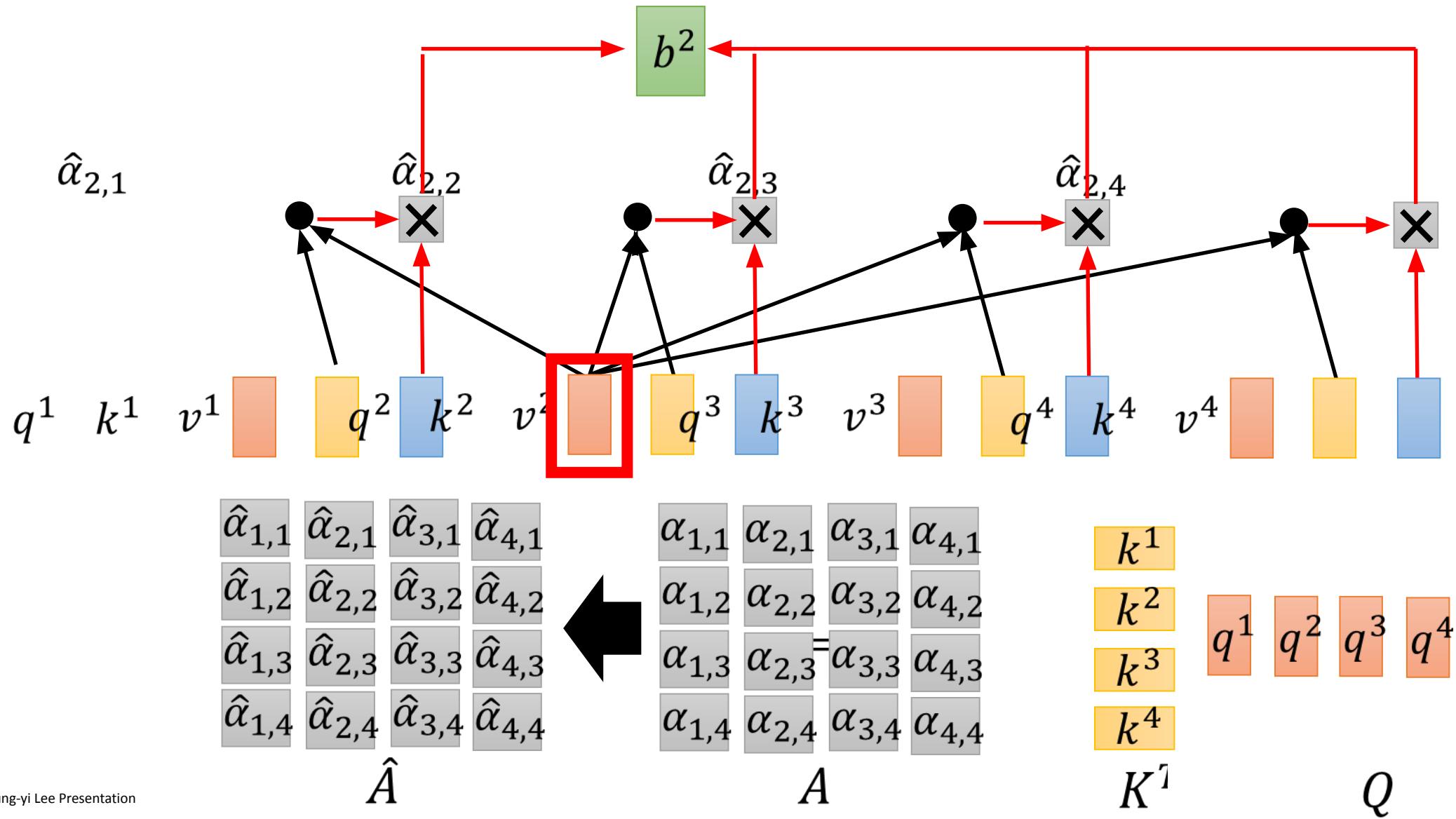


Self-Attention



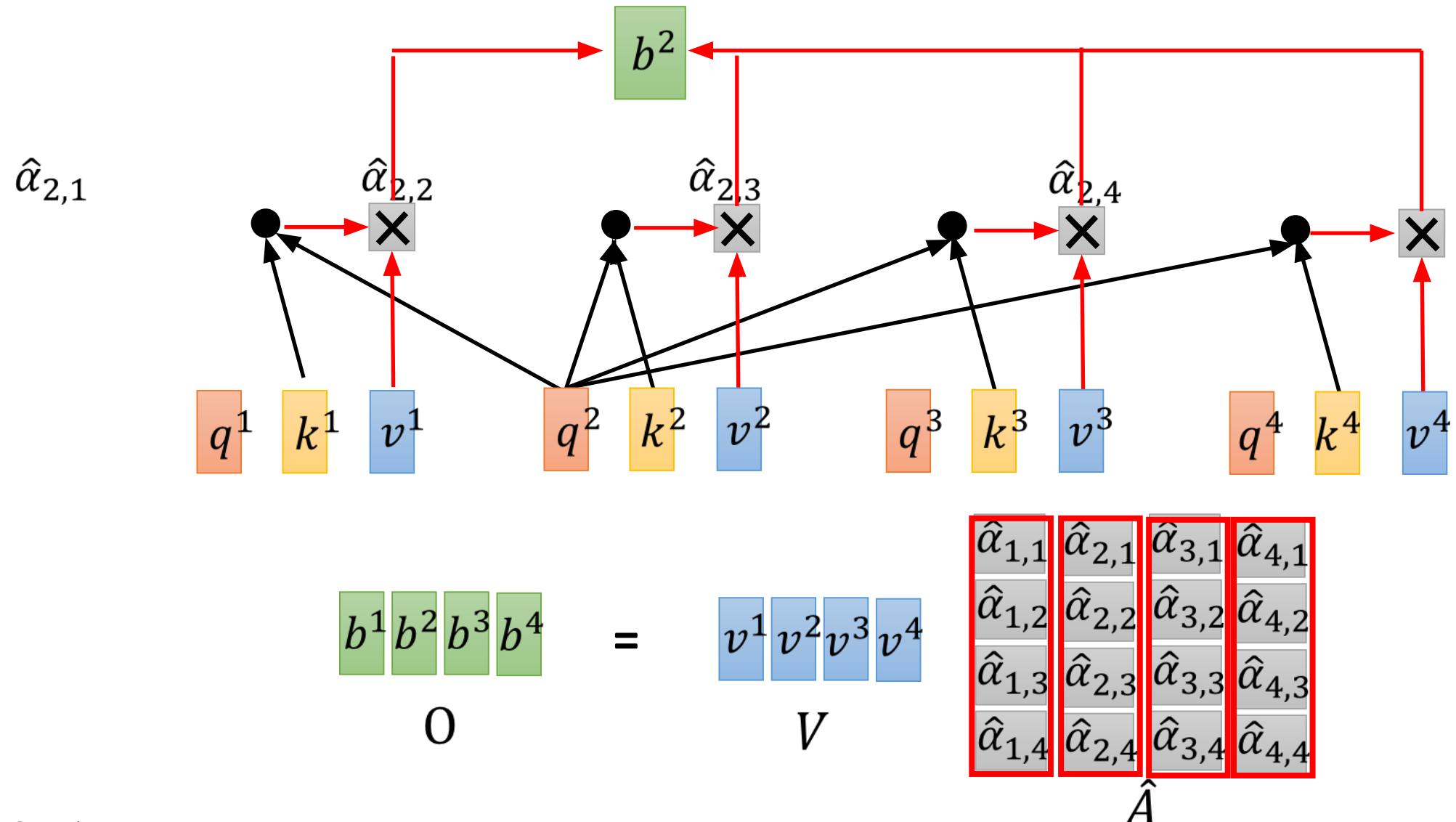
Self-Attention

$$b^2 = \sum_i \hat{\alpha}_{2,i} v^i$$

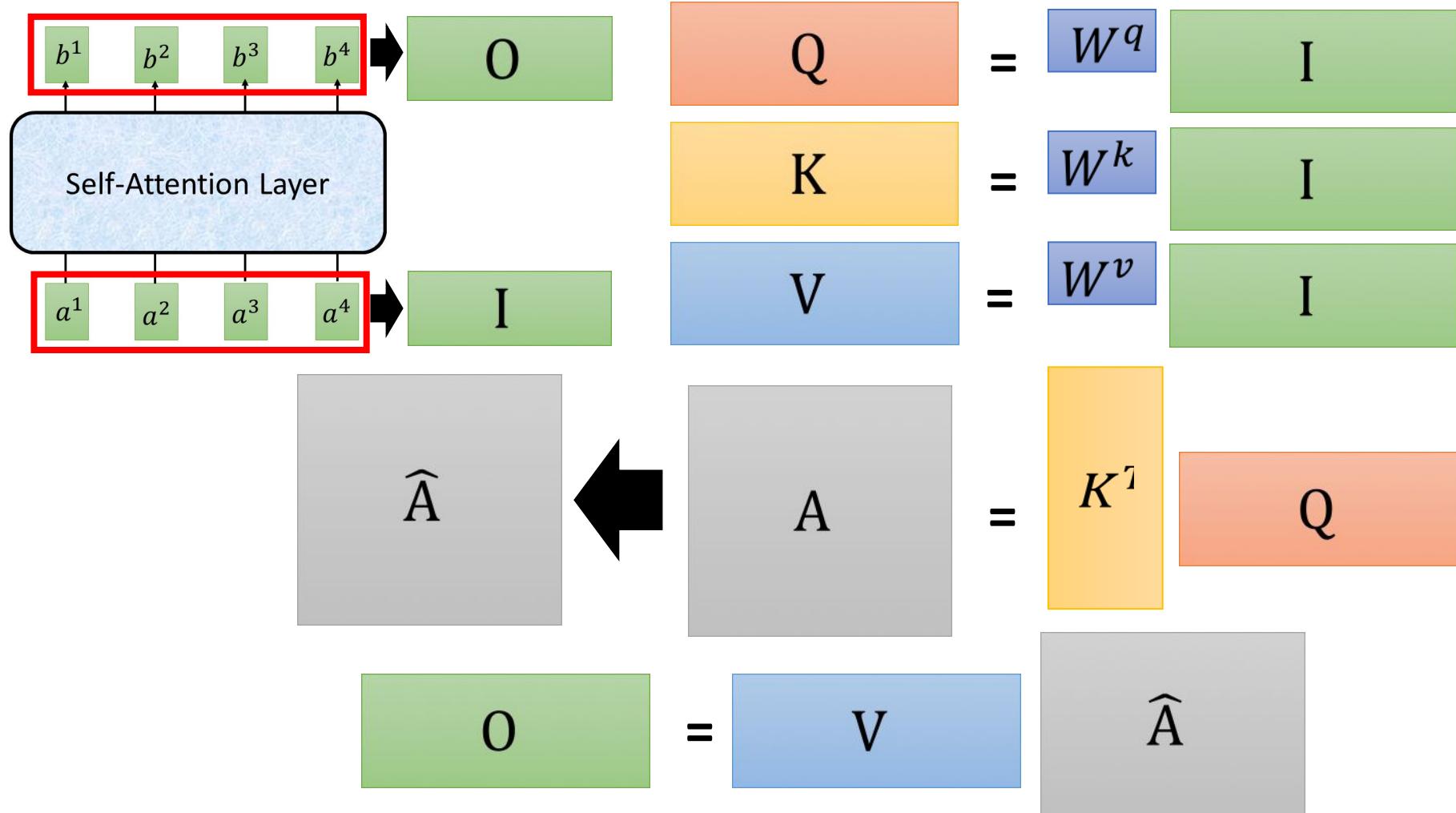


Self-Attention

$$b^2 = \sum_i \hat{\alpha}_{2,i} v^i$$

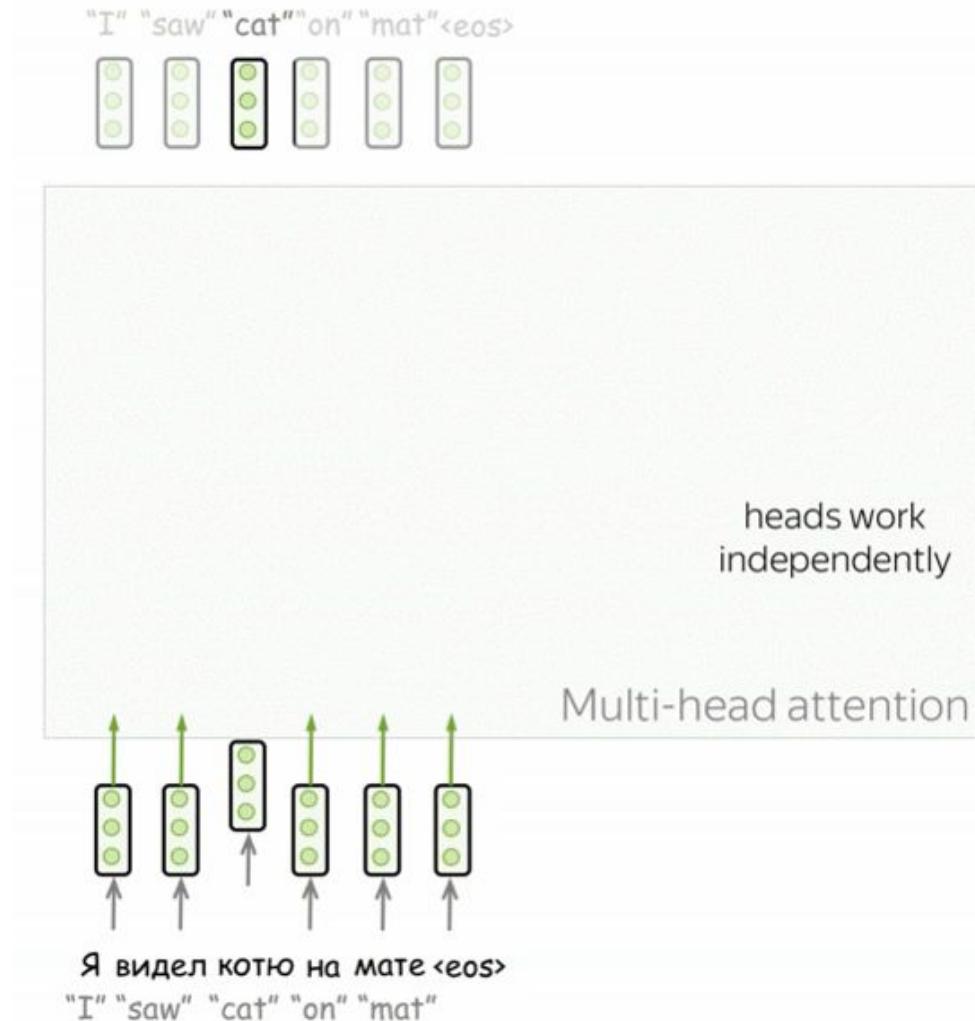


Self-Attention

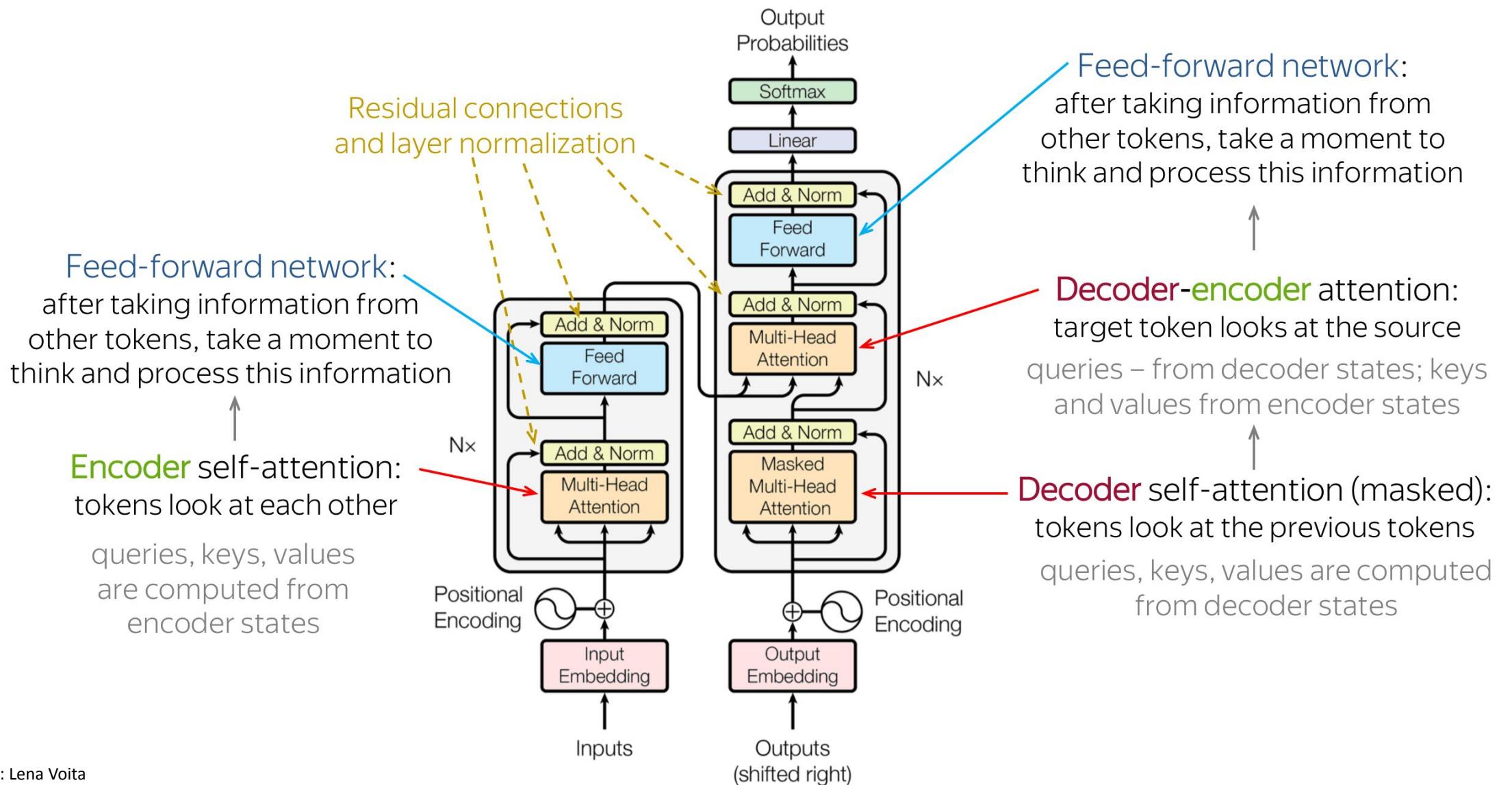


It's a bunch of matrix multiplication, which can be accelerated by GPU

Multi-head Attention



Transformer



Flow of Information

Consider the following sentences and their French translations:

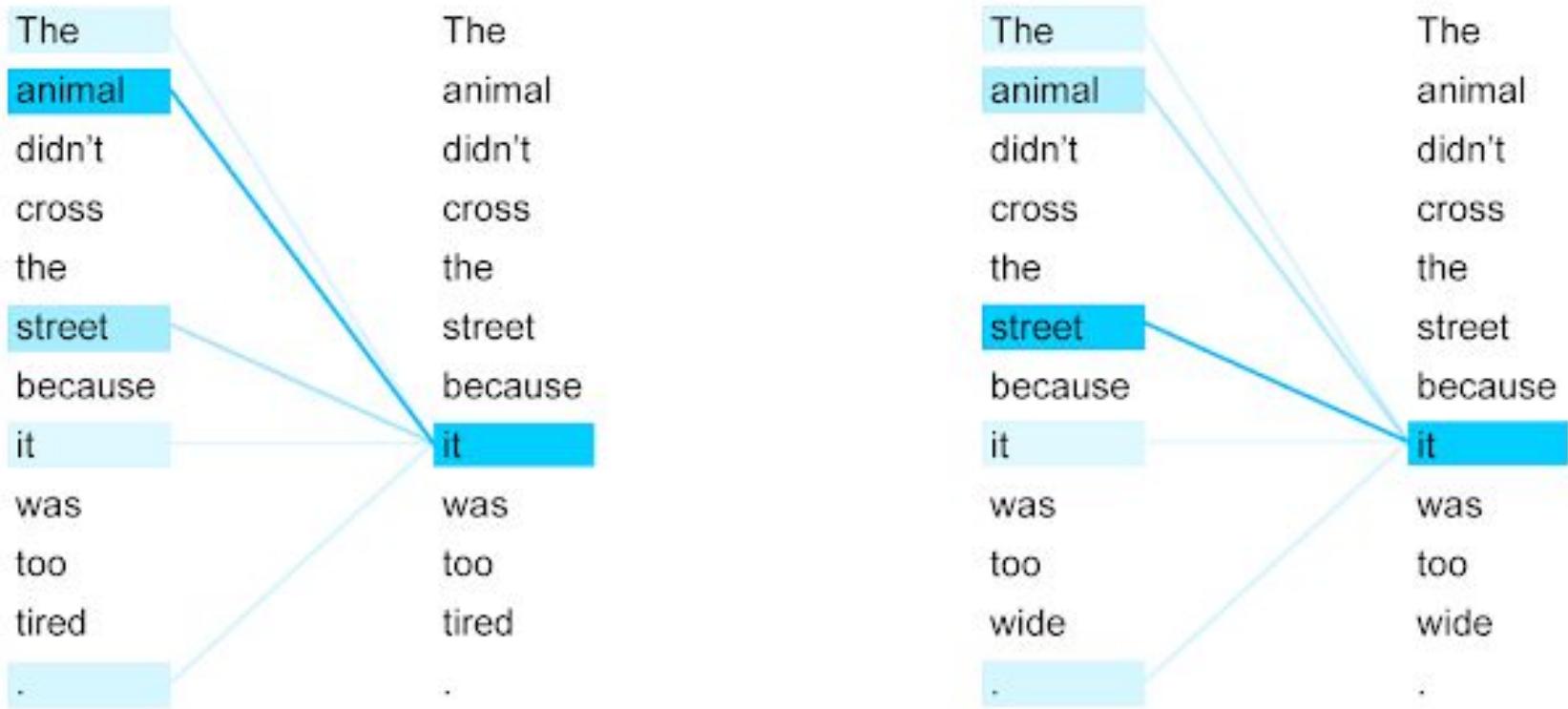
The animal didn't cross the street because it was too tired.
L'animal n'a pas traversé la rue parce qu'il était trop fatigué.

The animal didn't cross the street because it was too wide.
L'animal n'a pas traversé la rue parce qu'elle était trop large.

(Coreference resolution)

Transformer translates both of these sentences to French correctly

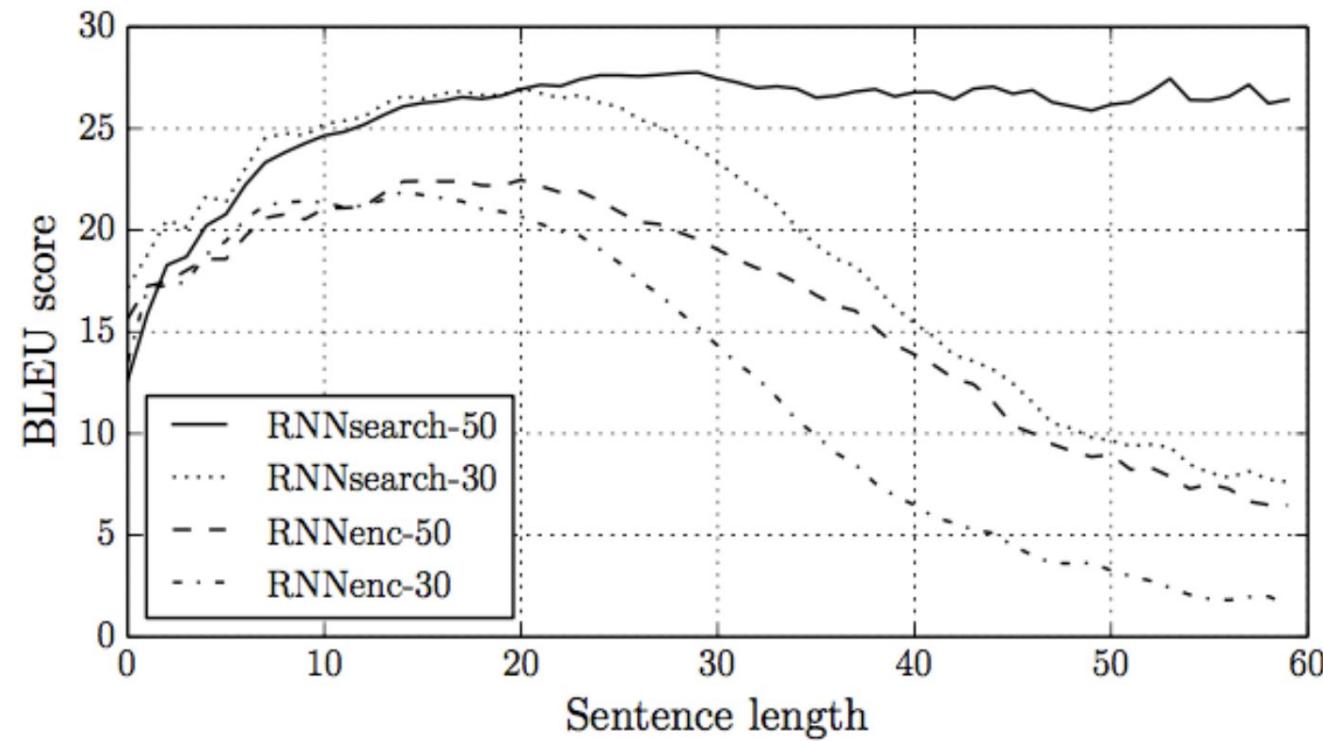
Attention Visualization



The encoder self-attention distribution for the word “it” from the 5th to the 6th layer of a Transformer trained on English to French translation.

Attention-Based Machine Translation

- The attention-based translation model does much better than the encoder/decoder model on long sentences.



GPT

Google's BERT

- Bidirectional Encoder Representations from Transformers (BERT)
- The BERT framework, a new language representation model from [Google AI](#), uses pre-training and fine-tuning to create state-of-the-art models for a wide range of tasks.
- These tasks include:
 - Question answering systems
 - Sentiment analysis
 - Language inference (http://nlpprogress.com/english/natural_language_inference.html)

BERT's Model Architecture

- BERT consists of a simple **stack of transformer blocks**
- This stack is *pre-trained* on a large general-domain corpus consisting of 800M words from English books (modern work, from unpublished authors), and 2.5B words of text from English Wikipedia articles (without markup).
- BERT uses a multi-layer bidirectional Transformer encoder.
- Its self-attention layer performs self-attention in both directions.

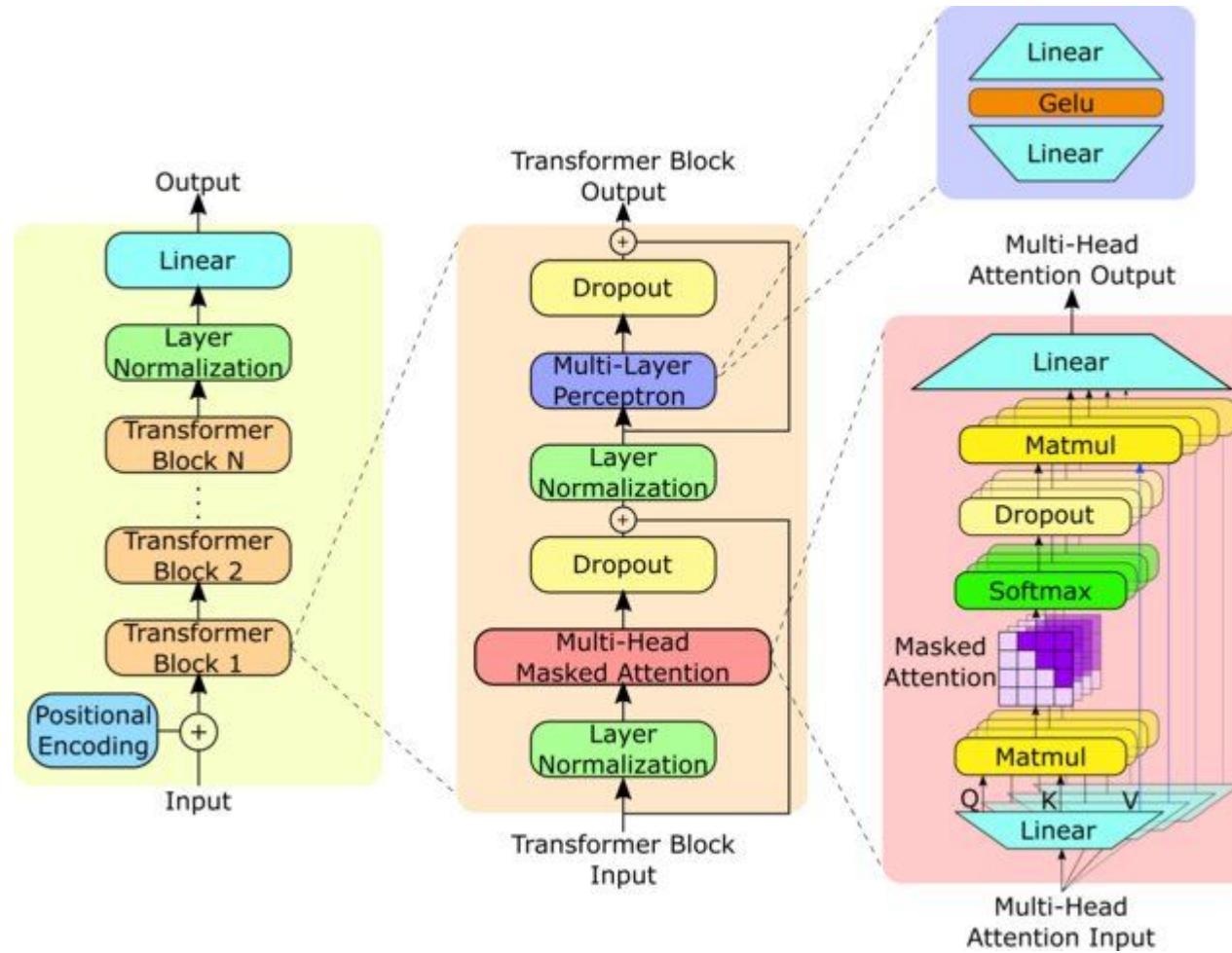
Google's BERT

- For pretraining, BERT uses the following hyper-parameters:
 - **Sequence length (single example):** 256
 - **Batch size:** 512
 - **Training steps:** 1,000,000 (Approximately 40 epochs)
 - **Optimizer:** Adam
 - **Learning rate:** 1e-4
 - **Learning rate schedule:** Warmup for 10,000 steps, then linear decay
 - **Dropout:** 0.1
 - **Activation function:** gelu (Gaussian Error Linear Unit)
 - **The training** took 4 days on 16 cloud TPUs (64 TPU chips).

Reference

- This lecture is based on the following resources + other web resources.
 - Chapter 10: <http://www.deeplearningbook.org/contents/rnn.html>
 - <https://arxiv.org/pdf/1703.03091.pdf>
 - https://lena-voita.github.io/nlp_course/seq2seq_and_attention.html
 - <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>
 - http://introtodeeplearning.com/materials/2018_6S191_Lecture2.pdf
 - RNN and LSTM slides: <http://bit.ly/2sO00ZC>
 - <https://medium.com/@jianqiangma/all-about-recurrent-neural-networks-9e5ae2936f6e>
 - <https://www.bioinf.jku.at/publications/older/2604.pdf>
 - CSC2535 2013: Advanced Machine Learning, Lecture 10 Recurrent neural networks by Geoffrey Hinton
 - <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>
 - <http://www.peterbloem.nl/blog/transformers>
 - <https://www.analyticsvidhya.com/blog/2019/06/understanding-transformers-nlp-state-of-the-art-models/>
 - <https://arxiv.org/pdf/1706.03762.pdf>
 - <https://web.stanford.edu/class/archive/cs/cs224n/cs224n.1184/lectures/lecture12.pdf>
 - <http://nlp.seas.harvard.edu/2018/04/03/attention.html>

GPT 2



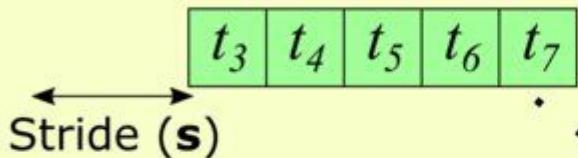
Striding

(a)

Training

t_1	t_2	t_3	t_4	t_5	t_6	...	t_{48}	t_{49}	t_{50}
-------	-------	-------	-------	-------	-------	-----	----------	----------	----------

t_1	t_2	t_3	t_4	t_5
-------	-------	-------	-------	-------



t_{44}	t_{45}	t_{46}	t_{47}	t_{48}
----------	----------	----------	----------	----------

t_{46}	t_{47}	t_{48}	t_{49}	t_{50}
----------	----------	----------	----------	----------

Window ($\mathbf{\omega}$)

(b)

t_1	t_2	t_3	t_4	t_5
-------	-------	-------	-------	-------

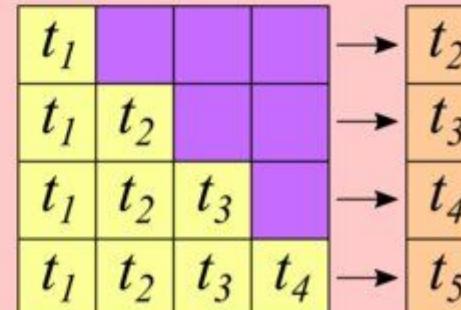
Training
Input

t_1	t_2	t_3	t_4
-------	-------	-------	-------

t_2	t_3	t_4	t_5
-------	-------	-------	-------

Training
Label

(c)



GPT 3

Bigger + fine-tuning (RL or supervised based on human feedback)

GPT 4

Bigger yet
8 submodels

Fine-tuning:

Reinforcement learning through human feedback RLHF -> reward agent
reward + “RL” -> fine-tuning
Also supervised learning to improve performance on rules and agents

Thank You

Appendix

Positional Encoding

- Transformer, all words in the input sentence are processed simultaneously – there is no inherent “word ordering” and thus no inherent position information.
- The authors of the Transformer propose adding a “positional encoding” to address this problem.
- The positional encoding makes it possible for the Transformer to use information about word order.
- The positional encoding uses a sequence of real-valued vectors that capture ordering information.
- Each word in a sentence is summed with a particular positional encoding vector based on its position within the sentence

<https://glassboxmedicine.com/2019/08/15/the-transformer-attention-is-all-you-need/>

Positional Encoding Vectors

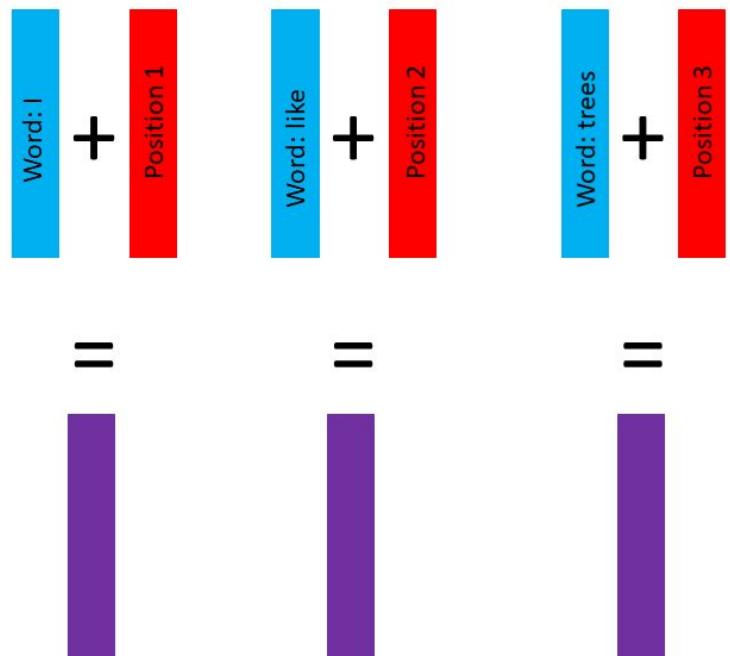
- **Option 1:** learning the positional encoding vectors (requires trainable parameters)
- **Option 2:** calculating the positional encoding vectors using an equation (requires no trainable parameters)
- They found the performance of both options was similar, so they went with option 2.

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

Positional Encoding

“I” “like” “trees”



COLOR SCHEME:

- word embedding (blue)
- positional encoding (red)
- final representation (purple)