

CS 4644-DL / 7643-A: LECTURE 13

DANFEI XU

Attention for Sequence Modeling

Attention is (Mostly) All you Need: Transformers

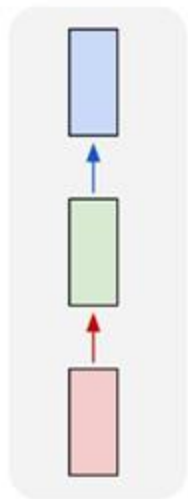
Administrative:

HW 3 due 10/20

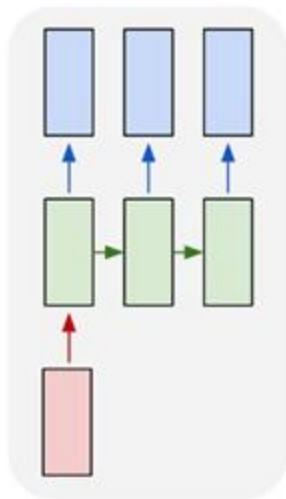
Milestone report due by 11/3: Need baseline results
(qualitative and quantitative)!

Recurrent Neural Networks: Process Sequences

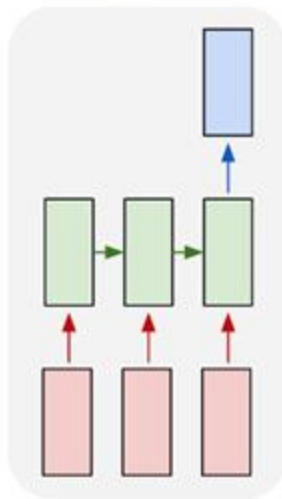
one to one



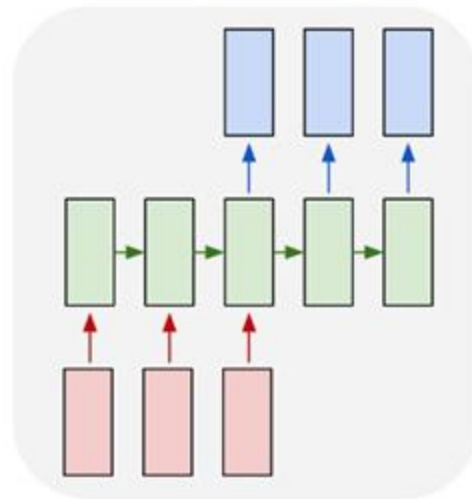
one to many



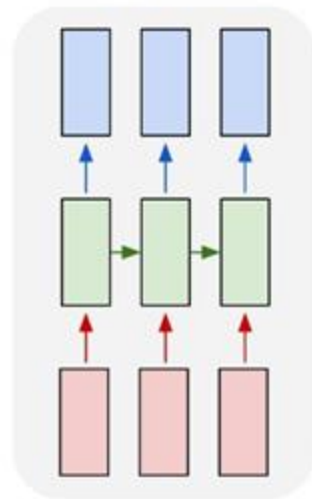
many to one



many to many



many to many



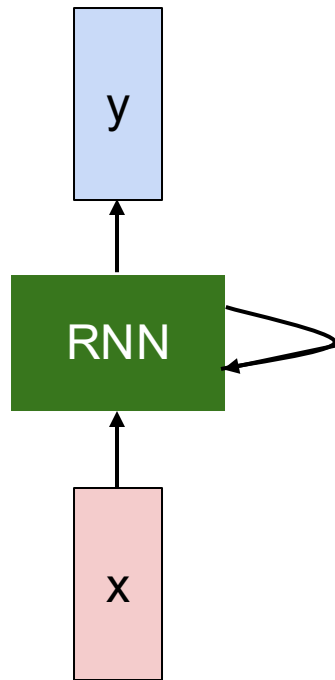
RNN hidden state update

We can process a sequence of vectors \mathbf{x} by applying a **recurrence formula** at every time step:

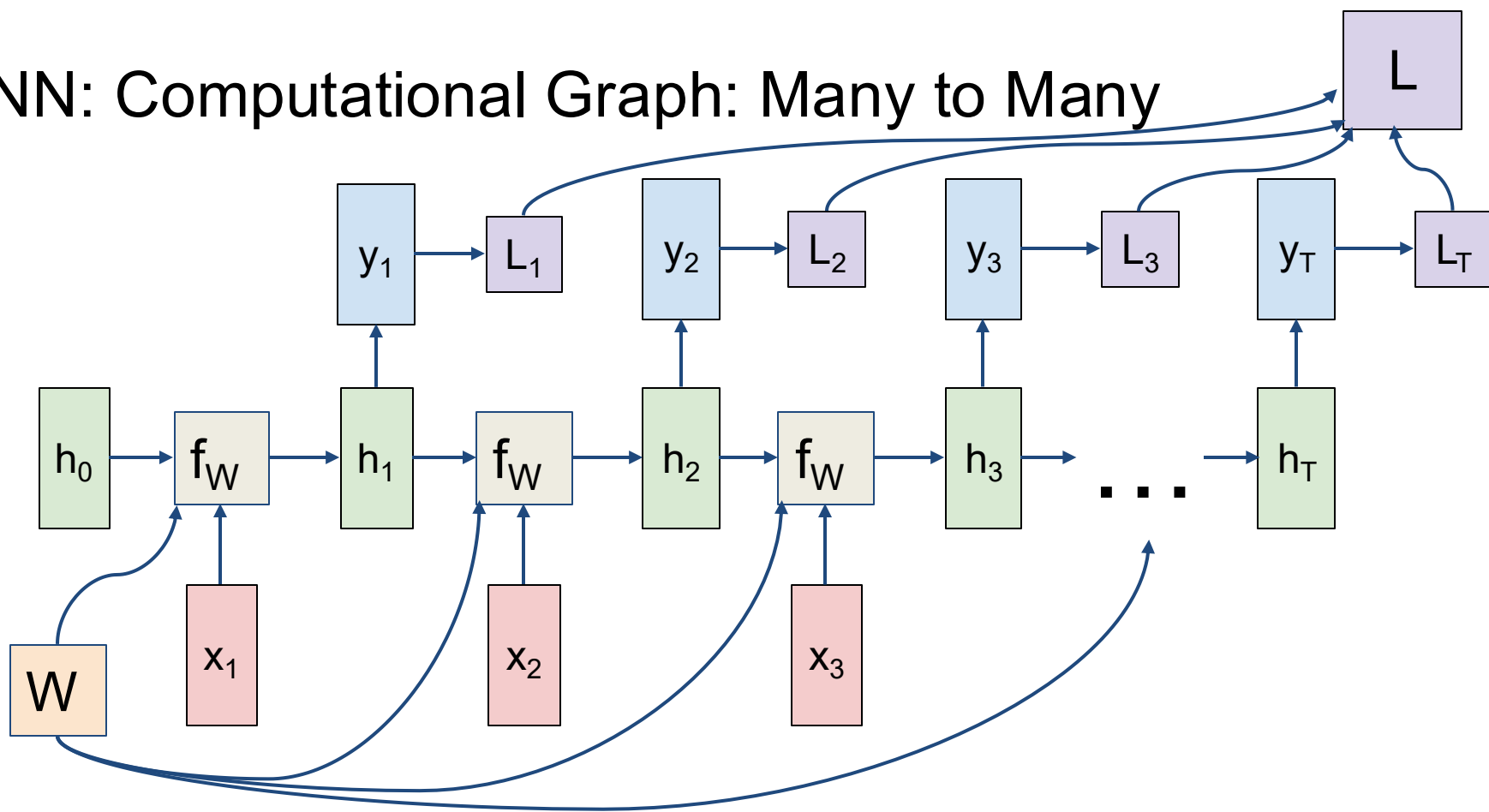
$$\boxed{h_t} = \boxed{f_W}(\boxed{h_{t-1}}, \boxed{x_t})$$

new state (vector) some function with parameters W old state (vector) input vector at some time step

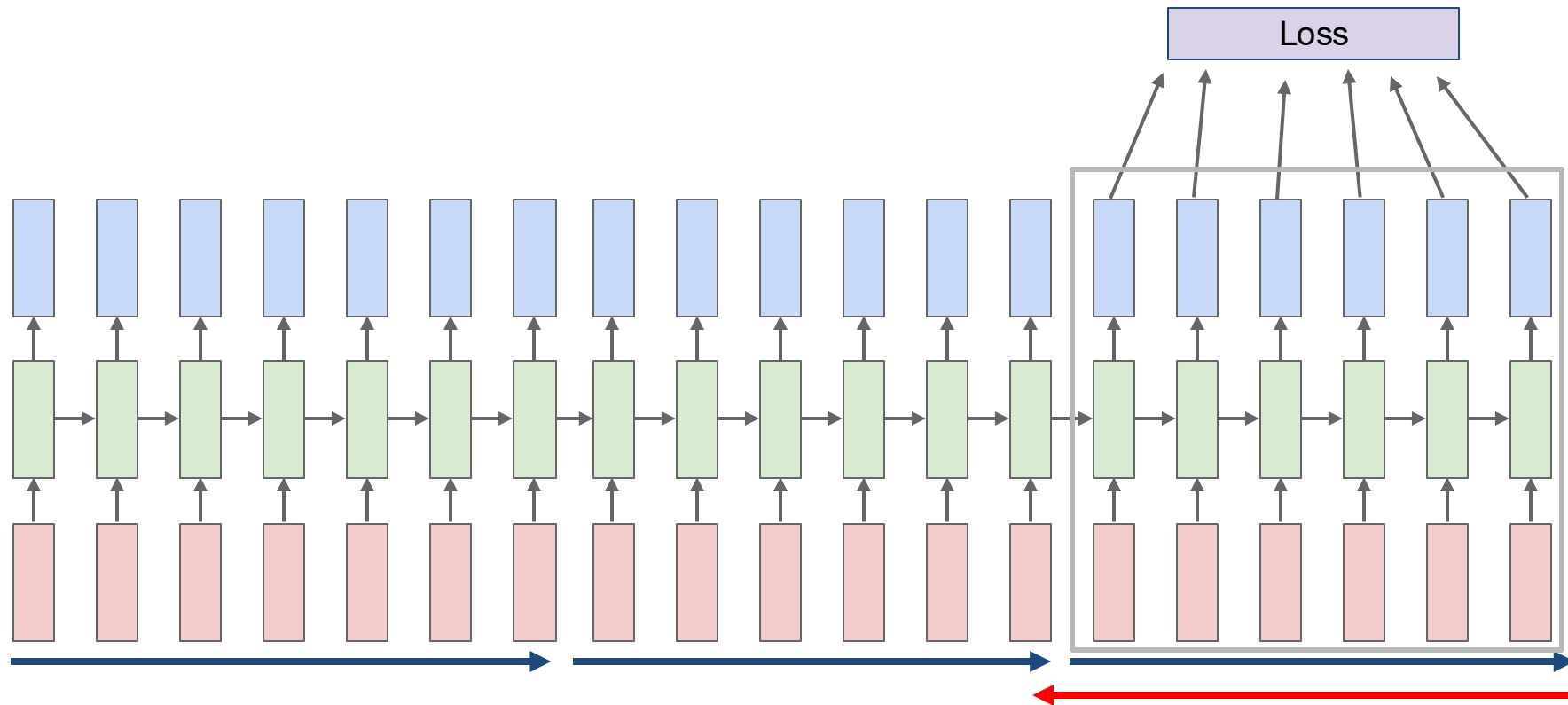
Can set initial state h_0 to all 0's



RNN: Computational Graph: Many to Many



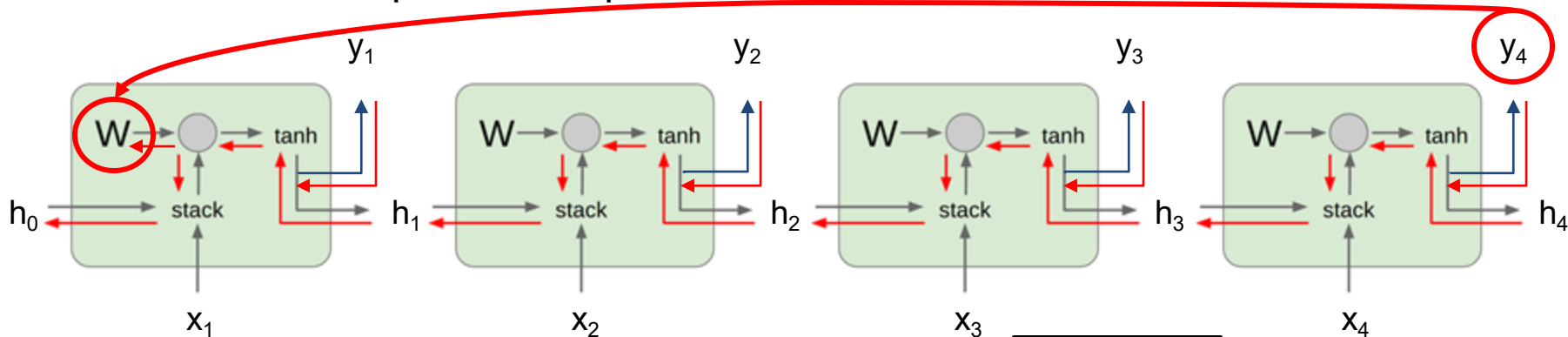
Truncated Backpropagation through time



Vanilla RNN Gradient Flow

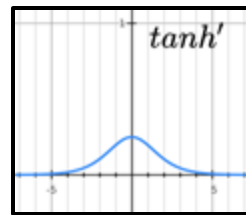
Gradients over multiple time steps:

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial L_t}{\partial W}$$

Always < 1
Vanishing gradients

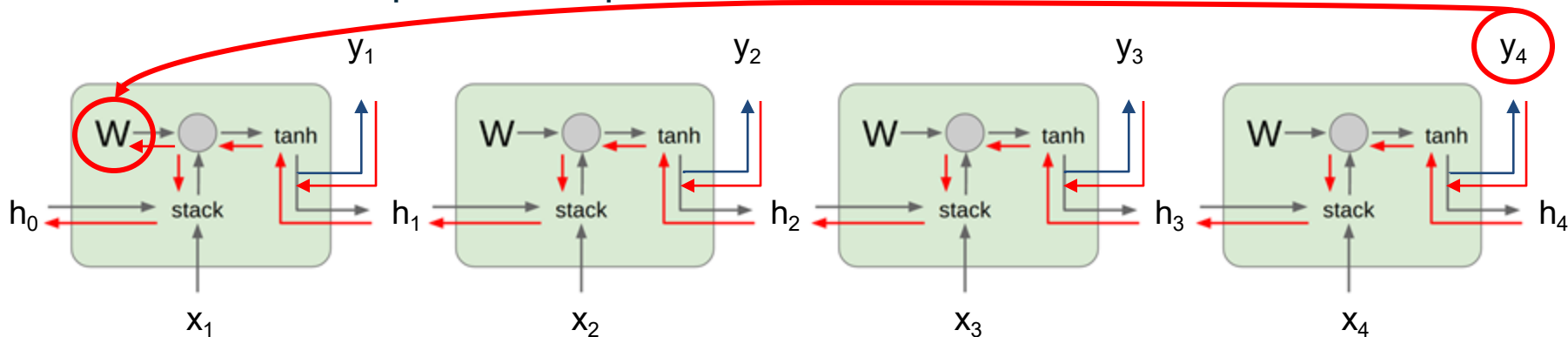


$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T} \left(\prod_{t=2}^T \tanh'(W_{hh} h_{t-1} + W_{xh} x_t) \right) W_{hh}^{T-1} \frac{\partial h_1}{\partial W}$$

Vanilla RNN Gradient Flow

Gradients over multiple time steps:

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



What if we assumed no non-linearity?

$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial L_t}{\partial W}$$

Largest eigen value > 1 :
Exploding gradients

$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial h_T} \boxed{W_{hh}^{T-1}} \frac{\partial h_1}{\partial W}$$

Largest eigen value < 1 :
Vanishing gradients

→ We need a new RNN architecture!

Long-Short Term Memory (Incomplete)

RNN directly updates h_t through **multiplying** with a weight matrix:

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$\frac{\partial h_t}{\partial h_{t-1}} = \tanh'(W_{hh}h_{t-1} + W_{xh}x_t)W_{hh}$$

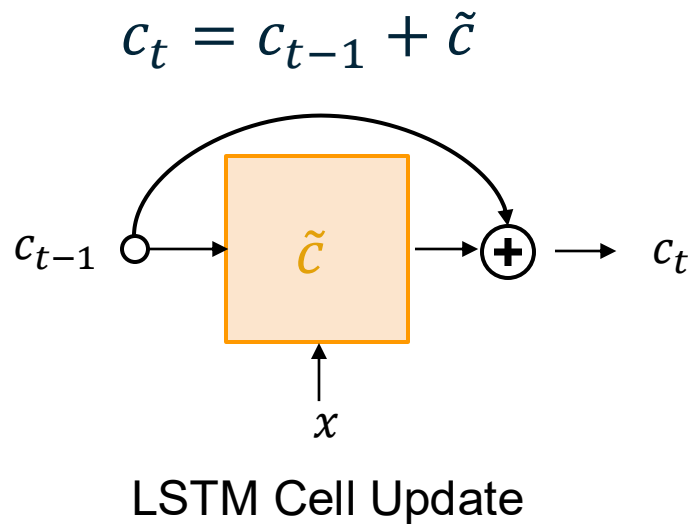
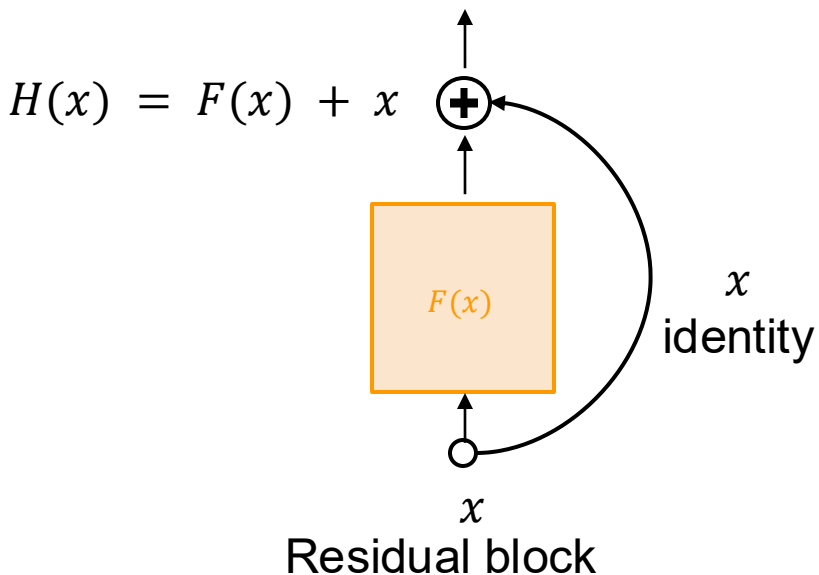
LSTM keeps track of a cell state c_t that's updated through **addition**

$$c_t = c_{t-1} + \tilde{c}$$
$$\tilde{c} = \tanh(W[h_{t-1}, x_t]), h_t = \tanh(c_t)$$

Gradient of cell state: $\frac{\partial c_t}{\partial c_{t-1}} = 1 + \frac{\partial \tilde{c}}{\partial c_{t-1}}$

This should look familiar ...
Residual connection!

Long-Short Term Memory (Incomplete)



$$\tilde{c} = \tanh(W[h_{t-1}, x_t]), h_t = \tanh(c_t)$$

Long Short Term Memory (LSTM)

Vanilla RNN

$$h_t = \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

LSTM

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$
$$c_t = f \odot c_{t-1} + i \odot g$$
$$h_t = o \odot \tanh(c_t)$$

Learn to control information flow from previous state to the next state

Long Short Term Memory (LSTM)

Vanilla RNN

$$h_t = \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

LSTM

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$
$$c_t = f \odot c_{t-1} + i \odot g$$
$$h_t = o \odot \tanh(c_t)$$

Learn to control information flow from previous state to the next state

Long Short Term Memory (LSTM)

Vanilla RNN

$$h_t = \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

Long-term memory c determines how much information should go into the hidden state h (short-term memory)

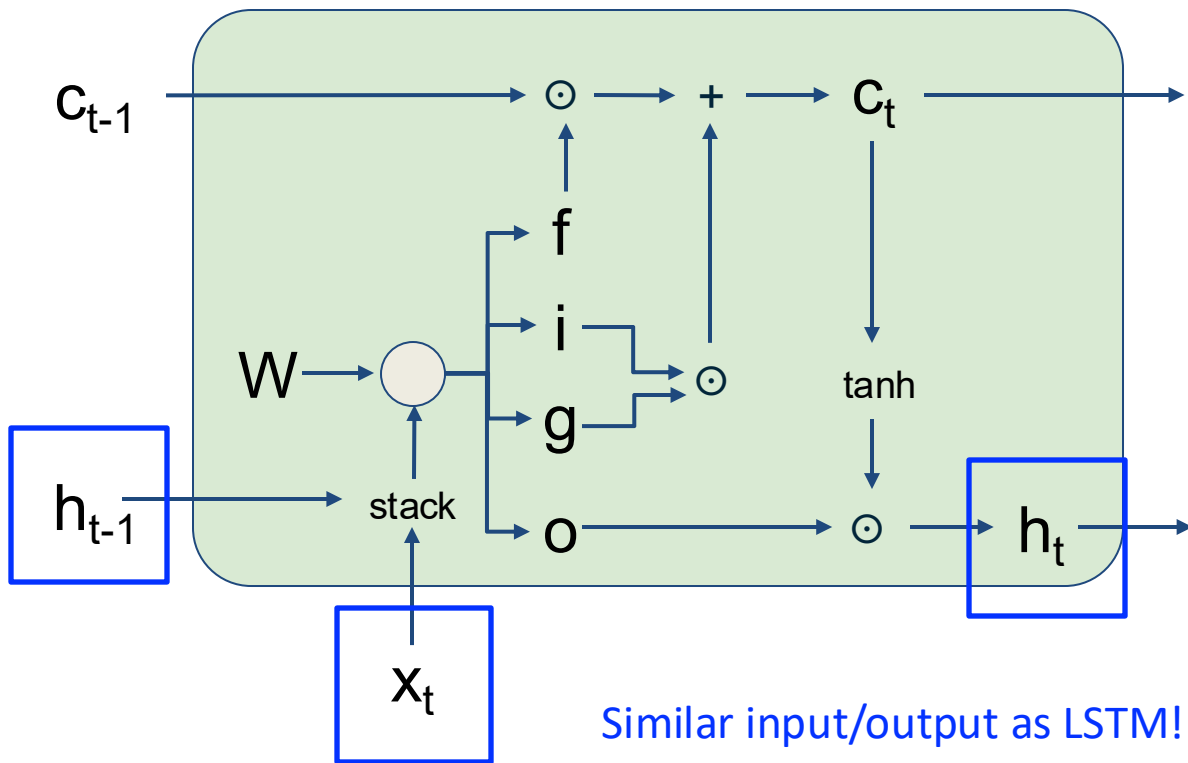
LSTM

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$
$$\begin{aligned} c_t &= f \odot c_{t-1} + i \odot g \\ h_t &= o \odot \tanh(c_t) \end{aligned}$$

Two “memory vectors”

Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]



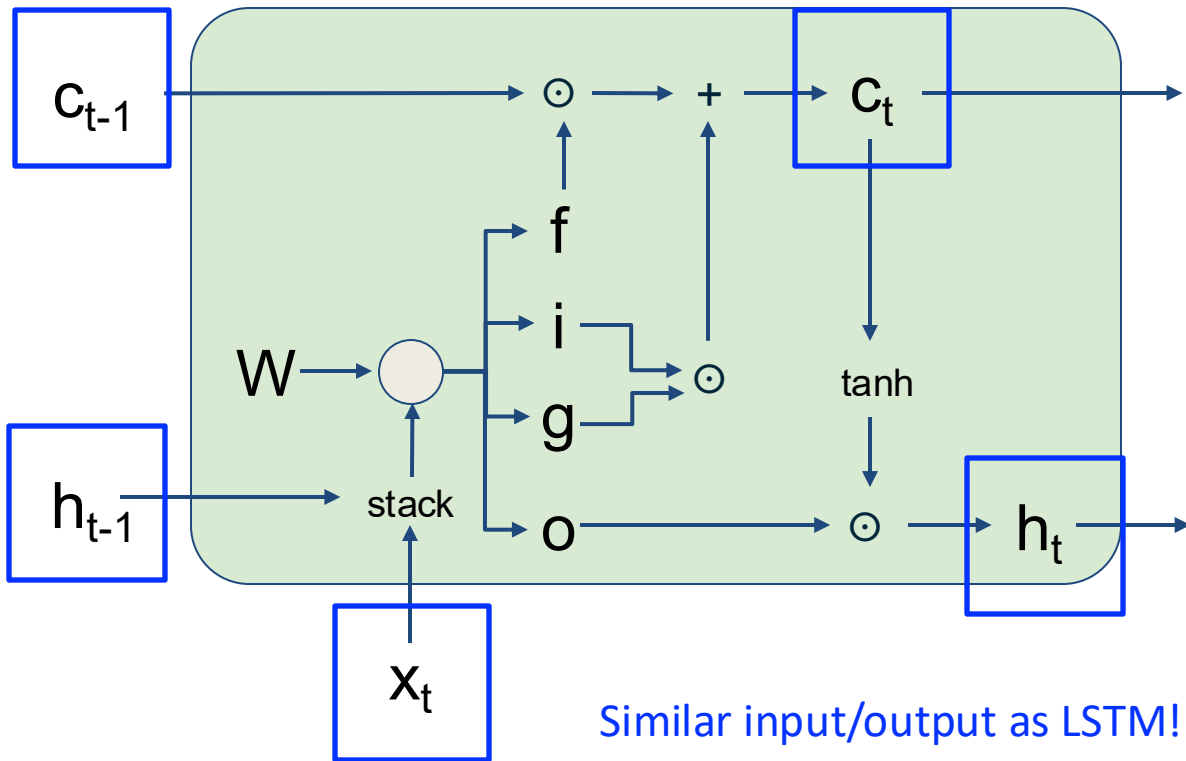
$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$
$$h_t = o_t \odot \tanh(c_t)$$

Similar input/output as LSTM!

Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]

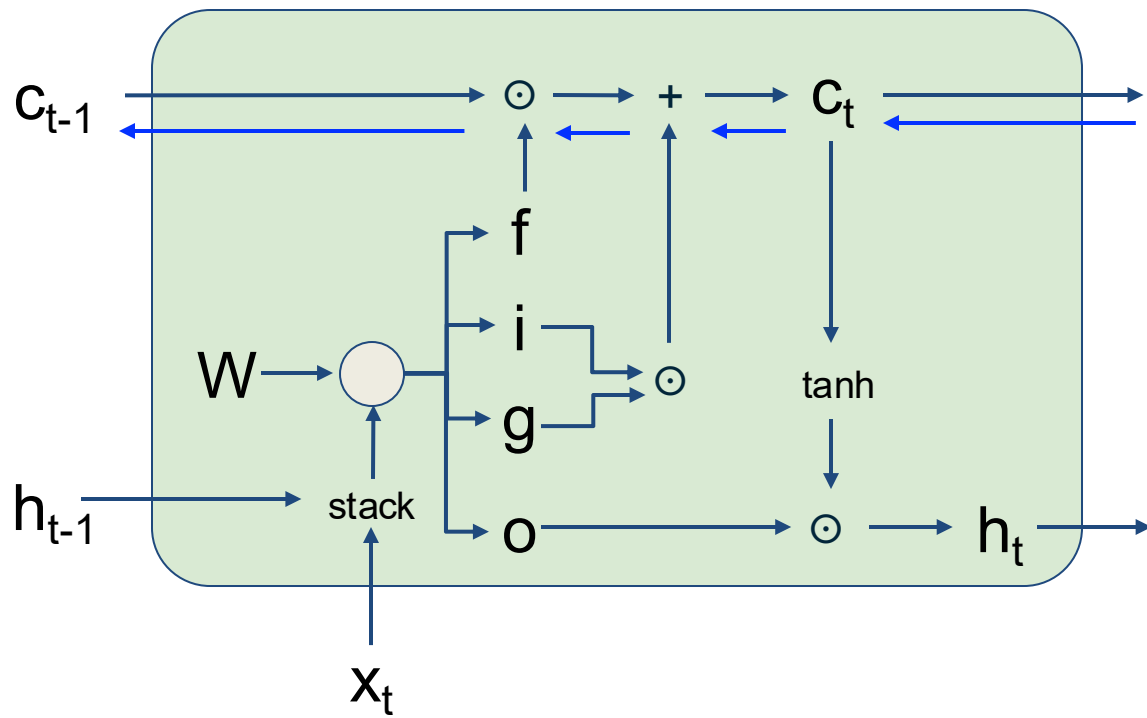
Keep long-term memory cell c in addition
to short term memory h



$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$
$$h_t = o_t \odot \tanh(c_t)$$

Long Short Term Memory (LSTM): Gradient Flow

[Hochreiter et al., 1997]



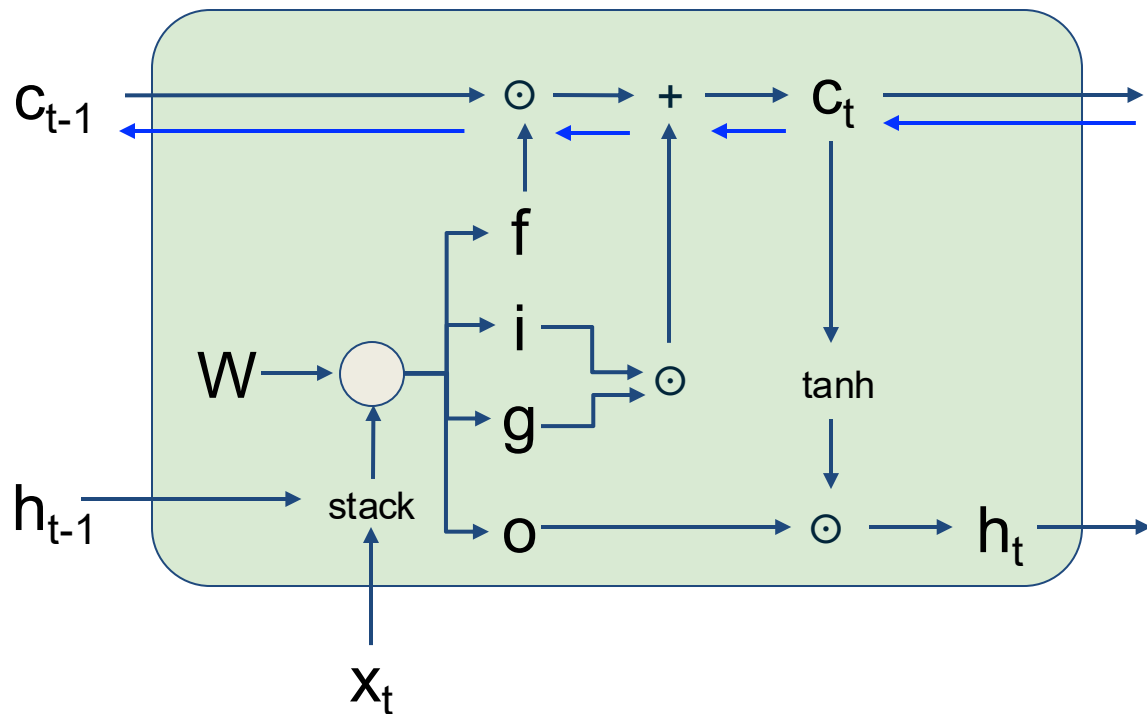
Backpropagation from c_t to c_{t-1}
only elementwise multiplication
by f (forget gate), no matrix
multiply by W

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$
$$h_t = o_t \odot \tanh(c_t)$$

$$\frac{\partial c_t}{\partial c_{t-1}} = ?$$

Long Short Term Memory (LSTM): Gradient Flow

[Hochreiter et al., 1997]



Backpropagation from c_t to c_{t-1}
only elementwise multiplication
by f (forget gate), no matrix
multiply by W

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$
$$h_t = o_t \odot \tanh(c_t)$$

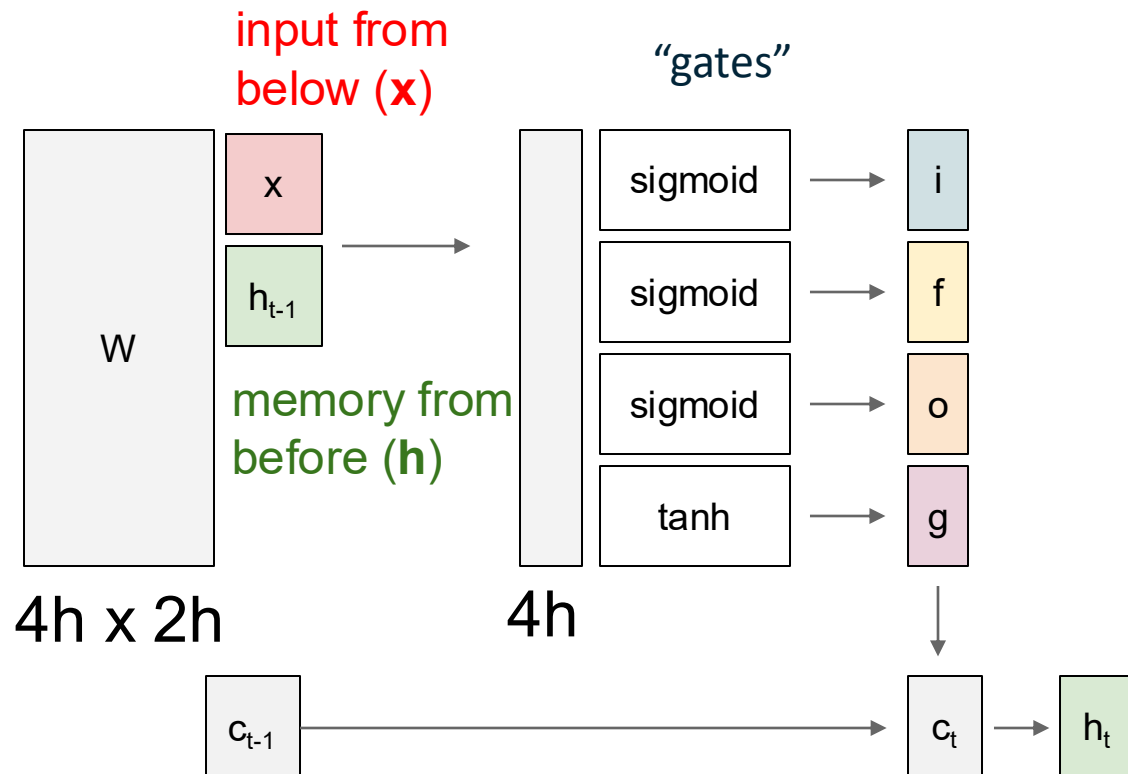
$$\frac{\partial c_t}{\partial c_{t-1}} = f_t + \dots \text{ (forget gate)}$$

Different each step!

When f_t is close to 1, it allows
gradient to flow back easily

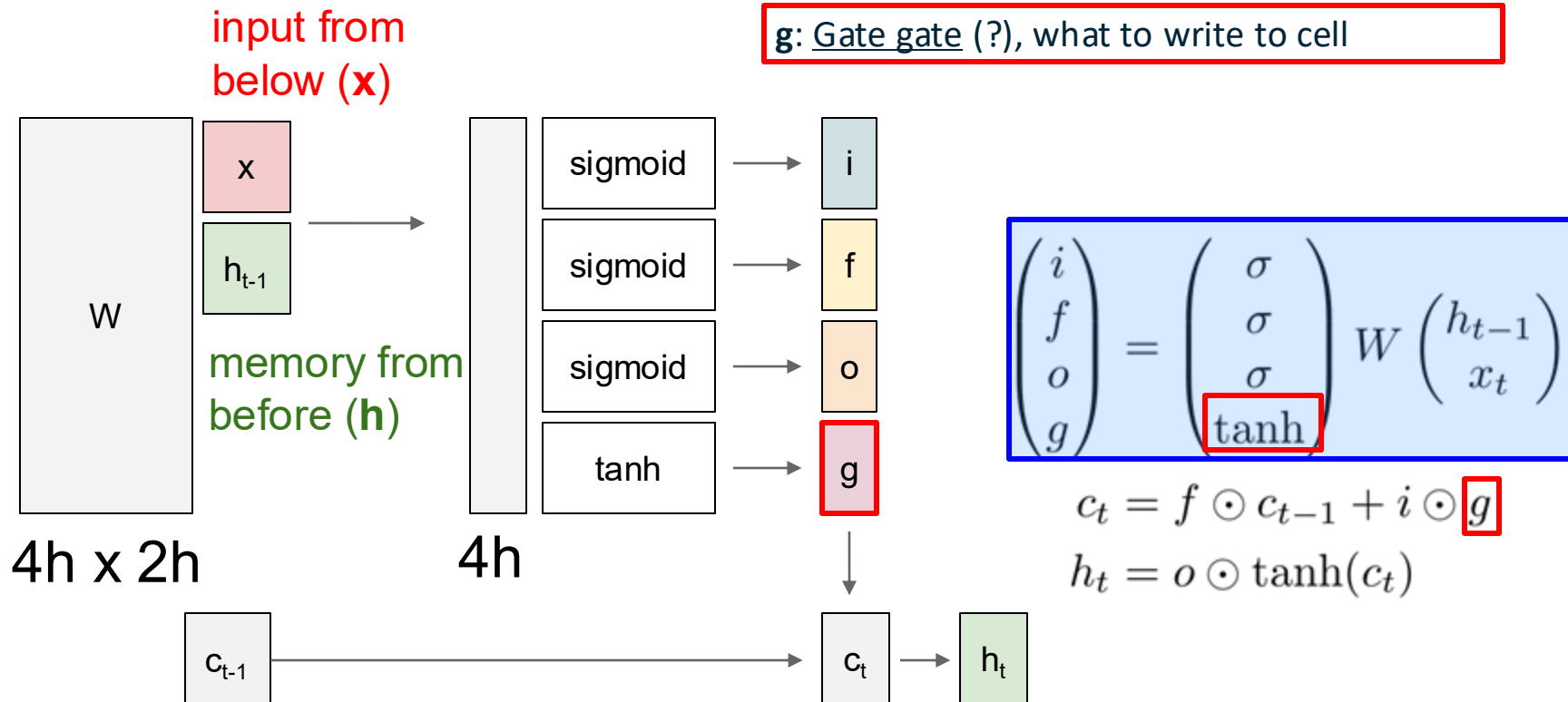
Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]



Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]

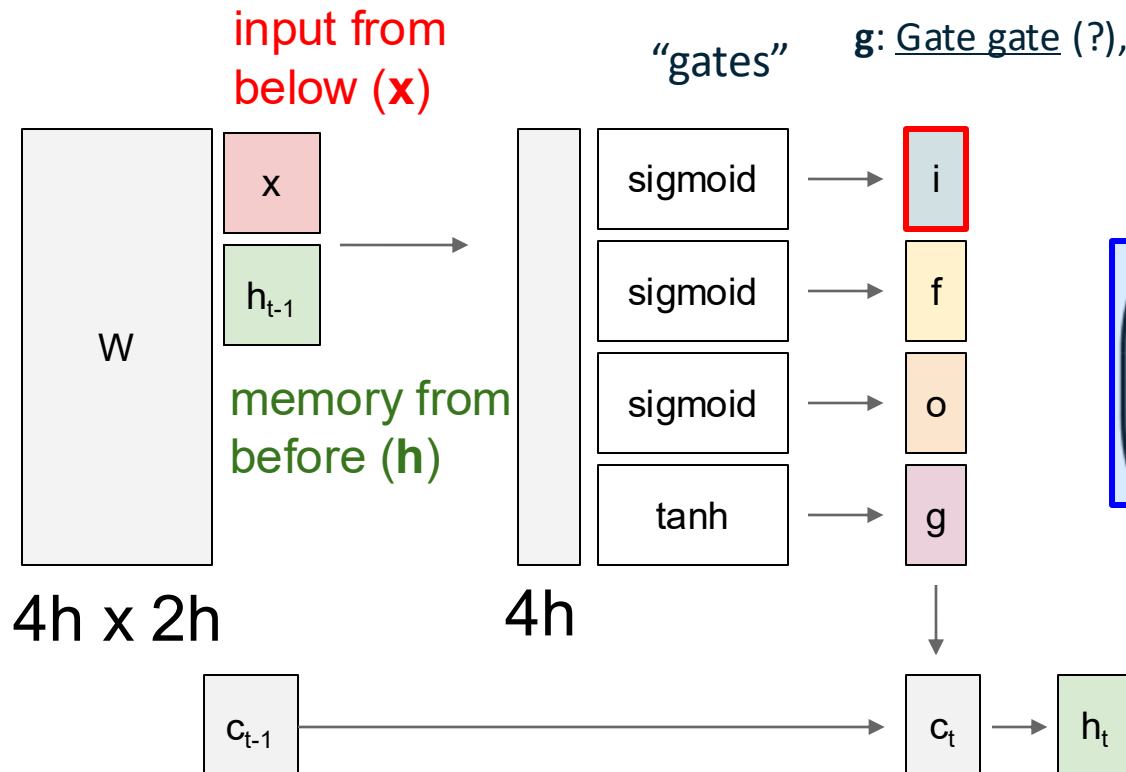


Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]

i: Input gate, whether to write to cell

g: Gate gate (?), what to write to cell



$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

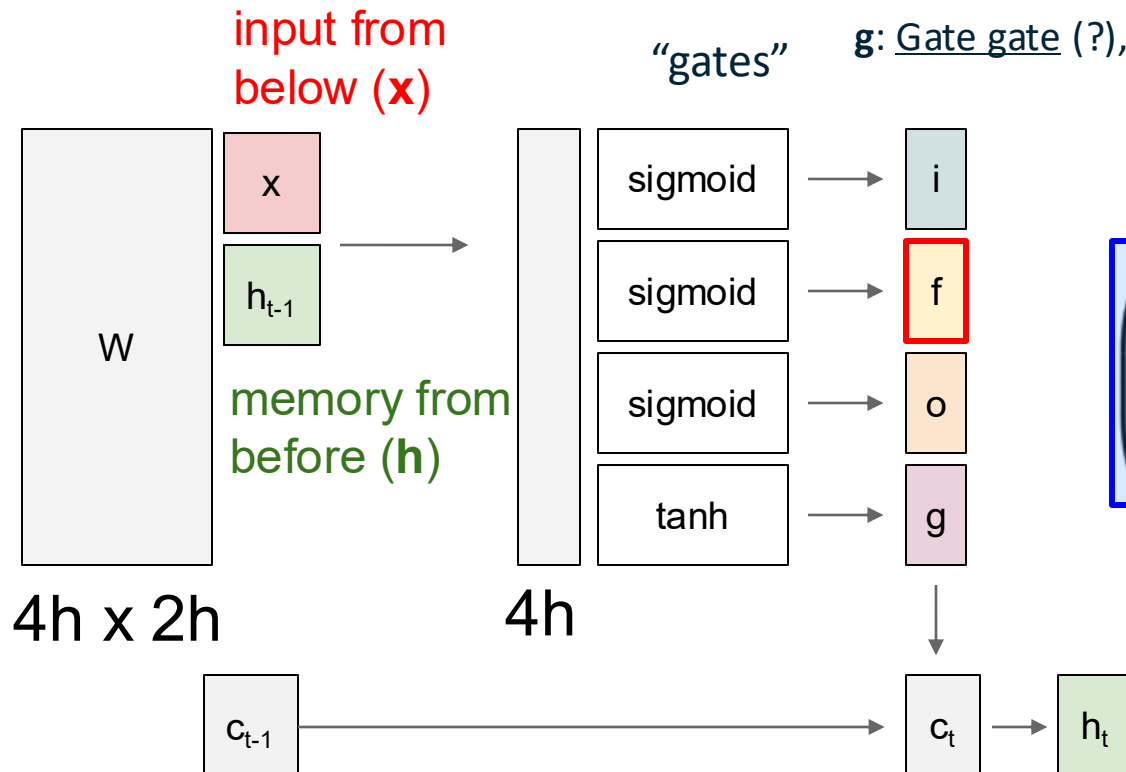
Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]

i: Input gate, whether to write to cell

f: Forget gate, whether to erase cell

g: Gate gate (?), what to write to cell



$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

Long Short Term Memory (LSTM)

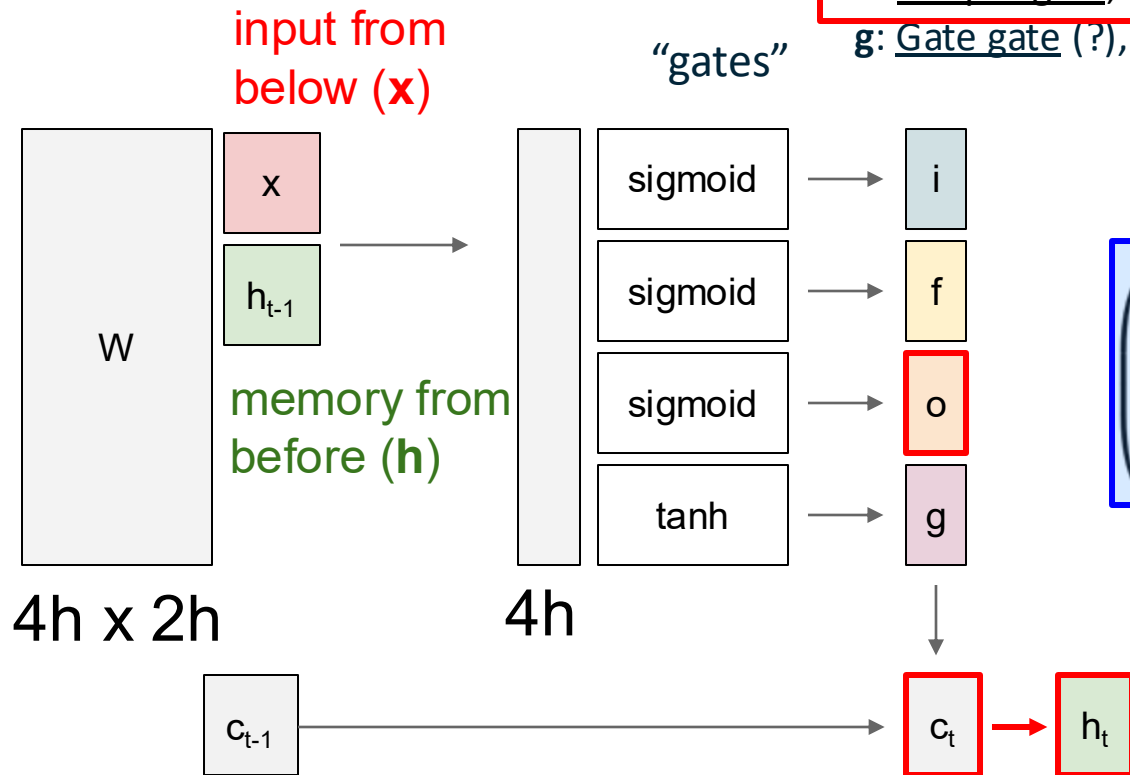
[Hochreiter et al., 1997]

i: Input gate, whether to write to cell

f: Forget gate, whether to erase cell

o: Output gate, how much to reveal cell

g: Gate gate (?), what to write to cell



$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

Do LSTMs solve the vanishing gradient problem?

The LSTM architecture makes it easier for the RNN to preserve information over many timesteps

- e.g. **if $f = 1$ and $i = 0$** , then the information of that cell is preserved indefinitely. Gradient flow back from cell c easily.
- By contrast, it's harder for vanilla RNN to learn a recurrent weight matrix W_h that preserves info in hidden state

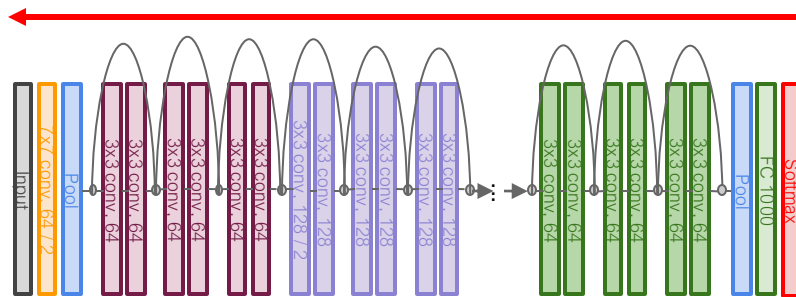
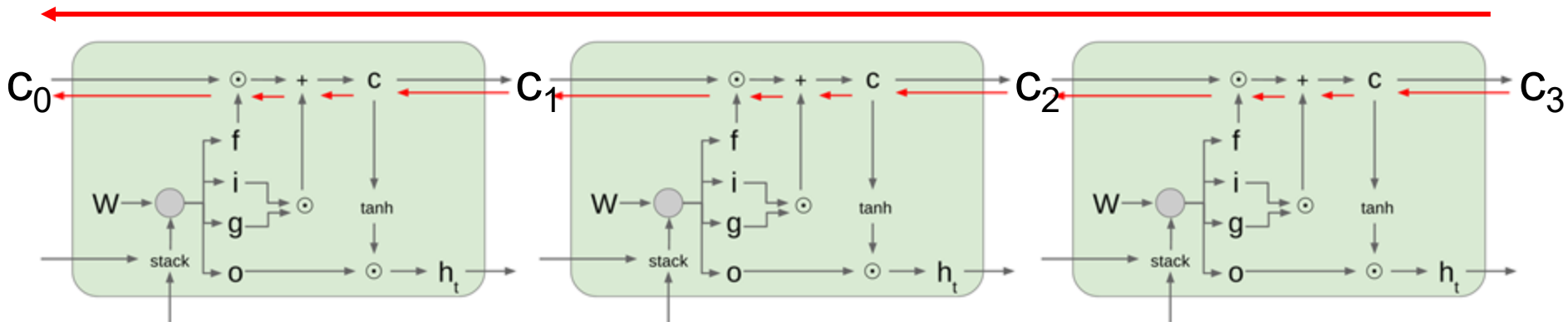
LSTM **doesn't guarantee** that there is no vanishing/exploding gradient, but it does provide an easier way for the model to learn long-distance dependencies.

It is possible to mitigate vanishing / exploding gradient by learning the correct f

Long Short Term Memory (LSTM): Gradient Flow

[Hochreiter et al., 1997]

Uninterrupted gradient flow!



Similar to ResNet!

Recommendations

- If you want to use RNN-like models, try LSTM
- Use variants like GRU if you want faster compute and less parameters
- New variants of RNNs are still active research topic. Example: RWKV (“Transformer-level performance but with RNN”)

Problem with Recurrent-style Models (RNN, LSTM, GRU, etc.)

Learning to memorize is still hard, especially for ultra-long sequences!

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$
$$c_t = f \odot c_{t-1} + i \odot g$$
$$h_t = o \odot \tanh(c_t)$$

Essentially trying to tune W such that the memory cell c can retain **important information** for **arbitrary future prediction problems**.

Example (Q&A):

[... (20-page long transcript)]. Q: What did the CEO say about their competitor company? ...

[... (same 20-page transcript)]. Q: How many times did the journalist use the word “interesting”? ...

Very difficult learning problem!

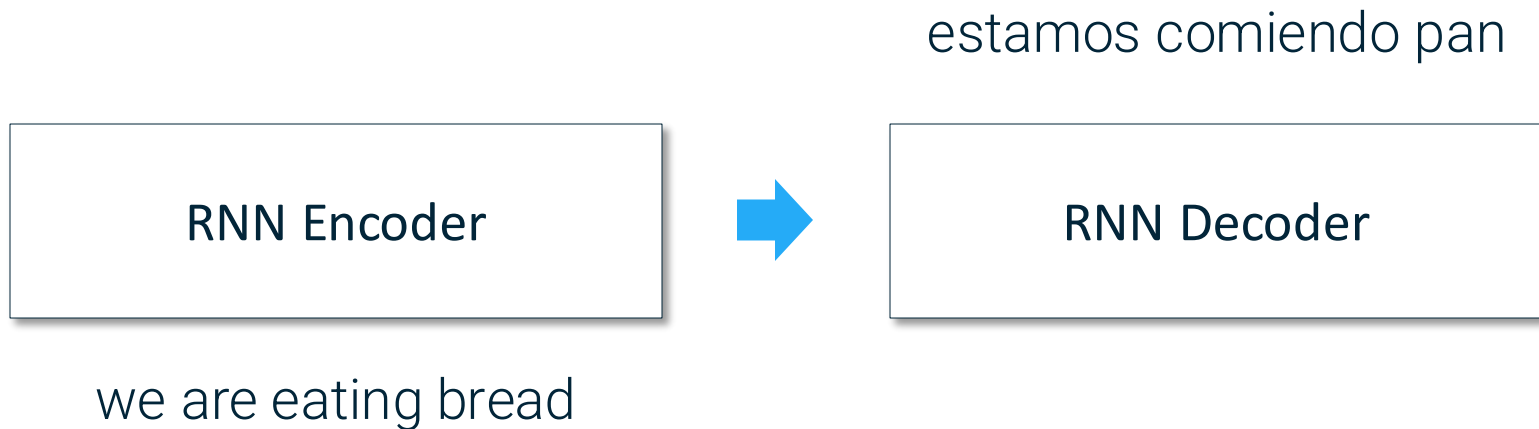
Attention Mechanism

(What memory? Just show me the sequence again)

Attention Mechanism

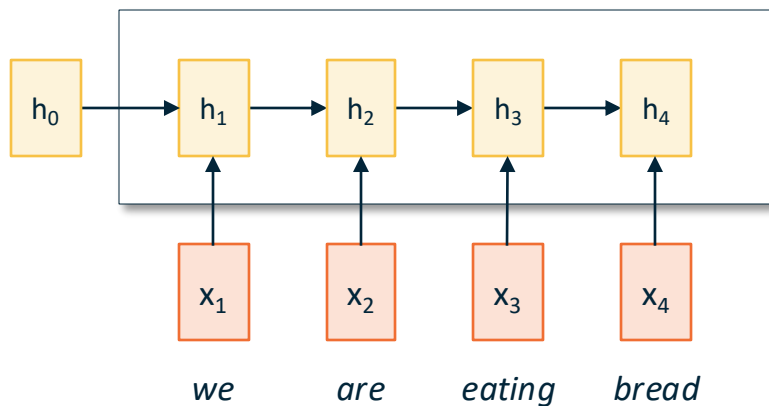


Example: Machine Translation



Machine Translation with RNNs

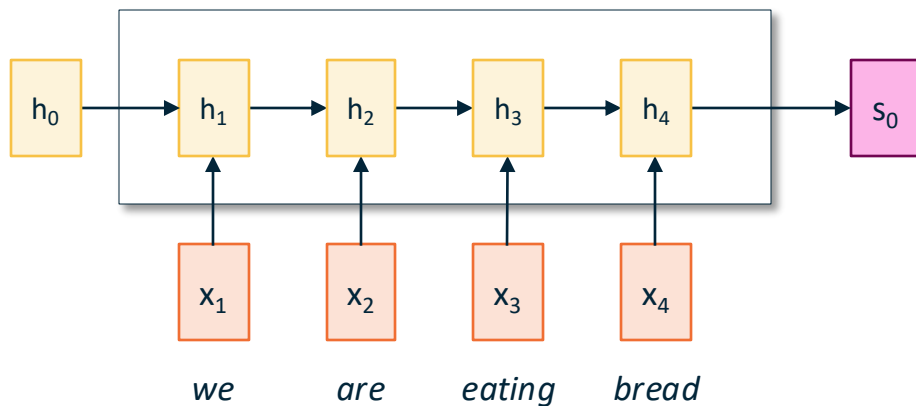
Encoder: $h_t = f_w(x_t, h_{t-1})$



Machine Translation with RNNs

Encoder: $h_t = f_W(x_t, h_{t-1})$

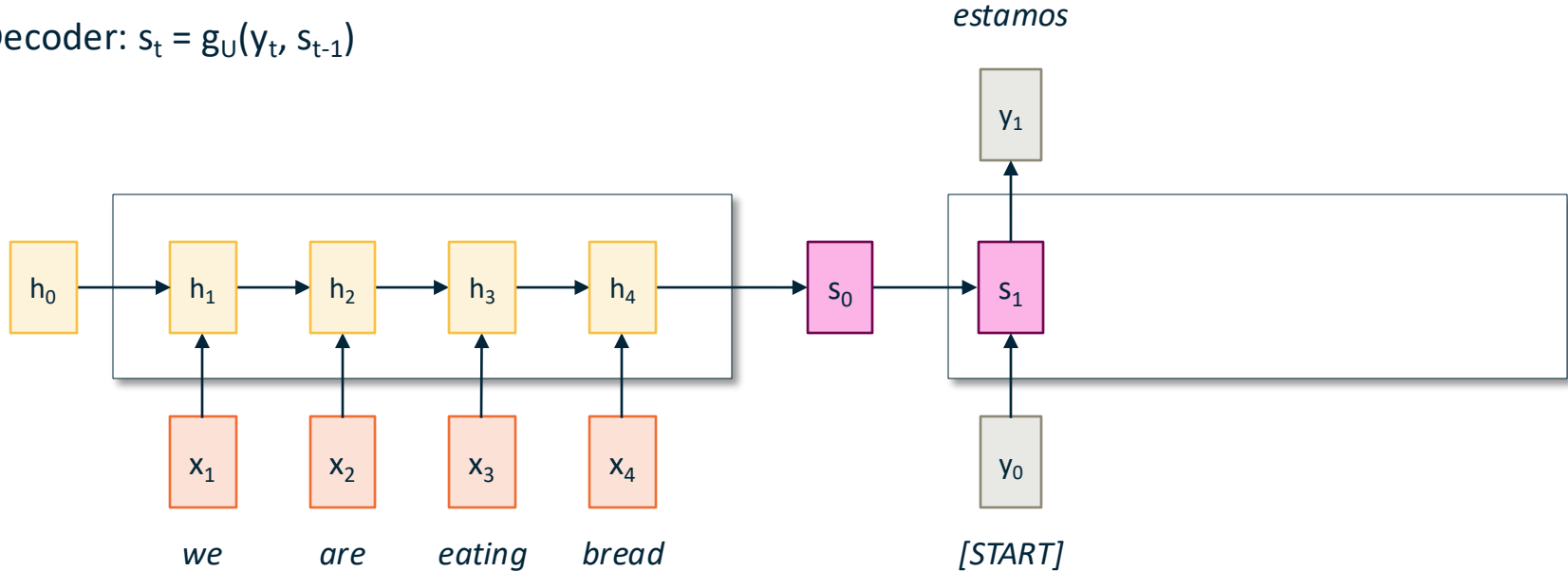
$$s_0 = h_4$$



Machine Translation with RNNs

Encoder: $h_t = f_W(x_t, h_{t-1})$

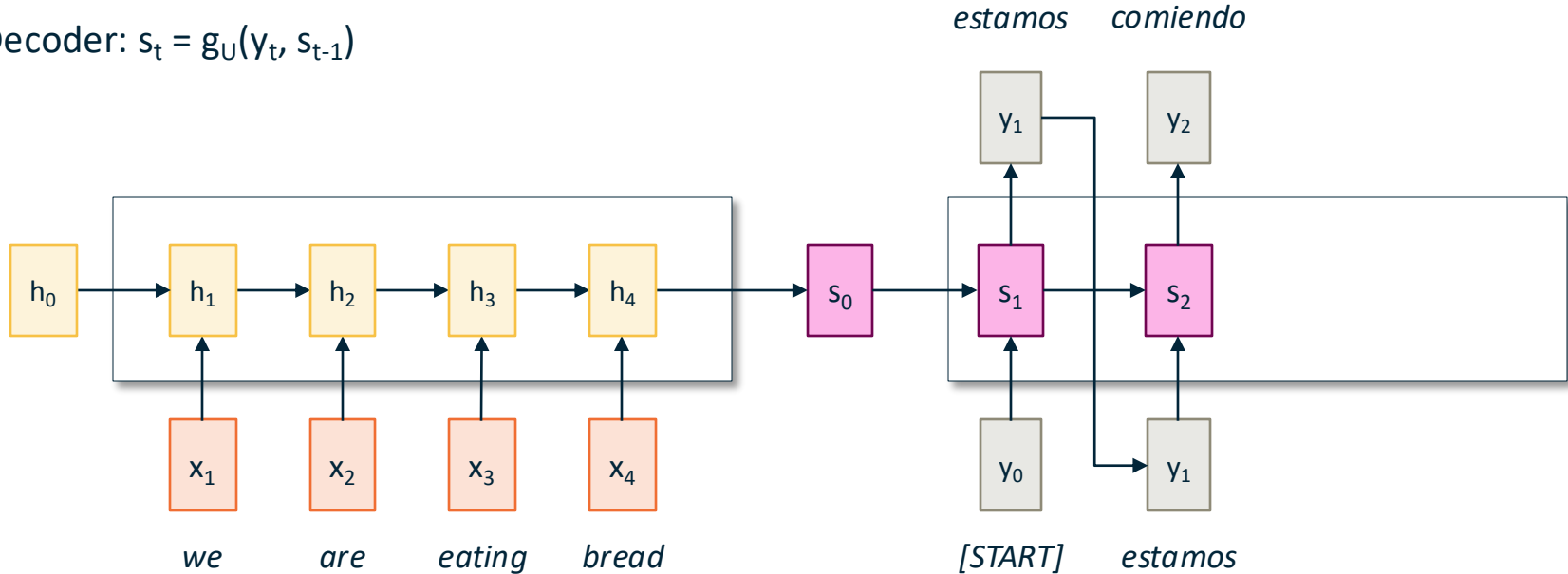
Decoder: $s_t = g_U(y_t, s_{t-1})$



Machine Translation with RNNs

Encoder: $h_t = f_W(x_t, h_{t-1})$

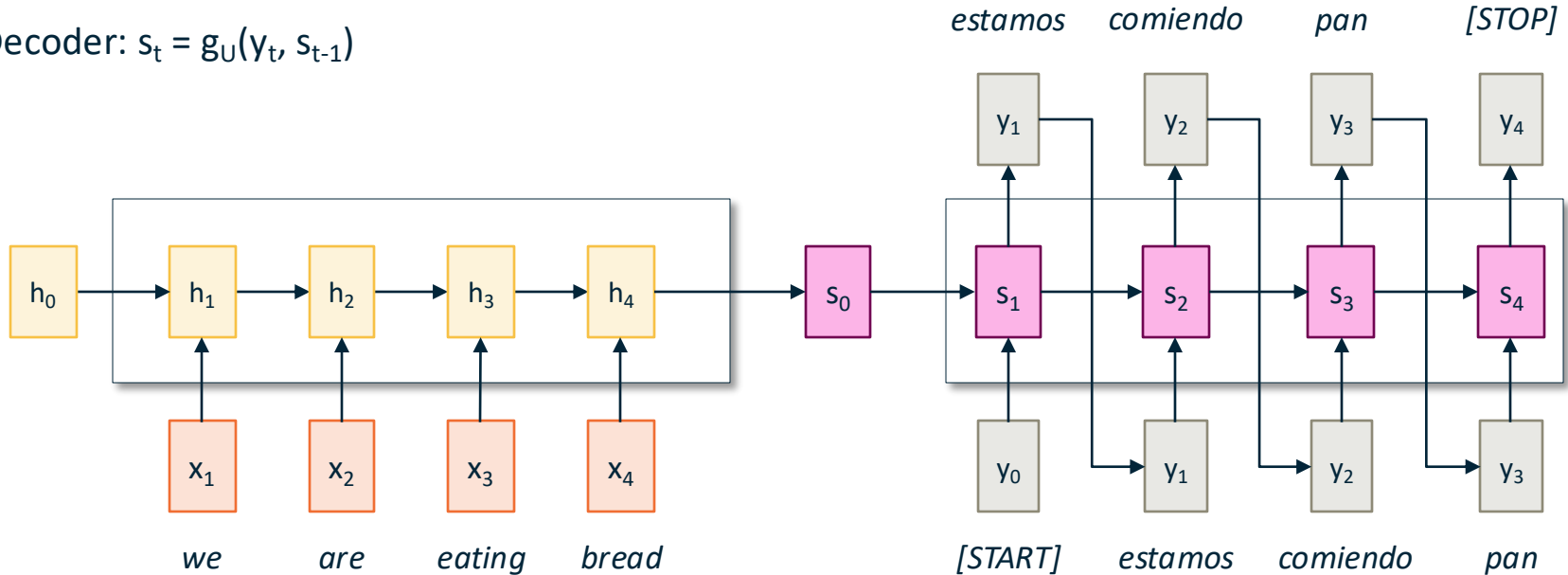
Decoder: $s_t = g_U(y_t, s_{t-1})$



Machine Translation with RNNs

Encoder: $h_t = f_W(x_t, h_{t-1})$

Decoder: $s_t = g_U(y_t, s_{t-1})$



Machine Translation with RNNs

Encoder: $h_t = f_W(x_t, h_{t-1})$

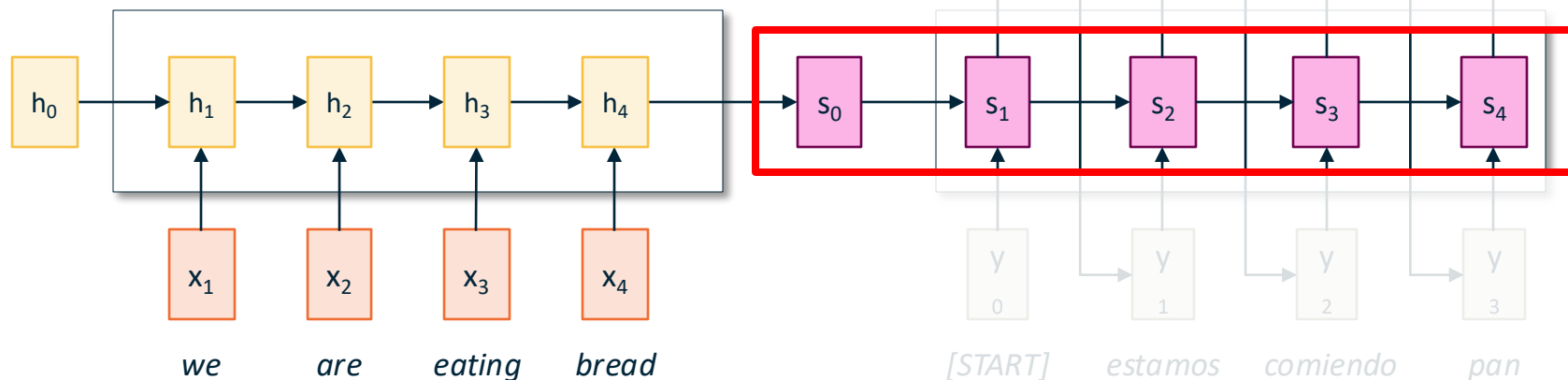
Decoder: $s_t = g_U(y_t, s_{t-1})$

Problem: s_i is used to both

(1) **encode input sequence**

(2) **maintain decoder state.**

Very difficult!



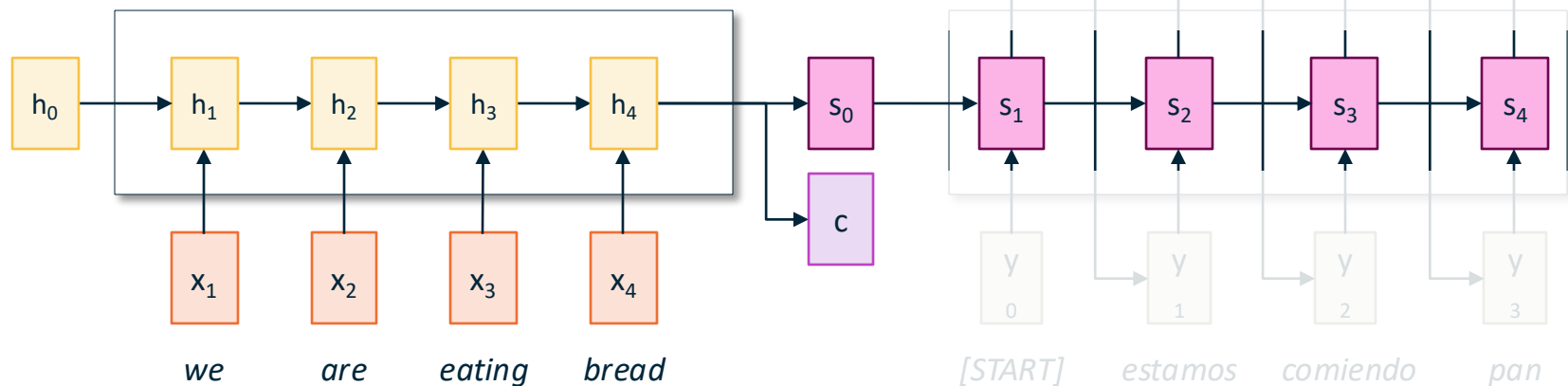
Machine Translation with RNNs

Encoder: $h_t = f_W(x_t, h_{t-1})$

Decoder: $s_t = g_U(y_t, s_{t-1}, \mathbf{c})$

Solution: add a context
vector $\mathbf{c} = h_4$ and
generate s_0 from h_4

\mathbf{c} doesn't change during
decoding



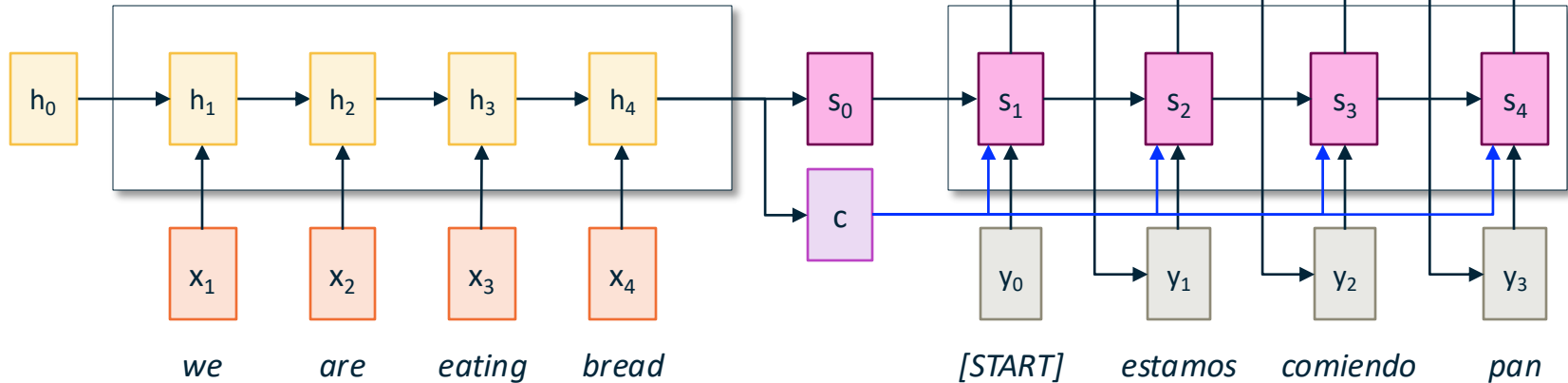
Machine Translation with RNNs

Encoder: $h_t = f_W(x_t, h_{t-1})$

Decoder: $s_t = g_U(y_t, s_{t-1}, \mathbf{c})$

Solution: add a context
vector $\mathbf{c} = h_4$ and
generate s_0 from h_4

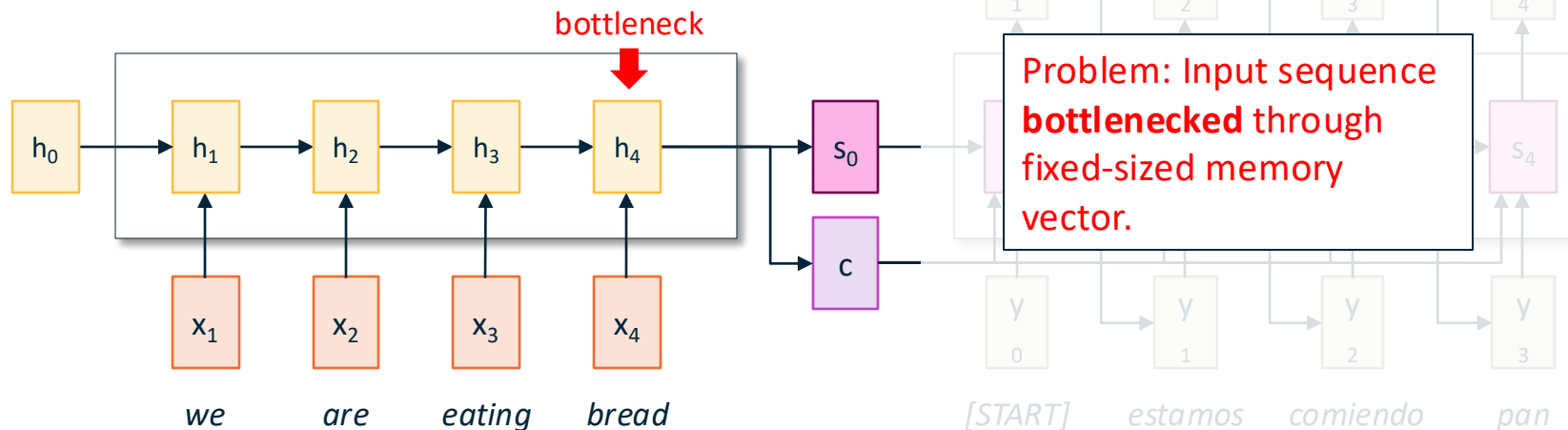
\mathbf{c} doesn't change during
decoding



Machine Translation with RNNs

Encoder: $h_t = f_W(x_t, h_{t-1})$

Decoder: $s_t = g_U(y_t, s_{t-1}, \mathbf{c})$



Problem with Recurrent-style Models (RNN, LSTM, GRU, etc.)

Learning to memorize is still hard, especially for ultra-long sequences!

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$
$$c_t = f \odot c_{t-1} + i \odot g$$
$$h_t = o \odot \tanh(c_t)$$

Essentially trying to tune W such that the memory cell c can retain **important information** for **arbitrary future prediction problems**.

Example (Q&A):

[... (20-page long transcript)]. Q: What did the CEO say about their competitor company? ...

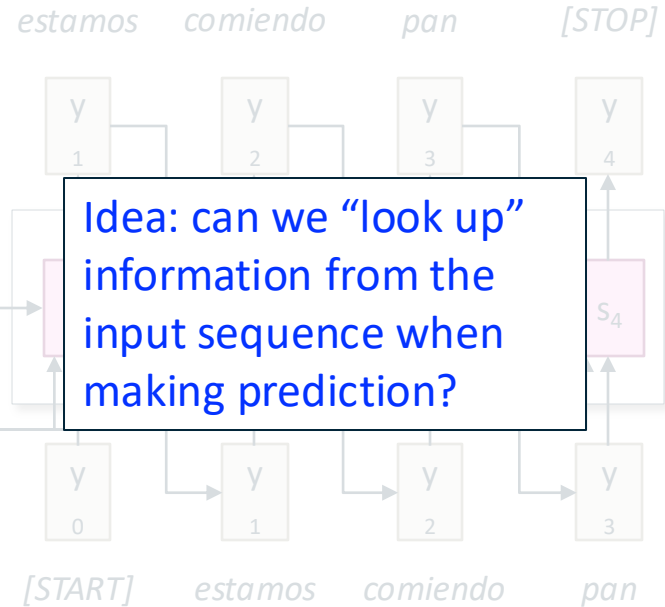
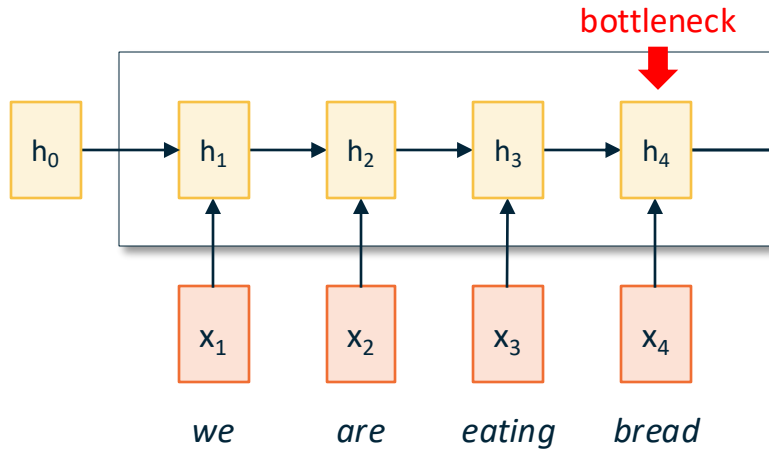
[... (same 20-page transcript)]. Q: How many times did the journalist use the word “interesting”? ...

Very difficult learning problem!

Machine Translation with RNNs

Encoder: $h_t = f_W(x_t, h_{t-1})$

Decoder: $s_t = g_U(y_t, s_{t-1}, \mathbf{c})$

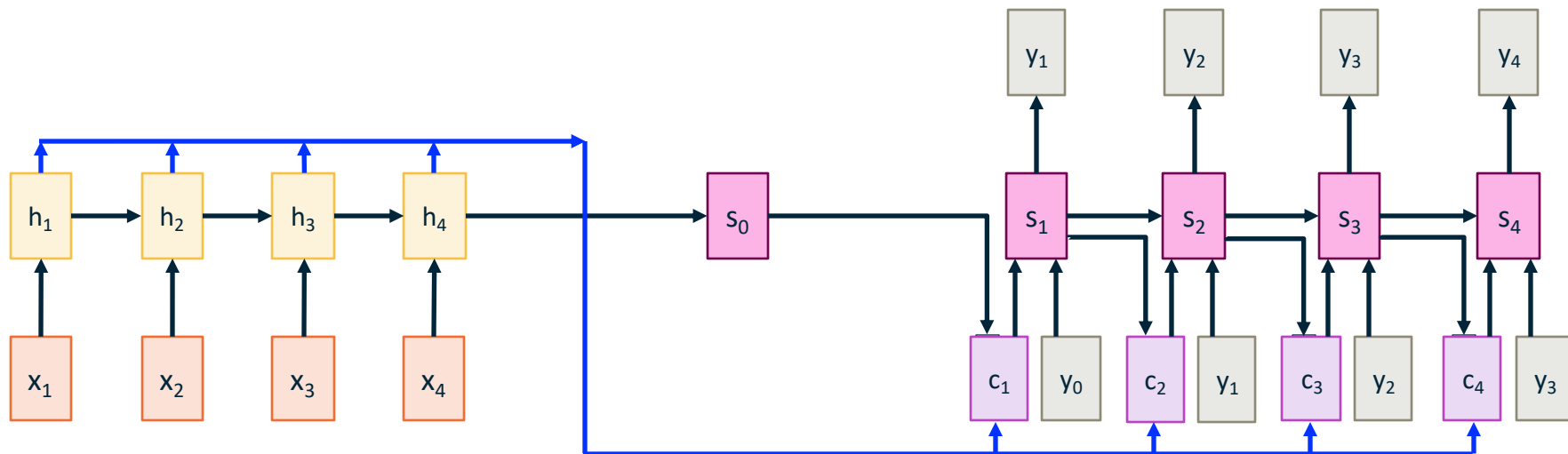


Machine Translation with RNNs **and Attention**

Conceptually, Attention is to **adaptively extract information** from input sequence based on the current decoding step

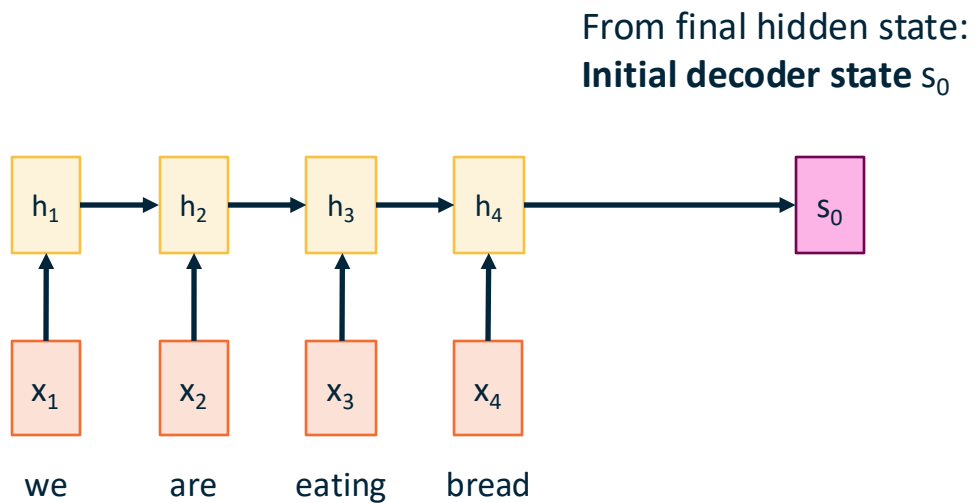
Machine Translation with RNNs **and Attention**

Conceptually, Attention is to **adaptively extract information** from input sequence based on the current decoding step



Goal: Adaptive context related to each prediction step

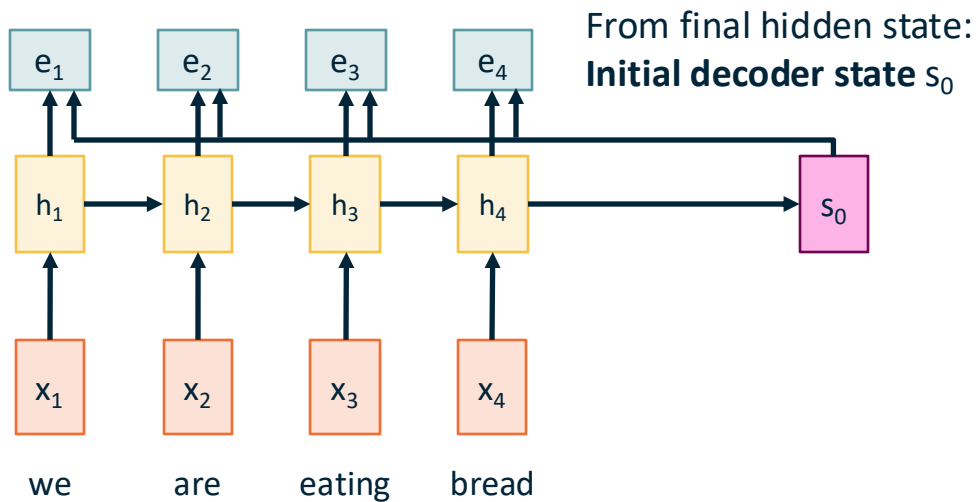
Machine Translation with RNNs **and Attention**



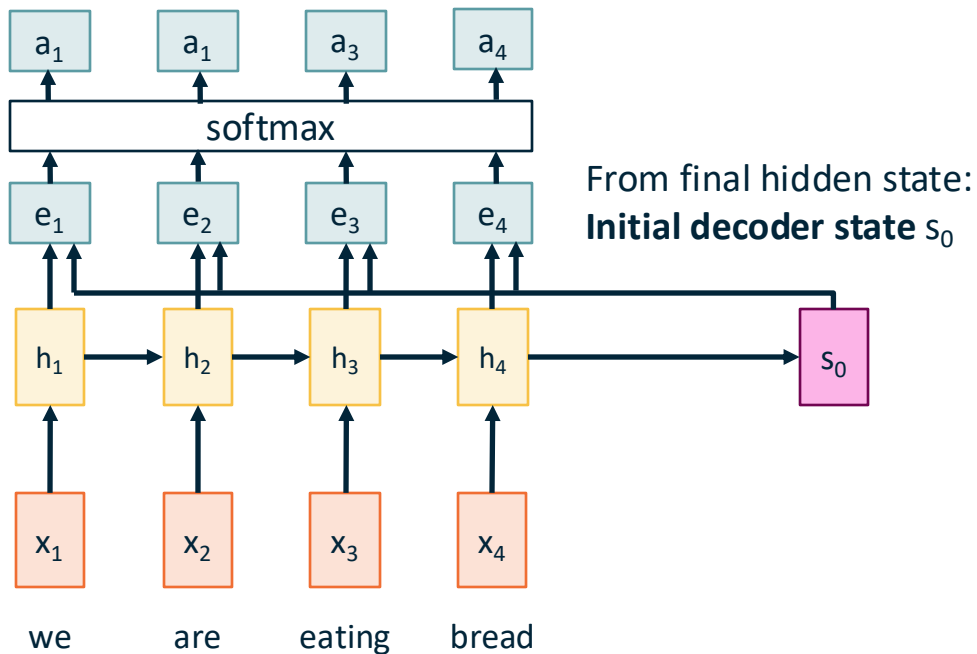
Machine Translation with RNNs **and Attention**

Compute **affinity scores**

$$e_{t,i} = f_{\text{att}}(s_{t-1}, h_i)$$



Machine Translation with RNNs **and Attention**



Compute **affinity scores**

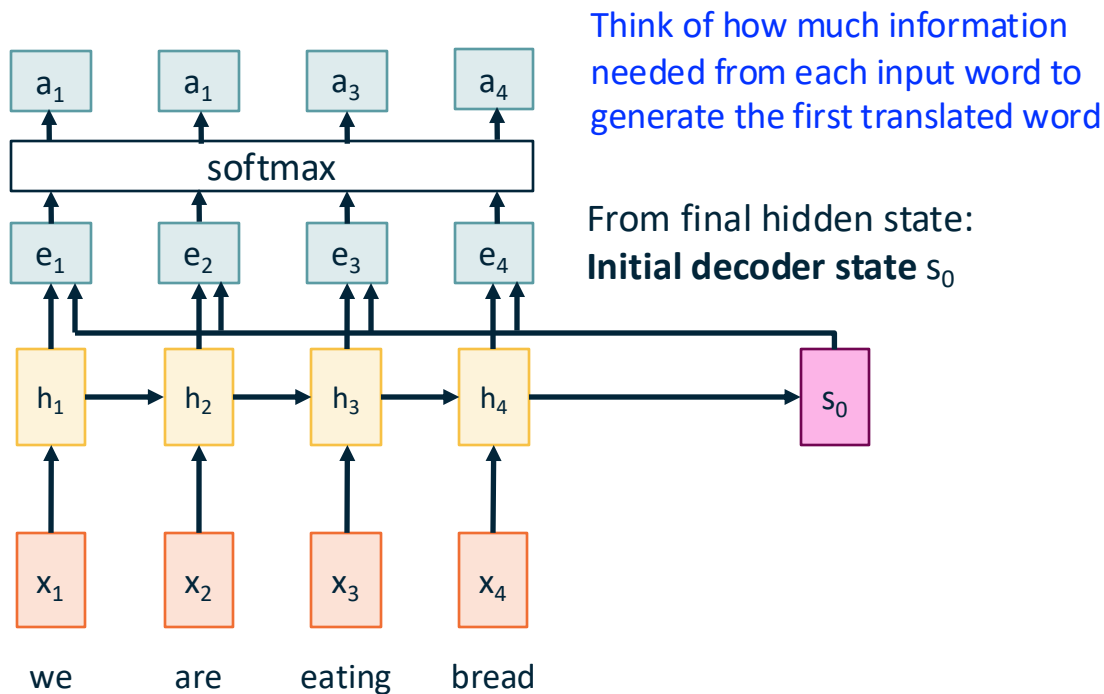
$$e_{t,i} = f_{\text{att}}(s_{t-1}, h_i)$$

Normalize to get

attention weights

$$0 < a_{t,i} < 1 \quad \sum_i a_{t,i} = 1$$

Machine Translation with RNNs **and Attention**



Compute **affinity scores**

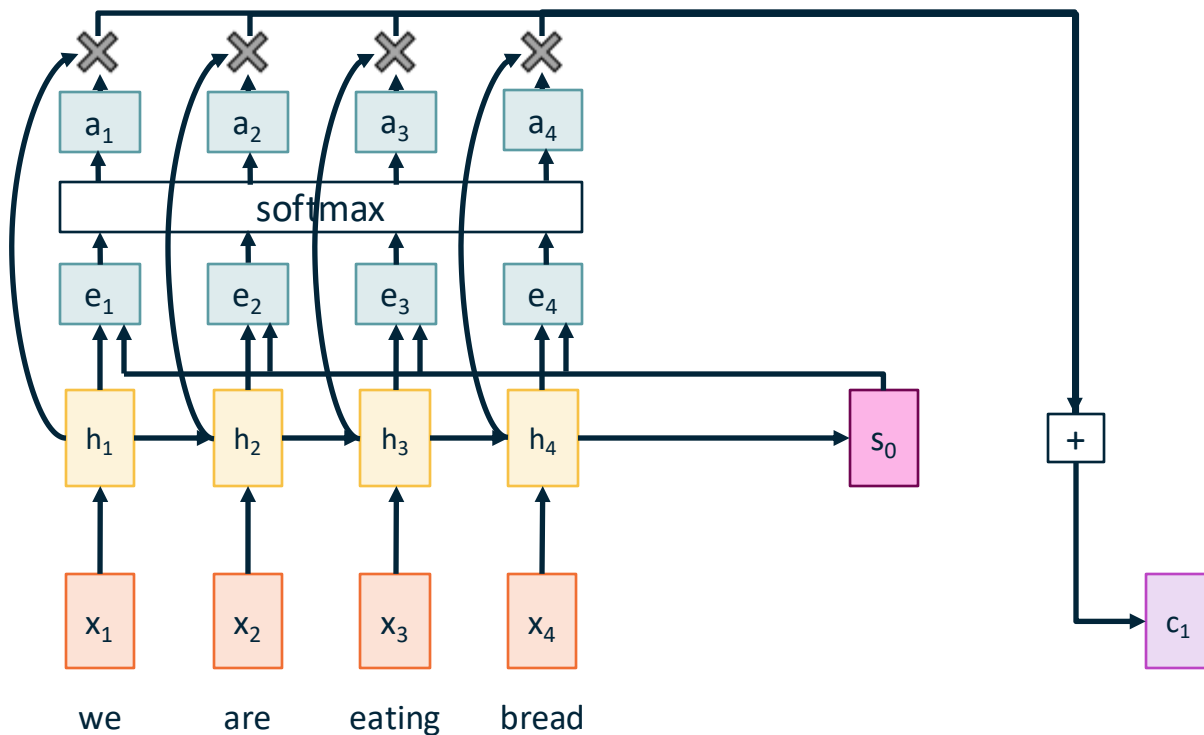
$$e_{t,i} = f_{\text{att}}(s_{t-1}, h_i)$$

Normalize to get

attention weights

$$0 < a_{t,i} < 1 \quad \sum_i a_{t,i} = 1$$

Machine Translation with RNNs **and Attention**



Compute **affinity scores**

$$e_{t,i} = f_{\text{att}}(s_{t-1}, h_i)$$

Normalize to get

attention weights

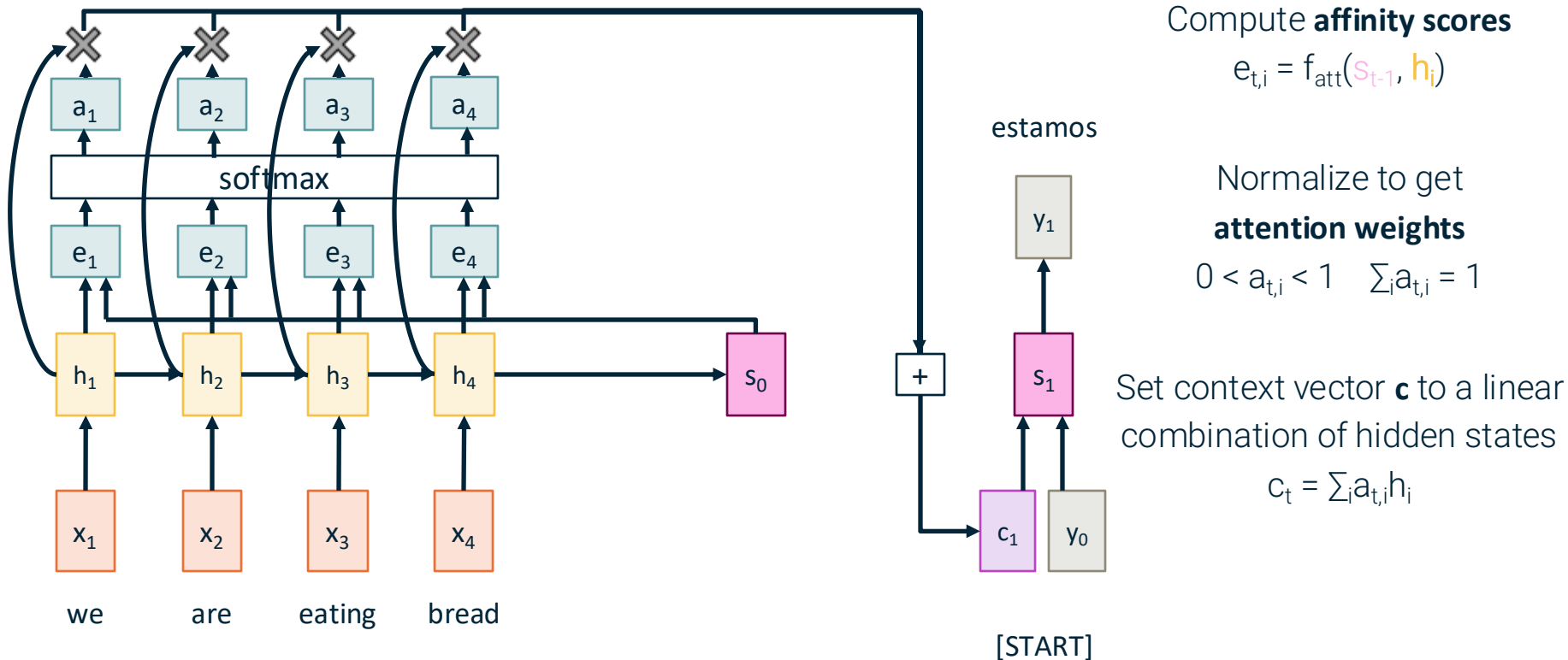
$$0 < a_{t,i} < 1 \quad \sum_i a_{t,i} = 1$$

Set context vector \mathbf{c} to a linear combination of hidden states

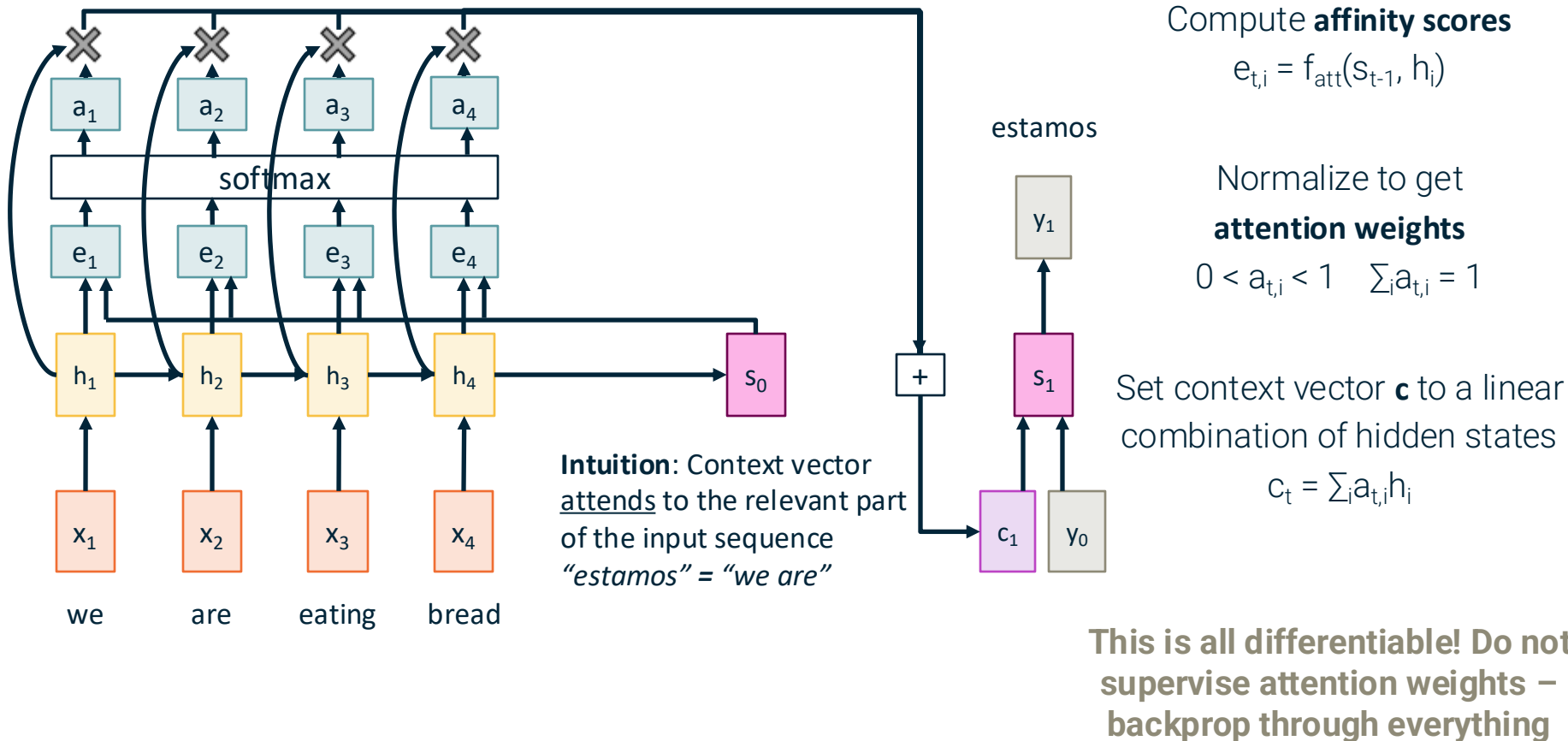
$$c_t = \sum_i a_{t,i} h_i$$

“Summarize the input sequence related to translating the t -th word”

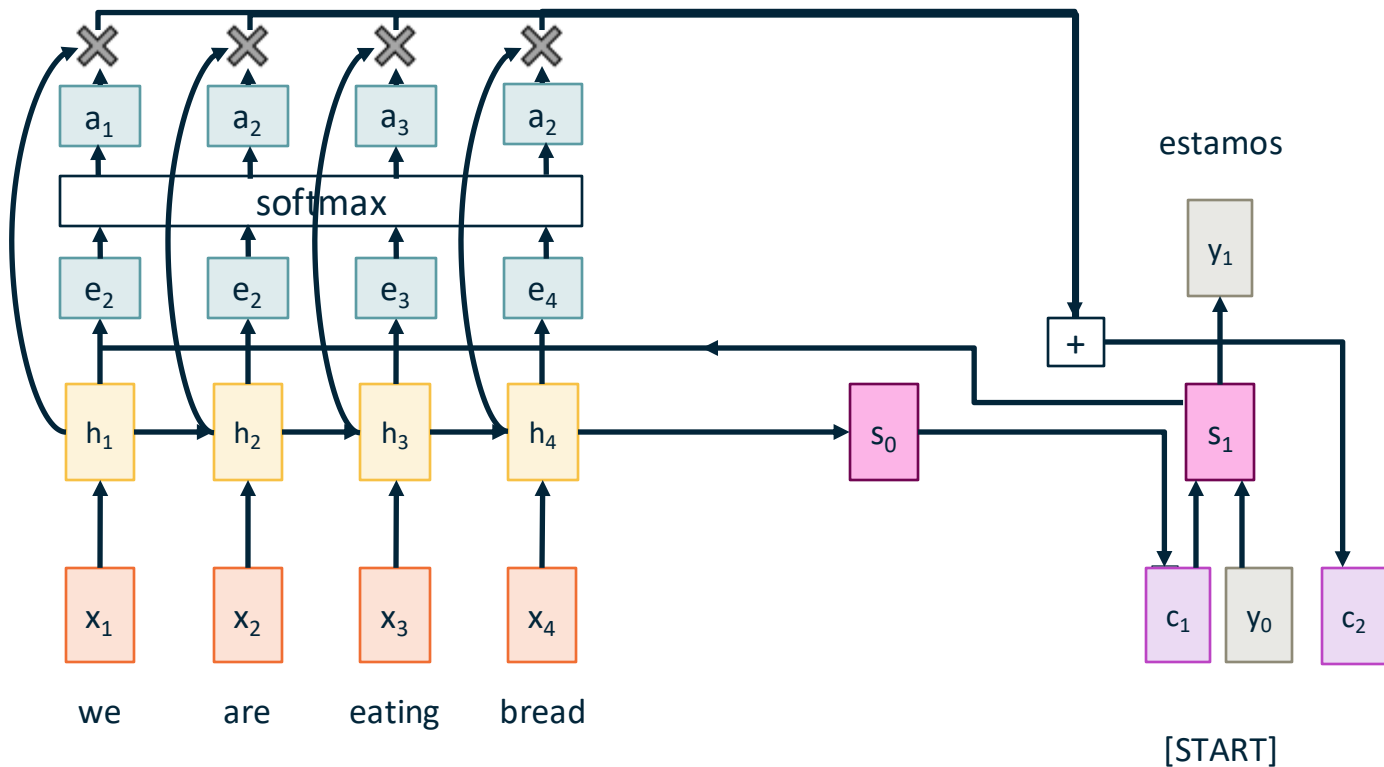
Machine Translation with RNNs **and Attention**



Machine Translation with RNNs **and Attention**

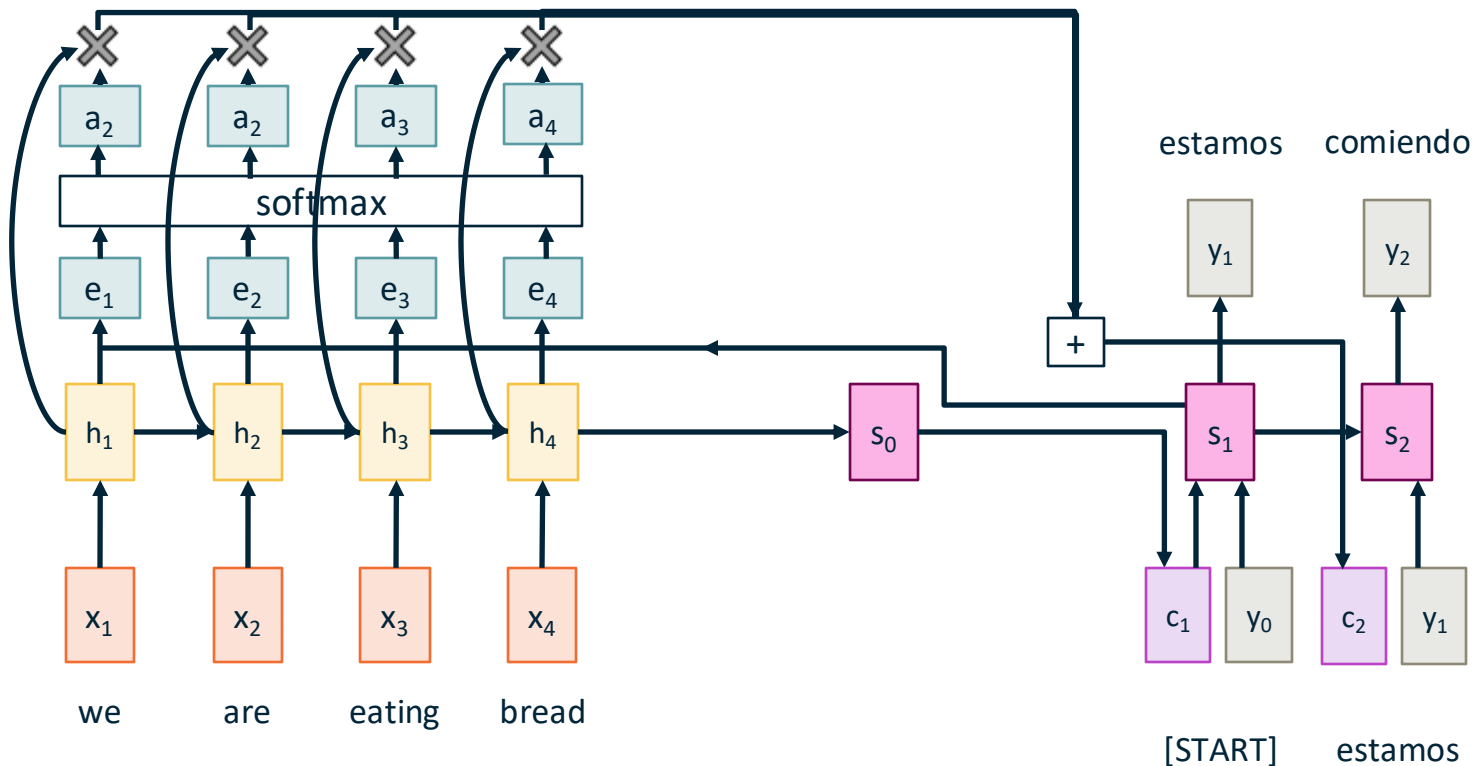


Machine Translation with RNNs **and Attention**



Repeat: Use s_1
to compute
attention and
get the new
context vector
 c_2

Machine Translation with RNNs **and Attention**

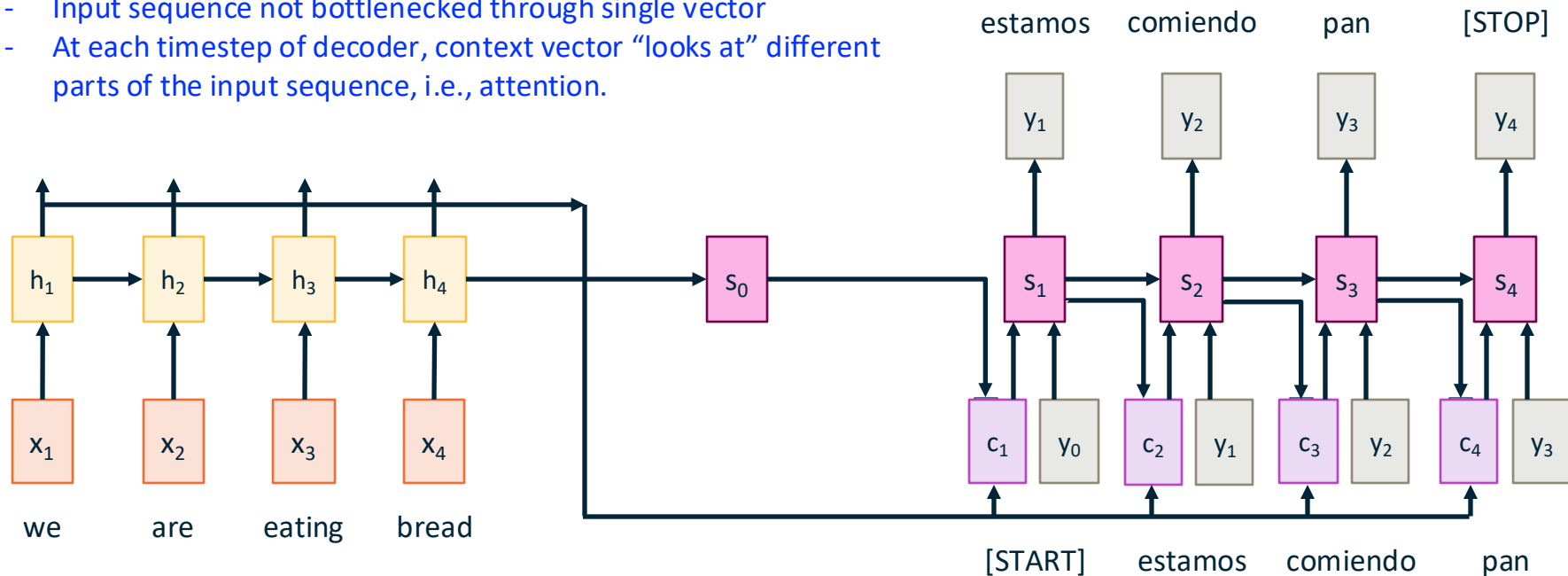


Repeat: Use s_1
to compute
attention and
get the new
context vector
 c_2
Use c_2 to
compute s_2, y_2

Machine Translation with RNNs **and Attention**

Use a different context vector in each timestep of decoder

- Input sequence not bottlenecked through single vector
- At each timestep of decoder, context vector “looks at” different parts of the input sequence, i.e., attention.



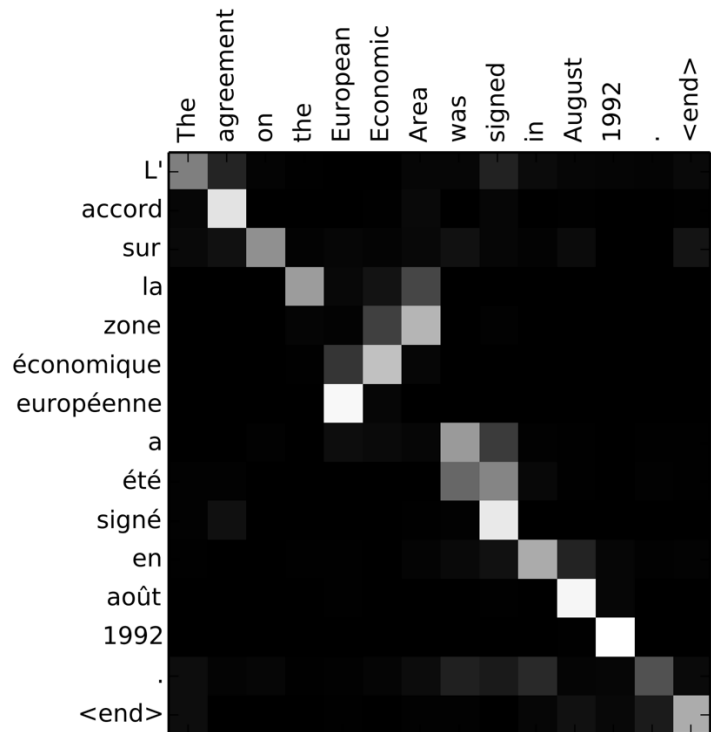
Machine Translation with RNNs **and Attention**

Example: English to French translation

Input: “The agreement on the European Economic Area was signed in August 1992.”

Output: “L’accord sur la zone économique européenne a été signé en août 1992.”

Visualize attention weights $a_{t,i}$



Machine Translation with RNNs **and Attention**

Example: English to French translation

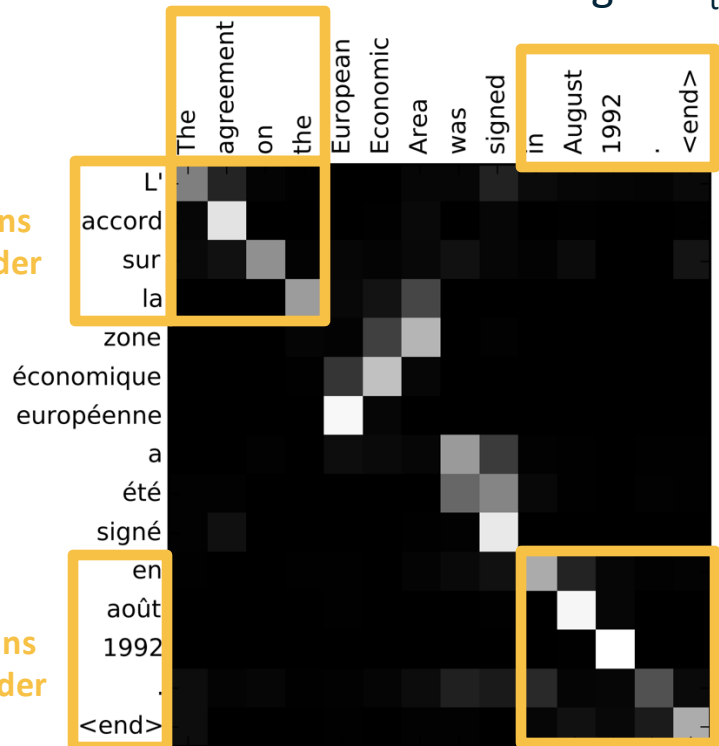
Input: “**The agreement on the** European Economic Area was signed **in August 1992.**”

Output: “**L’accord sur la** zone économique européenne a été signé **en août 1992.**”

Diagonal attention means words correspond in order

Diagonal attention means words correspond in order

Visualize attention weights $a_{t,i}$



Machine Translation with RNNs **and Attention**

Example: English to French translation

Input: “**The agreement on the** European Economic Area was signed **in August 1992.**”

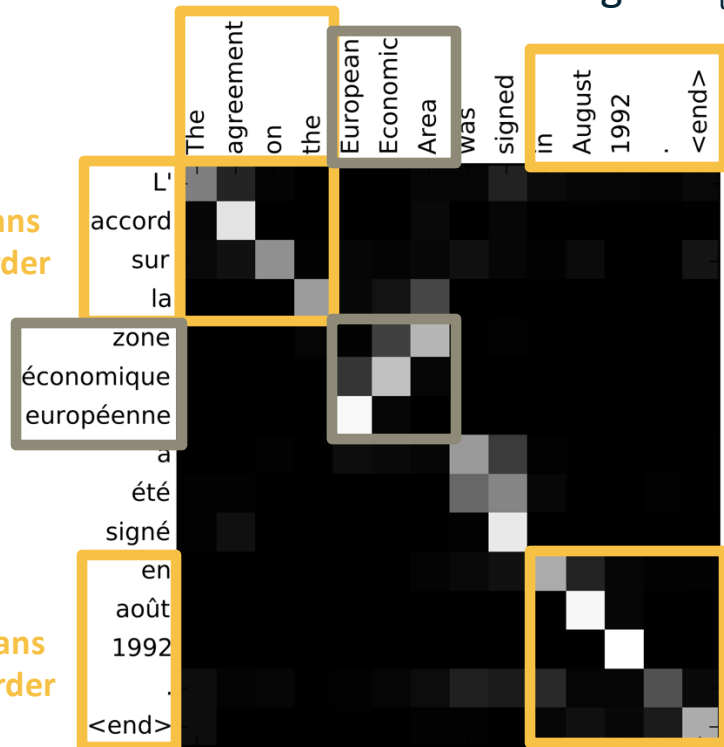
Output: “**L'accord sur la zone** économique européenne a été signé **en août 1992.**”

Diagonal attention means words correspond in order

Attention figures out different word orders

Diagonal attention means words correspond in order

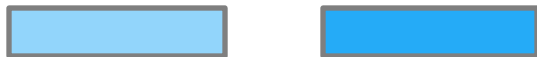
Visualize attention weights $a_{t,i}$



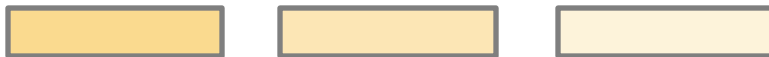
So ... What is Attention?

The attention mechanism in allows a model to **focus on specific parts of an input sequence** when processing or generating output, **dynamically weighting the importance of different elements**.

A set of query vectors:



A set of input vectors:



An attention function:

$$f(\text{yellow bar}, \text{light blue bar}) = 0.7$$

An attention map:

0.7	0.3	0.0
0.8	0.1	0.1

An aggregation function:

$$\text{orange bar} = 0.7 * \text{dark yellow bar} + 0.3 * \text{medium yellow bar} + 0.0 * \text{light yellow bar}$$

Aside: Differentiable Neural Computer/ Neural Turing Machine

Article | Published: 12 October 2016

Hybrid computing using a neural network with dynamic external memory

[Alex Graves](#) , [Greg Wayne](#) , [Malcolm Reynolds](#), [Tim Harley](#), [Ivo Danihelka](#), [Agnieszka Grabska-Barwińska](#), [Sergio Gómez Colmenarejo](#), [Edward Grefenstette](#), [Tiago Ramalho](#), [John Agapiou](#), [Adrià Puigdomènech Badia](#), [Karl Moritz Hermann](#), [Yori Zwols](#), [Georg Ostrovski](#), [Adam Cain](#), [Helen King](#), [Christopher Summerfield](#), [Phil Blunsom](#), [Koray Kavukcuoglu](#) & [Demis Hassabis](#)

[Nature](#) **538**, 471–476 (2016) | [Cite this article](#)

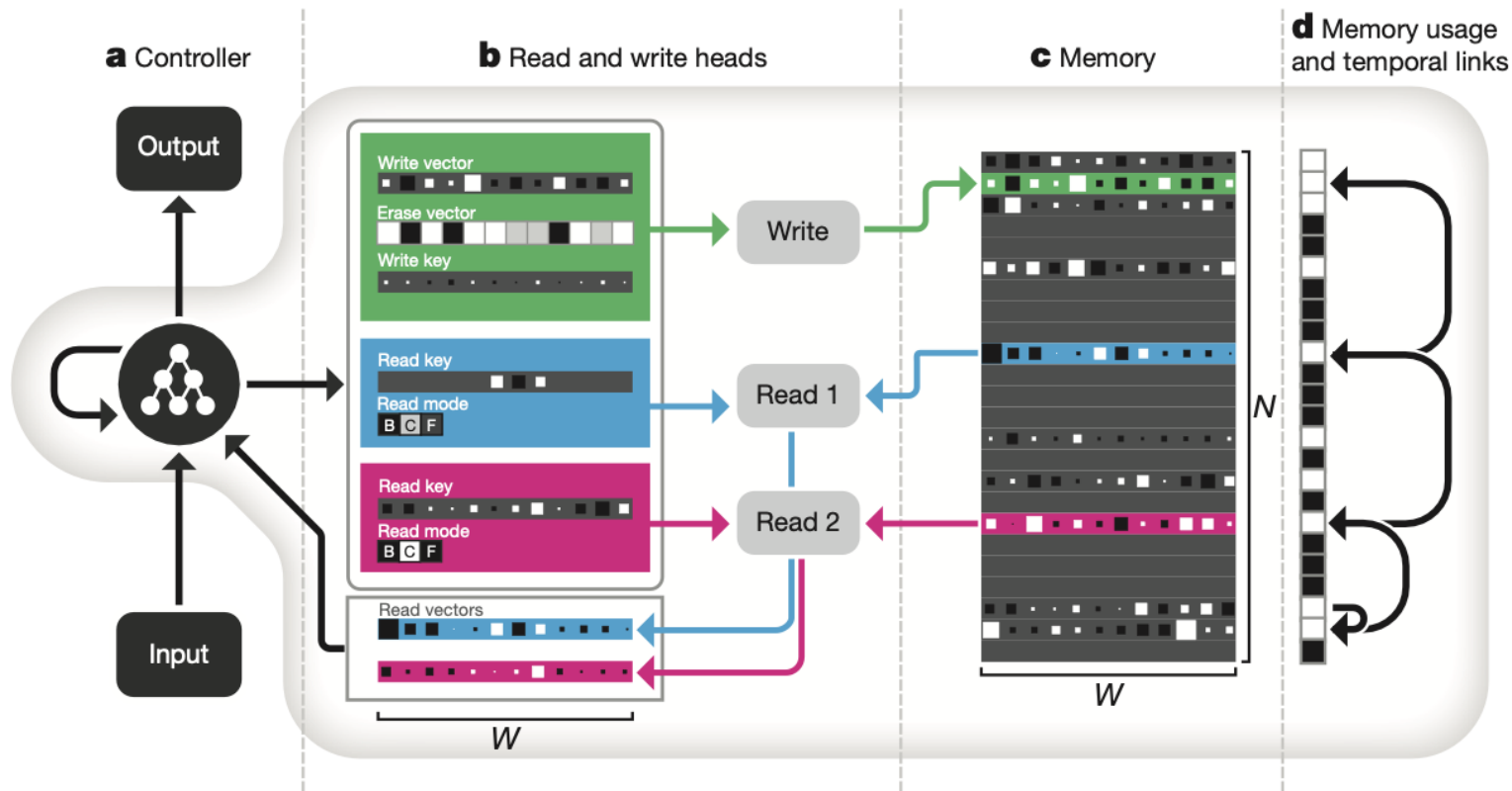
130k Accesses | **1030** Citations | **1270** Altmetric | [Metrics](#)

Can we make a neural network that behaves like a computer?
Can we make a computer-like mechanism fully-differentiable
and parameterized?

Aside: Differentiable Neural Computer/ Neural Turing Machine

Artificial neural networks are remarkably adept at sensory processing, sequence learning and reinforcement learning, but are limited in their ability to represent variables and data structures and to store data over long timescales, owing to the lack of an external memory. Here we introduce a machine learning model called a differentiable neural computer (DNC), which consists of a neural network that can read from and write to an external memory matrix, analogous to the random-access memory in a conventional computer. Like a conventional computer, it can use its memory to represent and manipulate complex data structures, but, like a neural network, it can learn to do so from data. When trained with supervised learning, we demonstrate that a DNC can successfully answer synthetic questions designed to emulate reasoning and inference problems in natural language. We show that it

Aside: Differentiable Neural Computer/ Neural Turing Machine



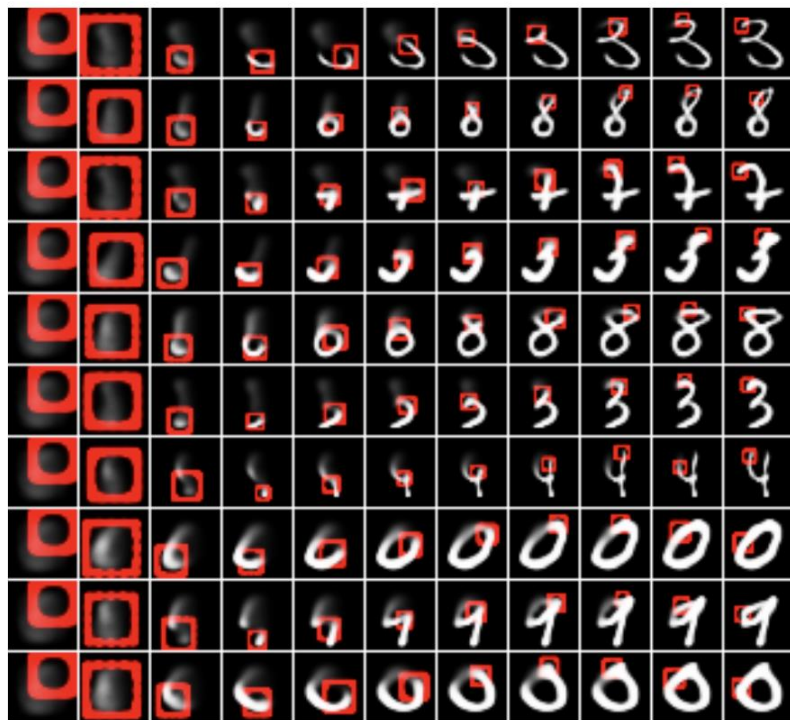
Aside: Attention for Generation

DRAW: A Recurrent Neural Network For Image Generation

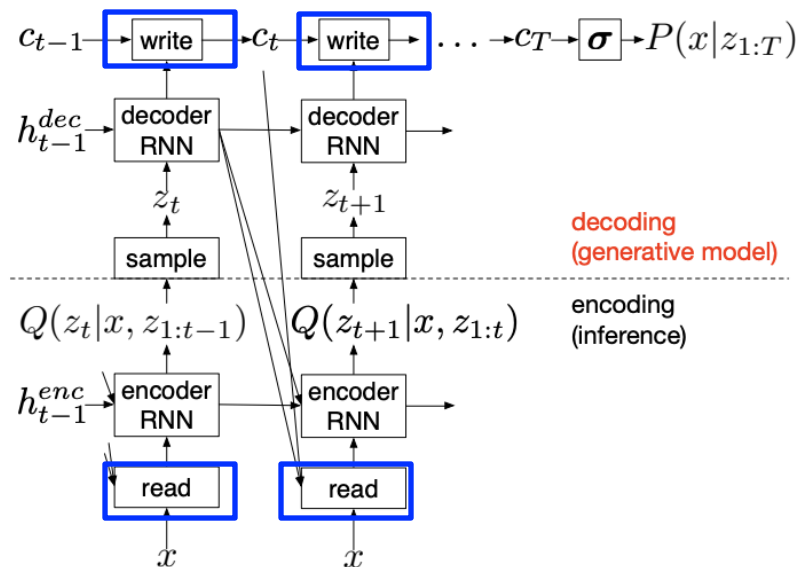
Karol Gregor
Ivo Danihelka
Alex Graves
Danilo Jimenez Rezende
Daan Wierstra
Google DeepMind

KAROLG@GOOGLE.COM
DANIHELKA@GOOGLE.COM
GRAVES@GOOGLE.COM
DANILOR@GOOGLE.COM
WIERSTRA@GOOGLE.COM

Aside: Attention for Generation



Time →



Optimize image-based attention (crop and blur) to learn to decide where to read and write (draw)

So ... What *is* Attention?

The attention mechanism in allows a model to **focus** on **specific parts of an input sequence** when processing or generating output, **dynamically weighting the importance of different elements**.

A set of query vectors:



A set of input vectors:



An attention function:

$$f(\text{yellow bar}, \text{light blue bar}) = 0.7$$

An attention map:

0.7	0.3	0.0
0.8	0.1	0.1

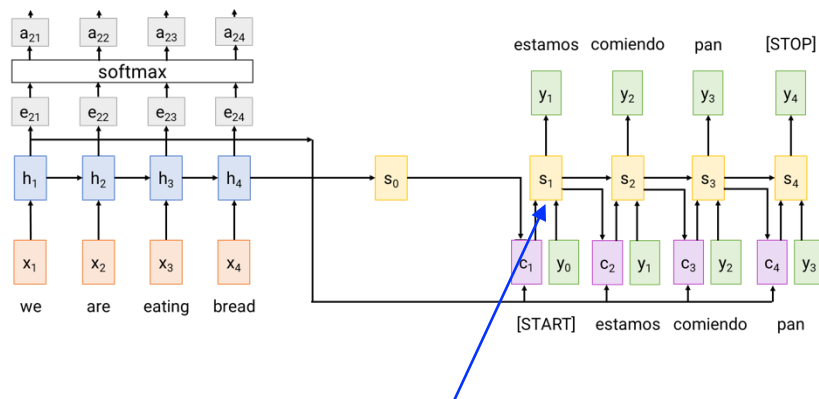
An aggregation function:

$$\text{orange bar} = 0.7 * \text{dark yellow bar} + 0.3 * \text{medium yellow bar} + 0.0 * \text{light yellow bar}$$

Attention Layer

Inputs:

Query vector: \mathbf{q} (Shape: D_0)



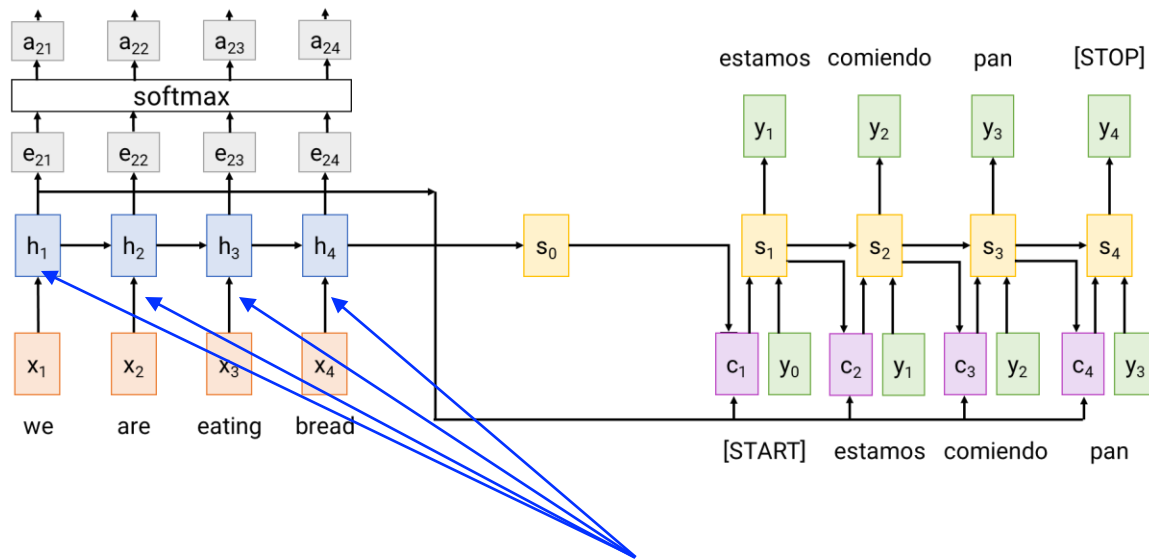
E.g., current decoding state

Attention Layer

Inputs:

Query vector: \mathbf{q} (Shape: D_Q)

Input vectors: \mathbf{X} (Shape: $N_X \times D_Q$)



E.g., all encoding states

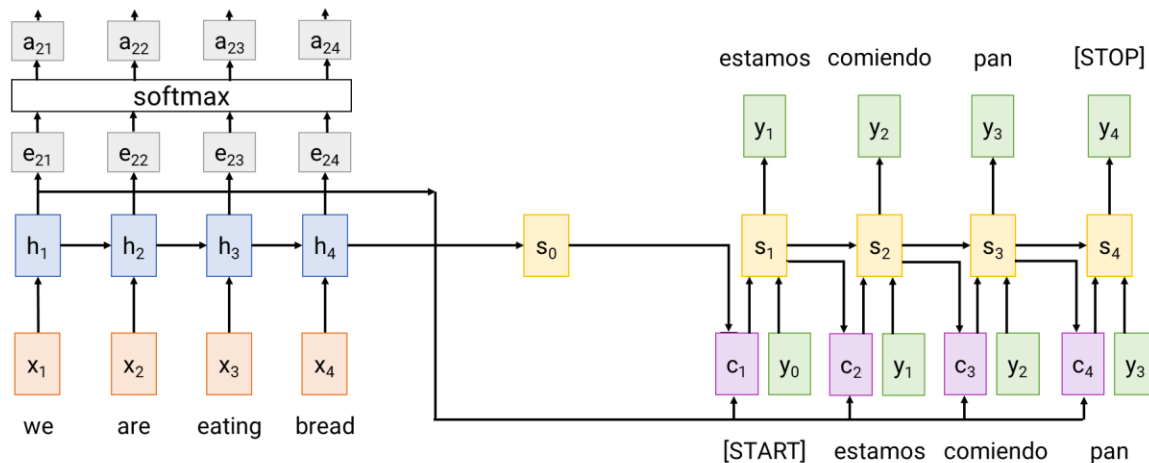
Attention Layer

Inputs:

Query vector: \mathbf{q} (Shape: D_Q)

Input vectors: \mathbf{X} (Shape: $N_X \times D_Q$)

Similarity function: f_{att}



Attention Layer

Inputs:

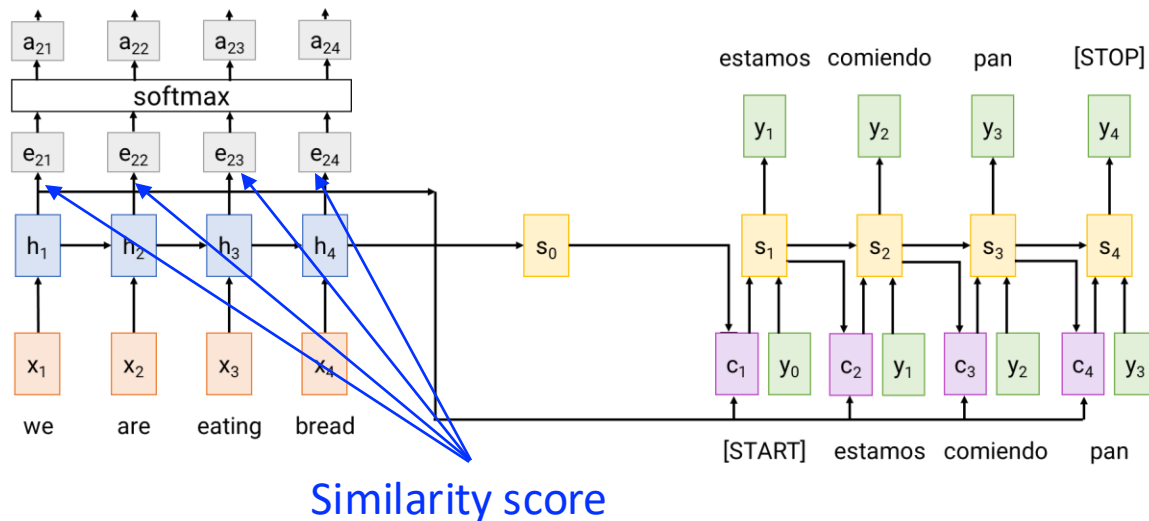
Query vector: \mathbf{q} (Shape: D_Q)

Input vectors: \mathbf{X} (Shape: $N_X \times D_Q$)

Similarity function: f_{att}

Computation:

Similarities: \mathbf{e} (Shape: N_X) $e_i = f_{\text{att}}(\mathbf{q}, \mathbf{x}_i)$



Attention Layer

Inputs:

Query vector: \mathbf{q} (Shape: D_Q)

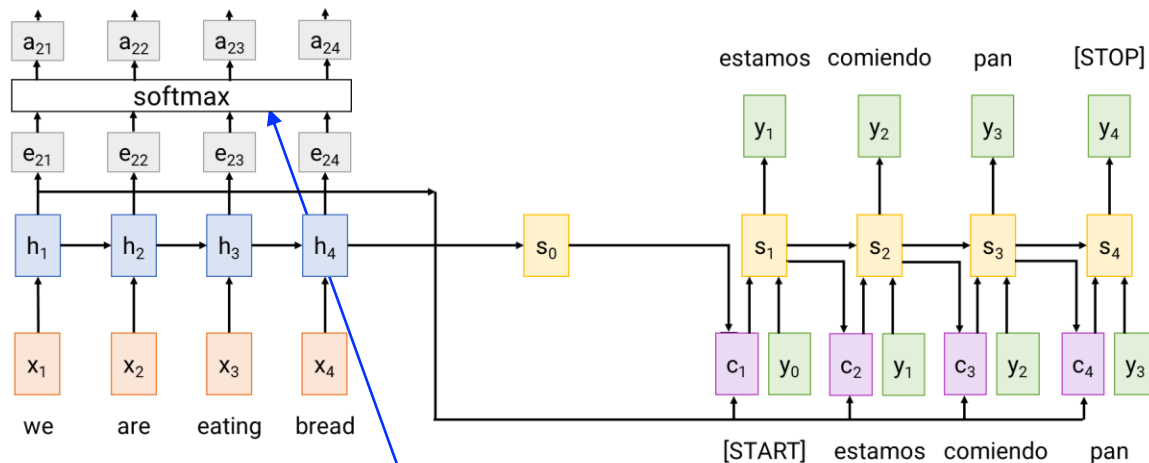
Input vectors: \mathbf{X} (Shape: $N_X \times D_Q$)

Similarity function: f_{att}

Computation:

Similarities: \mathbf{e} (Shape: N_X) $e_i = f_{\text{att}}(\mathbf{q}, \mathbf{x}_i)$

Attention weights: $\mathbf{a} = \text{softmax}(\mathbf{e})$ (Shape: N_X)



normalize attention weights

Attention Layer

Inputs:

Query vector: \mathbf{q} (Shape: D_Q)

Input vectors: \mathbf{X} (Shape: $N_X \times D_Q$)

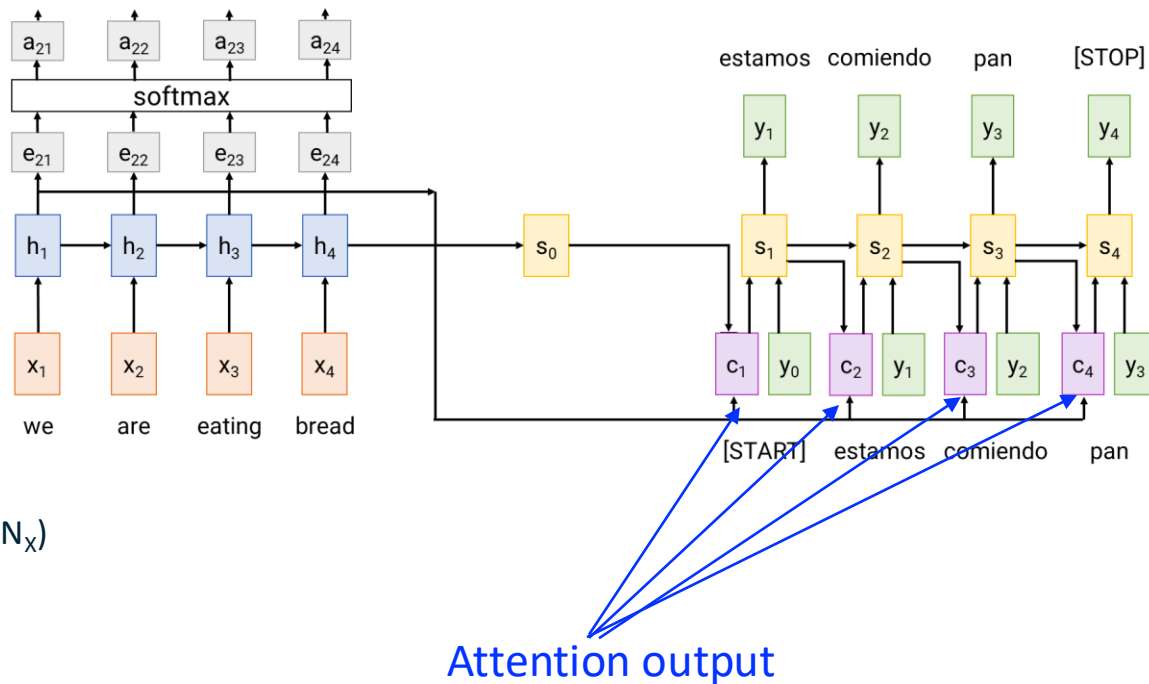
Similarity function: f_{att}

Computation:

Similarities: \mathbf{e} (Shape: N_X) $e_i = f_{\text{att}}(\mathbf{q}, \mathbf{x}_i)$

Attention weights: $\mathbf{a} = \text{softmax}(\mathbf{e})$ (Shape: N_X)

Output vector: $\mathbf{y} = \sum_i a_i \mathbf{x}_i$ (Shape: D_X)



Attention Layer

Inputs:

Query vector: \mathbf{q} (Shape: D_Q)

Input vectors: \mathbf{X} (Shape: $N_X \times D_O$)

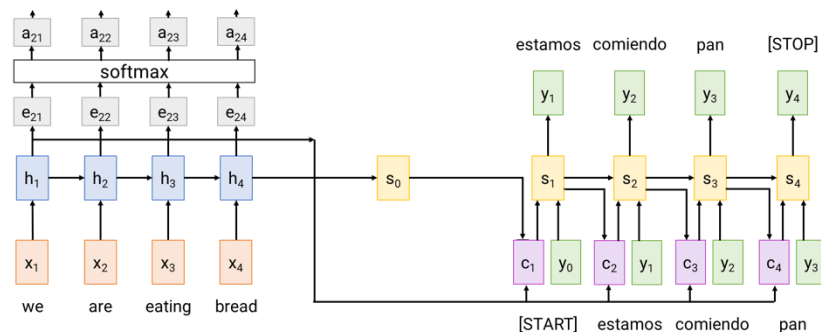
Similarity function: dot product

Computation:

Similarities: \mathbf{e} (Shape: N_X) $e_i = \mathbf{q} \cdot \mathbf{X}_i$

Attention weights: $\mathbf{a} = \text{softmax}(\mathbf{e})$ (Shape: N_X)

Output vector: $\mathbf{y} = \sum_i a_i \mathbf{X}_i$ (Shape: D_X)



Changes:

- Use dot product for similarity

Attention Layer

Inputs:

Query vector: \mathbf{q} (Shape: D_Q)

Input vectors: \mathbf{X} (Shape: $N_X \times D_Q$)

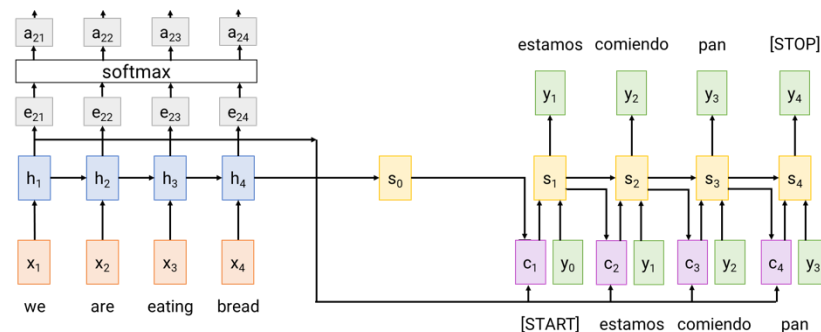
Similarity function: scaled dot product

Computation:

Similarities: \mathbf{e} (Shape: N_X) $e_i = \mathbf{q} \cdot \mathbf{X}_i / \sqrt{D_Q}$

Attention weights: $\mathbf{a} = \text{softmax}(\mathbf{e})$ (Shape: N_X)

Output vector: $\mathbf{y} = \sum_i a_i \mathbf{X}_i$ (Shape: D_X)



Changes:

- Use **scaled** dot product for similarity to account for the input size D_Q

Attention Layer

Inputs:

Query vectors: \mathbf{Q} (Shape: $N_Q \times D_Q$)

Input vectors: \mathbf{X} (Shape: $N_X \times D_Q$)

Computation:

Similarities: $E = \mathbf{QX}^T / \sqrt{D_Q}$ (Shape: $N_Q \times N_X$)

Attention matrix: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_Q \times N_X$)

Output vectors: $Y = AX$ (Shape: $N_Q \times D_X$) $Y_i = \sum_j A_{i,j} X_j$

Changes:

- Use **matrix multiplication** for similarity (many-to-many dot product)
- Multiple **query** vectors

Attention Layer

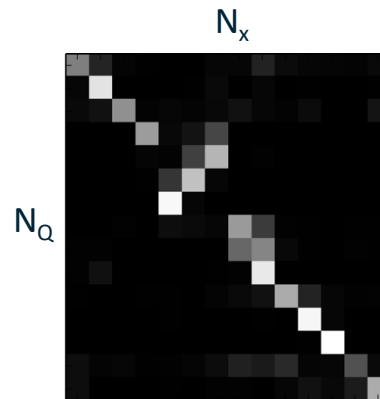
Inputs:

Query vectors: \mathbf{Q} (Shape: $N_Q \times D_Q$)

Input vectors: \mathbf{X} (Shape: $N_X \times D_X$)

Attention matrix (A)

Each row sums up to 1



Computation:

Similarities: $E = \mathbf{Q}\mathbf{X}^T / \sqrt{D_Q}$ (Shape: $N_Q \times N_X$)

Attention matrix: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_Q \times N_X$)

Output vectors: $\mathbf{Y} = \mathbf{A}\mathbf{X}$ (Shape: $N_Q \times D_X$) $Y_i = \sum_j A_{i,j} X_j$

Changes:

- Use **matrix multiplication** for similarity (many-to-many dot product)
- Multiple **query** vectors

Attention Layer

Inputs:

Query vectors: \mathbf{Q} (Shape: $N_Q \times D_Q$)

Input vectors: \mathbf{X} (Shape: $N_X \times D_Q$)

Computation:

Similarities: $E = \mathbf{QX}^T / \sqrt{D_Q}$ (Shape: $N_Q \times N_X$)

Attention matrix: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_Q \times N_X$)

Output vectors: $Y = AX$ (Shape: $N_Q \times D_X$) $Y_i = \sum_j A_{i,j} X_j$

Problem: use the same set of input vectors to compute both affinity and output.

Ideally, similarity/affinity should be based on some compact “signature vectors” that are efficient to compute.

Think of a key-value storage --- you typically use metadata (key) to look up an item (value), instead of using the item itself!

Changes:

- Use **matrix multiplication** for similarity (many-to-many dot product)
- Multiple **query** vectors

Attention Layer

Inputs:

Query vectors: \mathbf{Q} (Shape: $N_Q \times D_Q$)

Input vectors: \mathbf{X} (Shape: $N_X \times D_X$)

Key matrix: \mathbf{W}_K (Shape: $D_X \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_X \times D_V$)

Computation:

Similarities: $E = \mathbf{QX}^T / \sqrt{D_Q}$ (Shape: $N_Q \times N_X$)

Attention matrix: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_Q \times N_X$)

Output vectors: $Y = AX$ (Shape: $N_Q \times D_X$) $Y_i = \sum_j A_{i,j} X_j$

Problem: use the same set of input vectors to compute both affinity and output

Solution: project input to two sets of vectors: Keys (K) and Values (V).

Q,K,V attention: Compute attention matrix using Queries (Q) and Keys (K). Then compute output using attention and Values (V).

Changes:

- Use **matrix multiplication** for similarity (many-to-many dot product)
- Multiple **query** vectors
- Separate **key** and **value**

Attention Layer

Inputs:

Query vectors: \mathbf{Q} (Shape: $N_Q \times D_Q$)

Input vectors: \mathbf{X} (Shape: $N_X \times D_X$)

Key matrix: \mathbf{W}_K (Shape: $D_X \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_X \times D_V$)

Computation:

Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_X \times D_Q$)

Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_X \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_Q \times N_X$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_Q \times N_X$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_Q \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

Problem: use the same set of input vectors to compute both affinity and output

Solution: project input to two sets of vectors: Keys (K) and Values (V).

Q,K,V attention: Compute attention matrix using Queries (Q) and Keys (K). Then compute output using attention and Values (V).

Changes:

- Use **matrix multiplication** for similarity (many-to-many dot product)
- Multiple **query** vectors
- Separate **key** and **value**

Attention Layer

Inputs:

Query vectors: \mathbf{Q} (Shape: $N_Q \times D_Q$)

Input vectors: \mathbf{X} (Shape: $N_X \times D_X$)

Key matrix: \mathbf{W}_K (Shape: $D_X \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_X \times D_V$)

Computation:

Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_X \times D_Q$)

Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_X \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_Q \times N_X$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_Q \times N_X$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_Q \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

X_1

X_2

X_3

Q_1

Q_2

Q_3

Q_4

Attention Layer

Inputs:

Query vectors: \mathbf{Q} (Shape: $N_Q \times D_Q$)

Input vectors: \mathbf{X} (Shape: $N_X \times D_X$)

Key matrix: \mathbf{W}_K (Shape: $D_X \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_X \times D_V$)

Computation:

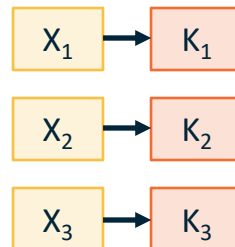
→ **Key vectors:** $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_X \times D_Q$)

Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_X \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_Q \times N_X$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_Q \times N_X$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_Q \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$



Attention Layer

Inputs:

Query vectors: \mathbf{Q} (Shape: $N_Q \times D_Q$)

Input vectors: \mathbf{X} (Shape: $N_X \times D_X$)

Key matrix: \mathbf{W}_K (Shape: $D_X \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_X \times D_V$)

Computation:

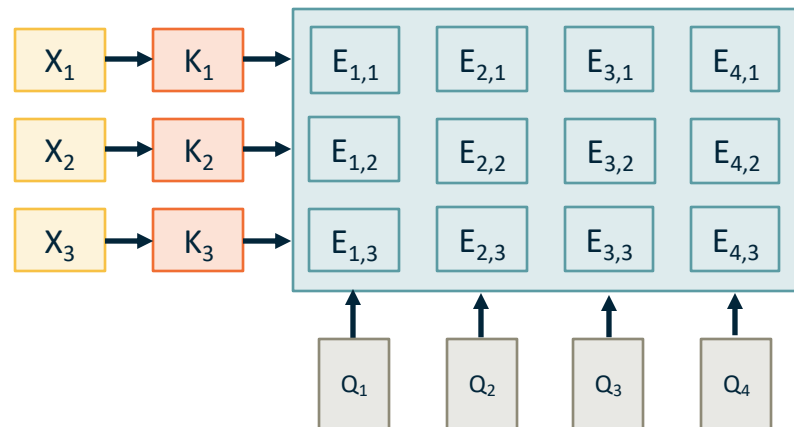
Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_X \times D_Q$)

Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_X \times D_V$)

→ Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_Q \times N_X$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_Q \times N_X$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_Q \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$



Attention Layer

Inputs:

Query vectors: \mathbf{Q} (Shape: $N_Q \times D_Q$)

Input vectors: \mathbf{X} (Shape: $N_X \times D_X$)

Key matrix: \mathbf{W}_K (Shape: $D_X \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_X \times D_V$)

Computation:

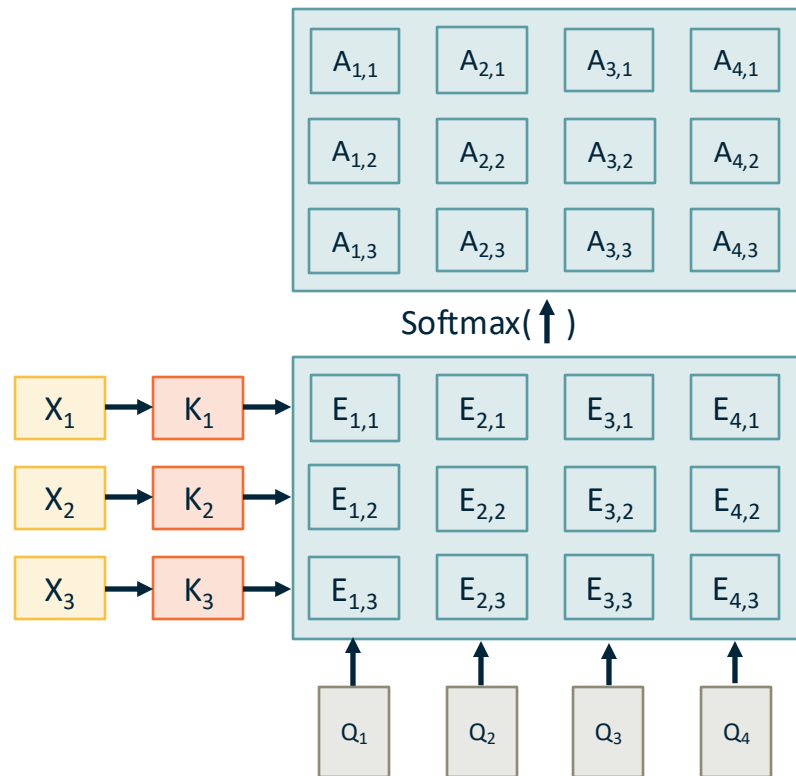
Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_X \times D_Q$)

Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_X \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_Q \times N_X$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

→ **Attention weights:** $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_Q \times N_X$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_Q \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$



Attention Layer

Inputs:

Query vectors: \mathbf{Q} (Shape: $N_Q \times D_Q$)

Input vectors: \mathbf{X} (Shape: $N_X \times D_X$)

Key matrix: \mathbf{W}_K (Shape: $D_X \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_X \times D_V$)

Computation:

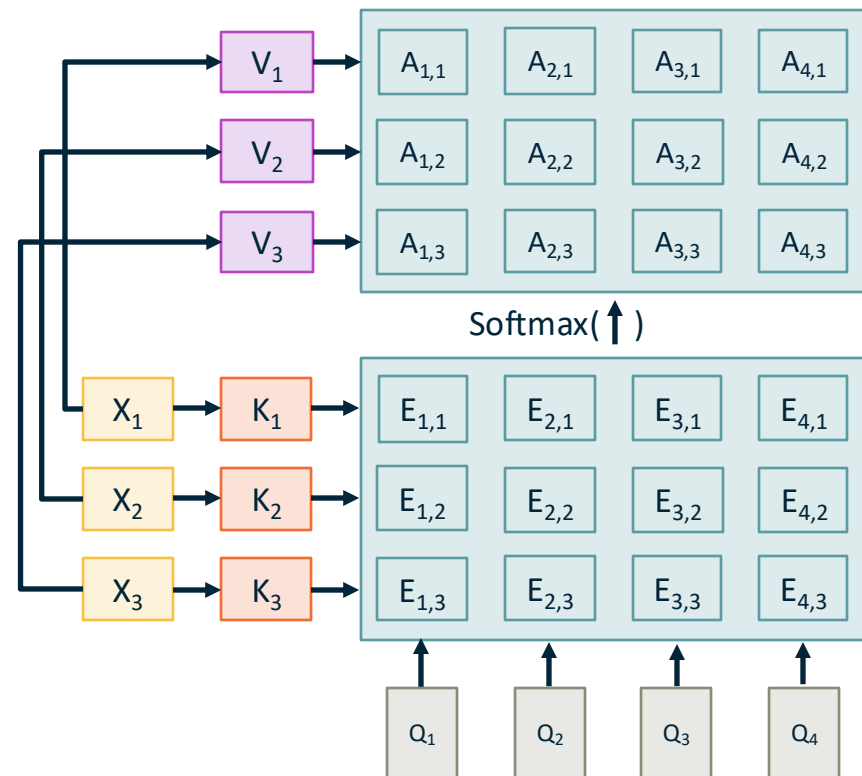
Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_X \times D_Q$)

→ Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_X \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_Q \times N_X$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_Q \times N_X$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_Q \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$



Attention Layer

Inputs:

Query vectors: \mathbf{Q} (Shape: $N_Q \times D_Q$)

Input vectors: \mathbf{X} (Shape: $N_X \times D_X$)

Key matrix: \mathbf{W}_K (Shape: $D_X \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_X \times D_V$)

Computation:

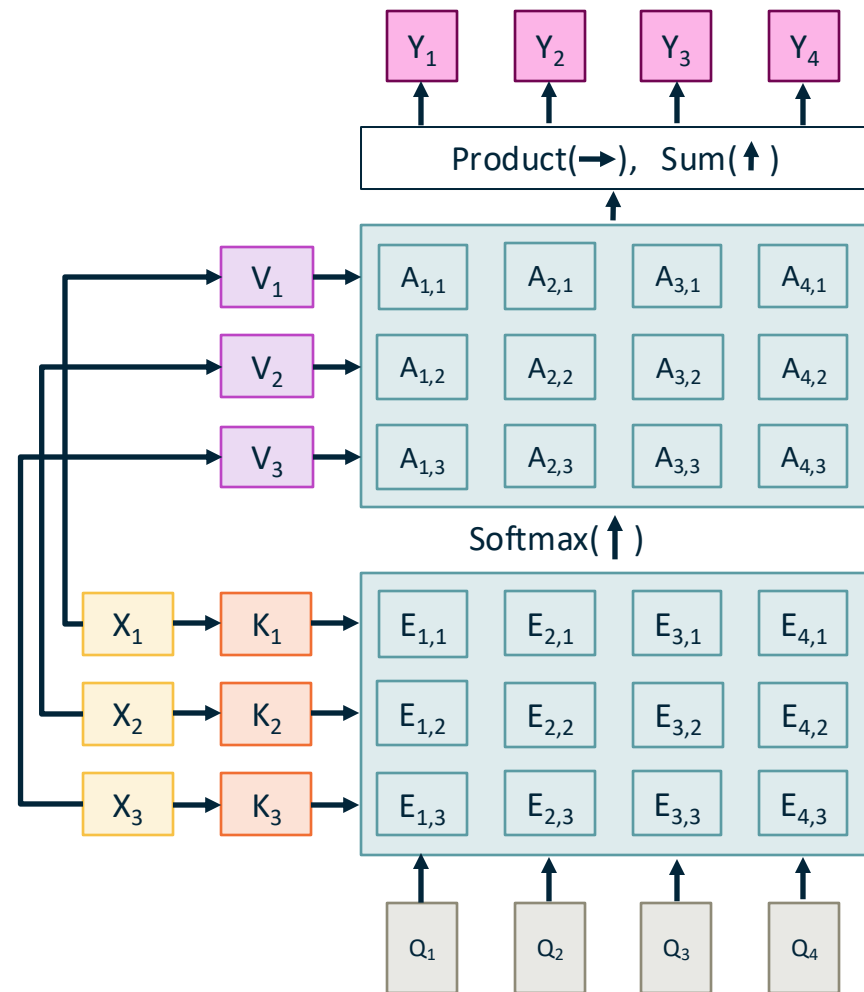
Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_X \times D_Q$)

Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_X \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_Q \times N_X$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_Q \times N_X$)

→ **Output vectors:** $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_Q \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$



Attention Layer

Inputs:

Query vectors: \mathbf{Q} (Shape: $N_Q \times D_Q$)

Input vectors: \mathbf{X} (Shape: $N_X \times D_X$)

Key matrix: \mathbf{W}_K (Shape: $D_X \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_X \times D_V$)

Computation:

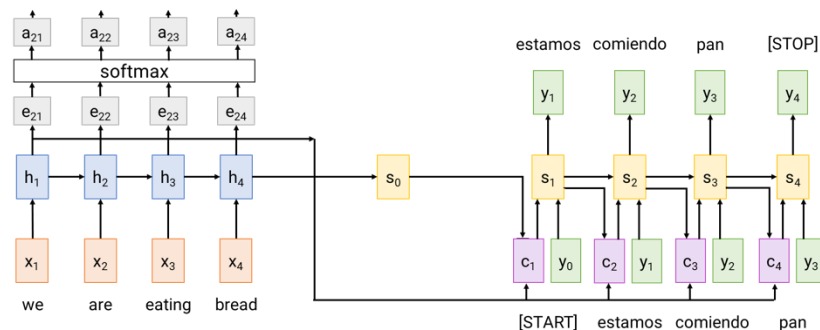
Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_X \times D_Q$)

Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_X \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_Q \times N_X$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_Q \times N_X$)

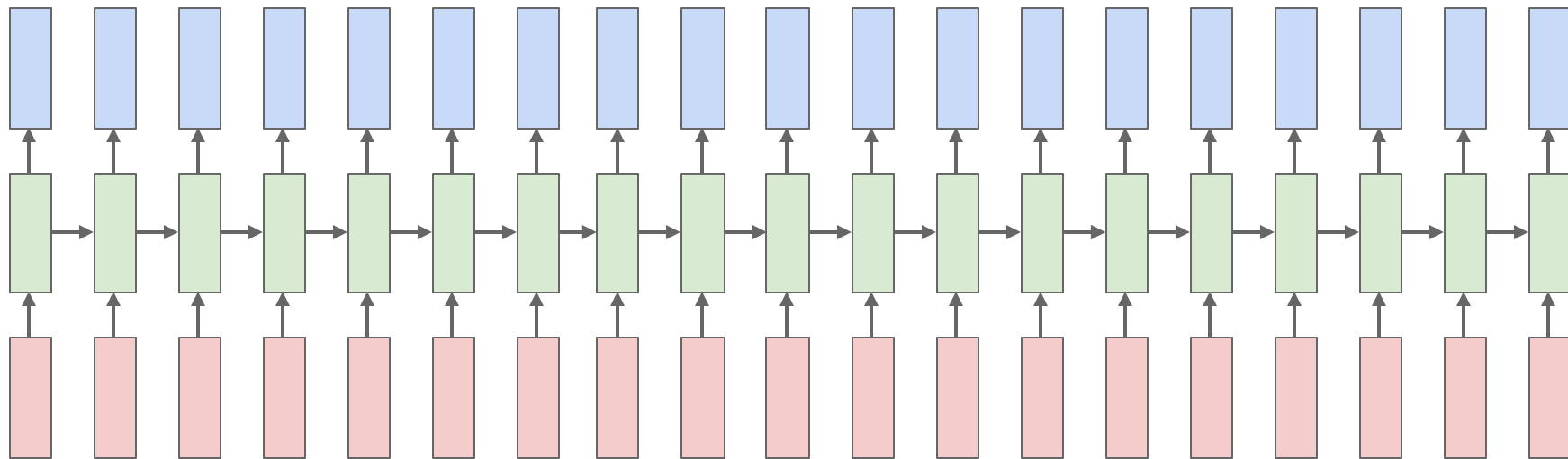
Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_Q \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$



Attention seems to be really powerful ...

Do we still need RNN?

RNN is bad at encoding long-range relationships!



Recurrent update can easily “forget” information

Attention Layer

Inputs:

Query vectors: \mathbf{Q} (Shape: $N_Q \times D_Q$)

Input vectors: \mathbf{X} (Shape: $N_X \times D_X$)

Key matrix: \mathbf{W}_K (Shape: $D_X \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_X \times D_V$)

Computation:

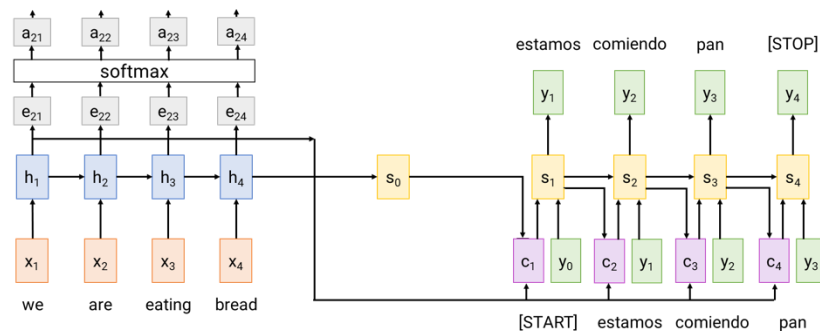
Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_X \times D_Q$)

Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_X \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_Q \times N_X$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_Q \times N_X$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_Q \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$



Attention seems to be really powerful ...

Do we still need RNN?

Can we use **only attention layers** to encode an entire sequence?

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
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

“The Transformer Paper”

Self-Attention Layer

Sequence encode -> use each input element as query!

Inputs:

Query vectors: \mathbf{Q} (Shape: $N_Q \times D_Q$)

Input vectors: \mathbf{X} (Shape: $N_X \times D_X$)

Key matrix: \mathbf{W}_K (Shape: $D_X \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_X \times D_V$)

Computation:

Query vectors: $\mathbf{Q} = \mathbf{XW}_Q$

Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_X \times D_Q$)

Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_X \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_X \times N_X$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_X \times N_X$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

Without recurrent-based encoding, still need to somehow represent inter-token connection in the input sequence.

Goal: encode the input sequence with only attention, without a recurrent network.

X_1

X_2

X_3

Self-Attention Layer

Sequence encode -> use each input element as query!

Inputs:

Query vectors: \mathbf{Q} (Shape: $N_Q \times D_Q$)

Input vectors: \mathbf{X} (Shape: $N_X \times D_X$)

Key matrix: \mathbf{W}_K (Shape: $D_X \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_X \times D_V$)

Query matrix: \mathbf{W}_Q (Shape: $D_X \times D_Q$)

Computation:

Query vectors: $\mathbf{Q} = \mathbf{XW}_Q$

Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_X \times D_Q$)

Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_X \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_X \times N_X$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_X \times N_X$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

Goal: encode the input sequence with only attention, without a recurrent network.

Encoding only -> no external queries

Use each element to query other elements in the same sequence, i.e., “self-attention”

X_1

X_2

X_3

Self-Attention Layer

Sequence encode -> use each input element as query!

Inputs:

Input vectors: \mathbf{X} (Shape: $N_x \times D_x$)

Key matrix: \mathbf{W}_K (Shape: $D_x \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_x \times D_V$)

Query matrix: \mathbf{W}_Q (Shape: $D_x \times D_Q$)

Computation:

→ **Query vectors:** $\mathbf{Q} = \mathbf{XW}_Q$

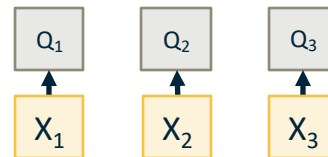
Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_x \times D_Q$)

Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_x \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_x \times N_x$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_x \times N_x$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_x \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$



Self-Attention Layer

Sequence encode -> use each input element as query!

Inputs:

Input vectors: \mathbf{X} (Shape: $N_x \times D_x$)

Key matrix: \mathbf{W}_K (Shape: $D_x \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_x \times D_V$)

Query matrix: \mathbf{W}_Q (Shape: $D_x \times D_Q$)

Computation:

Query vectors: $\mathbf{Q} = \mathbf{XW}_Q$

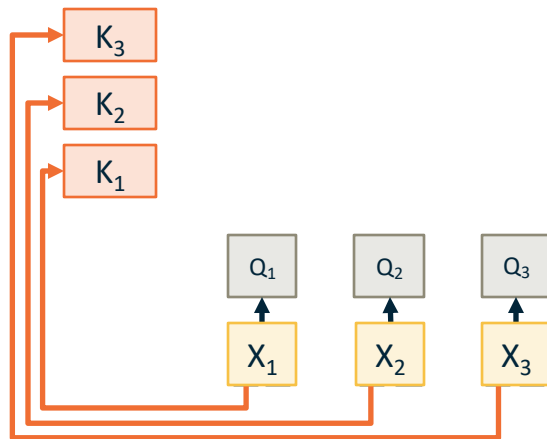
→ **Key vectors:** $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_x \times D_Q$)

Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_x \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_x \times N_x$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_x \times N_x$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_x \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$



Self-Attention Layer

Sequence encode -> use each input element as query!

Inputs:

Input vectors: \mathbf{X} (Shape: $N_x \times D_x$)

Key matrix: \mathbf{W}_K (Shape: $D_x \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_x \times D_V$)

Query matrix: \mathbf{W}_Q (Shape: $D_x \times D_Q$)

Computation:

Query vectors: $\mathbf{Q} = \mathbf{XW}_Q$

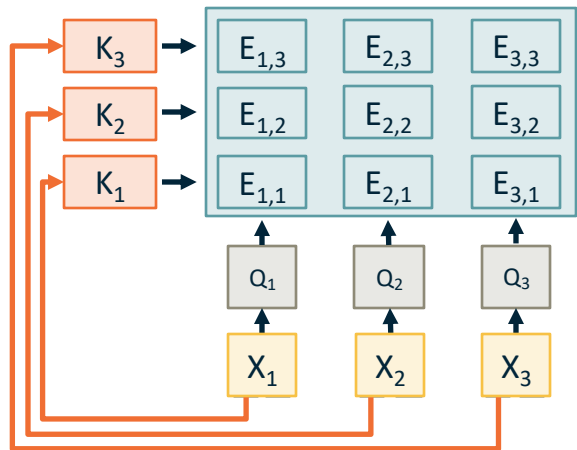
Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_x \times D_Q$)

Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_x \times D_V$)

→ **Similarities:** $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_x \times N_x$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_x \times N_x$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_x \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$



Self-Attention Layer

Sequence encode -> use each input element as query!

Inputs:

Input vectors: \mathbf{X} (Shape: $N_x \times D_x$)

Key matrix: \mathbf{W}_K (Shape: $D_x \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_x \times D_V$)

Query matrix: \mathbf{W}_Q (Shape: $D_x \times D_Q$)

Computation:

Query vectors: $\mathbf{Q} = \mathbf{XW}_Q$

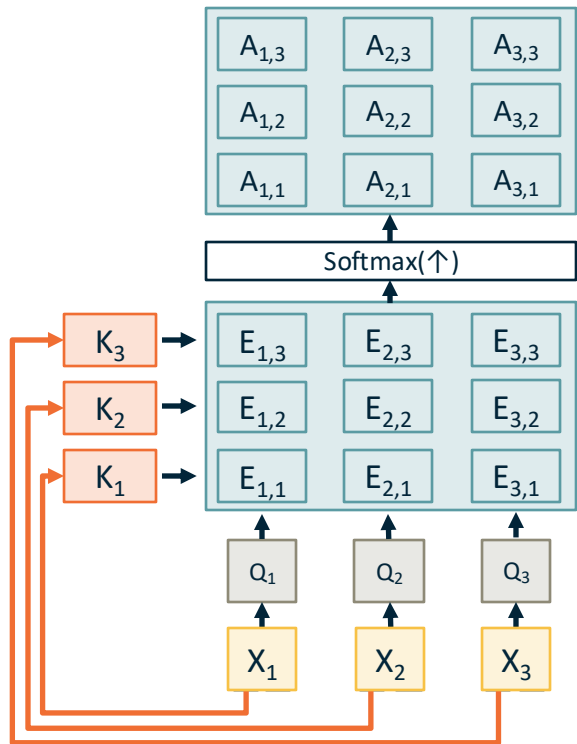
Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_x \times D_Q$)

Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_x \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_x \times N_x$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

→ **Attention weights:** $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_x \times N_x$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_x \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$



Self-Attention Layer

Sequence encode -> use each input element as query!

Inputs:

Input vectors: \mathbf{X} (Shape: $N_x \times D_x$)

Key matrix: \mathbf{W}_K (Shape: $D_x \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_x \times D_V$)

Query matrix: \mathbf{W}_Q (Shape: $D_x \times D_Q$)

Computation:

Query vectors: $\mathbf{Q} = \mathbf{XW}_Q$

Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_x \times D_Q$)

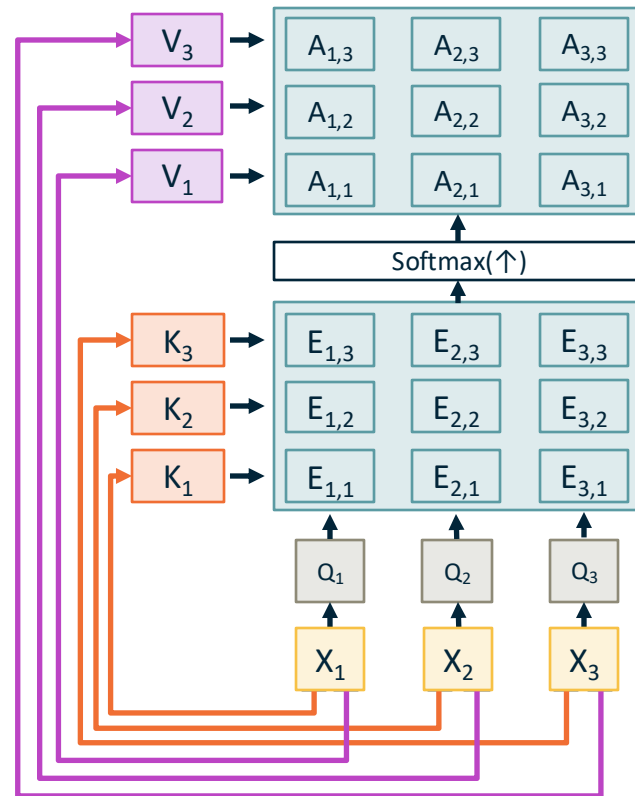
→ **Value vectors:** $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_x \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_x \times N_x$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_x \times N_x$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_x \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

Q,K,V are all generated from X!



Self-Attention Layer

Sequence encode \rightarrow use each input element as query!

Inputs:

Input vectors: \mathbf{X} (Shape: $N_X \times D_X$)

Key matrix: \mathbf{W}_K (Shape: $D_X \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_X \times D_V$)

Query matrix: \mathbf{W}_Q (Shape: $D_X \times D_Q$)

Computation:

Query vectors: $\mathbf{Q} = \mathbf{XW}_Q$

Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_X \times D_Q$)

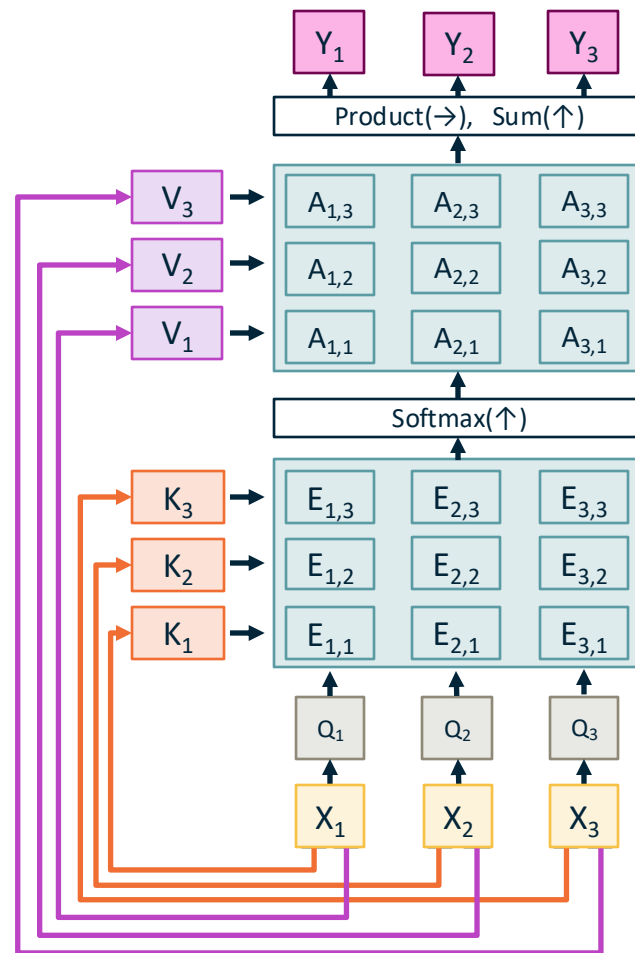
Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_X \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_X \times N_X$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_X \times N_X$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

$\mathbf{Q}, \mathbf{K}, \mathbf{V}$ are all generated from \mathbf{X} !



Self-Attention Layer

Sequence encode -> use each input element as query!

Inputs:

Input vectors: \mathbf{X} (Shape: $N_x \times D_x$)

Key matrix: \mathbf{W}_K (Shape: $D_x \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_x \times D_V$)

Query matrix: \mathbf{W}_Q (Shape: $D_x \times D_Q$)

Q: Can we use self-attention to encode an input with specific sequential ordering?

Computation:

Query vectors: $\mathbf{Q} = \mathbf{XW}_Q$

Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_x \times D_Q$)

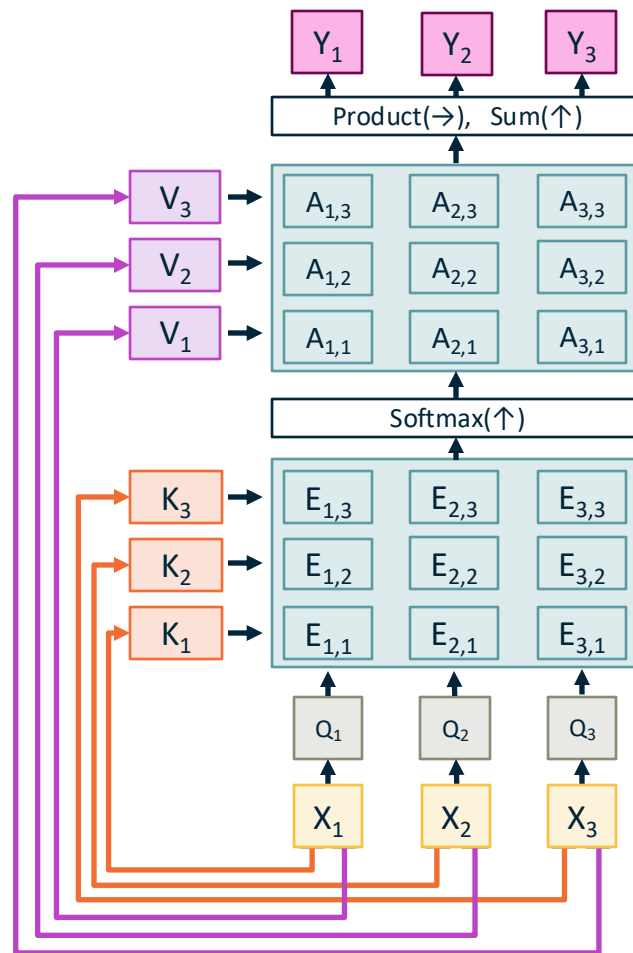
Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_x \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_x \times N_x$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_x \times N_x$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_x \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

Q,K,V are all generated from X!



Self-Attention Layer

Inputs:

Input vectors: \mathbf{X} (Shape: $N_x \times D_x$)

Key matrix: \mathbf{W}_K (Shape: $D_x \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_x \times D_V$)

Query matrix: \mathbf{W}_Q (Shape: $D_x \times D_Q$)

Computation:

Query vectors: $\mathbf{Q} = \mathbf{XW}_Q$

Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_x \times D_Q$)

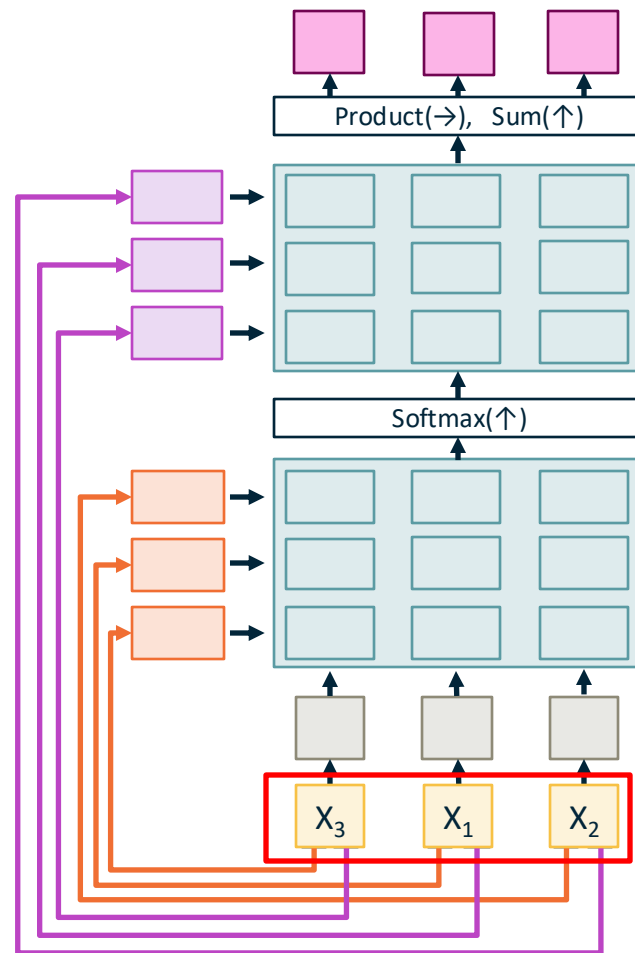
Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_x \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_x \times N_x$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_x \times N_x$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_x \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

Consider **permuting**
the input vectors:



Self-Attention Layer

Inputs:

Input vectors: \mathbf{X} (Shape: $N_x \times D_x$)

Key matrix: \mathbf{W}_K (Shape: $D_x \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_x \times D_V$)

Query matrix: \mathbf{W}_Q (Shape: $D_x \times D_Q$)

Computation:

Query vectors: $\mathbf{Q} = \mathbf{XW}_Q$

Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_x \times D_Q$)

Value Vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_x \times D_V$)

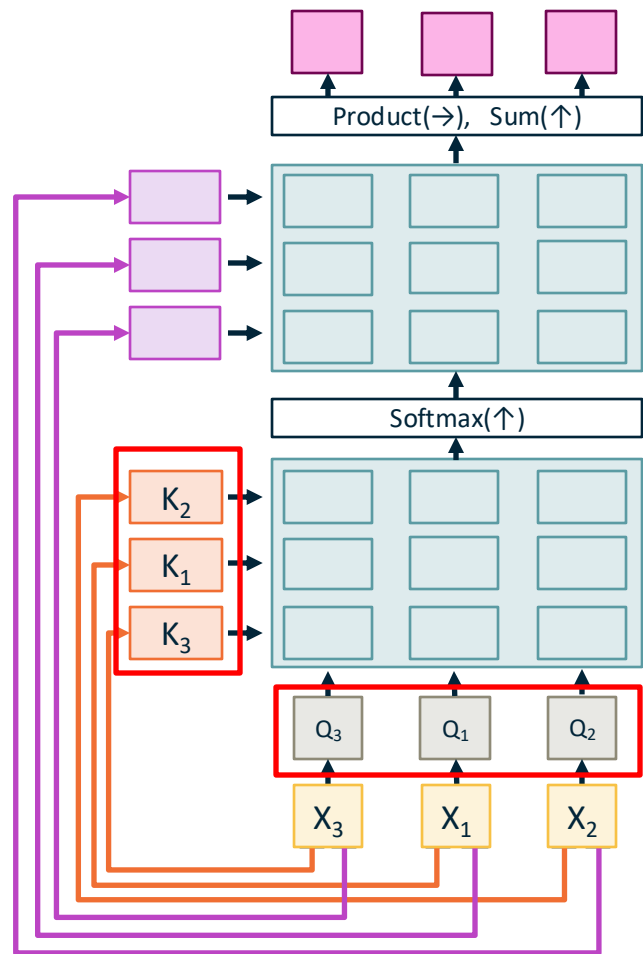
Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_x \times N_x$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_x \times N_x$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_x \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

Consider **permuting**
the input vectors:

Queries and Keys will be
the same, but permuted



Self-Attention Layer

Inputs:

Input vectors: \mathbf{X} (Shape: $N_X \times D_X$)

Key matrix: \mathbf{W}_K (Shape: $D_X \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_X \times D_V$)

Query matrix: \mathbf{W}_Q (Shape: $D_X \times D_Q$)

Computation:

Query vectors: $\mathbf{Q} = \mathbf{XW}_Q$

Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_X \times D_Q$)

Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_X \times D_V$)

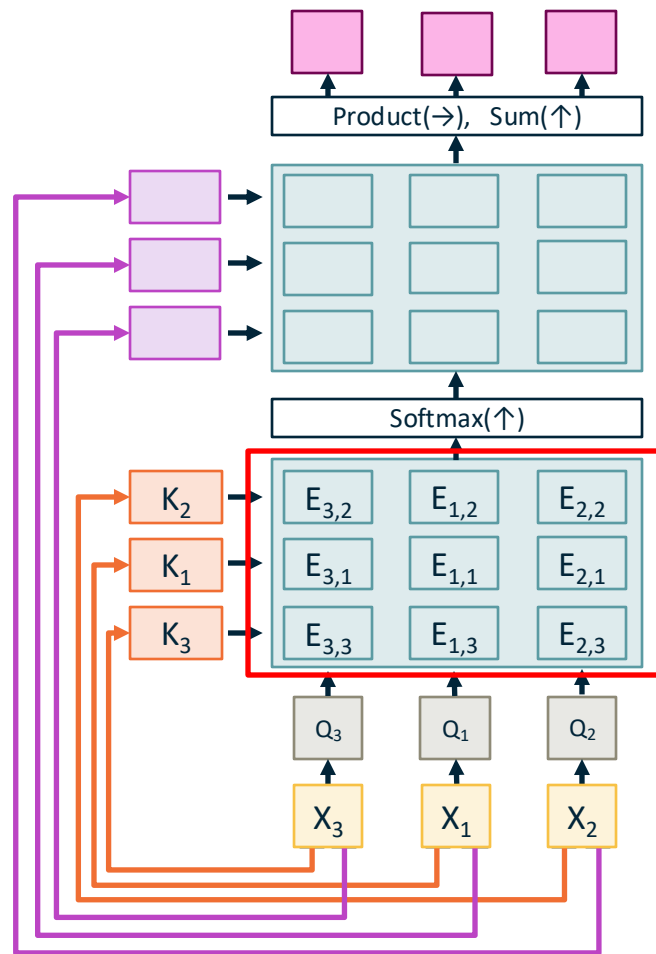
Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_X \times N_X$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_X \times N_X$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

Consider **permuting**
the input vectors:

Similarities will be the
same, but permuted



Self-Attention Layer

Inputs:

Input vectors: \mathbf{X} (Shape: $N_x \times D_x$)

Key matrix: \mathbf{W}_K (Shape: $D_x \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_x \times D_V$)

Query matrix: \mathbf{W}_Q (Shape: $D_x \times D_Q$)

Computation:

Query vectors: $\mathbf{Q} = \mathbf{XW}_Q$

Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_x \times D_Q$)

Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_x \times D_V$)

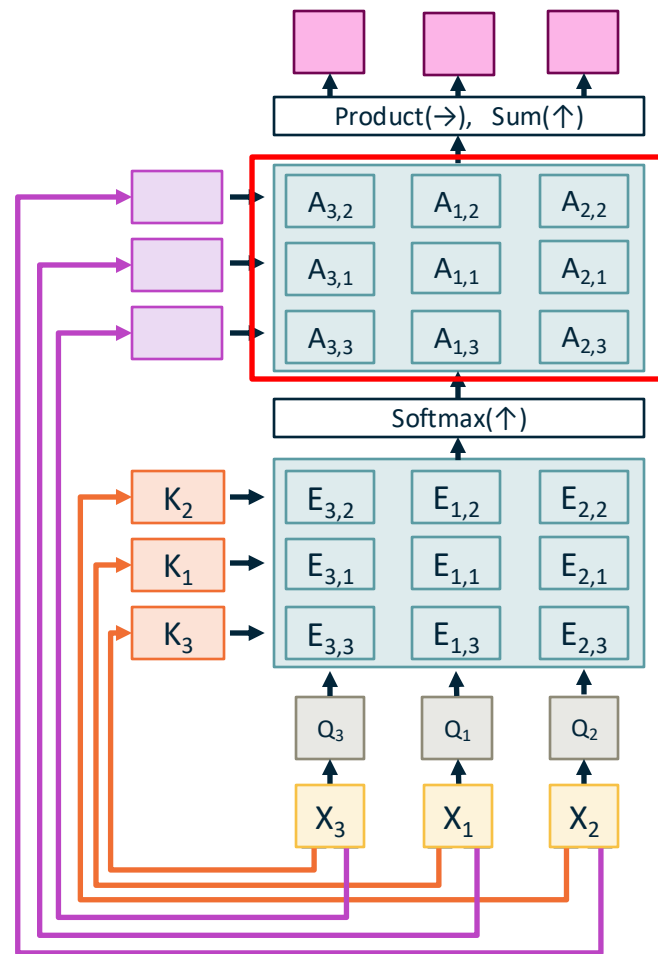
Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_x \times N_x$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_x \times N_x$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_x \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

Consider **permuting**
the input vectors:

Attention weights will be
the same, but permuted



Self-Attention Layer

Inputs:

Input vectors: \mathbf{X} (Shape: $N_X \times D_X$)

Key matrix: \mathbf{W}_K (Shape: $D_X \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_X \times D_V$)

Query matrix: \mathbf{W}_Q (Shape: $D_X \times D_Q$)

Computation:

Query vectors: $\mathbf{Q} = \mathbf{XW}_Q$

Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_X \times D_Q$)

Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_X \times D_V$)

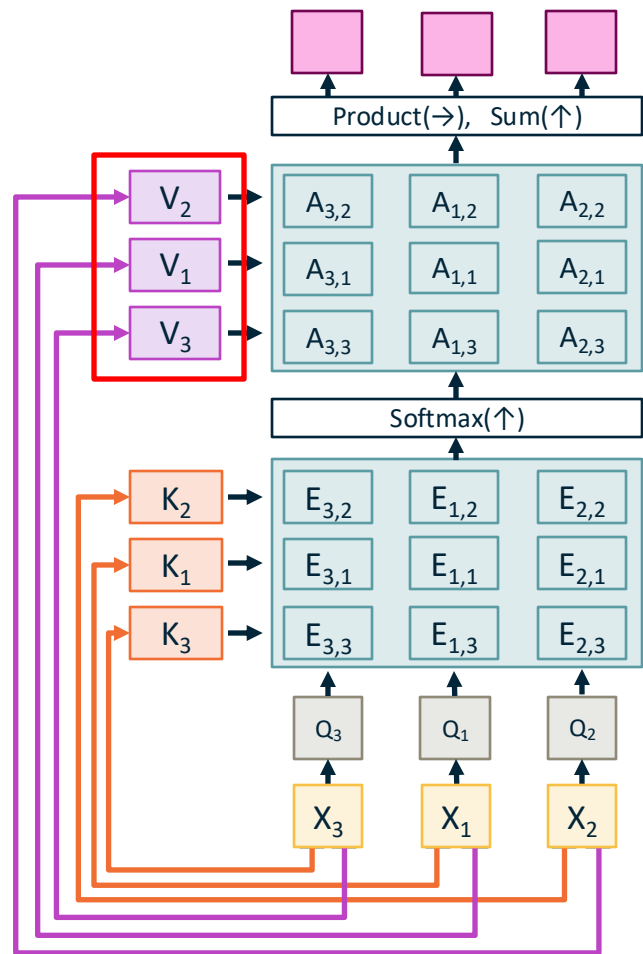
Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_X \times N_X$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_X \times N_X$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

Consider **permuting**
the input vectors:

Values will be the
same, but permuted



Self-Attention Layer

Inputs:

Input vectors: \mathbf{X} (Shape: $N_X \times D_X$)

Key matrix: \mathbf{W}_K (Shape: $D_X \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_X \times D_V$)

Query matrix: \mathbf{W}_Q (Shape: $D_X \times D_Q$)

Computation:

Query vectors: $\mathbf{Q} = \mathbf{XW}_Q$

Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_X \times D_Q$)

Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_X \times D_V$)

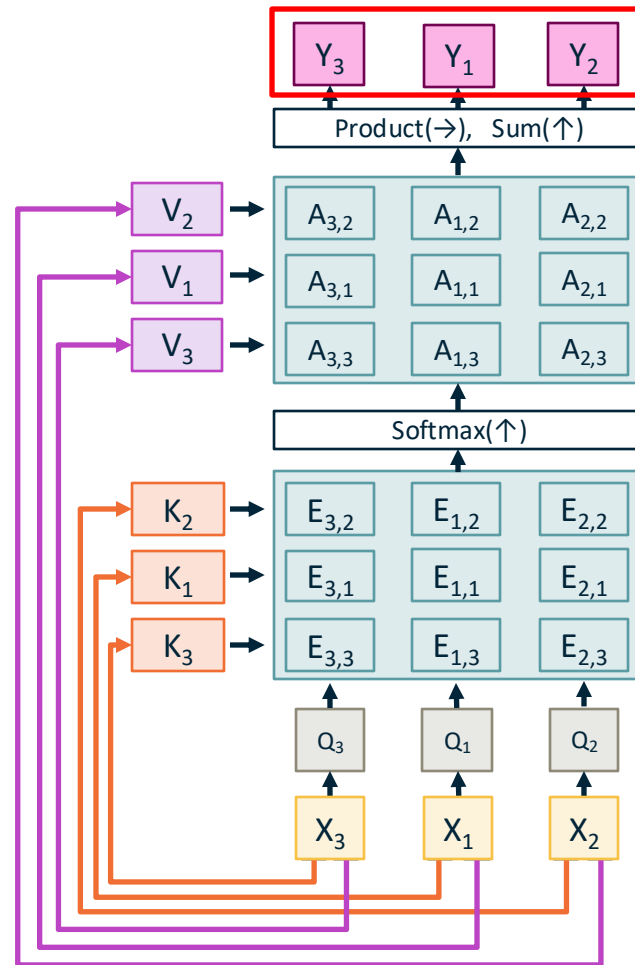
Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_X \times N_X$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_X \times N_X$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

Consider **permuting**
the input vectors:

Outputs will be the
same, but permuted



Self-Attention Layer

Inputs:

Input vectors: \mathbf{X} (Shape: $N_x \times D_x$)

Key matrix: \mathbf{W}_K (Shape: $D_x \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_x \times D_V$)

Query matrix: \mathbf{W}_Q (Shape: $D_x \times D_Q$)

Computation:

Query vectors: $\mathbf{Q} = \mathbf{XW}_Q$

Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_x \times D_Q$)

Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_x \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_x \times N_x$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

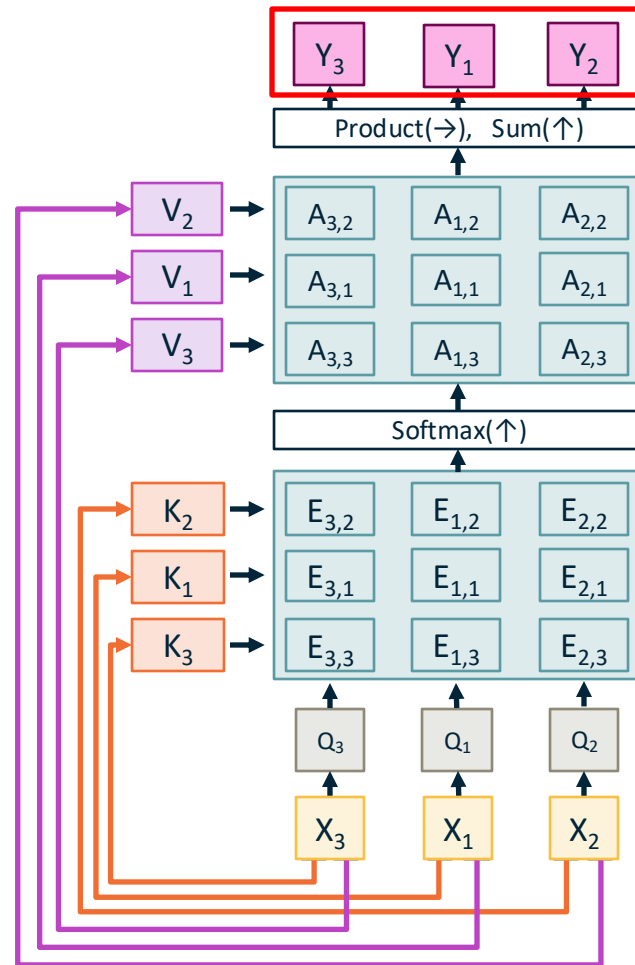
Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_x \times N_x$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_x \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

Consider **permuting**
the input vectors:

Outputs will be the
same, but permuted

Self-attention layer is
Permutation Equivariant
 $f(s(x)) = s(f(x))$



Self-Attention Layer

Inputs:

Input vectors: \mathbf{X} (Shape: $N_x \times D_x$)

Key matrix: \mathbf{W}_K (Shape: $D_x \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_x \times D_V$)

Query matrix: \mathbf{W}_Q (Shape: $D_x \times D_Q$)

Computation:

Query vectors: $\mathbf{Q} = \mathbf{XW}_Q$

Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_x \times D_Q$)

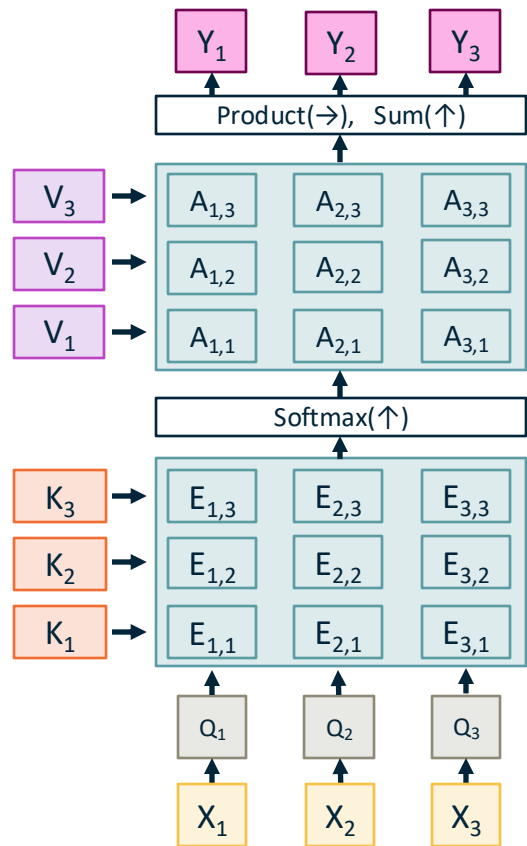
Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_x \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_x \times N_x$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_x \times N_x$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_x \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

Self attention doesn't "know" the order of the vectors it is processing! Not good for sequence encoding.



Self-Attention Layer

Inputs:

Input vectors: \mathbf{X} (Shape: $N_X \times D_X$)

Key matrix: \mathbf{W}_K (Shape: $D_X \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_X \times D_V$)

Query matrix: \mathbf{W}_Q (Shape: $D_X \times D_Q$)

Computation:

Query vectors: $\mathbf{Q} = \mathbf{XW}_Q$

Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_X \times D_Q$)

Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_X \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_X \times N_X$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

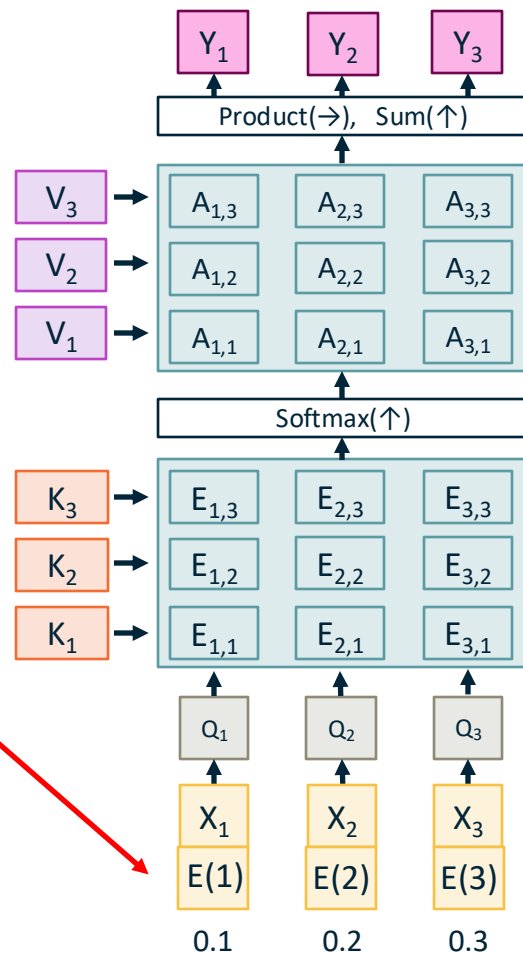
Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_X \times N_X$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

In order to make processing position-aware, concatenate input with **positional encoding** \mathbf{E}

$E(i)$ encodes the position of the i -th element in a sequence

$E()$ can be a simple function (e.g., linear or sin functions) or a learned lookup table.



Aside: Positional Encoding (PE) for Self-Attention

Motivation: Maintain the order of input data since attention mechanisms are permutation invariant. PEs are shared across all input sequences.

Linear Positional Encoding: $PE(pos) = a \cdot pos + b$.

Problem: encoding increases with the sequence length, causing gradient problem for long sequences.

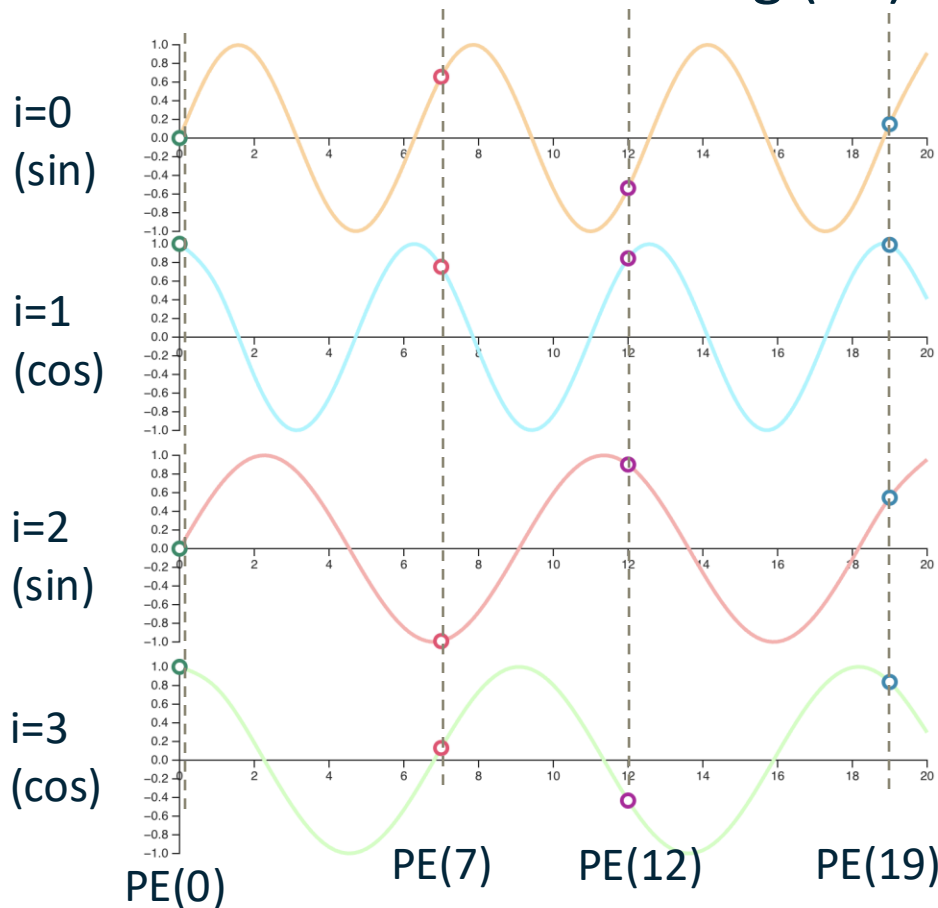
Sin/cos Positional Encoding (Default):

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

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

PE for each dimension (i) repeats periodically, combine different waveforms at each dimension to get a unique embedding.

Aside: Positional Encoding (PE) for Self-Attention



p0	p1	p2	p3	
0.000	0.657	-0.537	0.150	i=0
1.000	0.754	0.844	0.989	i=1
0.000	-0.992	0.901	0.547	i=2
1.000	0.130	-0.433	0.837	i=3

Positional Encoding

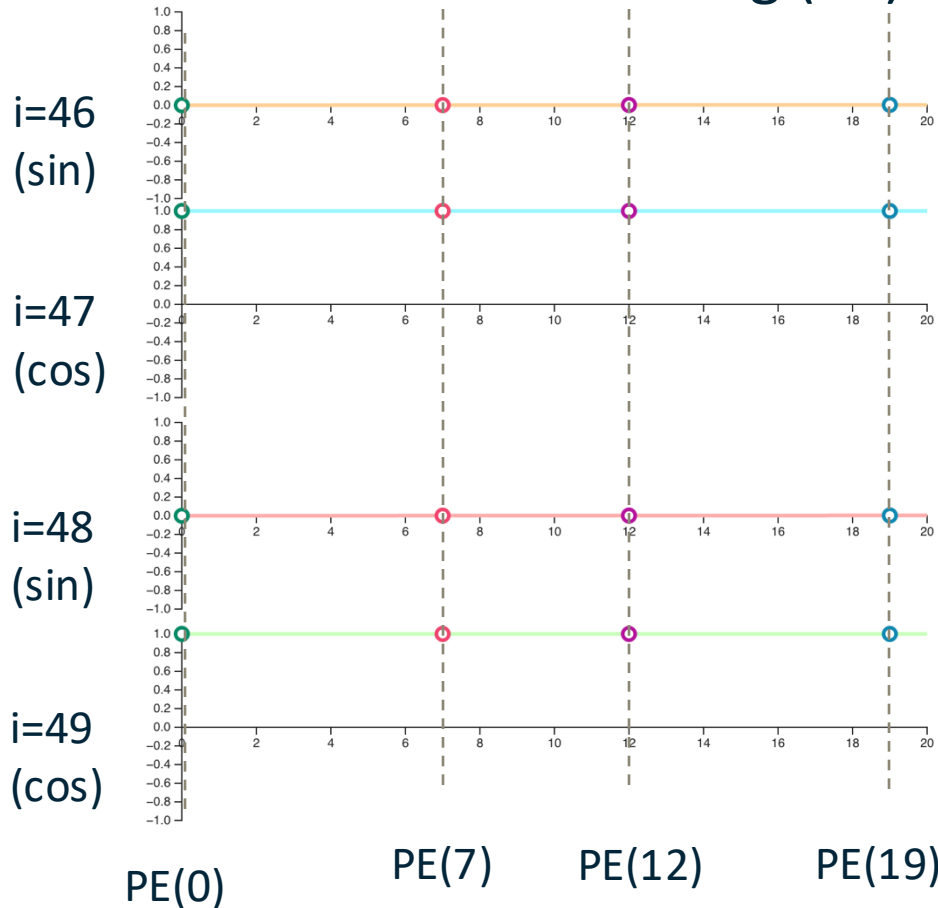
$$PE_{(pos, 2i)} = \sin\left(\frac{pos}{10000^{2i/d_{\text{model}}}}\right)$$

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

i determines the sin/cos period

$d_{\text{model}}=50$

Aside: Positional Encoding (PE) for Self-Attention



p0	p1	p2	p3	
0.000	0.001	0.003	0.004	$i=46$
1.000	1.000	1.000	1.000	$i=47$
0.000	0.001	0.002	0.003	$i=48$
1.000	1.000	1.000	1.000	$i=49$

Important!
indexes has
changed

Positional Encoding

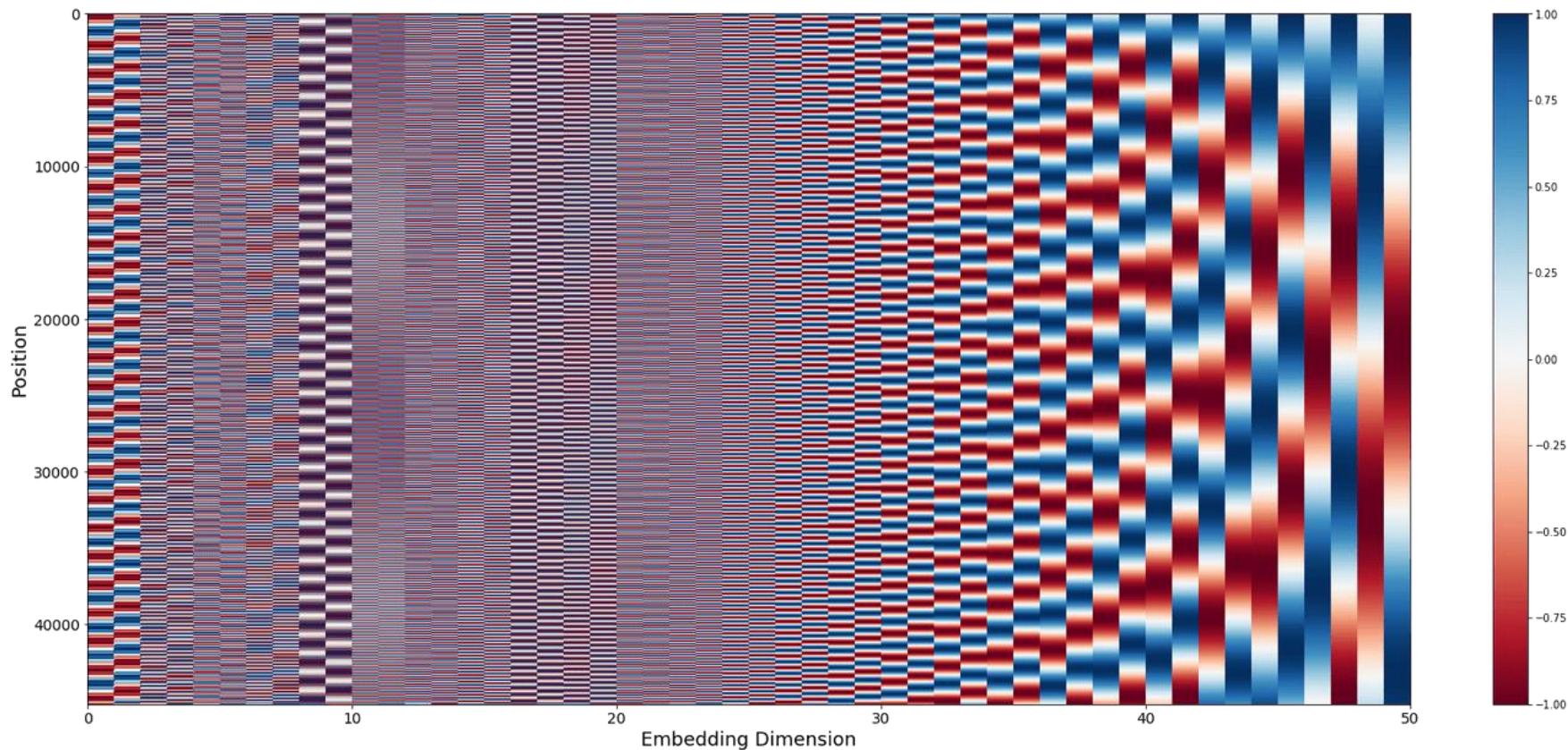
$$PE_{(pos, 2i)} = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

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

i determines the sin/cos period

$d_{model}=50$

Aside: Positional Encoding (PE) for Self-Attention



Aside: Positional Encoding (PE) for Self-Attention

Motivation: Maintain the order of input data since attention mechanisms are permutation invariant. PEs are shared across all input sequences.

Linear Positional Encoding: $PE(pos) = a \cdot pos + b$.

Problem: encoding increases with the sequence length, causing gradient problem for long sequences.

Sin/cos Positional Encoding (Default):

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

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

PE for each dimension (i) repeats periodically, combine different waveforms at each dimension to get a unique embedding.

Learned Positional Encoding: $PE_{\theta}(pos, i)$.

Learn the most suitable position embedding for the training set.

Masked Self-Attention Layer

Inputs:

Input vectors: \mathbf{X} (Shape: $N_x \times D_x$)

Key matrix: \mathbf{W}_K (Shape: $D_x \times D_o$)

Value matrix: W_v (Shape: $D_x \times D_v$)

Query matrix: W_Q (Shape: $D_x \times D_Q$)

Computation:

Query vectors: $Q = XW_o$

Key vectors: $K = XW_K$ (Shape: $N_x \times D_Q$)

Value vectors: $\mathbf{V} = \mathbf{XW}_v$ (Shape: $N_x \times D_v$)

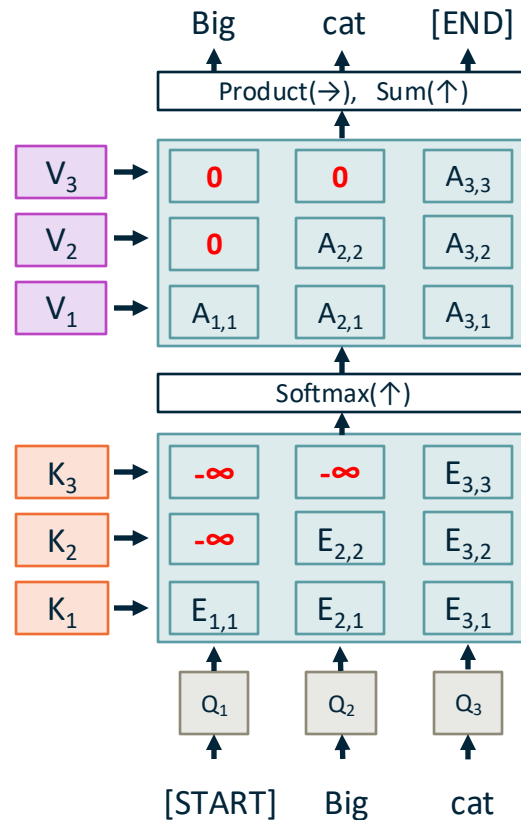
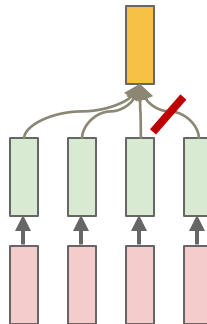
Similarities: $E = QK^T$ (Shape: $N_x \times N_x$) $E_{i,j} = Q_i \cdot K_j / \text{sqrt}(D_Q)$

Attention weights: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_x \times N_x$)

Output vectors: $Y = AV$ (Shape: $N_x \times D_v$) $Y_i = \sum_j A_{i,j} v_j$

Don't let vectors "look ahead" in the sequence

Used for sequence decoding
(predict next word)



Multi-headed Self-Attention Layer

Inputs:

Input vectors: \mathbf{X} (Shape: $N_x \times D_x$)

Key matrix: \mathbf{W}_K (Shape: $D_x \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_x \times D_V$)

Query matrix: \mathbf{W}_Q (Shape: $D_x \times D_Q$)

Computation:

Query vectors: $\mathbf{Q} = \mathbf{XW}_Q$

Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_x \times D_Q$)

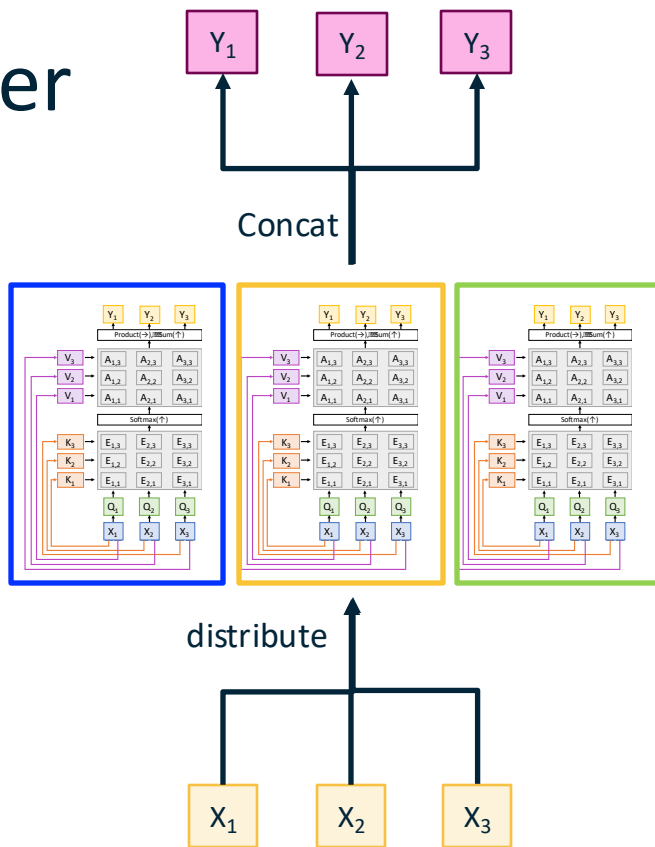
Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_x \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_x \times N_x$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_x \times N_x$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_x \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

Use H independent
“Attention Heads” in
parallel



Multi-headed Self-Attention Layer

Inputs:

Input vectors: \mathbf{X} (Shape: $N_X \times D_X$)

Key matrix: \mathbf{W}_K (Shape: $D_X \times D_Q$)

Value matrix: \mathbf{W}_V (Shape: $D_X \times D_V$)

Query matrix: \mathbf{W}_Q (Shape: $D_X \times D_Q$)

Computation:

Query vectors: $\mathbf{Q} = \mathbf{XW}_Q$

Key vectors: $\mathbf{K} = \mathbf{XW}_K$ (Shape: $N_X \times D_Q$)

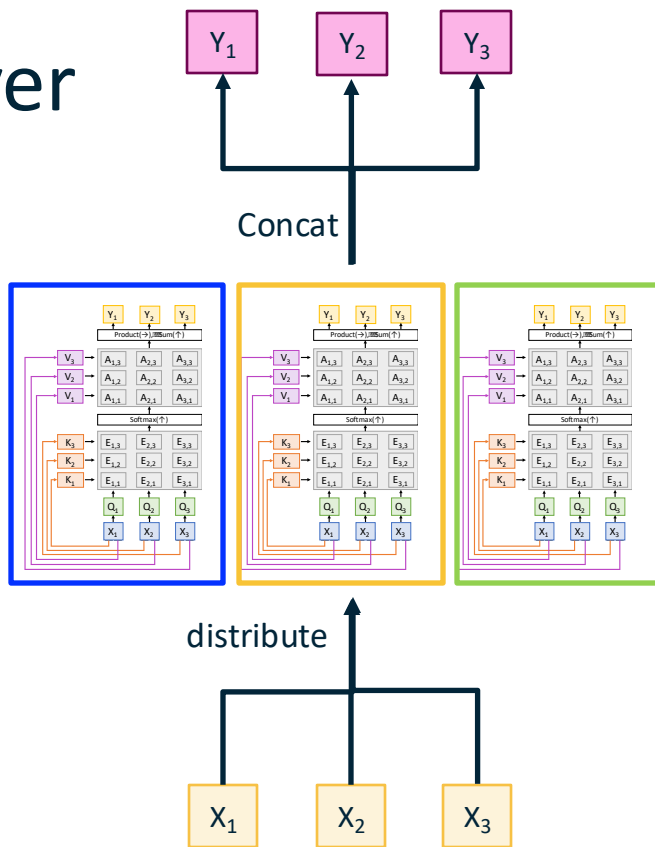
Value vectors: $\mathbf{V} = \mathbf{XW}_V$ (Shape: $N_X \times D_V$)

Similarities: $\mathbf{E} = \mathbf{QK}^T$ (Shape: $N_X \times N_X$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

Attention weights: $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$ (Shape: $N_X \times N_X$)

Output vectors: $\mathbf{Y} = \mathbf{AV}$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

Use H independent
“Attention Heads” in
parallel



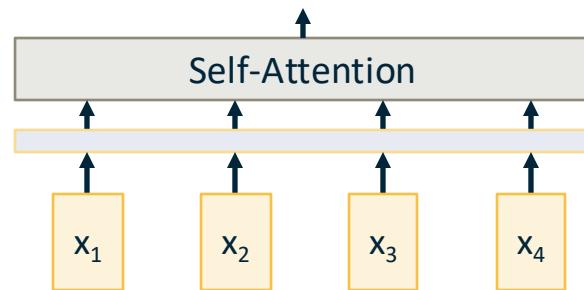
Highly parallelizable: Can compute attentions for all
input element from all head in parallel!

The Transformer Block



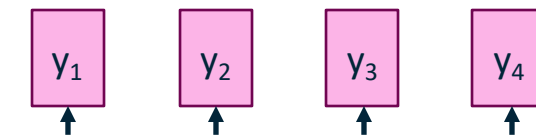
The Transformer Block

All vectors interact
with each other

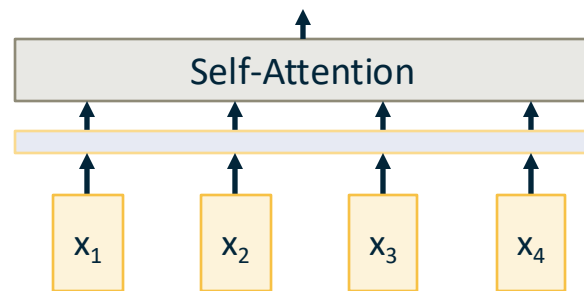


The Transformer Block

MLP (same copy)
independently on each
vector

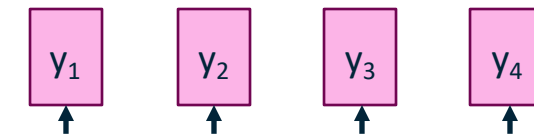


All vectors interact
with each other



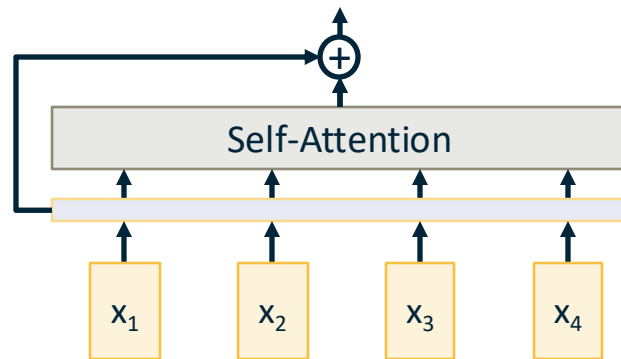
The Transformer Block

MLP (same copy)
independently on each
vector



Residual connection

All vectors interact
with each other



The Transformer Block

Recall **Layer Normalization**:

Given h_1, \dots, h_N (shape: D)

scale: γ (shape: D)

shift: β (shape: D)

$\mu_i = (1/D) \sum_j h_{i,j}$ (scalar)

$\sigma_i = (\sum_j (h_{i,j} - \mu_i)^2)^{1/2}$ (scalar)

$z_i = (h_i - \mu_i) / \sigma_i$ (shape: D)

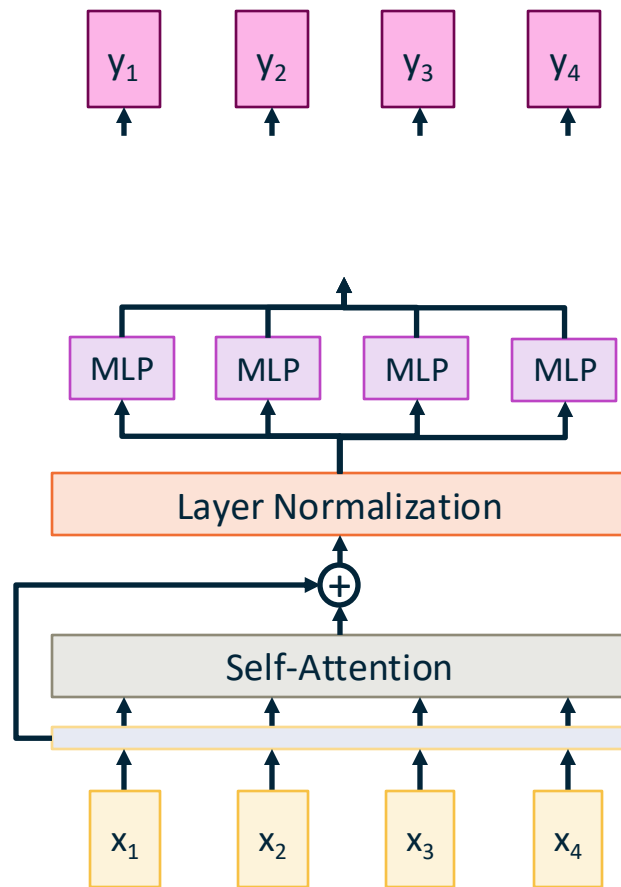
$y_i = \gamma * z_i + \beta$ (shape: D)

Applied **per dimension**, not
across the sequence

MLP (same copy)
independently on each
vector

Residual connection

All vectors interact
with each other



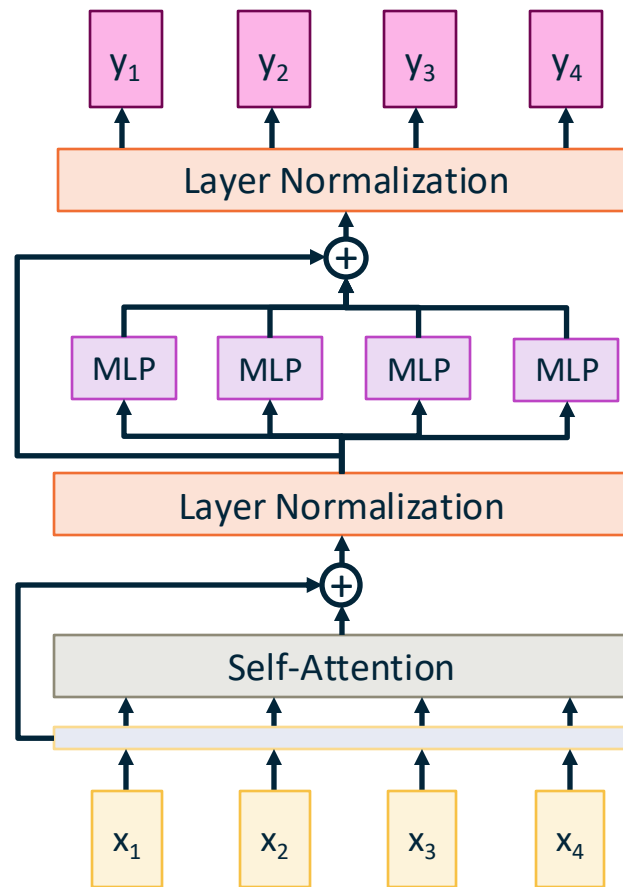
The Transformer Block

Residual connection

MLP (same copy)
independently on each
vector

Residual connection

All vectors interact
with each other



The Transformer Block

Transformer Block:

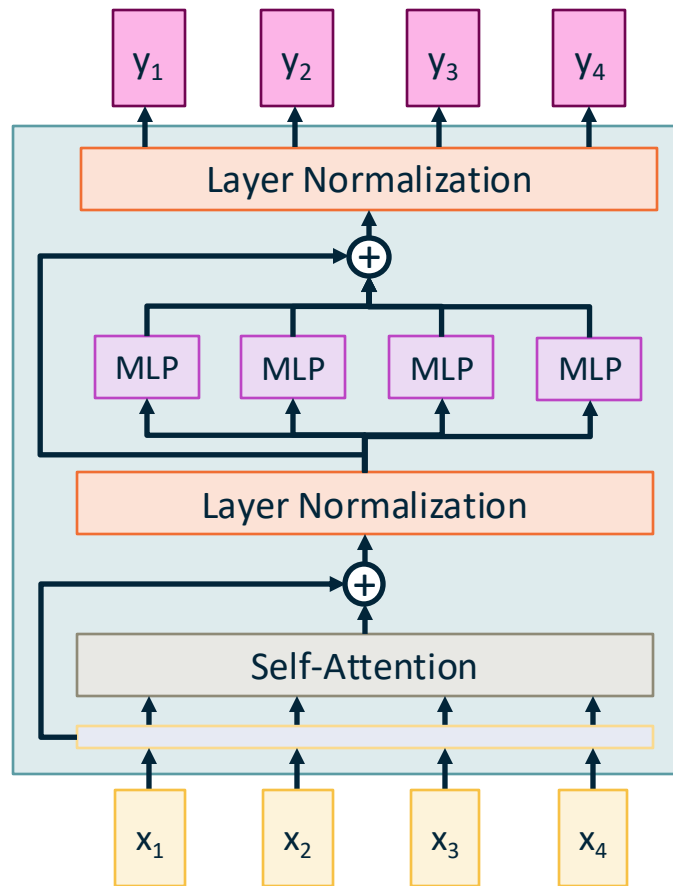
Input: Set of vectors x

Output: Set of vectors y

Self-attention is the only
interaction among vectors!

Layer norm and MLP work
independently per vector

Highly scalable, highly
parallelizable



The Transformer

Transformer Block:

Input: Set of vectors x

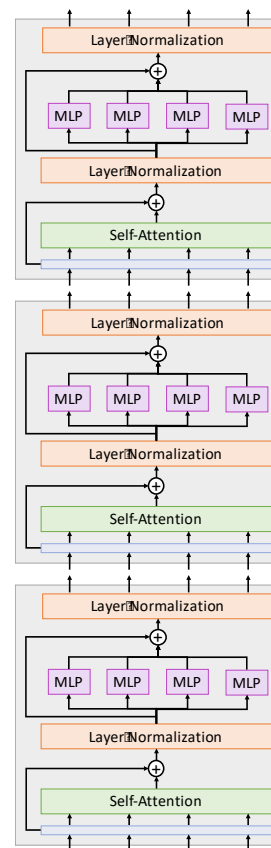
Output: Set of vectors y

Self-attention is the only
interaction among vectors!

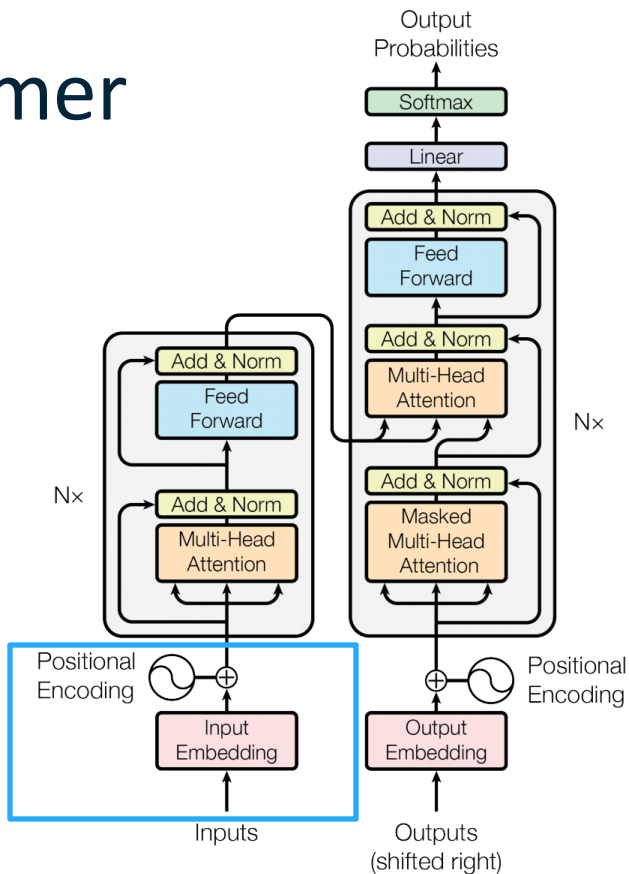
Layer norm and MLP work
independently per vector

Highly scalable, highly
parallelizable

A **Transformer** is a sequence
of transformer blocks

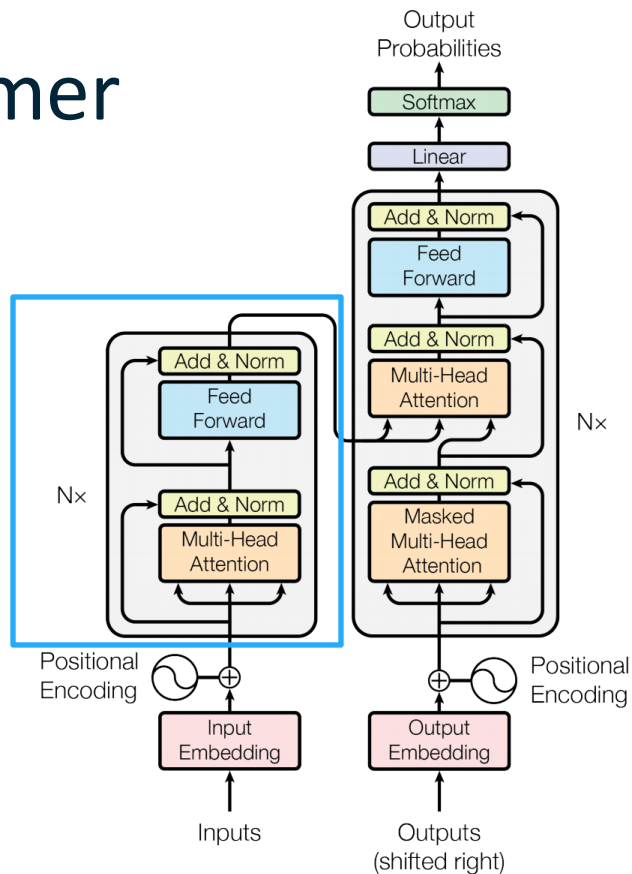


The Transformer



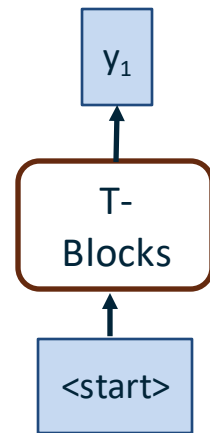
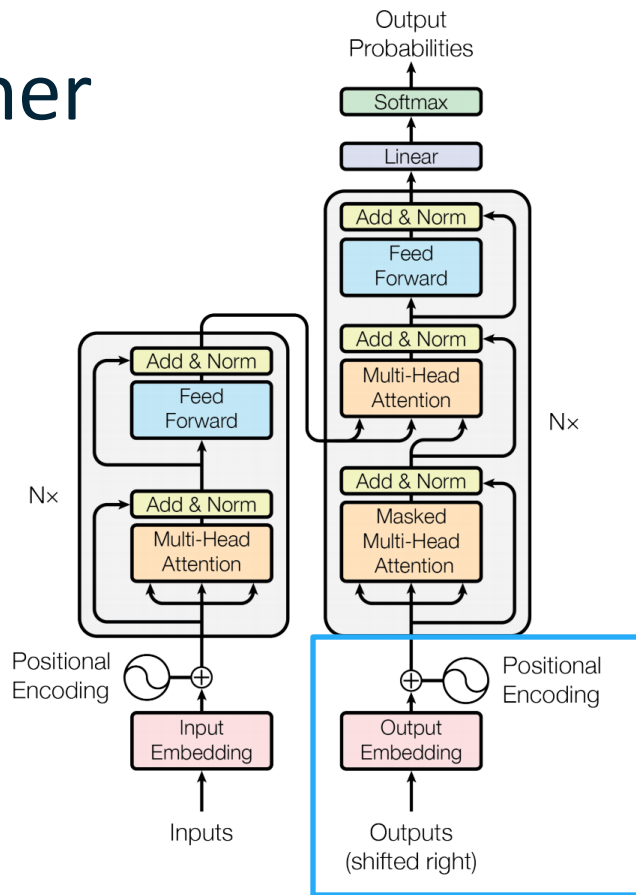
Encoder-Decoder

The Transformer



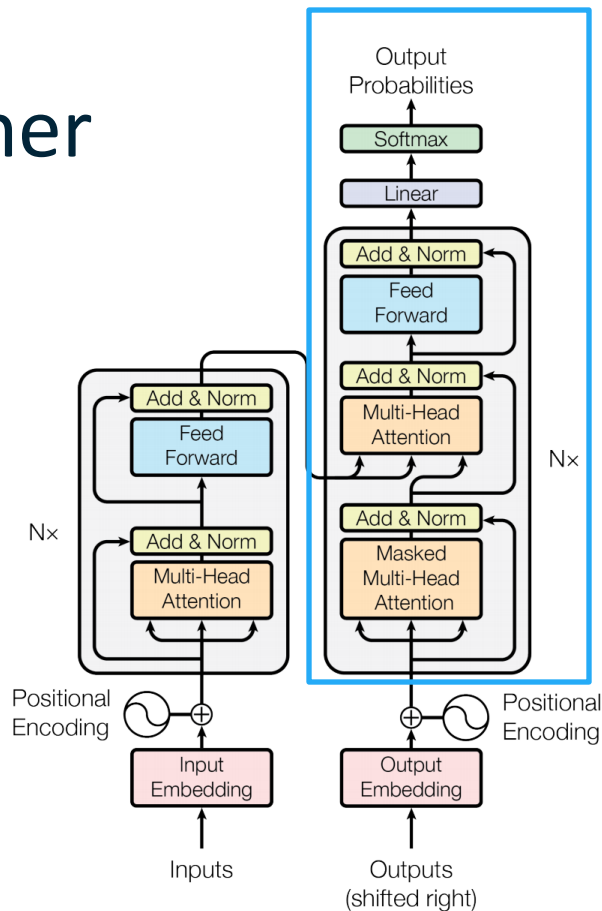
Encoder-Decoder

The Transformer

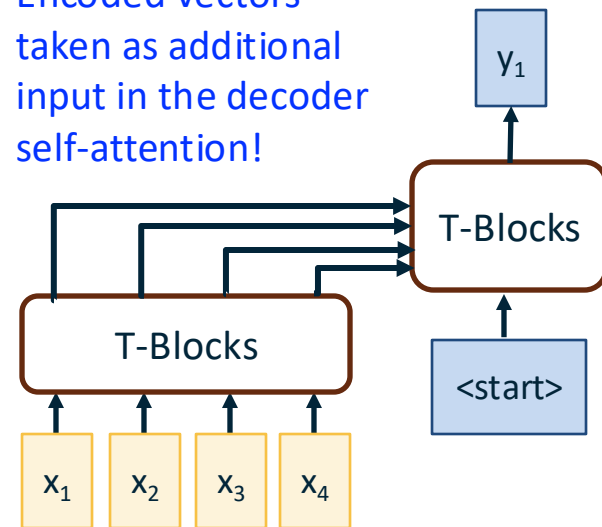


Encoder-Decoder

The Transformer

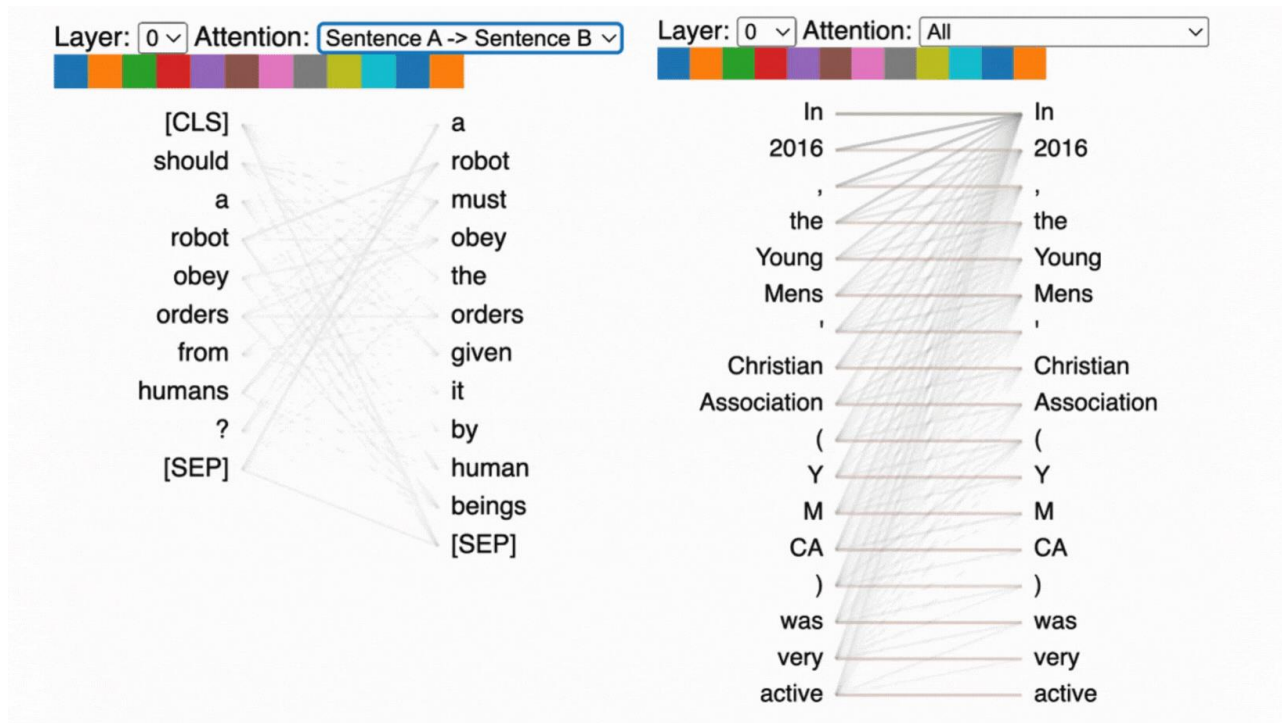


Encoded vectors
taken as additional
input in the decoder
self-attention!



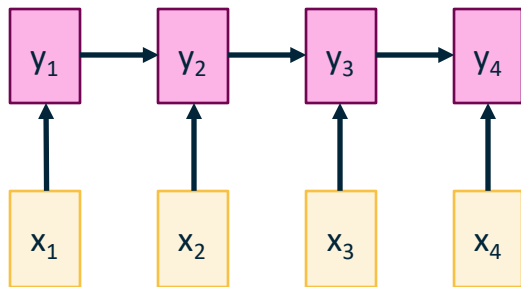
Encoder-Decoder

Visualizing Transformer Attentions



Three Ways of Processing Sequences

Recurrent Neural Network



Works on **Ordered Sequences**

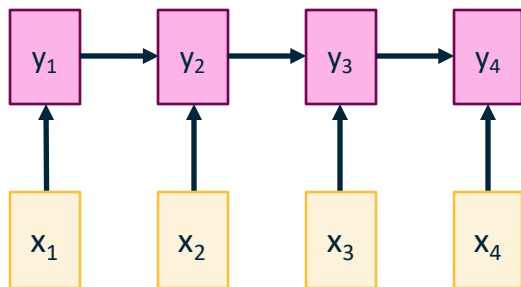
(+) Natural sequential processing:
“sees” the input sequence in its
original ordering

(-) Forgetful: difficult to handle
long-range dependencies.

(-) Not parallelizable: need to
compute hidden states sequentially

Three Ways of Processing Sequences

Recurrent Neural Network



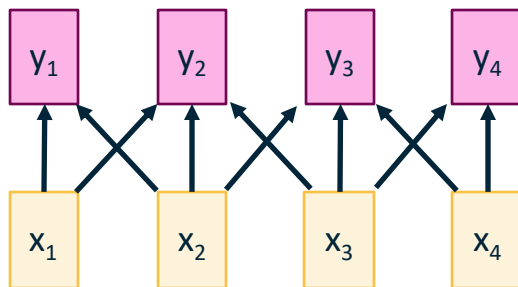
Works on **Ordered Sequences**

(+) **Natural sequential processing:**
“sees” the input sequence in its
original ordering

(-) **Forgetful:** difficult to handle
long-range dependencies.

(-) **Not parallelizable:** need to
compute hidden states sequentially

1D Convolution



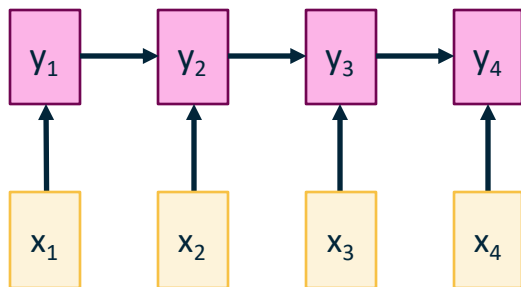
Works on **Multidimensional Grids**

(-) **Bad at long sequences:** Need to
stack many conv layers for outputs
to “see” the whole sequence

(+) **Highly parallel:** Each output can
be computed in parallel

Three Ways of Processing Sequences

Recurrent Neural Network



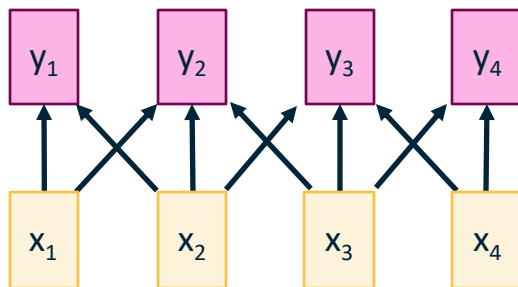
Works on **Ordered Sequences**

(+) **Natural sequential processing:** “sees” the input sequence in its original ordering

(-) **Forgetful:** difficult to handle long-range dependencies.

(-) **Not parallelizable:** need to compute hidden states sequentially

1D Convolution

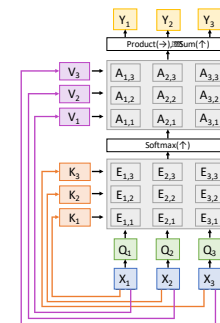


Works on **Multidimensional Grids**

(-) **Bad at long sequences:** Need to stack many conv layers for outputs to “see” the whole sequence

(+) **Highly parallel:** Each output can be computed in parallel

Self-Attention



Works on **Sets of Vectors**

(+) **Good at long sequences:** after one self-attention layer, each output “sees” all inputs!

(+) **Highly parallel:** Each output can be computed in parallel

(-) **Very memory intensive**

(-) **Requires positional encoding**

Some Recent Advances in Transformers

- **Compute Efficiency:** KV Cache, Grouped-query Attention
- **System-level optimization:** Paged Attention
- **Beyond language data**

Transformer Inference is Expensive!

Q: What's the $O()$ complexity of generating a length- N sequence with vanilla RNN?

A: $O(N)$: you only process each generated token / item once

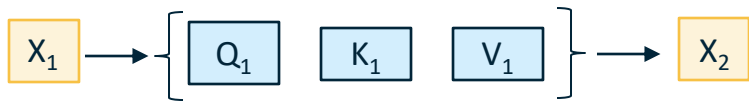
Q: What's the $O()$ complexity of generating a length- N sequence with Transformer?

A: $O(N^2)$. Generating a new token needs to attend to all existing tokens. Bad for long sequences!

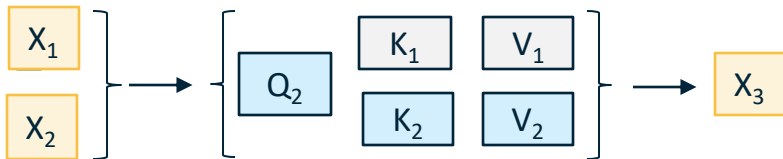
Recall: $\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{KV^T}{\sqrt{d_k}}\right)Q$

Observation: Only the query (Q) changes. No need to recompute K and V for the already-generated sequence!

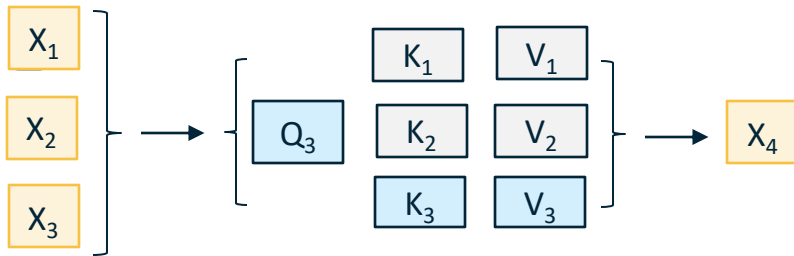
KV Cache: Trading GPU Memory for Efficiency



With KV Cache, complexity goes from quadratic to linear! $O(N)$

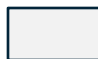



For ~ 1000 tokens, without KV cache requires 300 million attention operations per layer!



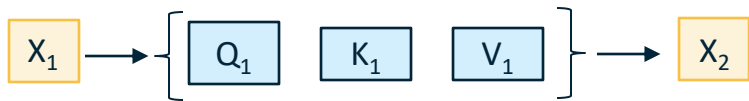
With KV cache: $\sim 500,000$ attention operations per layer.

600x speedup!

 =cached

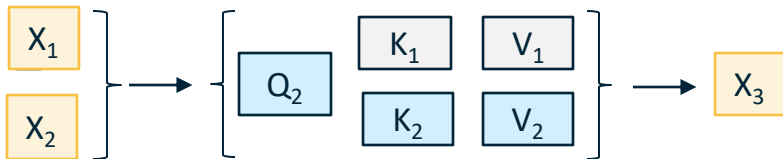
 =new

KV Cache: Trading GPU Memory for Efficiency

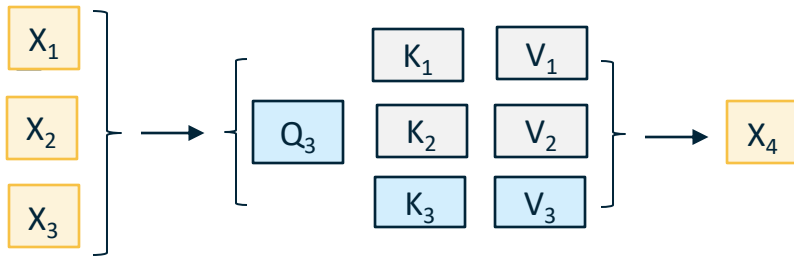


But we need to cache the KV for all of the intermediate attention layers

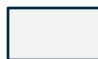
...




For a typical model: 32 layers X 4096 embedding dimension x 2 vectors



Takes **1GB of memory** to KV cache a **N=2048** sequence!

 =cached

 =new

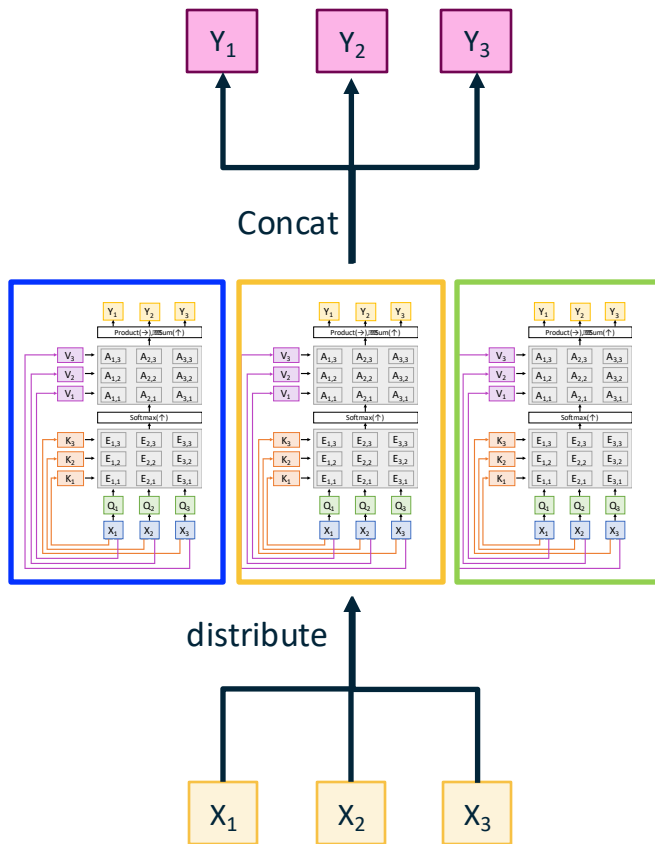
Optimizing KV Cache: Multi-Query Attention (MQA)

Recall Multi-headed Attention (MHA):

- Similar to having multiple kernels per layer in ConvNet, have multiple attention heads in each attention layer.
- Each head learns to attend to different patterns (syntax, semantics, long-range dependencies, etc.).
- The diversity comes from having separate initialization for each Q, K, V head.

Problem: For 32-head MHA, 32 K matrices + 32 V matrices = **64 matrices per layer**.

Big overhead for KV Cache!



Optimizing KV Cache: Multi-Query Attention (MQA)

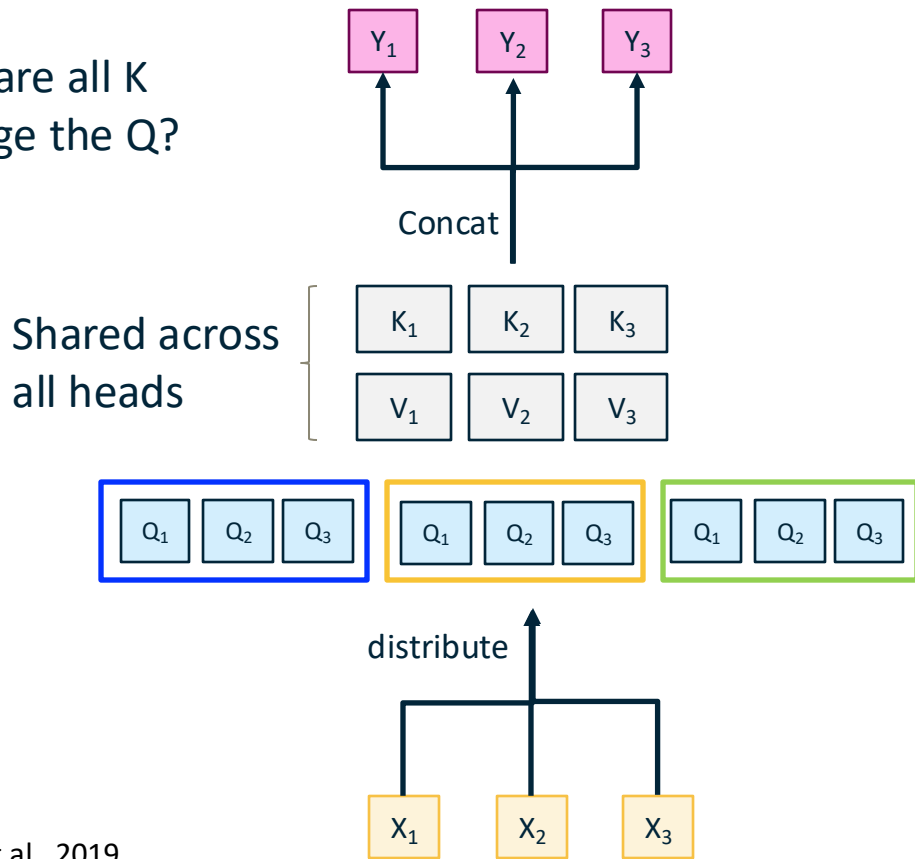
Multi-Query Attention (MQA): What if we share all K and V vectors across all heads and only change the Q?

Diversity comes from different Q.

1/32 KV cache size for 32-head MHA.

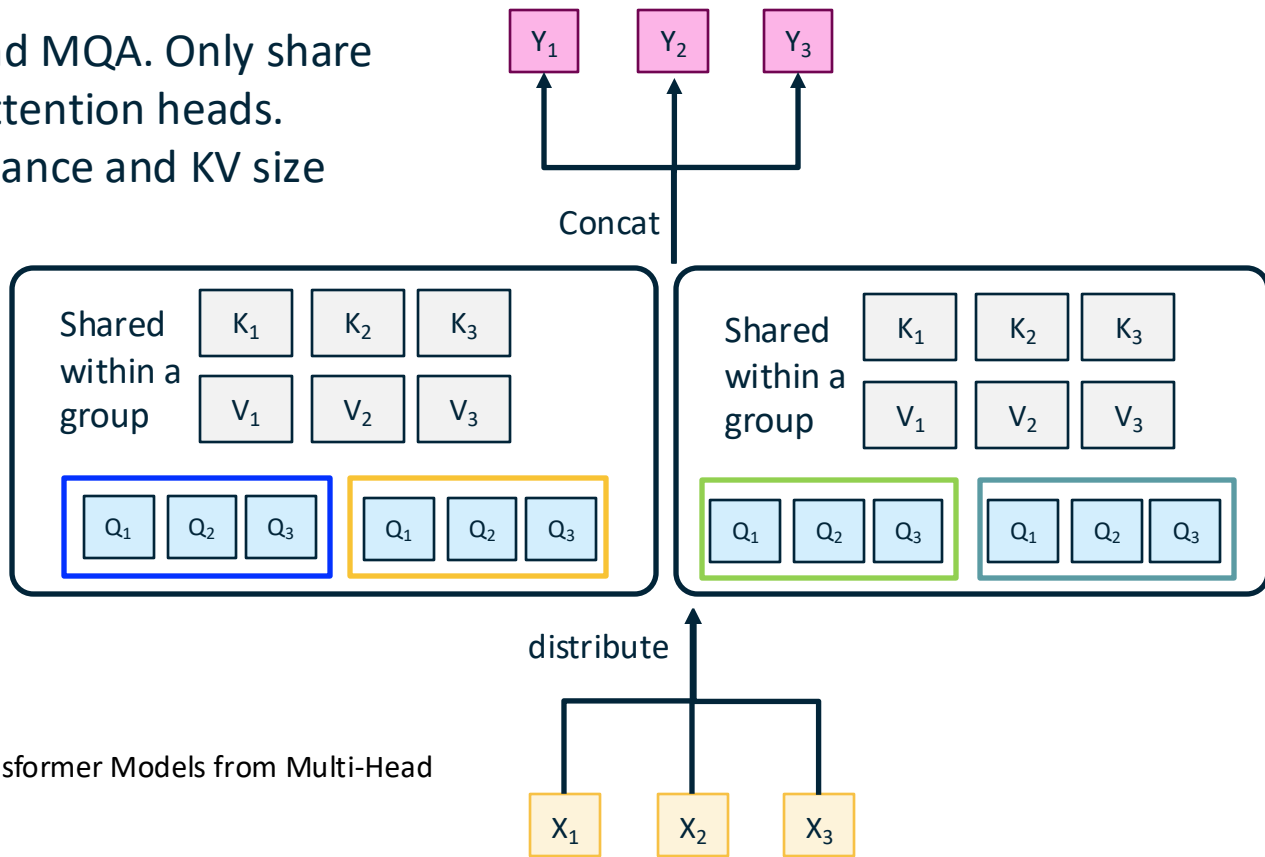
Not caching Q anyways

Slightly degraded performance (1%-2%),
but massive save on KV cache!



Optimizing KV Cache: Group-Query Attention (GQA)

GQA: In-between MHA and MQA. Only share KV between a subset of attention heads.
Balance between performance and KV size



GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints, Ainslie et al., 2023

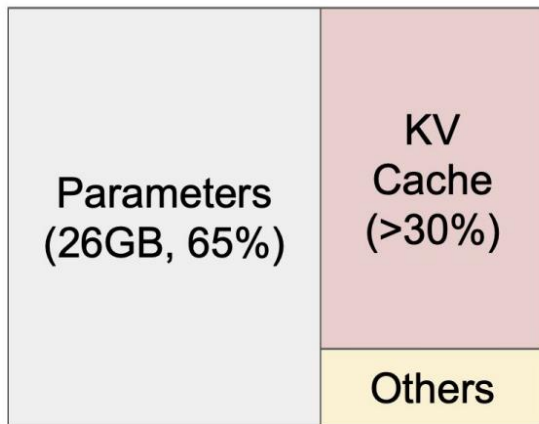
Can we do even better?

Need to go deeper than PyTorch/Python level ...

Serving Transformers / LLMs on Modern GPUs

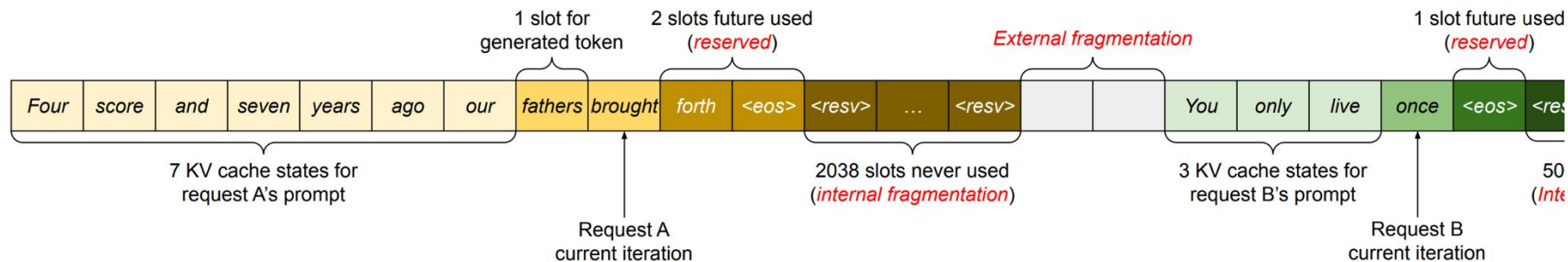
KV cache size in GPU memory determines:

- How long of a sequence a system could generate
- How many requests a system could respond to in a batch (parallel inference instead of sequential)



NVIDIA A100 40GB

Naïve KV Caching is Wasteful

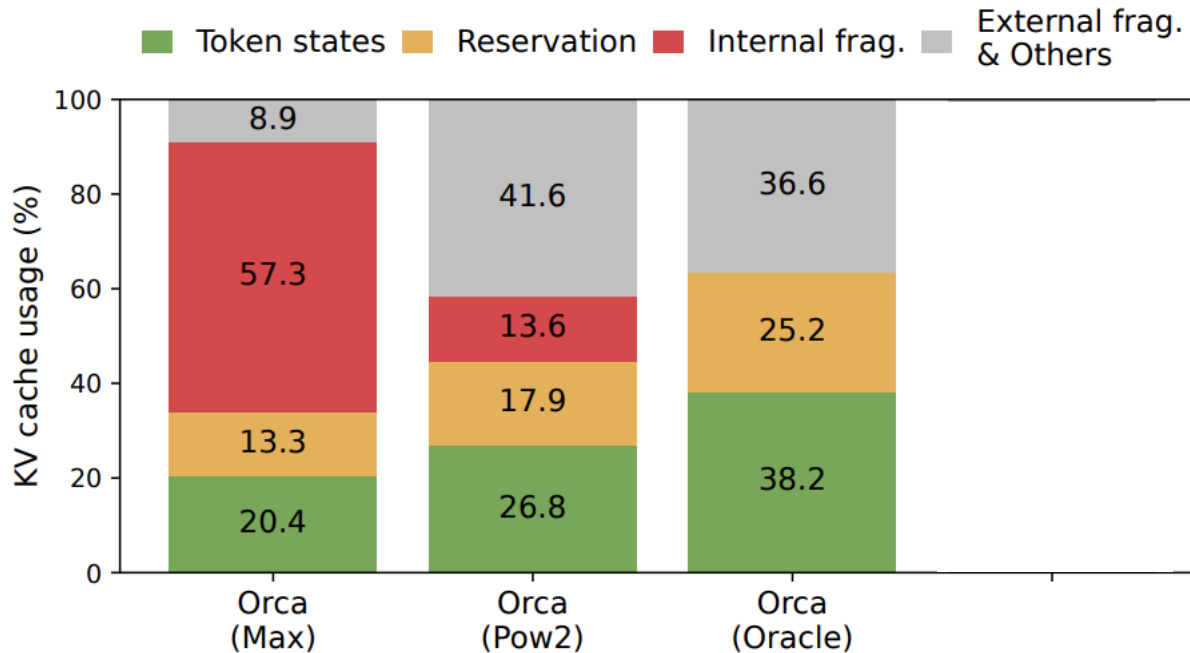


Internal fragmentation: memory reserved for a request but left unused because the final output length is shorter than the maximum length

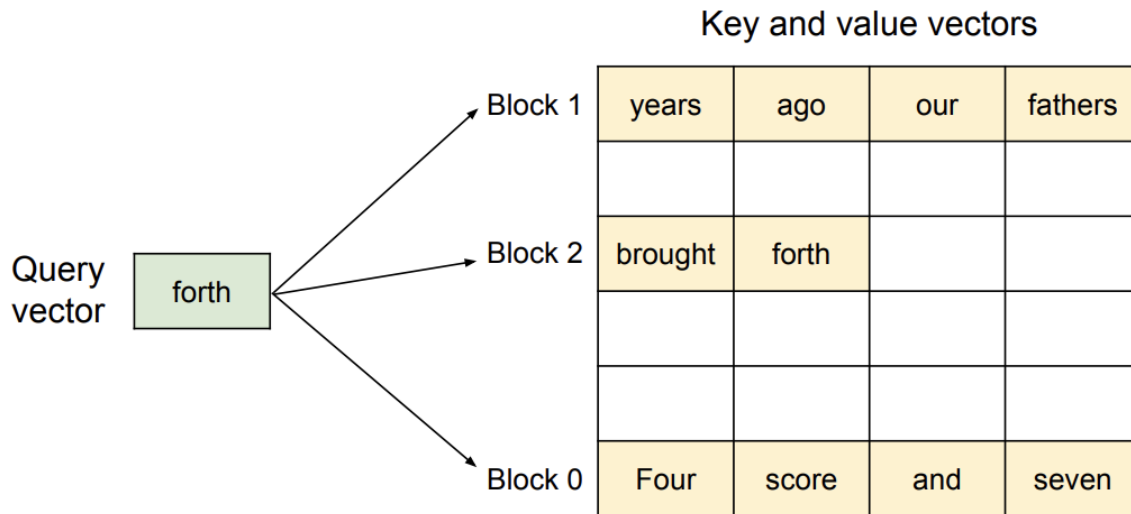
Reservation: memory currently unused but may be used by the same request

External fragmentation: memory gaps between requested allocations (A and B) --- different requests can have different maximum length

Naïve KV Caching is Wasteful

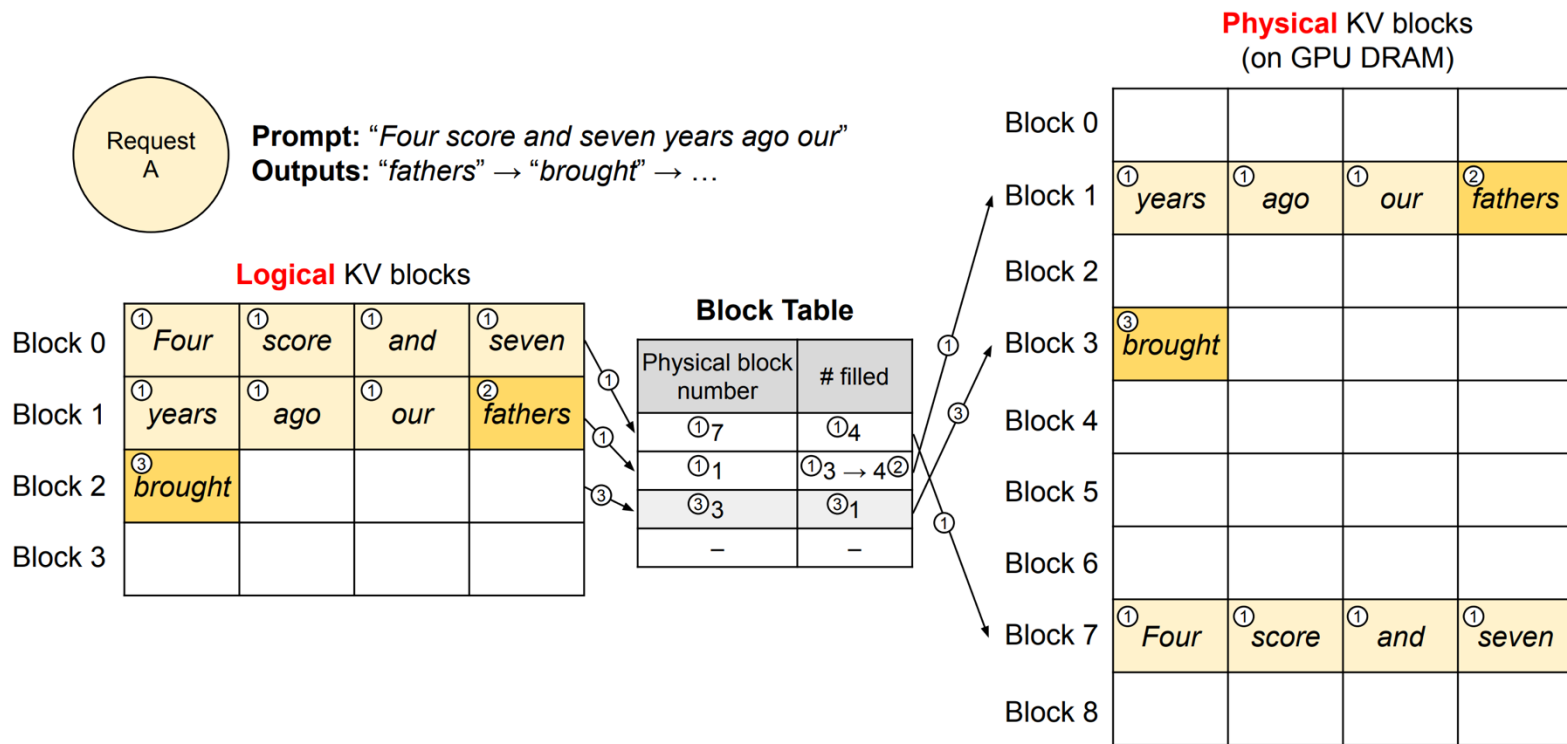


Idea: PagedAttention

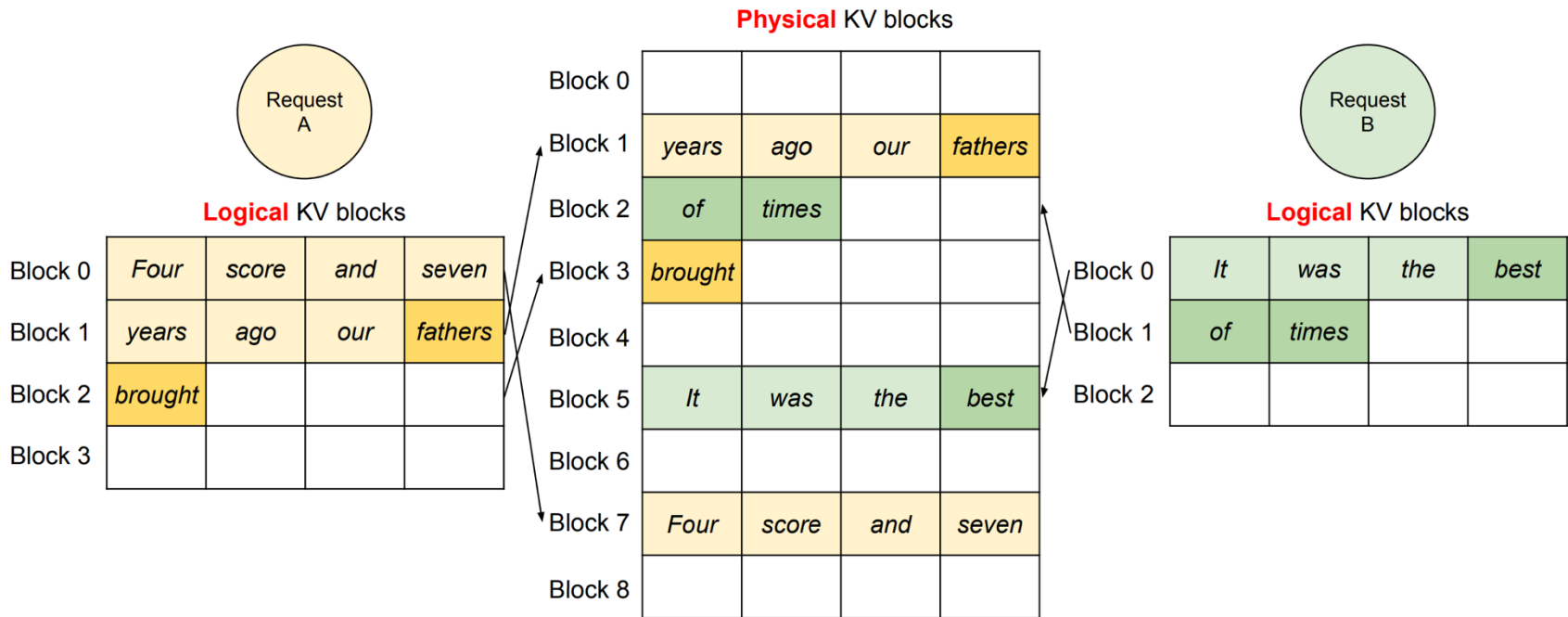


- Store KV vectors in **non-contiguous, fixed-sized blocks** in the memory
- Compute attention across all relevant blocks

Logical vs. Physical KV Blocks



Serving Multiple Requests



System: vLLM

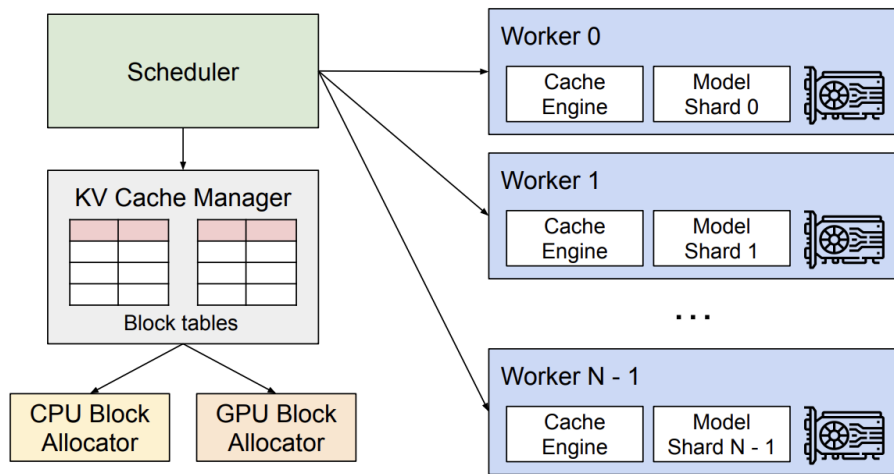
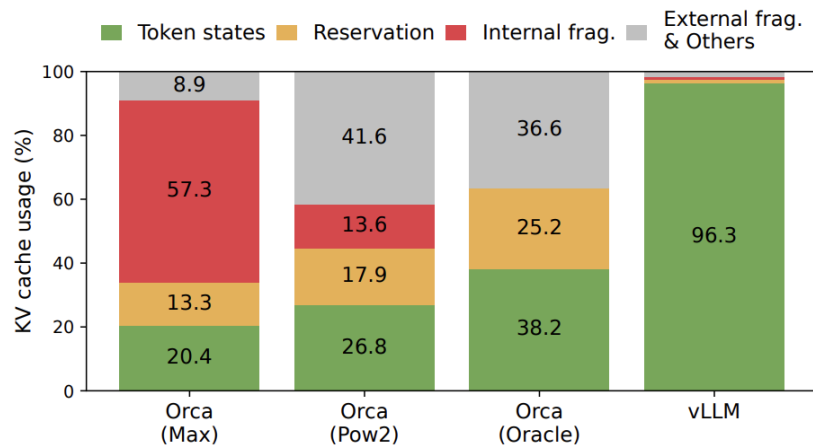


Figure 4. vLLM system overview.



System: vLLM

vLLM is an autoregressive-style transformer serving system that:

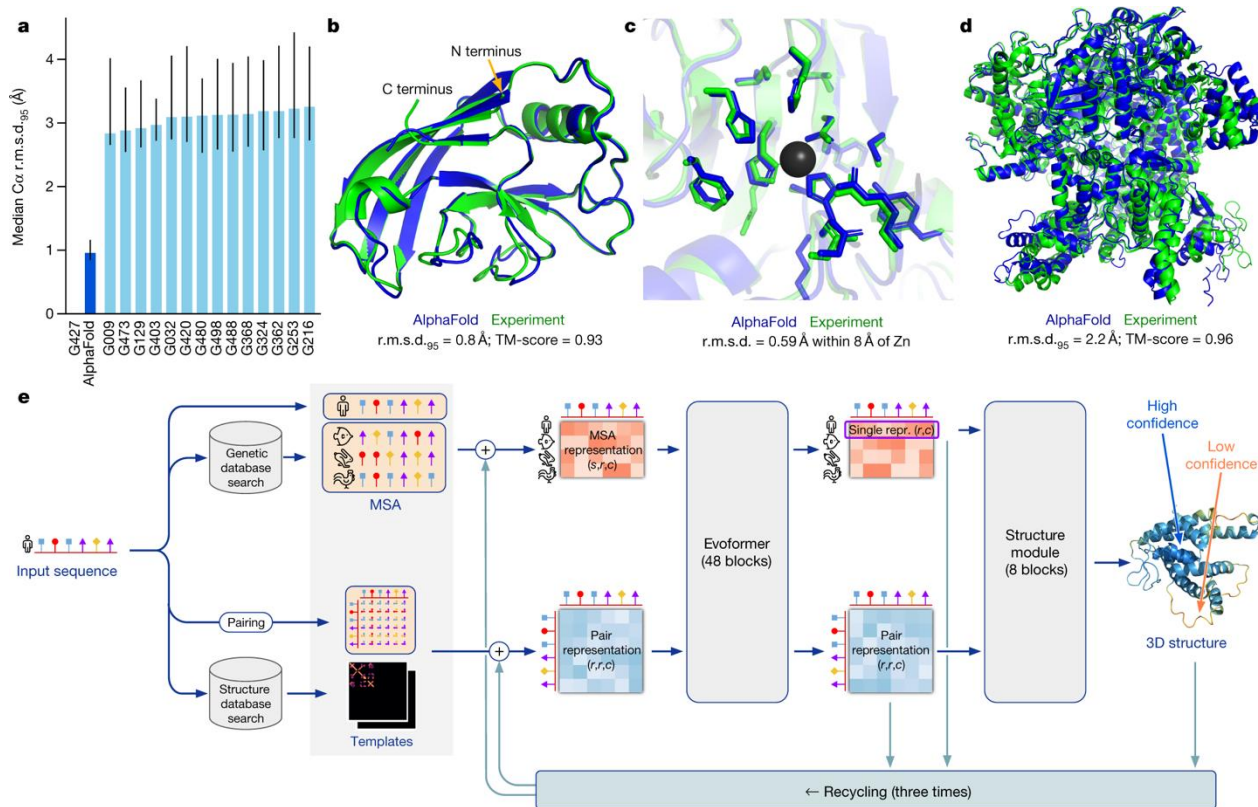
Breaks KV cache into small fixed-size pages (like OS virtual memory) instead of large contiguous blocks per sequence

Allocates pages on-demand as sequences grow, eliminating wasted pre-allocated memory

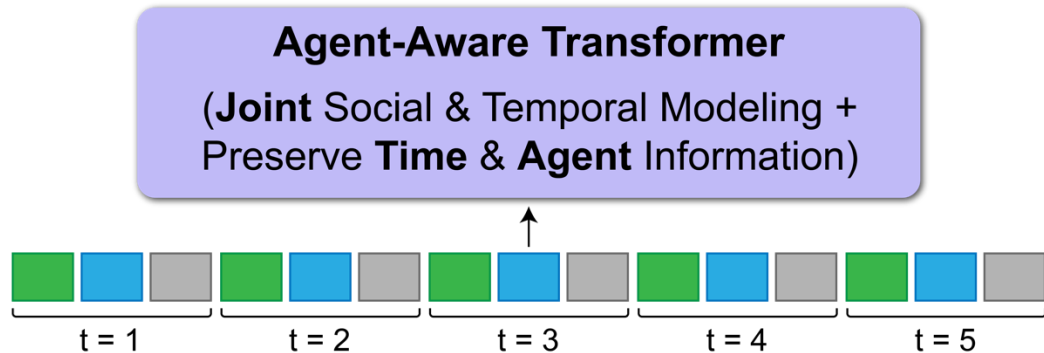
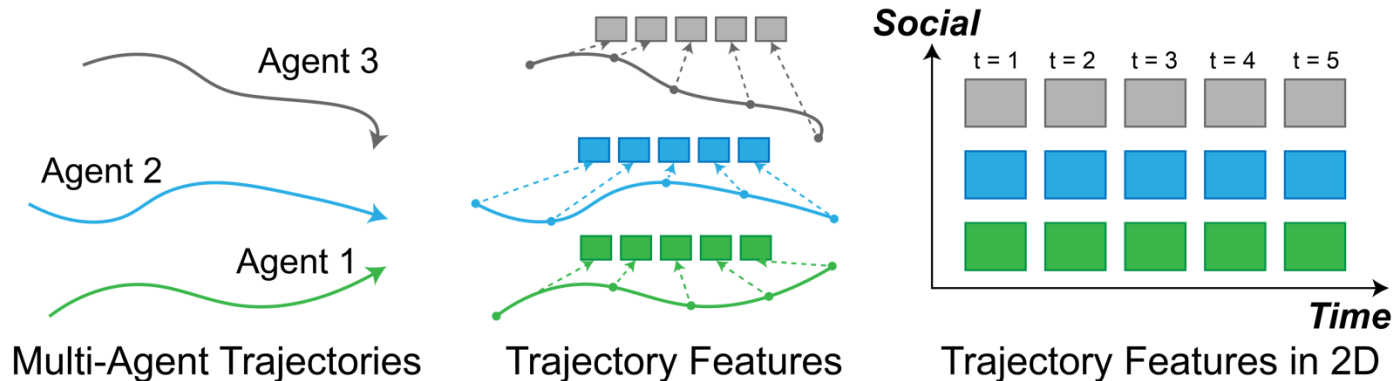
Shares pages between sequences with common prefixes (same prompt, beam search branches, etc.)

Can Attention/Transformers be used from
more than text processing?

Encoding/Decoding Protein Structures (AlphaFold)



Predicting Multi-agent Behaviors



ViLBERT: A Visolinguistic Transformer



pop artist performs at the
festival in a city.

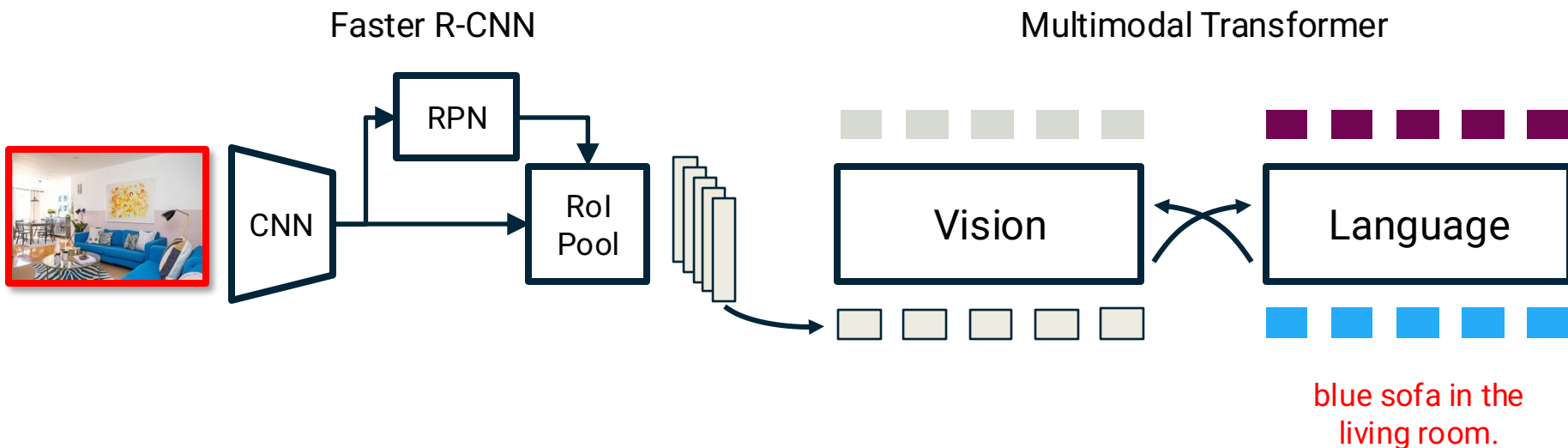


a worker helps to clear
the debris.



blue sofa in the
living room.

ViLBERT: A Visolinguistic Transformer



What about for just image inputs? Without Convolution?

Preprint. Under review.

AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

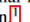
Alexey Dosovitskiy^{*,†}, Lucas Beyer^{*}, Alexander Kolesnikov^{*}, Dirk Weissenborn^{*},
Xiaohua Zhai^{*}, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer,
Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, Neil Houlsby^{*,†}

^{*}equal technical contribution, [†]equal advising

Google Research, Brain Team

{adosovitskiy, neilhoulby}@google.com

ABSTRACT

While the Transformer architecture has become the de-facto standard for natural language processing tasks, its applications to computer vision remain limited. In vision, attention is either applied in conjunction with convolutional networks, or used to replace certain components of convolutional networks while keeping their overall structure in place. We show that this reliance on CNNs is not necessary and a pure transformer applied directly to sequences of image patches can perform very well on image classification tasks. When pre-trained on large amounts of data and transferred to multiple mid-sized or small image recognition benchmarks (ImageNet, CIFAR-100, VTAB, etc.), Vision Transformer (ViT) attains excellent results compared to state-of-the-art convolutional networks while requiring substantially fewer computational resources to train. 

[cs.CV] 22 Oct 2020

Slide progression inspired by Soheil Feizi



How should we
“tokenize” images?

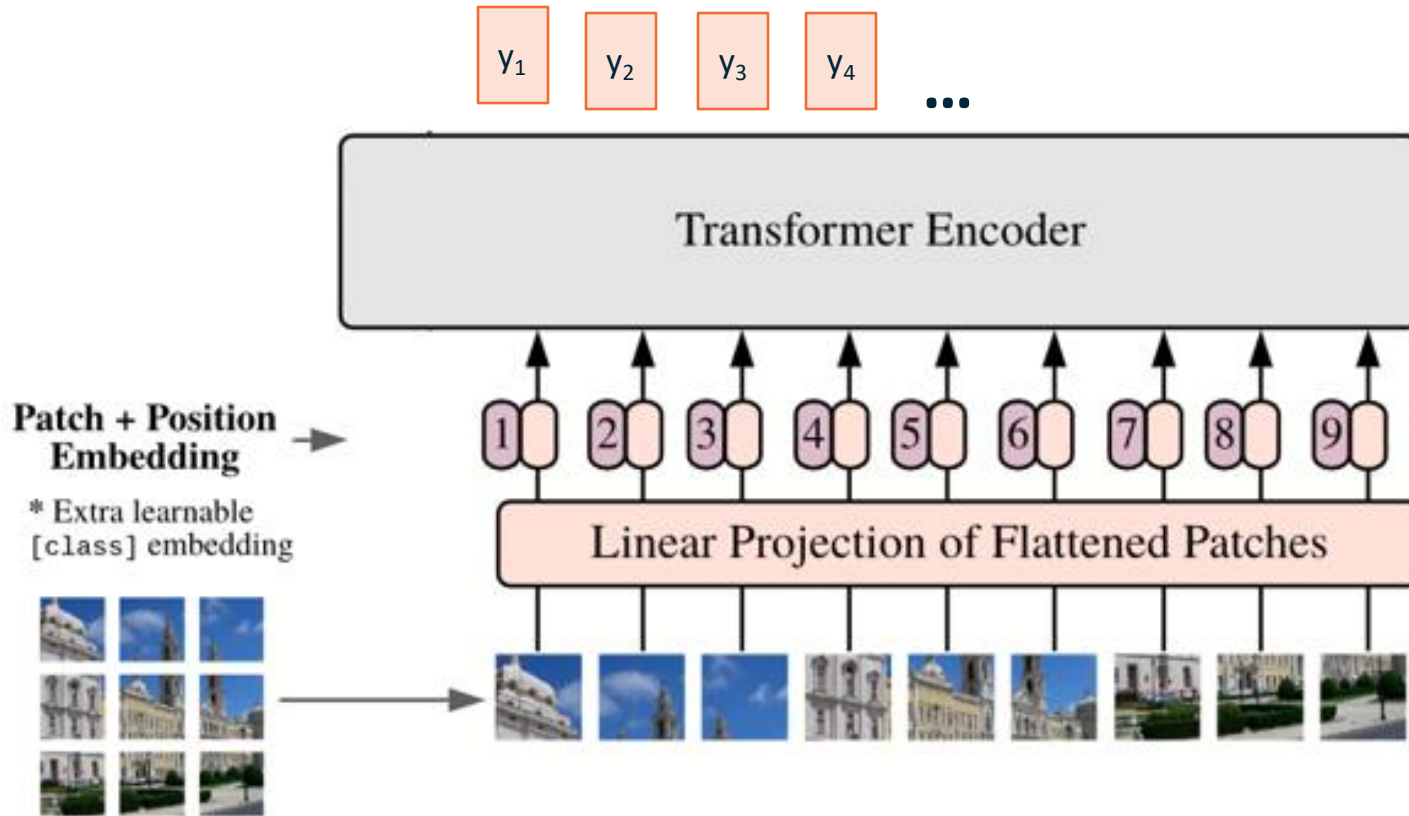
y_1 y_2 y_3 y_4

Self-Attention

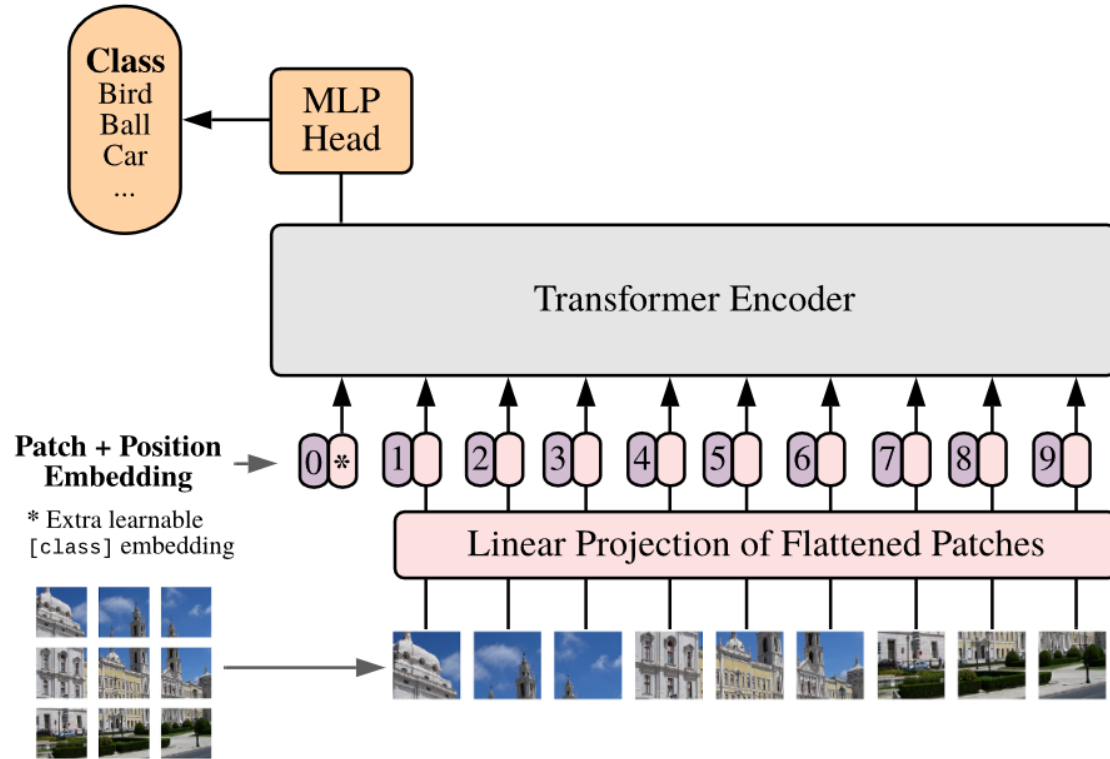
x_1 x_2 x_3 x_4

- **Pixels?** Too computationally intensive $O(n^2)$!
- **Patches!**

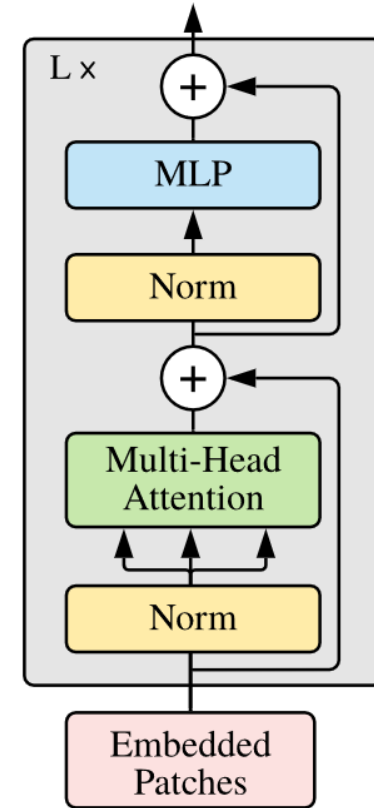
- How do we do classification?



Vision Transformer (ViT)



Transformer Encoder



Vision Transformer (ViT)

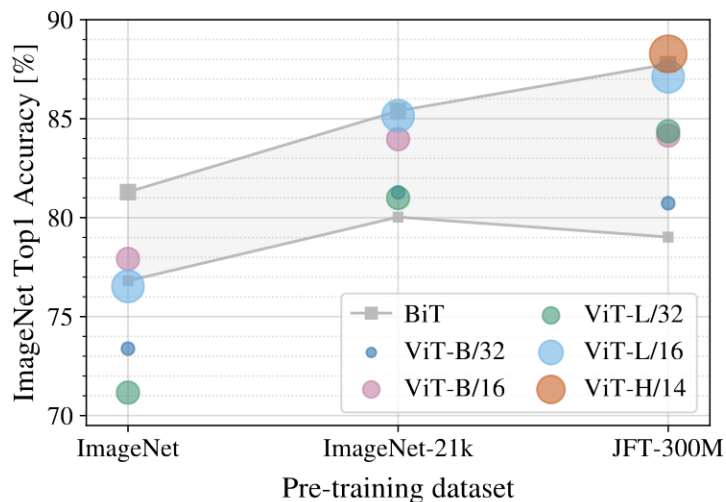


Figure 3: Transfer to ImageNet. While large ViT models perform worse than BiT ResNets (shaded area) when pre-trained on small datasets, they shine when pre-trained on larger datasets. Similarly, larger ViT variants overtake smaller ones as the dataset grows.

When trained on mid-sized datasets such as ImageNet, such models yield modest accuracies of a few percentage points below ResNets of comparable size.

Why?

Lacks some of the inductive biases:

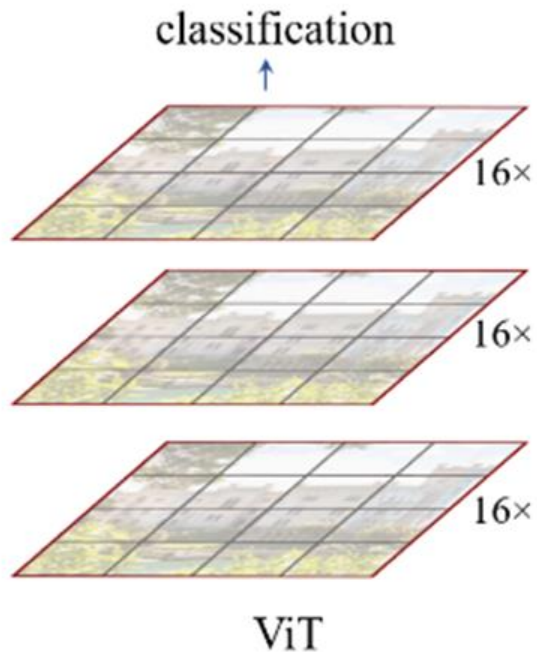
- Spatial locality
- Translation equivariance

Model	Layers	Hidden size D	MLP size	Heads	Params
ViT-Base	12	768	3072	12	86M
ViT-Large	24	1024	4096	16	307M
ViT-Huge	32	1280	5120	16	632M

However, the picture changes if the models are trained on larger datasets (14M-300M images). We find that large scale training trumps inductive bias.

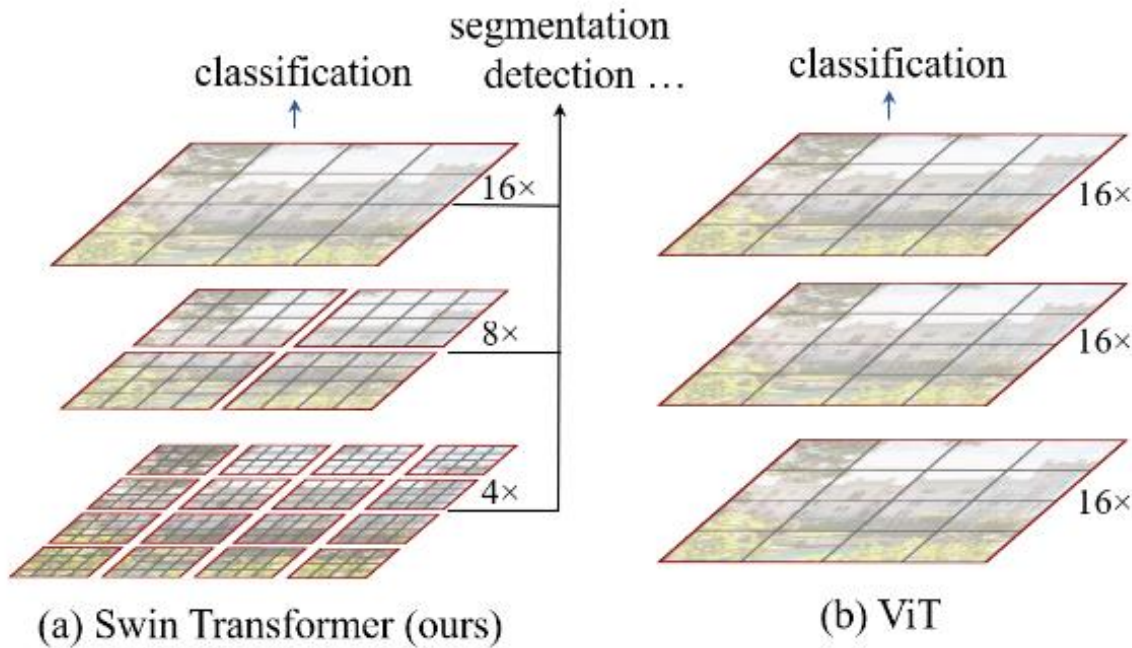
Table 1: Details of Vision Transformer model variants.

	Ours-JFT (ViT-H/14)	Ours-JFT (ViT-L/16)	Ours-I21K (ViT-L/16)	BiT-L (ResNet152x4)	Noisy Student (EfficientNet-L2)
ImageNet	88.55 ± 0.04	87.76 ± 0.03	85.30 ± 0.02	87.54 ± 0.02	88.4/88.5*
ImageNet ReaL	90.72 ± 0.05	90.54 ± 0.03	88.62 ± 0.05	90.54	90.55
CIFAR-10	99.50 ± 0.06	99.42 ± 0.03	99.15 ± 0.03	99.37 ± 0.06	—
CIFAR-100	94.55 ± 0.04	93.90 ± 0.05	93.25 ± 0.05	93.51 ± 0.08	—
Oxford-IIIT Pets	97.56 ± 0.03	97.32 ± 0.11	94.67 ± 0.15	96.62 ± 0.23	—
Oxford Flowers-102	99.68 ± 0.02	99.74 ± 0.00	99.61 ± 0.02	99.63 ± 0.03	—
VTAB (19 tasks)	77.63 ± 0.23	76.28 ± 0.46	72.72 ± 0.21	76.29 ± 1.70	—
TPUv3-core-days	2.5k	0.68k	0.23k	9.9k	12.3k



**What is
wrong with
this?**

$O((HW)^2)$ attention complexity!

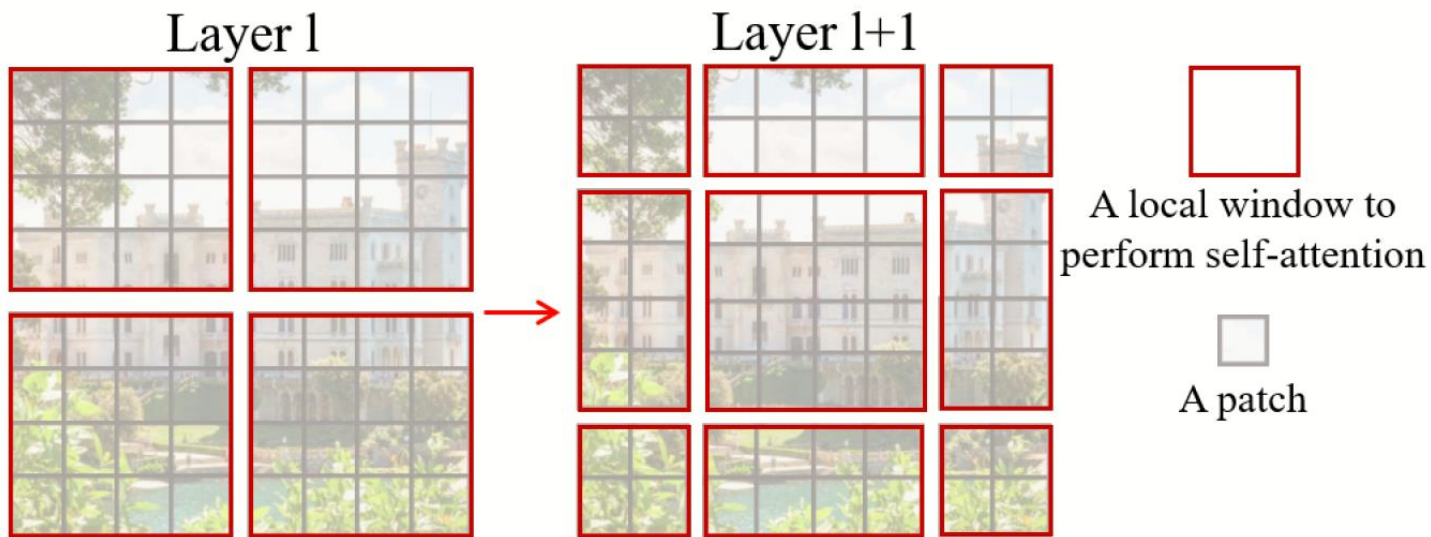


Swin Transformers key ideas:

- Feature pyramid!
- Self-attention only among patches within non-overlapping windows
- Shift windows across patches throughout the model
(Shifted **W**INdow)

Swin Transformer: Hierarchical Vision Transformer using Shifted Windows

Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, Baining Guo

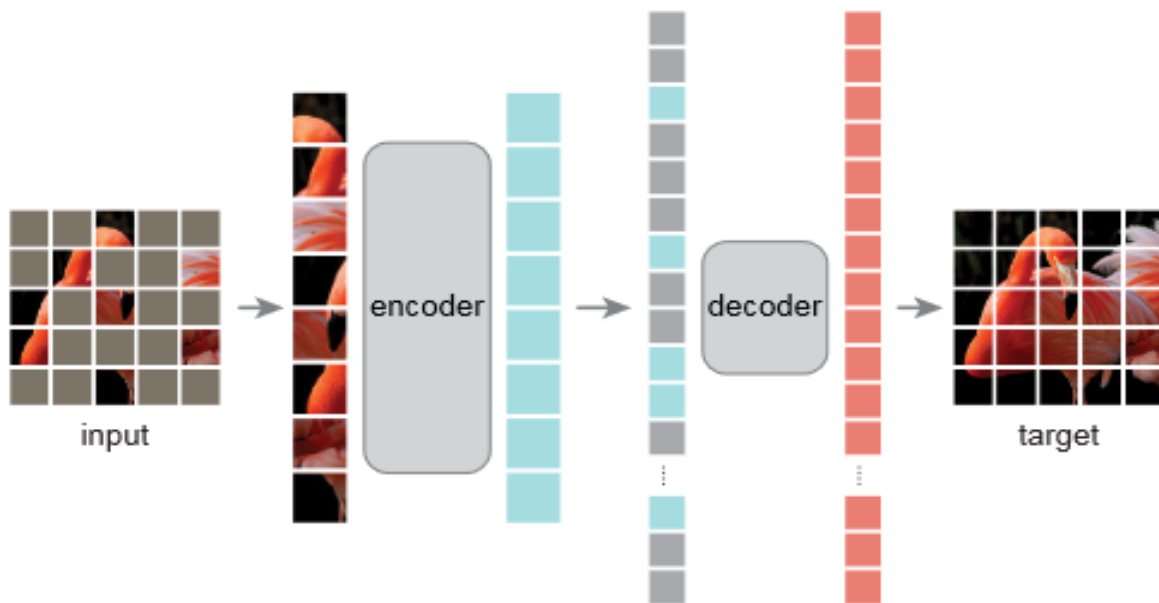


$O(HW)$ attention complexity
(number of windows is constant)

Swin Transformer: Hierarchical Vision Transformer using Shifted Windows

Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, Baining Guo

Shifted Window Attention



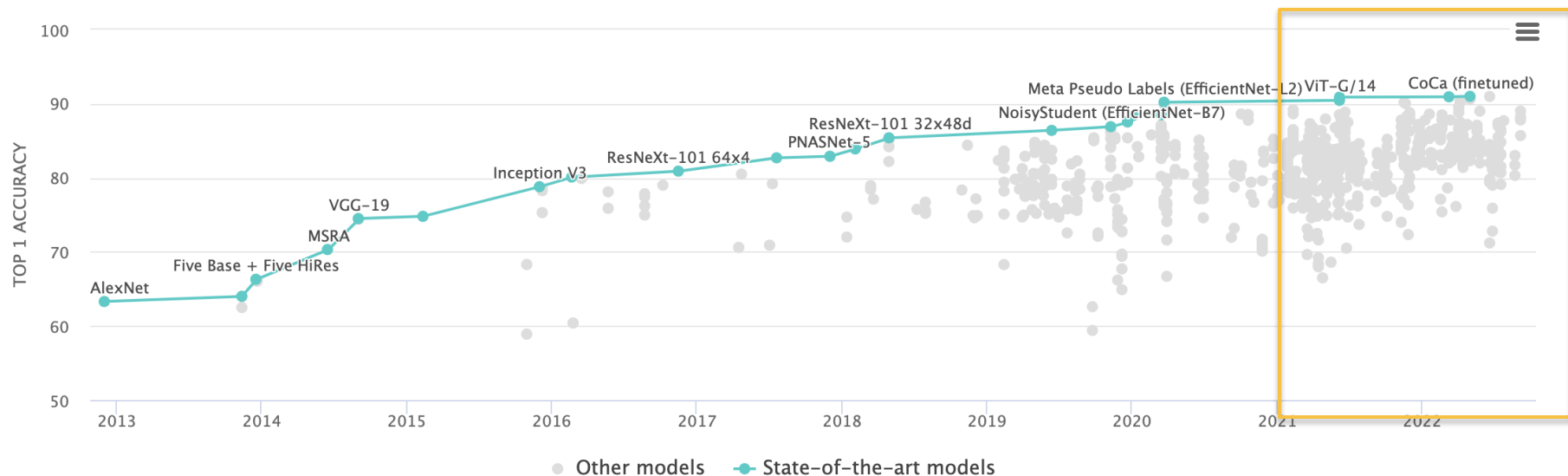
How can we learn
unsupervised
representations?

He et al., Masked Autoencoders Are Scalable Vision Learners

Masked Autoencoders

<https://paperswithcode.com/sota/instance-segmentation-on-coco>

ViT: Vision Transformer



Generally more expensive to train and execute than ConvNets-based models

Formal Algorithms for Transformers

Mary Phuong¹ and Marcus Hutter¹

¹DeepMind

This document aims to be a self-contained, mathematically precise overview of transformer architectures and algorithms (*not* results). It covers what transformers are, how they are trained, what they are used for, their key architectural components, and a preview of the most prominent models. The reader is assumed to be familiar with basic ML terminology and simpler neural network architectures such as MLPs.

Keywords: formal algorithms, pseudocode, transformers, attention, encoder, decoder, BERT, GPT, Gopher, tokenization, training, inference.

Contents

1	Introduction	1
2	Motivation	1
3	Transformers and Typical Tasks	3
4	Tokenization: How Text is Represented	4
5	Architectural Components	4
6	Transformer Architectures	7
7	Transformer Training and Inference	8
8	Practical Considerations	9
A	References	9
B	List of Notation	16

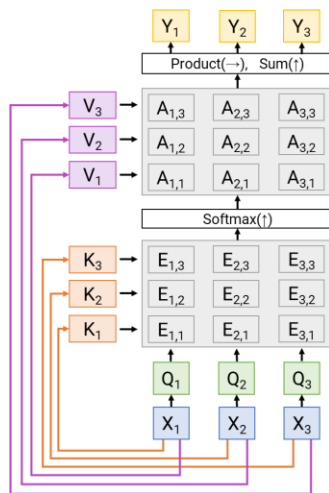
A famous colleague once sent an actually very well-written paper he was quite proud of to a famous complexity theorist. His answer: “I can’t find a theorem in the paper. I have no idea what this

plete, precise and compact overview of transformer architectures and formal algorithms (but *not* results). It covers what Transformers are (Section 6), how they are trained (Section 7), what they’re used for (Section 3), their key architectural components (Section 5), tokenization (Section 4), and a preview of practical considerations (Section 8) and the most prominent models.

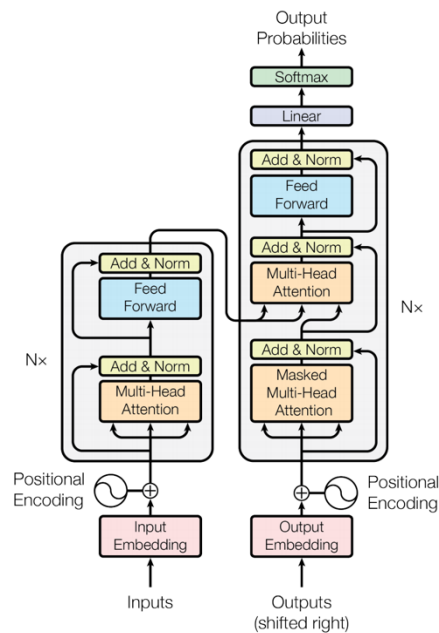
The essentially complete pseudocode is about 50 lines, compared to thousands of lines of actual real source code. We believe these formal algorithms will be useful for theoreticians who require compact, complete, and precise formulations, experimental researchers interested in implementing a Transformer from scratch, and

Summary

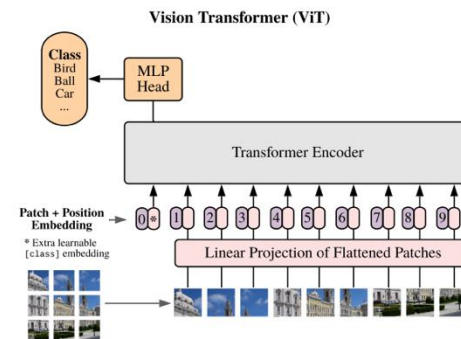
Self-Attention



Transformer Model



Beyond Language



Next Lecture (Virtual):

How to Train Your Large Language Models



Siddharth Karamcheti
(Incoming GT Prof.)