# CS 4644-DL / 7643-A: LECTURE 13
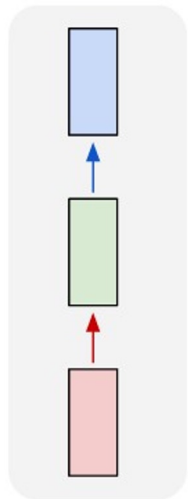# DANFEI XU

Attention for Sequence Modeling

Attention is (Mostly) All you Need: Transformers
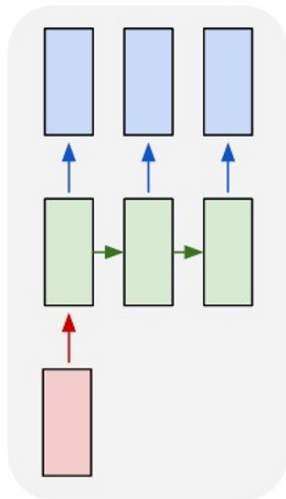
## Administrative:

- HW2 due 10/05 11:59pm + 48hr grace period.
- No class this Thu (10/05)
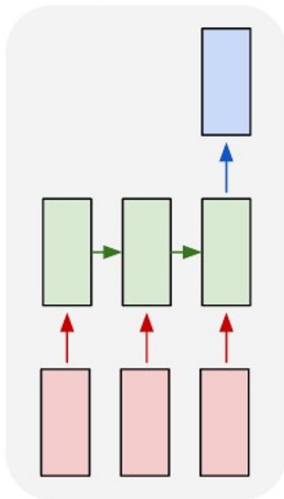
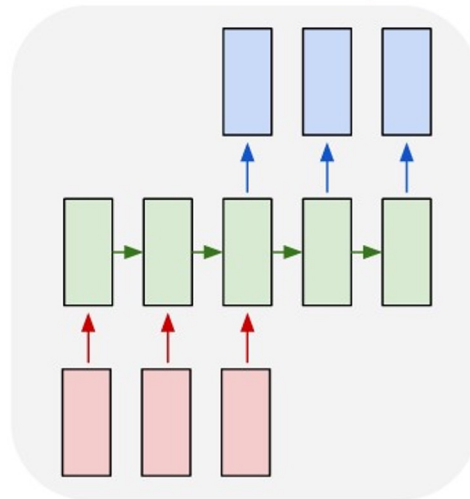# Recurrent Neural Networks: Process Sequences



one to one | one to many | many to one | many to many | many to many

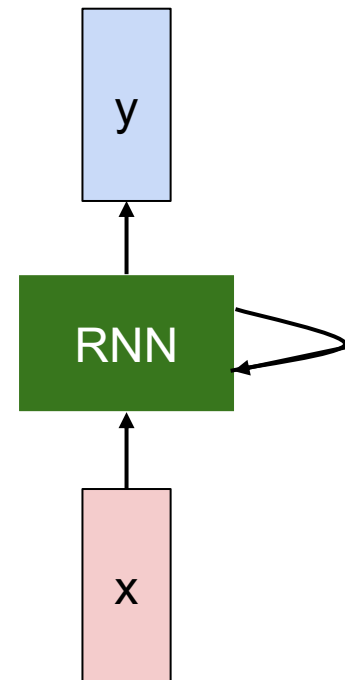# RNN hidden state update

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

$$h_t = f_W(h_{t-1}, x_t)$$
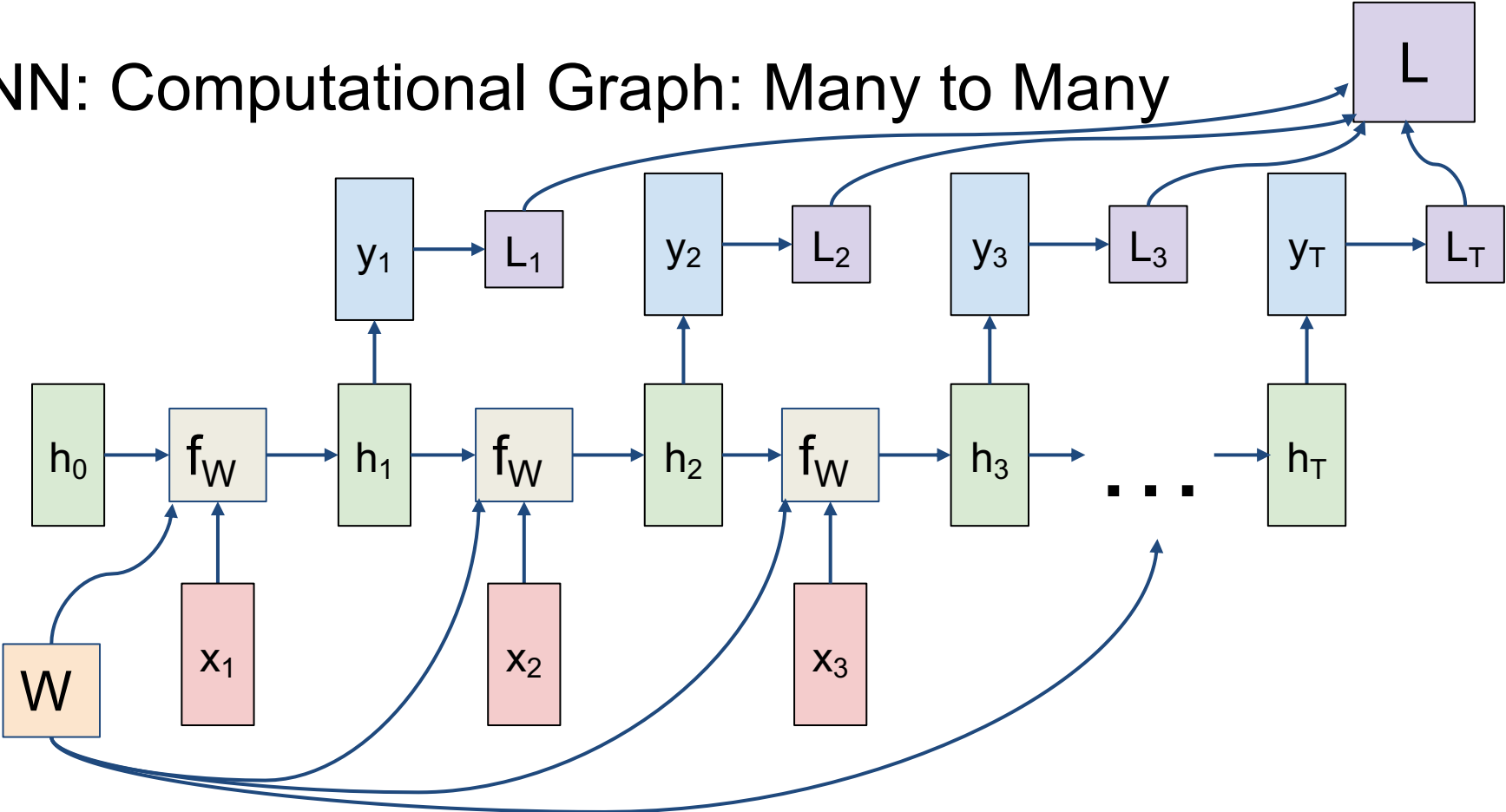
new state (vector)

some function with parameters W

old state (vector)

input vector at some time step

Can set initial state $h_0$ to all 0's

y

RNN

x

RNN: Computational Graph: Many to Many

# **Truncated** Backpropagation through time

# Vanilla RNN Gradient Flow

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

Gradients over multiple time steps:



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

Always < 1
**Vanishing gradients**

$tanh'$

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

# Vanilla RNN Gradient Flow

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

Gradients over multiple time steps:



What if we assumed no non-linearity?
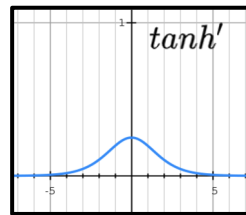
$$\frac{\partial L}{\partial W} = \sum_{t=1}^{T} \frac{\partial L_t}{\partial W}$$

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 (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

Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation 1997

# Long Short Term Memory (LSTM)

**Vanilla RNN**

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

**LSTM**

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

$$c_t = f \odot c_{t-1} + i \odot g$$
$$h_t = o \odot \tanh(c_t)$$

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

Two "memory vectors"

# Long Short Term Memory (LSTM)
*[Hochreiter et al., 1997]*

input from
below (**x**)

x

$h_{t-1}$

W

memory from
before (**h**)

4h x 2h

$c_{t-1}$

cell memory (c)

# Long Short Term Memory (LSTM)
*[Hochreiter et al., 1997]*

input from
below (**x**)

"gates"

x

$h_{t-1}$

W

memory from
before (**h**)

sigmoid $\longrightarrow$ i

sigmoid $\longrightarrow$ f

sigmoid $\longrightarrow$ o

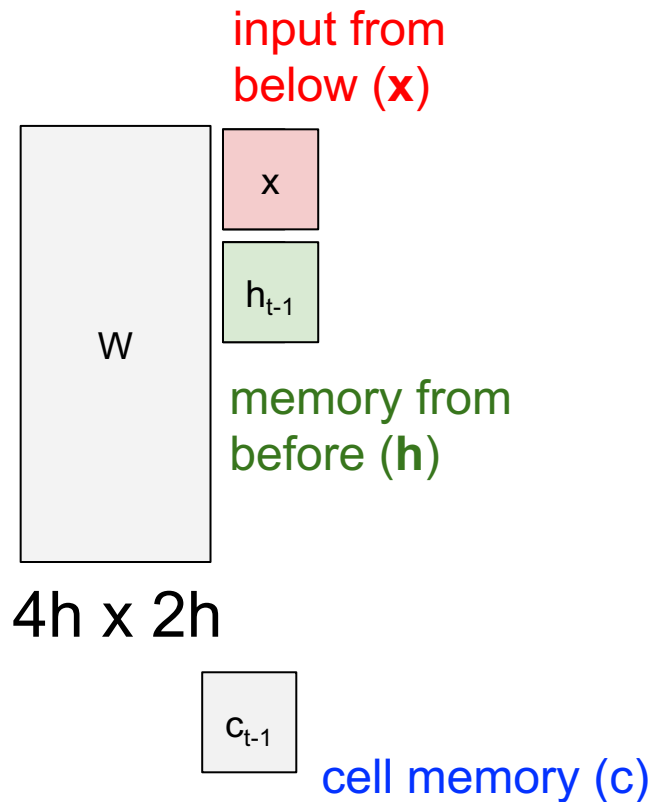tanh $\longrightarrow$ g

4h x 2h

4h

$c_{t-1}$

cell memory (c)

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

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

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

# Long Short Term Memory (LSTM)
*[Hochreiter et al., 1997]*

input from
below (**x**)

"gates"



memory from
before (**h**)

4h x 2h

4h

cell memory (c)

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

$$c_t = f \odot c_{t-1} + i \odot g$$
$$h_t = o \odot \tanh(c_t)$$

# Long Short Term Memory (LSTM)
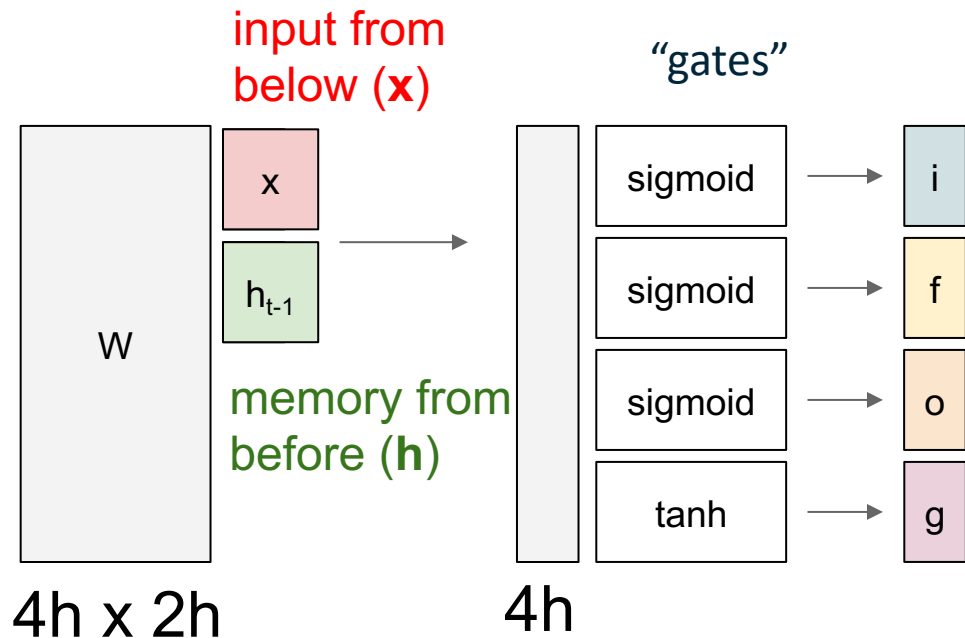*[Hochreiter et al., 1997]*

input from
below (**x**)

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

| | |
|---|---|
| W | x |
| | $h_{t-1}$ |

memory from
before (**h**)

4h x 2h

4h

sigmoid ⟶ i

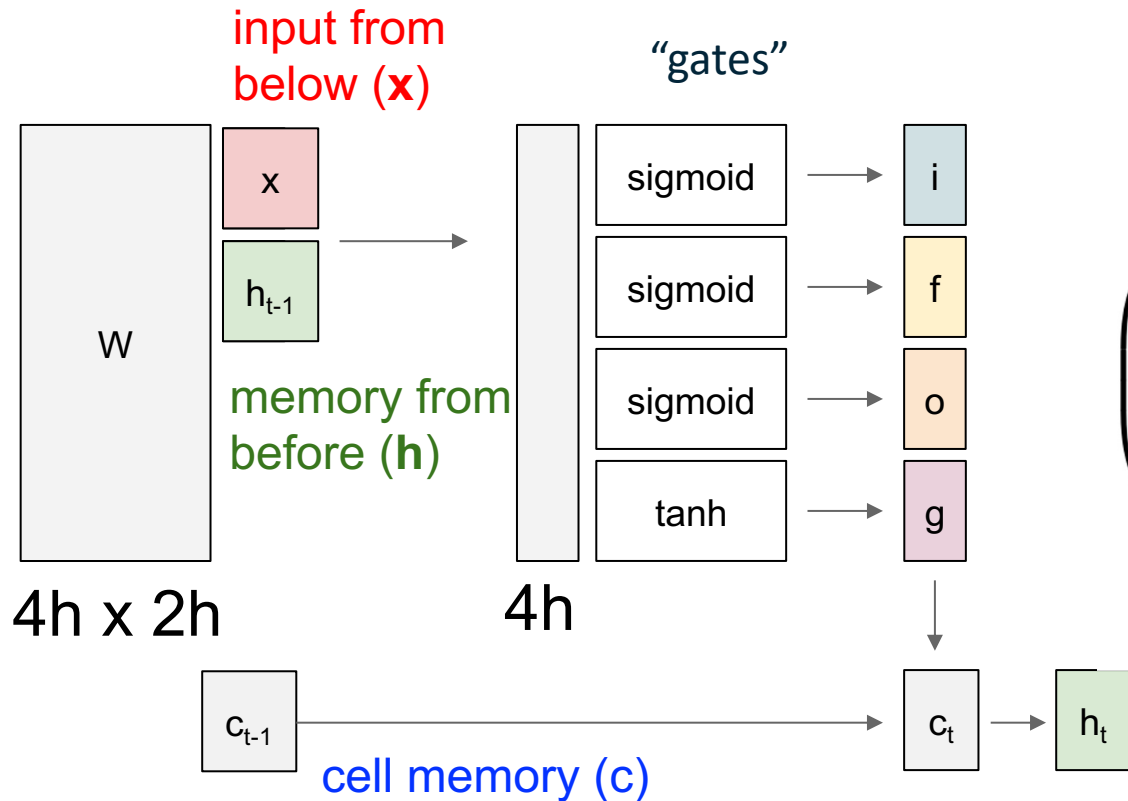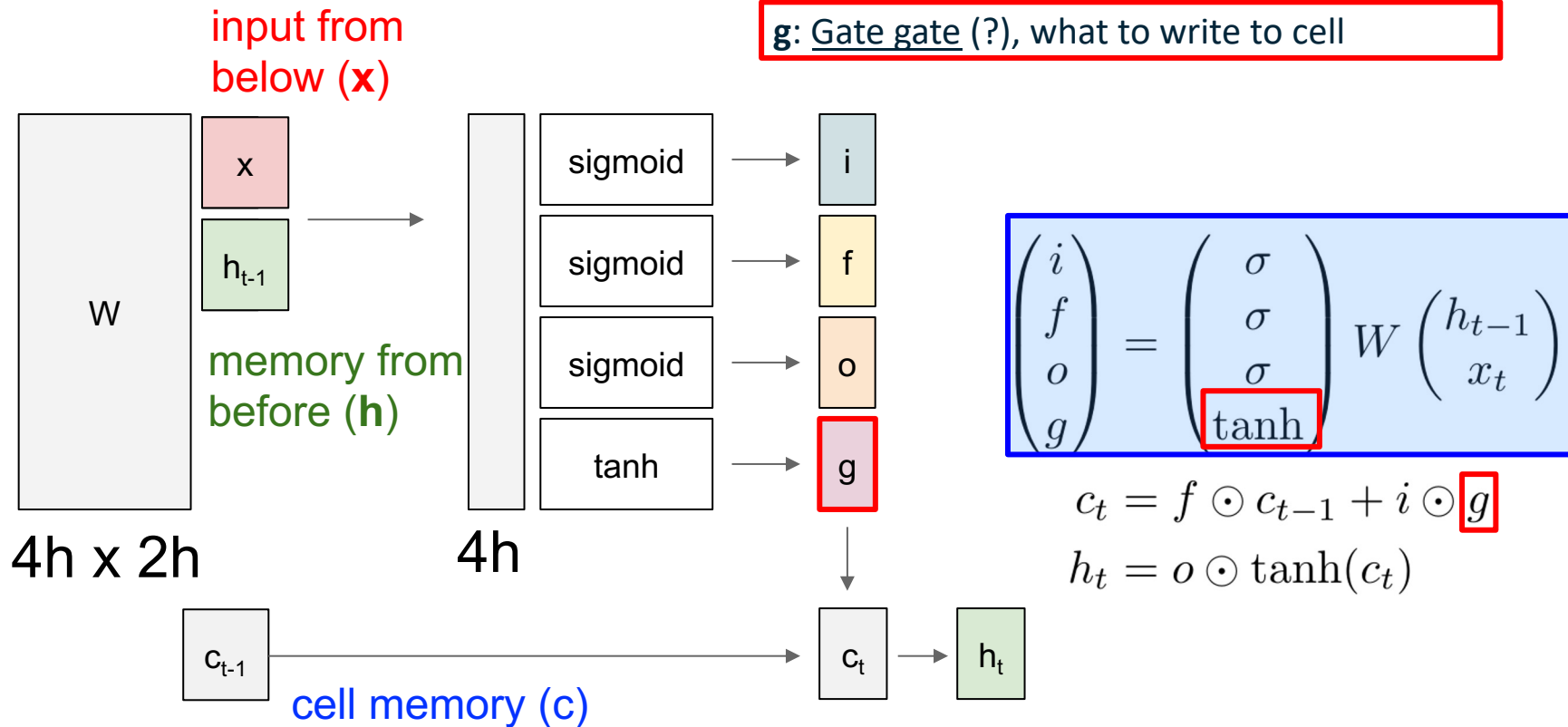sigmoid ⟶ f

sigmoid ⟶ o

tanh ⟶ g

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

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

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

$c_{t-1}$ ⟶ $c_t$ ⟶ $h_t$

cell memory (c)

# Long Short Term Memory (LSTM)
*[Hochreiter et al., 1997]*

**i**: Input gate, whether to write to cell

input from below (**x**)

"gates"

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



memory from before (**h**)

4h x 2h

4h

cell memory (c)

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

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

$$h_t = o \odot \tanh(c_t)$$
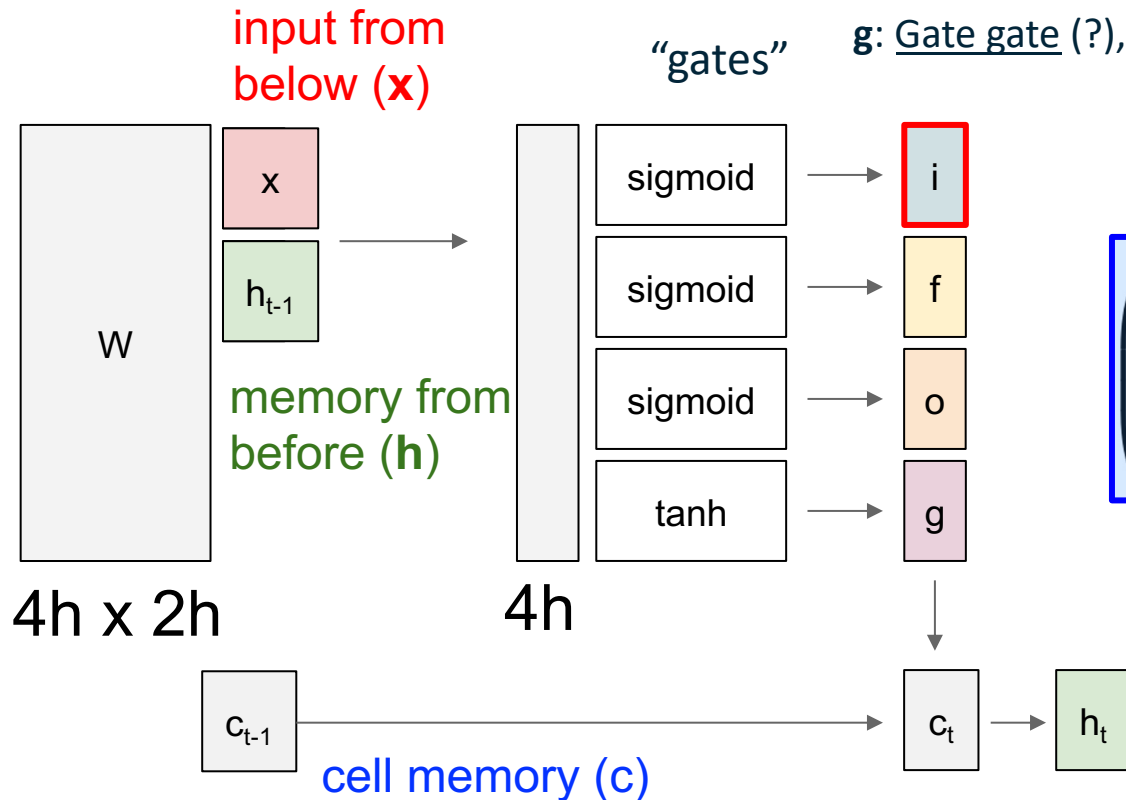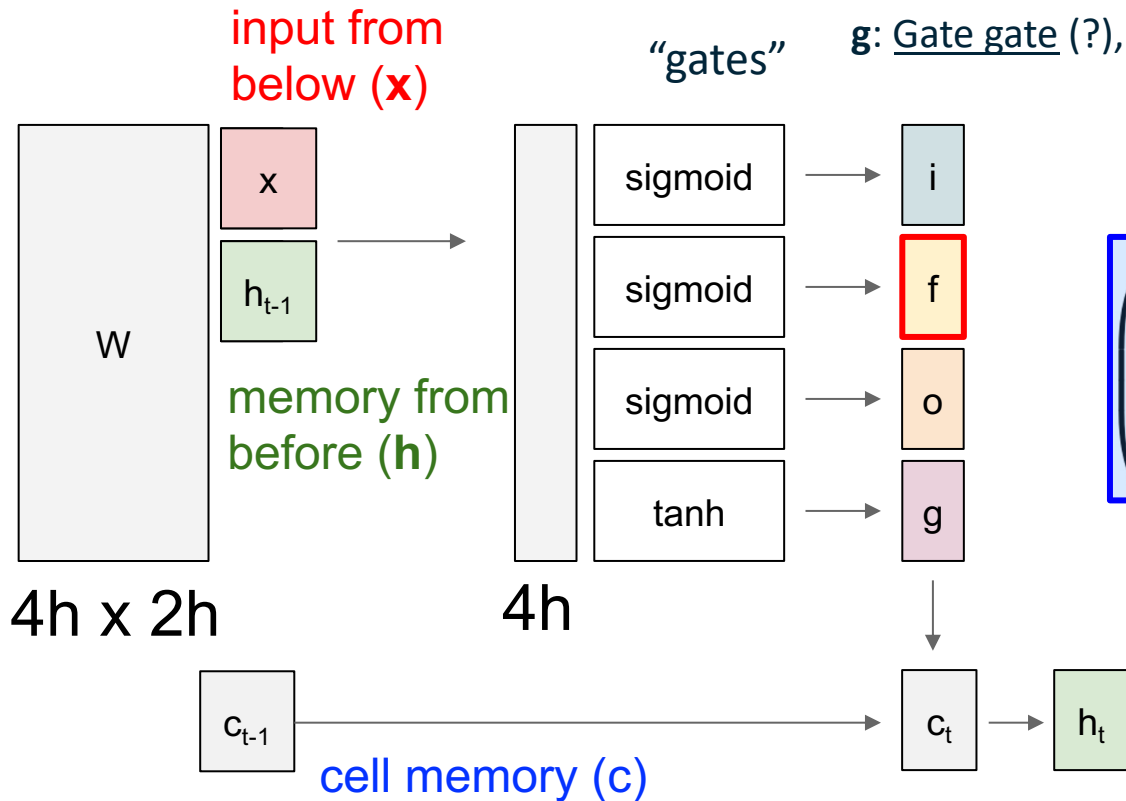
# Long Short Term Memory (LSTM)
*[Hochreiter et al., 1997]*

**i**: <u>Input gate</u>, whether to write to cell

**f**: <u>Forget gate</u>, whether to erase cell

**g**: <u>Gate gate</u> (?), what to write to cell

input from below (**x**)

"gates"



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

$$c_t = \boxed{f} \odot c_{t-1} + i \odot g$$
$$h_t = o \odot \tanh(c_t)$$

memory from before (**h**)

4h x 2h        4h

cell memory (c)

# Long Short Term Memory (LSTM)
*[Hochreiter et al., 1997]*

**i**: <u>Input gate</u>, whether to write to cell
**f**: Forget gate, whether to erase cell
**o**: <u>Output gate</u>, how much to reveal cell
**g**: <u>Gate gate</u> (?), what to write to cell

input from
below (**x**)

"gates"



memory from
before (**h**)

4h x 2h

4h

cell memory (c)

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

$$c_t = f \odot c_{t-1} + i \odot g$$
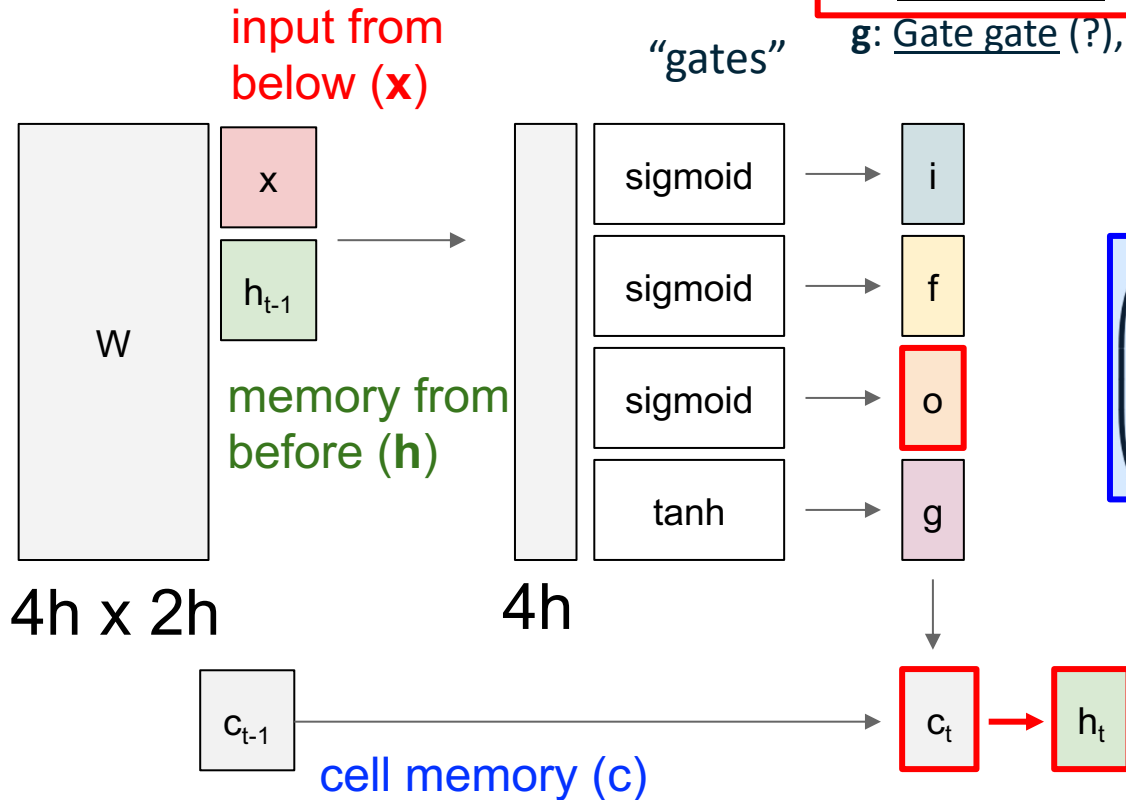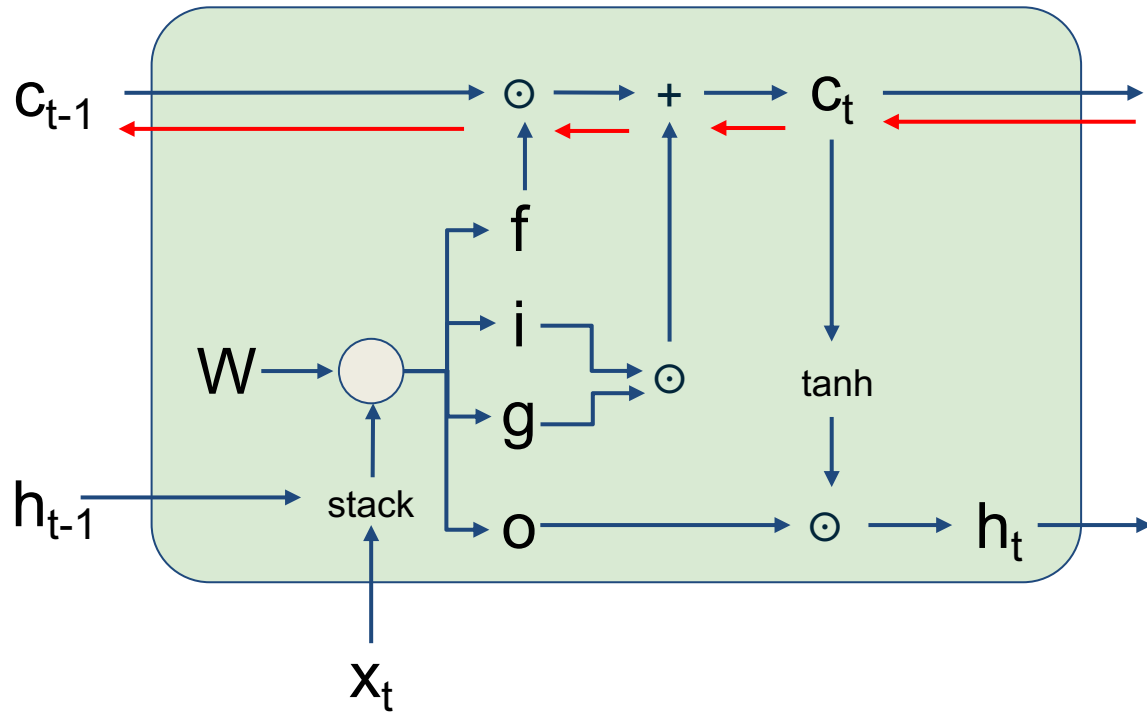$$h_t = o \odot \tanh(c_t)$$

# Long Short Term Memory (LSTM): Gradient Flow
*[Hochreiter et al., 1997]*

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



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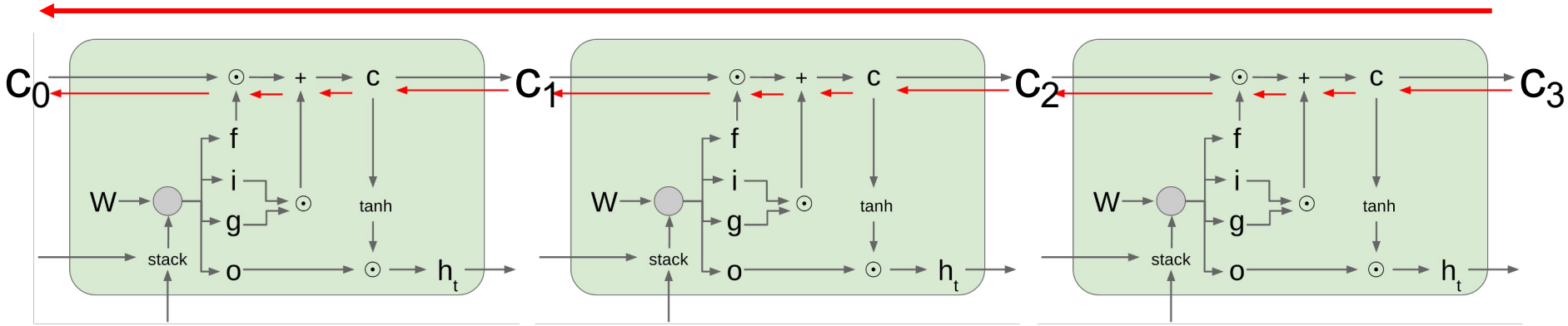
$$c_t = f \odot c_{t-1} + i \odot g$$

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

# Long Short Term Memory (LSTM): Gradient Flow
*[Hochreiter et al., 1997]*

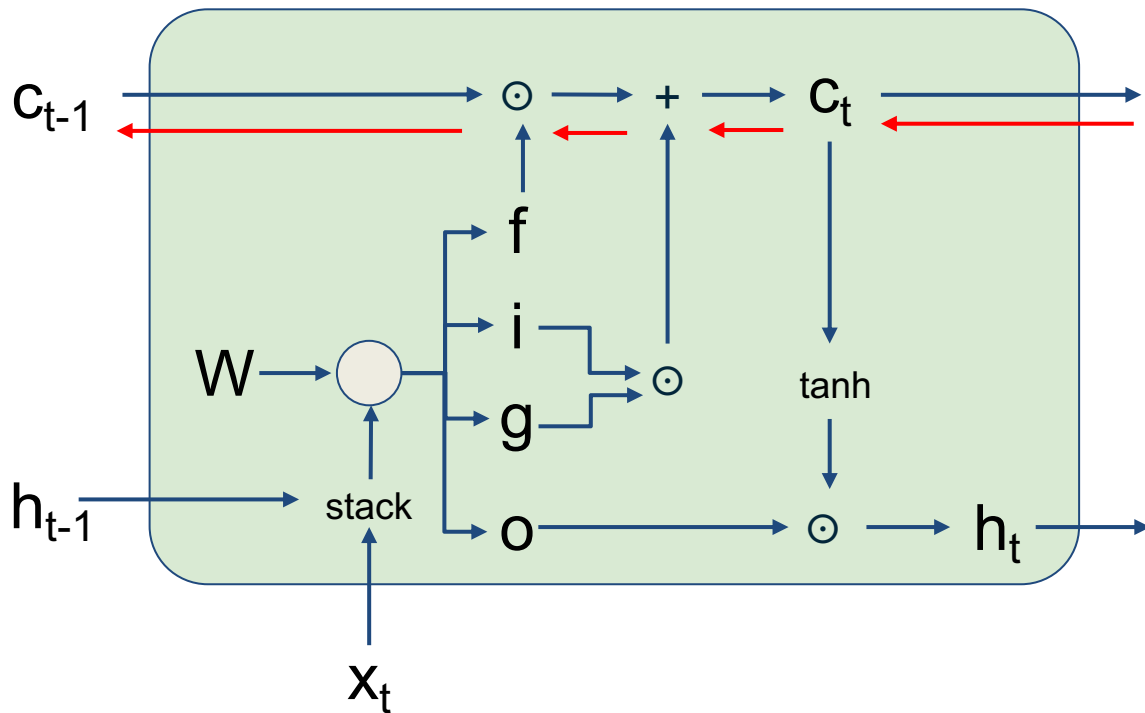## Uninterrupted gradient flow!



Notice that the gradient contains the **f** gate's vector of activations
- allows better control of gradients values, using suitable parameter updates of the forget gate.

The hidden state is emitted from *c* with an output gate (o), instead of recurrent multiplication with a weight vector.

# Long Short Term Memory (LSTM): Gradient Flow
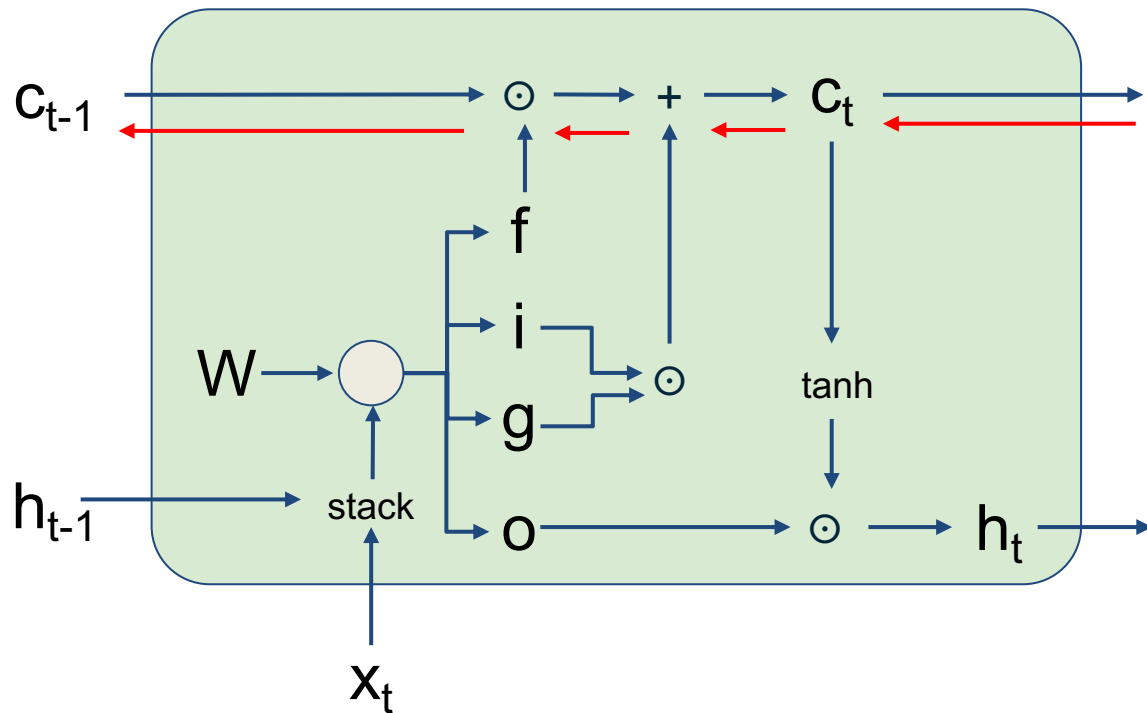*[Hochreiter et al., 1997]*



Q: What if f = 1 and i = 0?

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

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

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

# Long Short Term Memory (LSTM): Gradient Flow
*[Hochreiter et al., 1997]*



Q: What if f = 1 and i = 0?

A: LSTM doesn't forget / take in new information!

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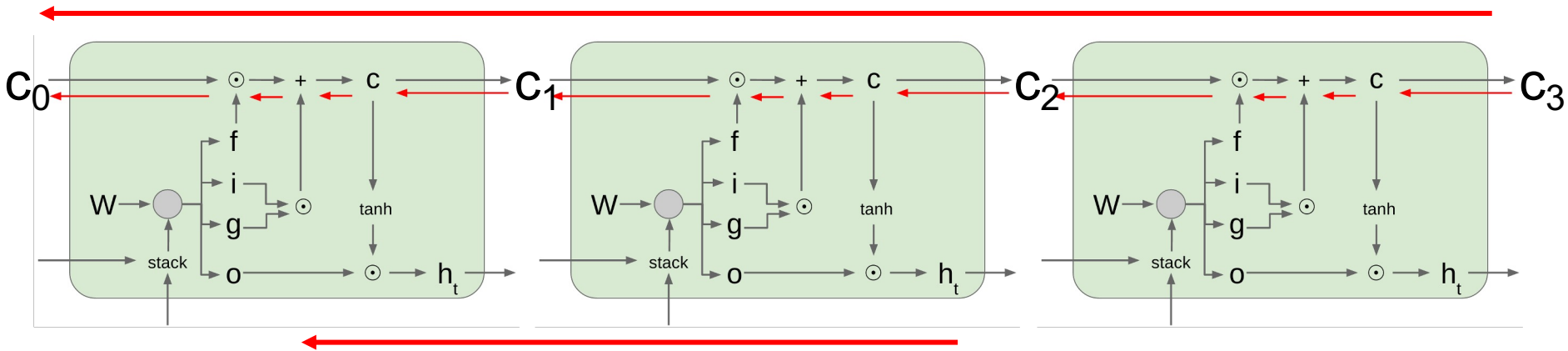
$$c_t = f \odot c_{t-1} + i \odot g$$
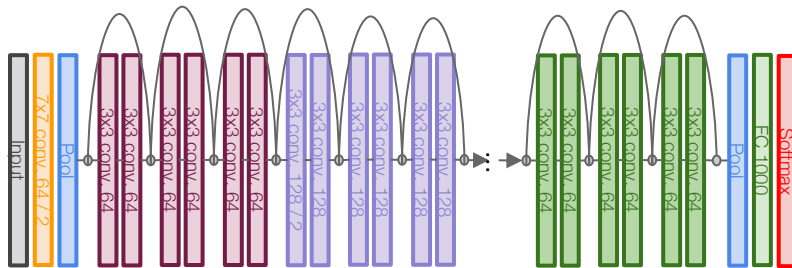
$$h_t = o \odot \tanh(c_t)$$

# Long Short Term Memory (LSTM): Gradient Flow
*[Hochreiter et al., 1997]*

## Uninterrupted gradient flow!



Similar to ResNet!

Possible to learn to set f = 1 and use i and g to learn "memory residual"

# Summary: LSTM

The LSTM architecture makes it easier for the RNN to preserve information over many timesteps
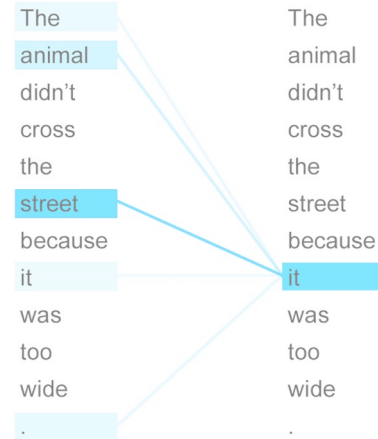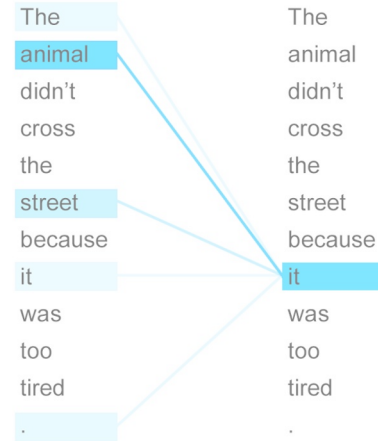- e.g. if the f = 1 and the i = 0, then the information of that cell is preserved indefinitely.
- By contrast, it's harder for vanilla RNN to learn a recurrent weight matrix Wh that preserves info in hidden state

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.

Possible to mitigate vanishing / exploding gradient by learning suitable i and f.

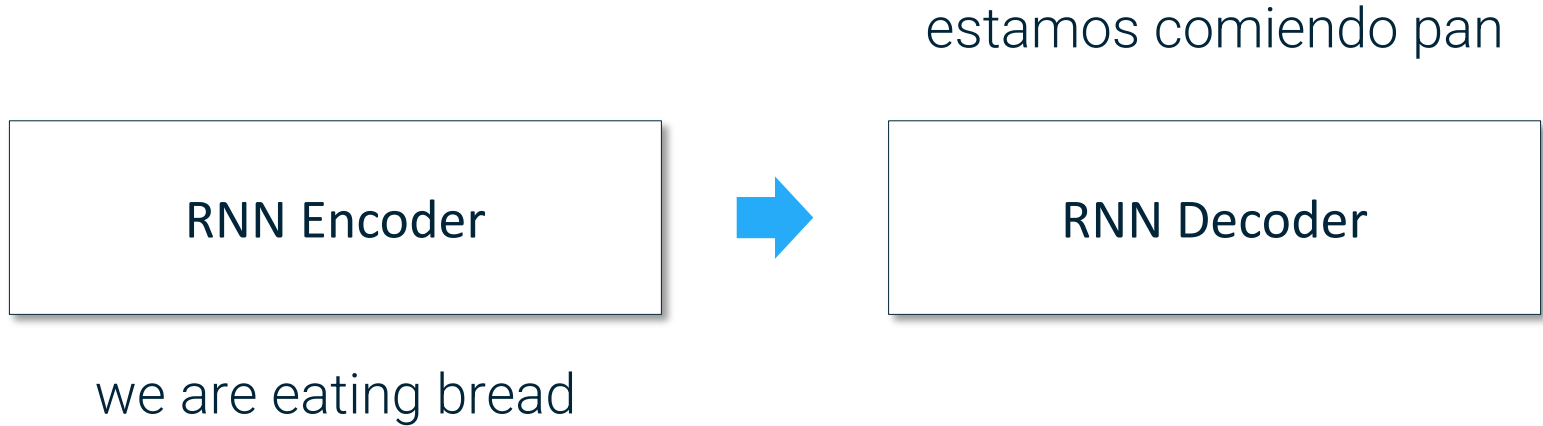**RNNs / LSTMs are still forgetful. Hard to represent a long sequence with a compact memory vector**

# Attention Mechanism

# Example: Machine Translation

estamos comiendo pan

| RNN Encoder | | RNN Decoder |

we are eating bread

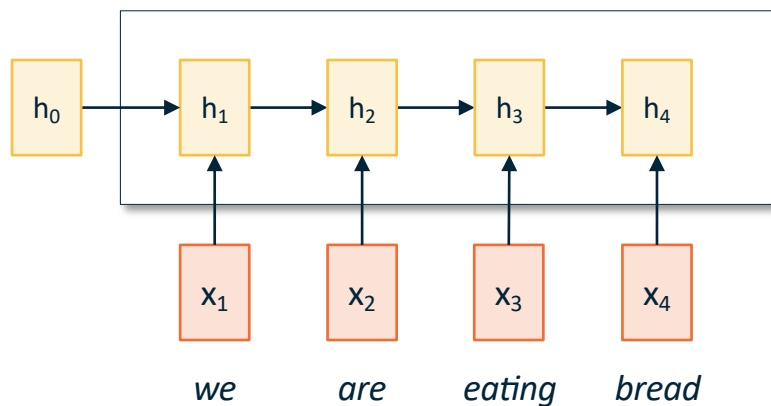# Machine Translation with RNNs

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

# Machine Translation with RNNs

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

$s_0 = h_4$



we    are    eating    bread

# Machine Translation with RNNs

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

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

# Machine Translation with RNNs

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

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

# Machine Translation with RNNs

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

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



Slide credit: Justin Johnson

# Machine Translation with RNNs

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

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

Problem: $s_i$ is used to encode input and maintain decoder state



Slide credit: Justin Johnson

# Machine Translation with RNNs

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

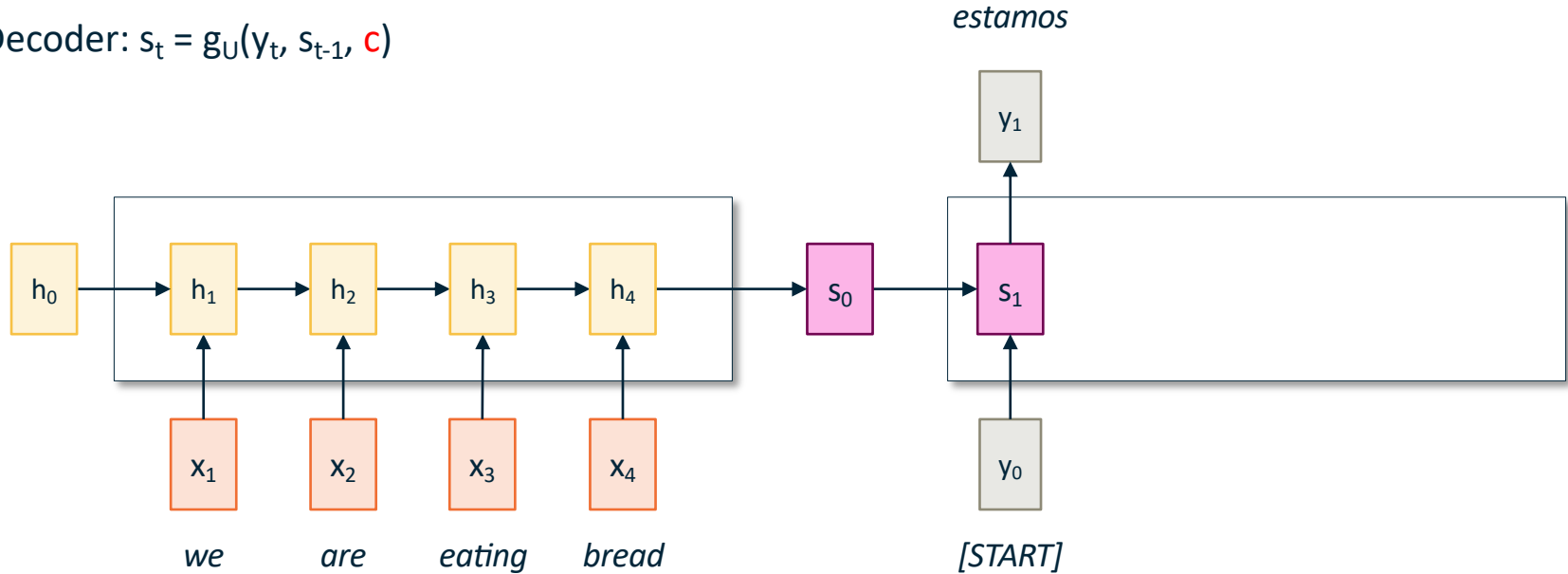Decoder: $s_t = g_U(y_t, s_{t-1}, c)$

Solution: add a context vector $c = h_4$ and generate $s_0$ from $h_4$

# Machine Translation with RNNs

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

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

Solution: add a context vector $c = h_4$ and generate $s_0$ from $h_4$

# Machine Translation with RNNs

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

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



Problem: Input sequence **bottlenecked** through fixed-sized memory vector.
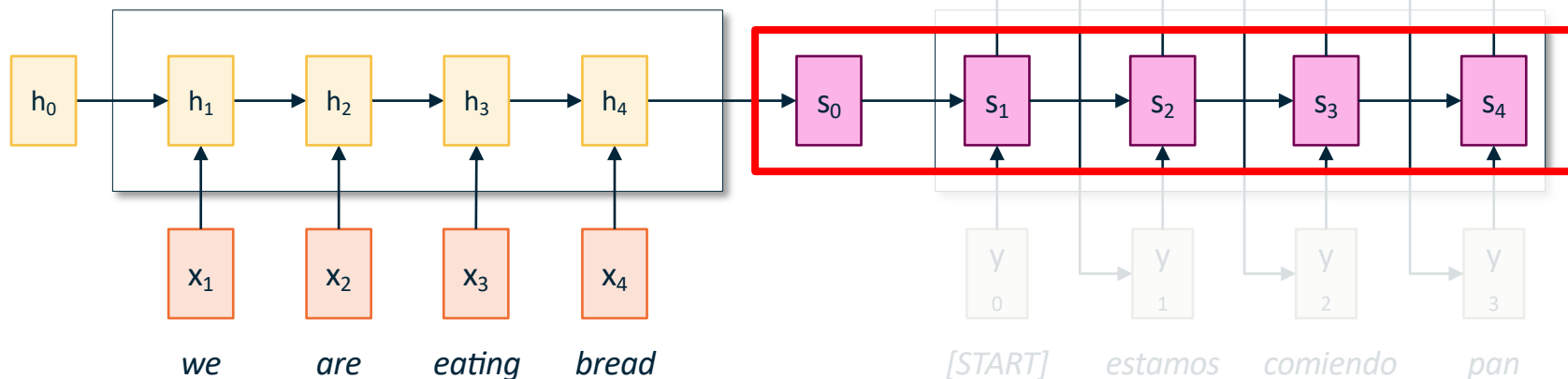
# Machine Translation with RNNs

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

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



bottleneck

Idea: can we "look up" information from the input sequence when making prediction?

# Machine Translation with RNNs **and Attention**

From final hidden state:
**Initial decoder state** $s_0$



we    are    eating    bread

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Slide credit: Justin Johnson

# Machine Translation with RNNs **and Attention**

Compute **affinity scores**

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

From final hidden state:
**Initial decoder state** $s_0$



Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

# Machine Translation with RNNs **and Attention**



Compute **affinity scores**
$$e_{t,i} = f_{att}(s_{t-1}, h_i) \qquad (f_{att} \text{ is an MLP})$$

Normalize to get
**attention weights**
$$0 < a_{t,i} < 1 \qquad \sum_i a_{t,i} = 1$$

From final hidden state:
**Initial decoder state** $s_0$

a₁ a₁ a₃ a₄

softmax

e₁ e₂ e₃ e₄

h₁ h₂ h₃ h₄ s₀

x₁ x₂ x₃ x₄

we    are    eating    bread

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

# Machine Translation with RNNs **and Attention**



Think of how much information needed from each input word to generate the first translated word

From final hidden state:
**Initial decoder state** $s_0$

Compute **affinity scores**

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

Normalize to get
**attention weights**

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

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

# Machine Translation with RNNs **and Attention**



Compute **affinity scores**

$e_{t,i} = f_{att}(s_{t-1}, h_i)$      ($f_{att}$ is an MLP)

Normalize to get
**attention weights**

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

Set context vector **c** to a linear
combination of hidden states

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

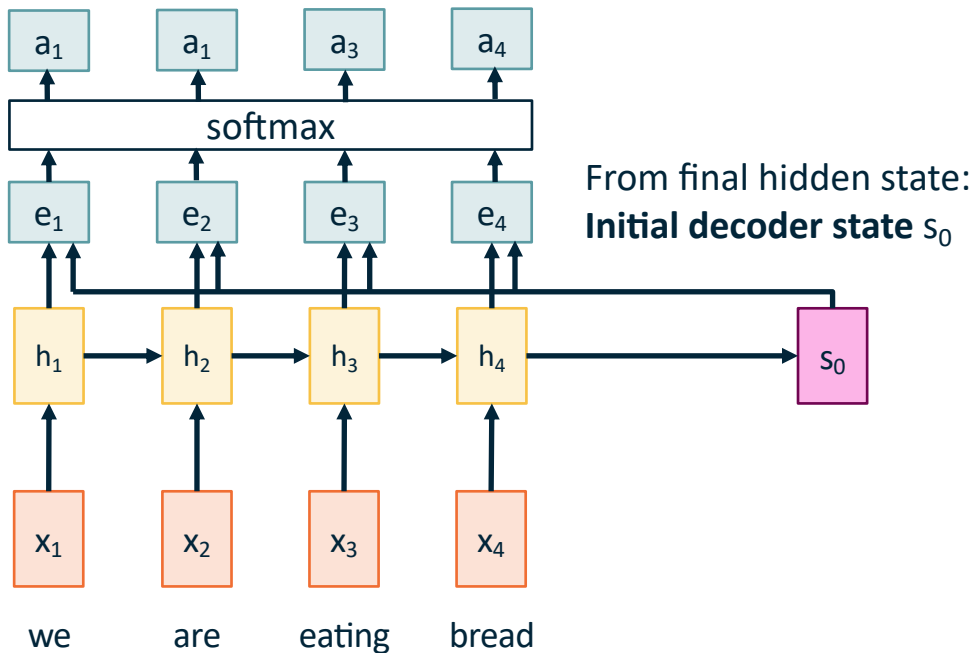Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

# Machine Translation with RNNs **and Attention**



Compute **affinity scores**
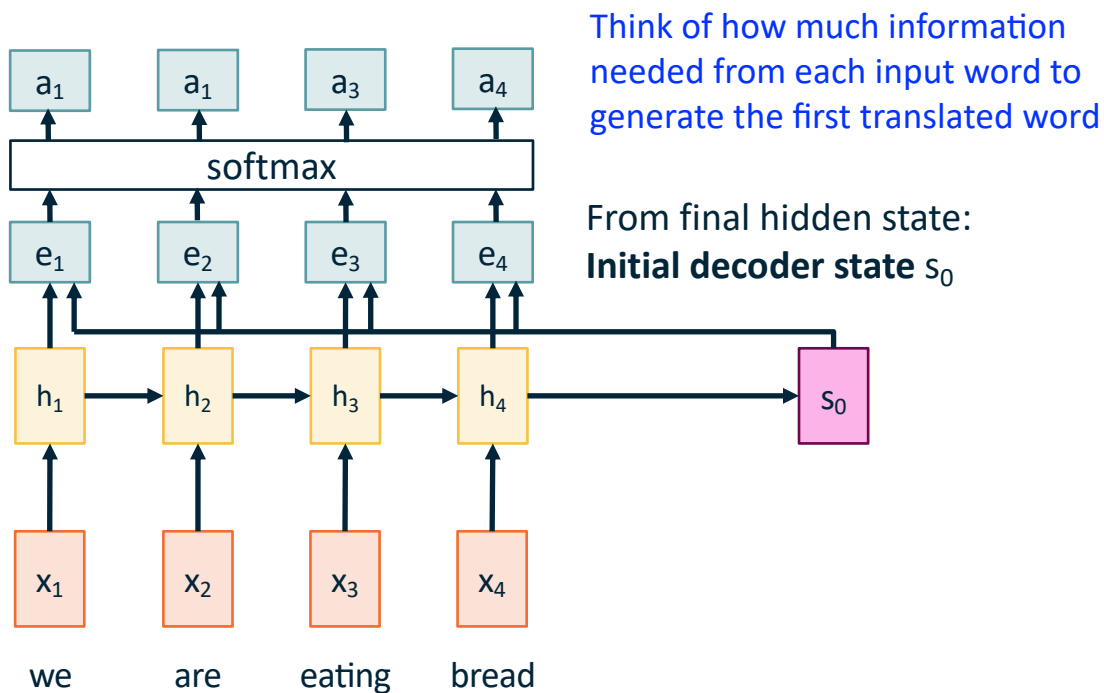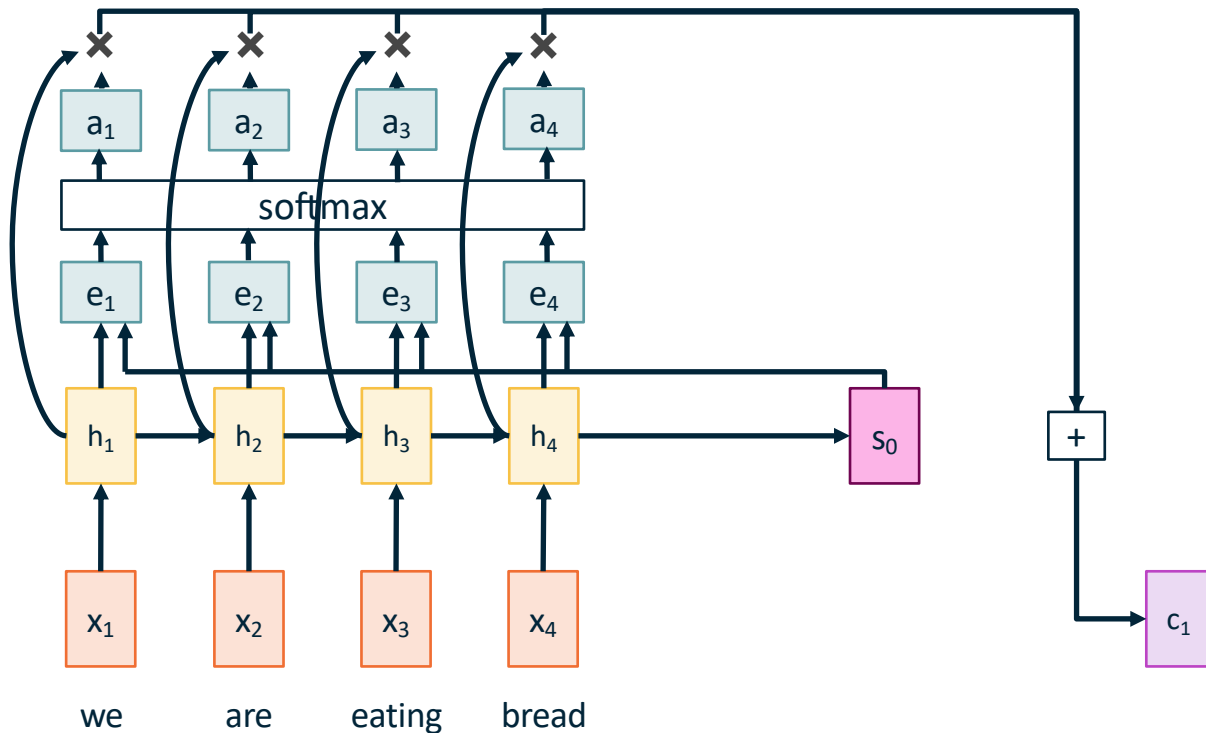
$$e_{t,i} = f_{att}(s_{t-1}, h_i) \qquad (f_{att} \text{ is an MLP})$$
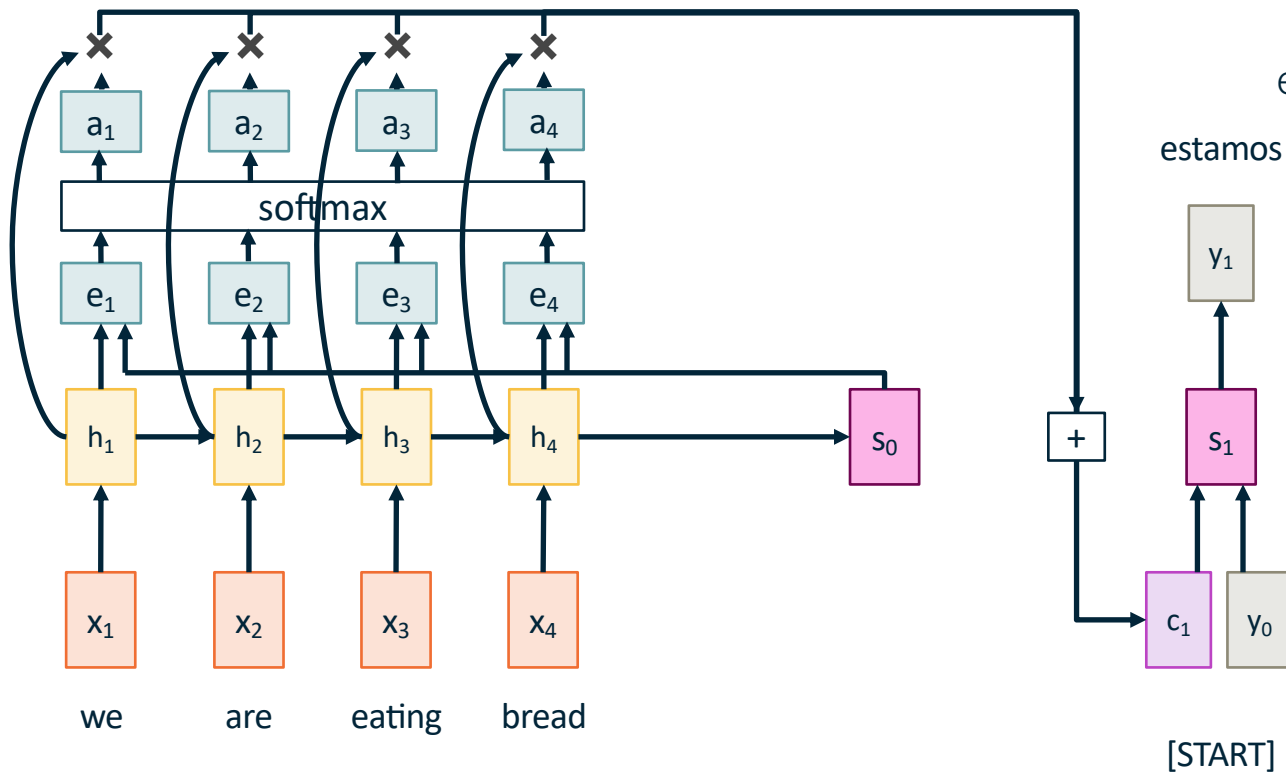
Normalize to get
**attention weights**

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

Set context vector **c** to a linear combination of hidden states

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

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

# Machine Translation with RNNs **and Attention**



Compute **affinity scores**

$e_{t,i} = f_{att}(s_{t-1}, h_i)$     ($f_{att}$ is an MLP)

Normalize to get
**attention weights**

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

Set context vector **c** to a linear
combination of hidden states

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

**Intuition**: Context vector
attends to the relevant part
of the input sequence
*"estamos" = "we are"*

**This is all differentiable! Do not
supervise attention weights –
backprop through everything**

# Machine Translation with RNNs **and Attention**



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

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Slide credit: Justin Johnson

# 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 $s_2$, $y_2$

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Slide credit: Justin Johnson

# 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.



Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Slide credit: Justin Johnson

# 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."

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Visualize attention weights $a_{t,i}$



Slide credit: Justin Johnson

# 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**."

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

Visualize attention weights $a_{t,i}$



Diagonal attention means words correspond in order

Diagonal attention means words correspond in order
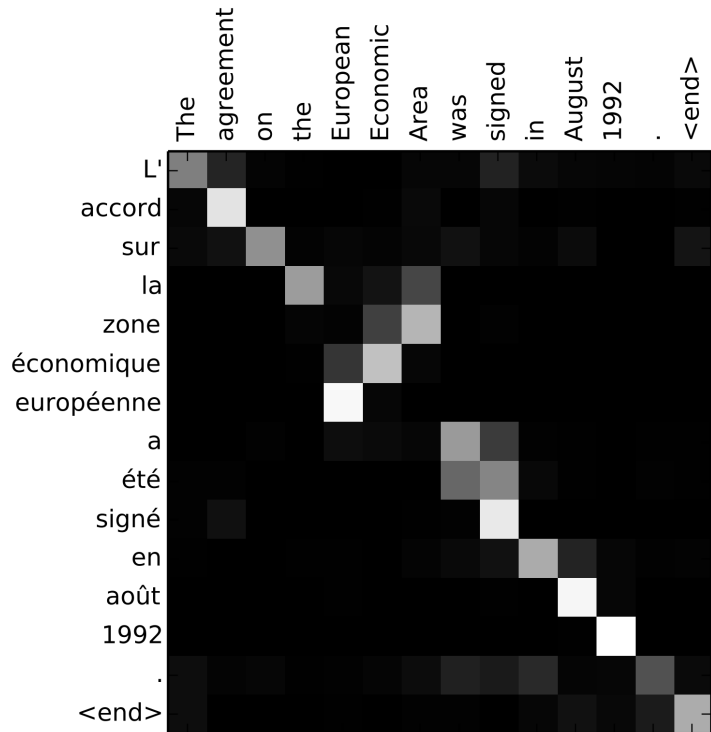
Slide credit: Justin Johnson

# Machine Translation with RNNs **and Attention**

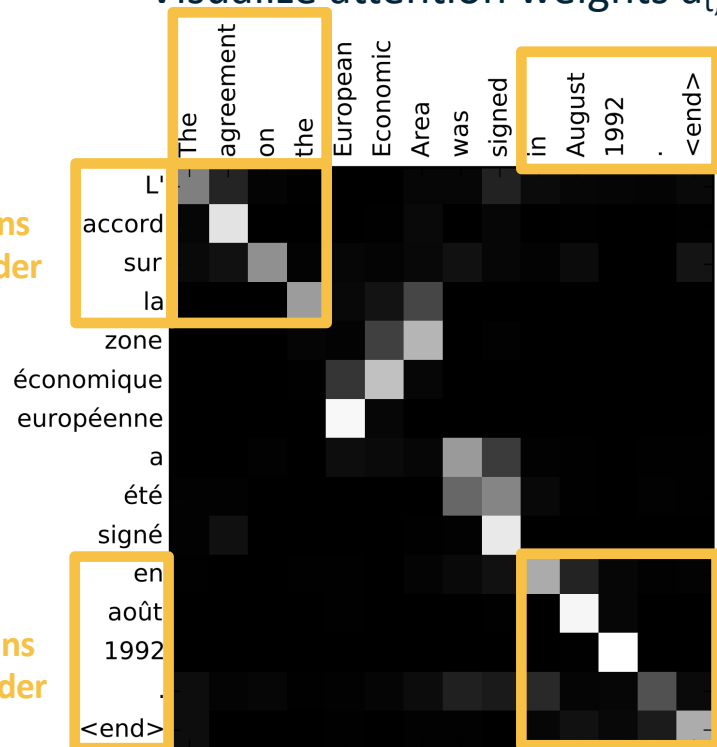**Example**: English to French translation

**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 $a_{t,i}$

Diagonal attention means words correspond in order

Attention figures out different word orders

Diagonal attention means words correspond in order

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

# Attention Layer

**Inputs**:
**State vector**: $s_i$ (Shape: $D_Q$)
**Hidden vectors**: $h_i$ (Shape: $N_X \times D_H$)
**Similarity function**: $f_{att}$



**Computation**:
**Similarities**: e (Shape: $N_X$)   $e_i = f_{att}(s_{t-1}, h_i)$
**Attention weights**: a = softmax(e)  (Shape: $N_X$)
**Output vector**: y = $\sum_i a_i h_i$   (Shape: $D_X$)

# Attention Layer

**Inputs**:
**Query vector**: $q$ (Shape: $D_Q$)
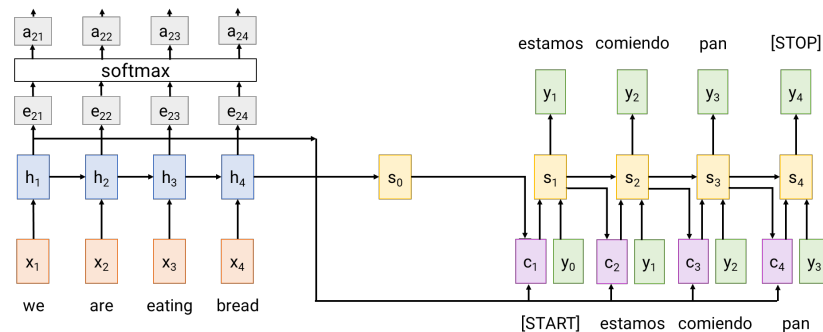**Input vectors**: $X$ (Shape: $N_X \times D_X$)
**Similarity function**: $f_{att}$



**Computation**:
**Similarities**: e (Shape: $N_X$)   $e_i = f_{att}(q, X_i)$
**Attention weights**: a = softmax(e)  (Shape: $N_X$)
**Output vector**: $y = \sum_i a_i X_i$    (Shape: $D_X$)

# Attention Layer

**Inputs**:
**Query vector**: $q$ (Shape: $D_Q$)
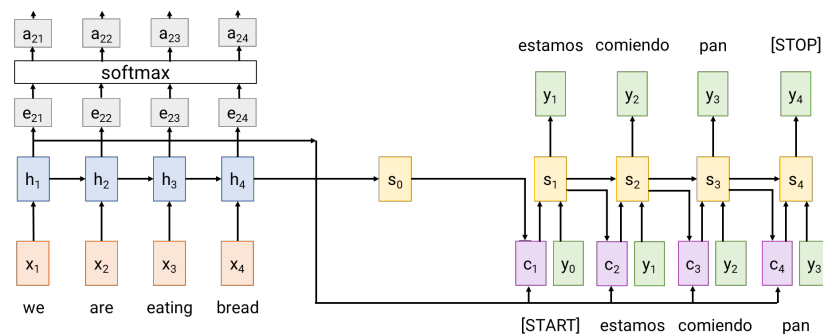**Input vectors**: $X$ (Shape: $N_X \times D_Q$)
**Similarity function**: dot product

**Computation**:
**Similarities**: e (Shape: $N_X$) $\quad e_i = q \cdot X_i$
**Attention weights**: a = softmax(e)  (Shape: $N_X$)
**Output vector**: y = $\sum_i a_i X_i$   (Shape: $D_X$)



Changes:
-   Use dot product for similarity

# Attention Layer

**Inputs**:
**Query vector**: $q$ (Shape: $D_Q$)
**Input vectors**: $X$ (Shape: $N_X$ x $D_Q$)
**Similarity function** scaled dot product



estamos  comiendo  pan  [STOP]

we    are    eating    bread

[START]  estamos  comiendo  pan

**Computation**:
**Similarities**: e (Shape: $N_X$)  $e_i = q \cdot X_i / \text{sqrt}(D_Q)$
**Attention weights**: a = softmax(e)  (Shape: $N_X$)
**Output vector**: y = $\sum_i a_i X_i$  (Shape: $D_X$)

Changes:
- Use **scaled** dot product for similarity

# Attention Layer

**Inputs**:
**Query vectors**: $Q$ (Shape: $N_Q \times D_Q$)
**Input vectors**: $X$ (Shape: $N_X \times D_Q$)

**Computation**:
**Similarities**: $E = QX^T / \sqrt{D_Q}$ (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 dot product for similarity
- Multiple **query** vectors

# Attention Layer

**Inputs**:
**Query vectors**: $Q$ (Shape: $N_Q$ x $D_Q$)
**Input vectors**: $X$ (Shape: $N_X$ x $D_Q$)

Attention matrix (A)
Each row sums up to 1

$N_X$

$N_Q$



**Computation**:
**Similarities**: $E = QX^T / sqrt(D_Q)$ (Shape: $N_Q$ x $N_X$)
**Attention matrix**: $A = softmax(E, dim=1)$ (Shape: $N_Q$ x $N_X$)
**Output vectors**: $Y = AX$ (Shape: $N_Q$ x $D_X$) $Y_i = \sum_j A_{i,j} X_j$

Changes:
- Use dot product for similarity
- Multiple **query** vectors

# Attention Layer

<u>Inputs</u>:
**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**: $E = \mathbf{QK}^T$ (Shape: $N_Q \times N_X$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j$ / sqrt($D_Q$)
**Attention weights**: $A$ = softmax($E$, dim=1) (Shape: $N_Q \times N_X$)
**Output vectors**: $Y = A\mathbf{V}$ (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 dot product for similarity
- Multiple **query** vectors
- Separate **key** and **value**

# Attention Layer

**Inputs**:
**Query vectors**: $Q$ (Shape: $N_Q \times D_Q$)
**Input vectors**: $X$ (Shape: $N_X \times D_X$)
**Key matrix**: $W_K$ (Shape: $D_X \times D_Q$)
**Value matrix**: $W_V$ (Shape: $D_X \times D_V$)

**Computation**:
**Key vectors**: $K = XW_K$ (Shape: $N_X \times D_Q$)
**Value vectors**: $V = XW_V$ (Shape: $N_X \times D_V$)
**Similarities**: $E = QK^T$ (Shape: $N_Q \times N_X$) $E_{i,j} = Q_i \cdot K_j / \sqrt{D_Q}$
**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_Q \times N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_Q \times D_V$) $Y_i = \sum_j A_{i,j} V_j$

$X_1$

$X_2$

$X_3$

$Q_1$  $Q_2$  $Q_3$  $Q_4$

# Attention Layer

**Inputs**:
**Query vectors**: $Q$ (Shape: $N_Q \times D_Q$)
**Input vectors**: $X$ (Shape: $N_X \times D_X$)
**Key matrix**: $W_K$ (Shape: $D_X \times D_Q$)
**Value matrix**: $W_V$ (Shape: $D_X \times D_V$)

**Computation**:
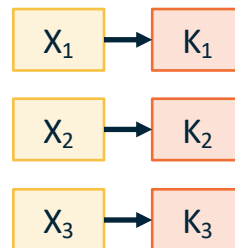**Key vectors**: $K = XW_K$ (Shape: $N_X \times D_Q$)
**Value vectors**: $V = XW_V$ (Shape: $N_X \times D_V$)
**Similarities**: $E = QK^T$ (Shape: $N_Q \times N_X$) $E_{i,j} = Q_i \cdot K_j / \text{sqrt}(D_Q)$
**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_Q \times N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_Q \times D_V$) $Y_i = \sum_j A_{i,j} V_j$

$X_1 \rightarrow K_1$

$X_2 \rightarrow K_2$

$X_3 \rightarrow K_3$

$Q_1 \quad Q_2 \quad Q_3 \quad Q_4$

# Attention Layer

**Inputs**:
**Query vectors**: $Q$ (Shape: $N_Q$ x $D_Q$)
**Input vectors**: $X$ (Shape: $N_X$ x $D_X$)
**Key matrix**: $W_K$ (Shape: $D_X$ x $D_Q$)
**Value matrix**: $W_V$ (Shape: $D_X$ x $D_V$)

**Computation**:
**Key vectors**: $K = XW_K$ (Shape: $N_X$ x $D_Q$)
**Value vectors**: $V = XW_V$ (Shape: $N_X$ x $D_V$)
**Similarities**: $E = QK^T$ (Shape: $N_Q$ x $N_X$) $E_{i,j} = Q_i \cdot K_j / sqrt(D_Q)$
**Attention weights**: $A = softmax(E, dim=1)$ (Shape: $N_Q$ x $N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_Q$ x $D_V$) $Y_i = \sum_j A_{i,j} V_j$

# 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: $D_X \times D_Q$)
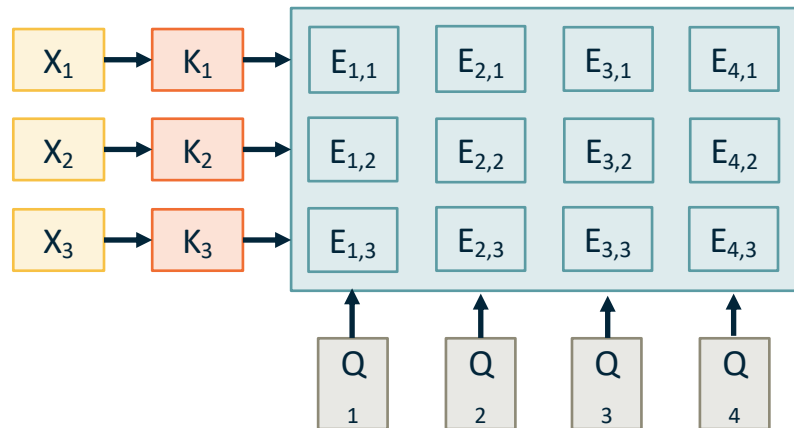**Value matrix**: $\mathbf{W_V}$ (Shape: $D_X \times D_V$)

**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**: $E = \mathbf{QK^T}$ (Shape: $N_Q \times N_X$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / sqrt(D_Q)$
**Attention weights**: $A = softmax(E, dim=1)$  (Shape: $N_Q \times N_X$)
**Output vectors**: $Y = A\mathbf{V}$ (Shape: $N_Q \times D_V$) $Y_i = \sum_j A_{i,j}\mathbf{V}_j$



Softmax( ↑ )

# Attention Layer

**Inputs**:
**Query vectors**: $Q$ (Shape: $N_Q \times D_Q$)
**Input vectors**: $X$ (Shape: $N_X \times D_X$)
**Key matrix**: $W_K$ (Shape: $D_X \times D_Q$)
**Value matrix**: $W_V$ (Shape: $D_X \times D_V$)

**Computation**:
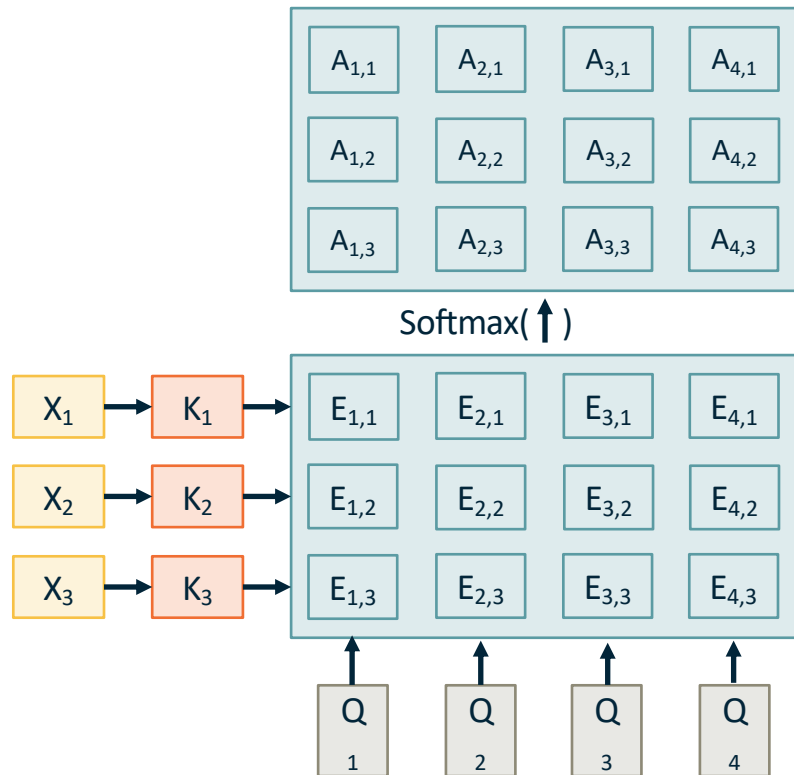**Key vectors**: $K = XW_K$ (Shape: $N_X \times D_Q$)
**Value vectors**: $V = XW_V$ (Shape: $N_X \times D_V$)
**Similarities**: $E = QK^T$ (Shape: $N_Q \times N_X$) $E_{i,j} = Q_i \cdot K_j / \sqrt{D_Q}$
**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_Q \times N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_Q \times D_V$) $Y_i = \sum_j A_{i,j} V_j$

| $V_1$ | $A_{1,1}$ | $A_{2,1}$ | $A_{3,1}$ | $A_{4,1}$ |
| $V_2$ | $A_{1,2}$ | $A_{2,2}$ | $A_{3,2}$ | $A_{4,2}$ |
| $V_3$ | $A_{1,3}$ | $A_{2,3}$ | $A_{3,3}$ | $A_{4,3}$ |

Softmax( ↑ )

| $X_1$ | $K_1$ | $E_{1,1}$ | $E_{2,1}$ | $E_{3,1}$ | $E_{4,1}$ |
| $X_2$ | $K_2$ | $E_{1,2}$ | $E_{2,2}$ | $E_{3,2}$ | $E_{4,2}$ |
| $X_3$ | $K_3$ | $E_{1,3}$ | $E_{2,3}$ | $E_{3,3}$ | $E_{4,3}$ |

| $Q$ 1 | $Q$ 2 | $Q$ 3 | $Q$ 4 |

# Attention Layer

**Inputs**:
**Query vectors**: $Q$ (Shape: $N_Q \times D_Q$)
**Input vectors**: $X$ (Shape: $N_X \times D_X$)
**Key matrix**: $W_K$ (Shape: $D_X \times D_Q$)
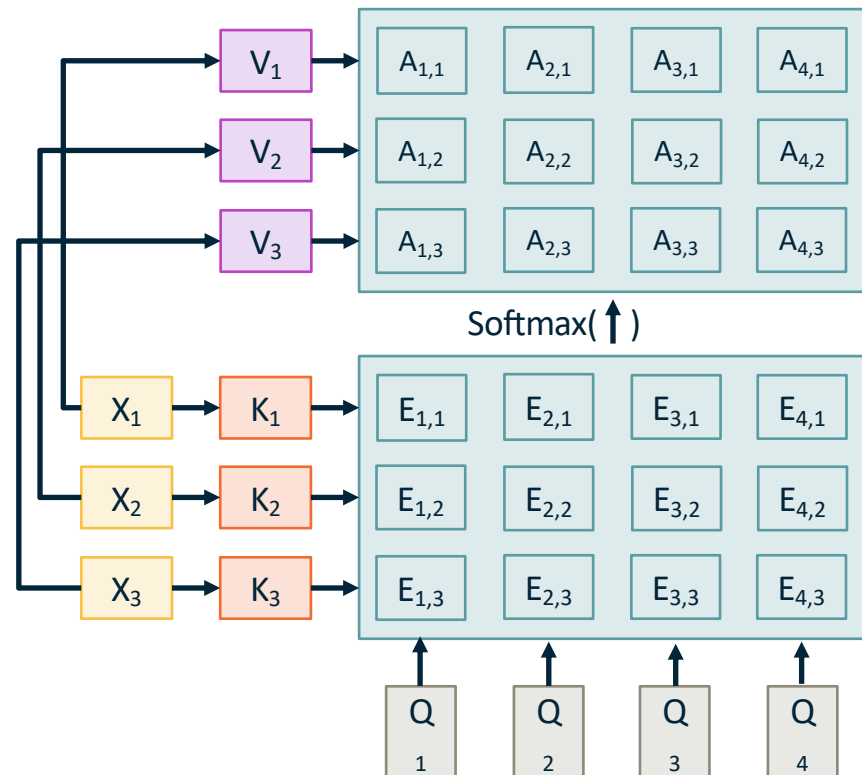**Value matrix**: $W_V$ (Shape: $D_X \times D_V$)

**Computation**:
**Key vectors**: $K = XW_K$  (Shape: $N_X \times D_Q$)
**Value vectors**: $V = XW_V$ (Shape: $N_X \times D_V$)
**Similarities**: $E = QK^T$ (Shape: $N_Q \times N_X$) $E_{i,j} = Q_i \cdot K_j / \sqrt{D_Q}$
**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$  (Shape: $N_Q \times N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_Q \times D_V$) $Y_i = \sum_j A_{i,j} V_j$

$Y_1$  $Y_2$  $Y_3$  $Y_4$

Product($\rightarrow$),  Sum($\uparrow$)

| $V_1$ | $A_{1,1}$ | $A_{2,1}$ | $A_{3,1}$ | $A_{4,1}$ |
| $V_2$ | $A_{1,2}$ | $A_{2,2}$ | $A_{3,2}$ | $A_{4,2}$ |
| $V_3$ | $A_{1,3}$ | $A_{2,3}$ | $A_{3,3}$ | $A_{4,3}$ |

Softmax($\uparrow$)

| $X_1$ | $K_1$ | $E_{1,1}$ | $E_{2,1}$ | $E_{3,1}$ | $E_{4,1}$ |
| $X_2$ | $K_2$ | $E_{1,2}$ | $E_{2,2}$ | $E_{3,2}$ | $E_{4,2}$ |
| $X_3$ | $K_3$ | $E_{1,3}$ | $E_{2,3}$ | $E_{3,3}$ | $E_{4,3}$ |

$Q_1$  $Q_2$  $Q_3$  $Q_4$

# 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: $D_X \times D_Q$)
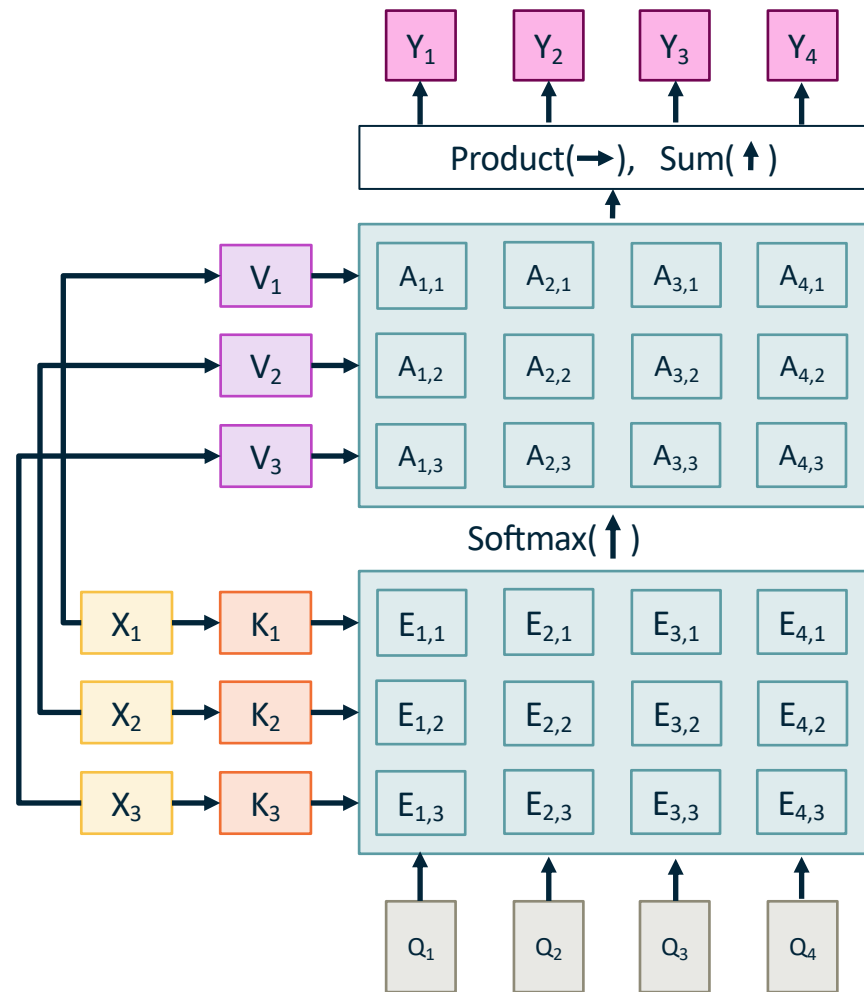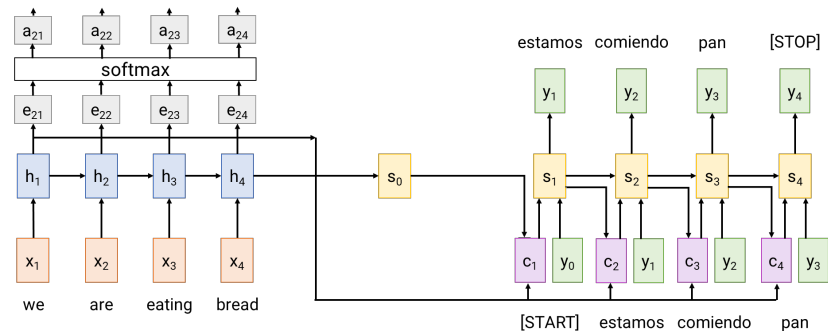**Value matrix**: $\mathbf{W_V}$ (Shape: $D_X \times D_V$)

**Computation**:
**Key vectors**: $\mathbf{K} = \mathbf{X}\mathbf{W_K}$ (Shape: $N_X \times D_Q$)
**Value vectors**: $\mathbf{V} = \mathbf{X}\mathbf{W_V}$ (Shape: $N_X \times D_V$)
**Similarities**: $E = \mathbf{Q}\mathbf{K}^T$ (Shape: $N_Q \times N_X$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \sqrt{D_Q}$
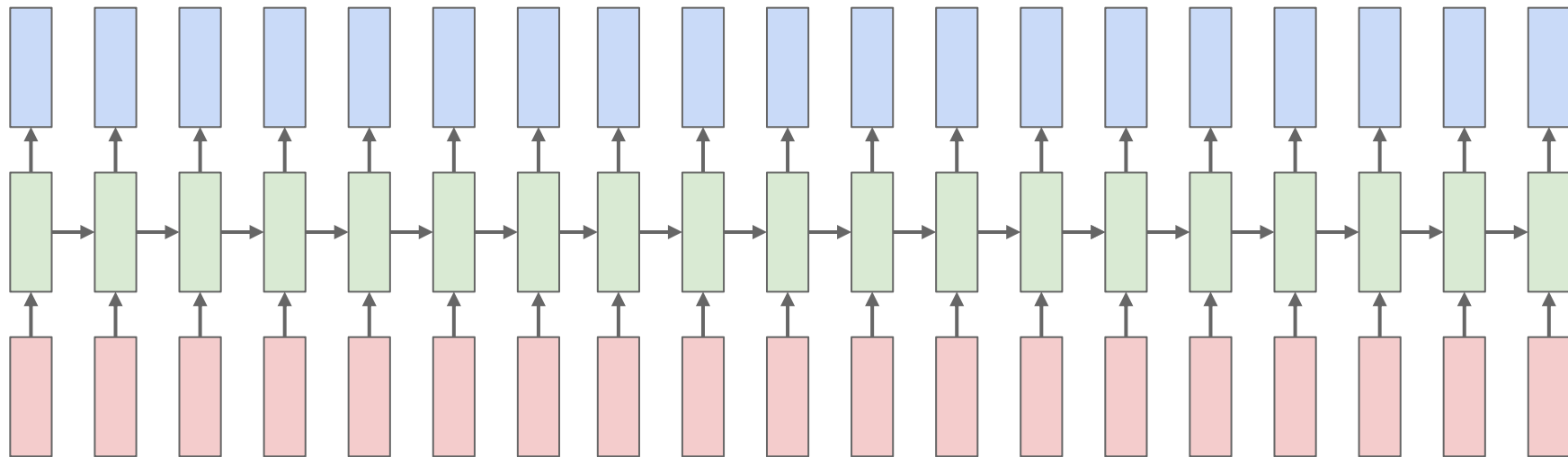**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_Q \times N_X$)
**Output vectors**: $Y = A\mathbf{V}$ (Shape: $N_Q \times D_V$) $Y_i = \sum_j A_{i,j}\mathbf{V}_j$



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**: $Q$ (Shape: $N_Q$ x $D_Q$)
**Input vectors**: $X$ (Shape: $N_X$ x $D_X$)
**Key matrix**: $W_K$ (Shape: $D_X$ x $D_Q$)
**Value matrix**: $W_V$ (Shape: $D_X$ x $D_V$)

**Computation**:
**Key vectors**: $K = XW_K$  (Shape: $N_X$ x $D_Q$)
**Value vectors**: $V = XW_V$ (Shape: $N_X$ x $D_V$)
**Similarities**: $E = QK^T$ (Shape: $N_Q$ x $N_X$) $E_{i,j} = Q_i \cdot K_j / \sqrt{D_Q}$
**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$  (Shape: $N_Q$ x $N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_Q$ x $D_V$) $Y_i = \sum_j A_{i,j} V_j$

Attention seems to be really powerful ...
Do we still need RNN?

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

# Self-Attention Layer

Sequence encode -> use each input element as query!

**Inputs**:
**Input vectors**: **X** (Shape: $N_X \times D_X$)
**Key matrix**: **W$_K$** (Shape: $D_X \times D_Q$)
**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$_Q$**
**Key vectors**: **K** = **XW$_K$** (Shape: $N_X \times D_Q$)
**Value vectors**: **V** = **XW$_V$** (Shape: $N_X \times D_V$)
**Similarities**: E = **QK$^T$** (Shape: $N_X \times N_X$) $E_{i,j}$ = **Q$_i$** · **K$_j$** / sqrt($D_Q$)
**Attention weights**: A = softmax(E, dim=1) (Shape: $N_X \times N_X$)
**Output vectors**: Y = A**V** (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j}$**V$_j$**

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!

<u>Inputs</u>:
**Input vectors**: **X** (Shape: $N_X$ x $D_X$)
**Key matrix**: **W**$_K$ (Shape: $D_X$ x $D_Q$)
**Value matrix**: **W**$_V$ (Shape: $D_X$ x $D_V$)
**Query matrix**: **W**$_Q$ (Shape: $D_X$ x $D_Q$)

<u>Computation</u>:
**Query vectors**: **Q** = **XW**$_Q$
**Key vectors**: **K** = **XW**$_K$  (Shape: $N_X$ x $D_Q$)
**Value vectors**: **V** = **XW**$_V$ (Shape: $N_X$ x $D_V$)
**Similarities**: E = **QK**$^T$ (Shape: $N_X$ x $N_X$) $E_{i,j}$ = **Q**$_i$ · **K**$_j$ / sqrt($D_Q$)
**Attention weights**: A = softmax(E, dim=1)  (Shape: $N_X$ x $N_X$)
**Output vectors**: Y = A**V** (Shape: $N_X$ x $D_V$) $Y_i$ = $\sum_j A_{i,j}$**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

$X_1$  $X_2$  $X_3$

# Self-Attention Layer

Sequence encode -> use each input element as query!

**Inputs**:
**Input vectors**: $X$ (Shape: $N_X \times D_X$)
**Key matrix**: $W_K$ (Shape: $D_X \times D_Q$)
**Value matrix**: $W_V$ (Shape: $D_X \times D_V$)
**Query matrix**: $W_Q$ (Shape: $D_X \times D_Q$)

**Computation**:
**Query vectors**: $Q = XW_Q$
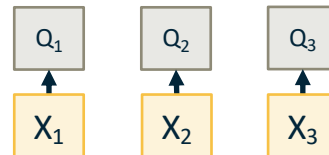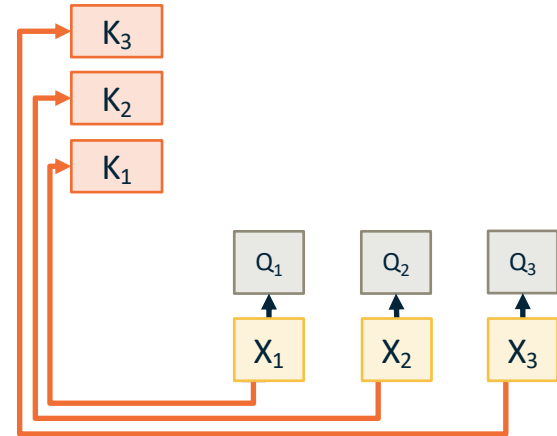**Key vectors**: $K = XW_K$ (Shape: $N_X \times D_Q$)
**Value vectors**: $V = XW_V$ (Shape: $N_X \times D_V$)
**Similarities**: $E = QK^T$ (Shape: $N_X \times N_X$) $E_{i,j} = Q_i \cdot K_j / sqrt(D_Q)$
**Attention weights**: $A = softmax(E, dim=1)$ (Shape: $N_X \times N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} V_j$

# Self-Attention Layer

Sequence encode -> use each input element as query!

**Inputs**:
**Input vectors**: $X$ (Shape: $N_X$ x $D_X$)
**Key matrix**: $W_K$ (Shape: $D_X$ x $D_Q$)
**Value matrix**: $W_V$ (Shape: $D_X$ x $D_V$)
**Query matrix**: $W_Q$ (Shape: $D_X$ x $D_Q$)

**Computation**:
**Query vectors**: $Q = XW_Q$
**Key vectors**: $K = XW_K$  (Shape: $N_X$ x $D_Q$)
**Value vectors**: $V = XW_V$ (Shape: $N_X$ x $D_V$)
**Similarities**: $E = QK^T$ (Shape: $N_X$ x $N_X$) $E_{i,j} = Q_i \cdot K_j / \sqrt{D_Q}$
**Attention weights**: $A = \mathrm{softmax}(E, \dim=1)$  (Shape: $N_X$ x $N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_X$ x $D_V$) $Y_i = \sum_j A_{i,j} V_j$

# Self-Attention Layer

Sequence encode -> use each input element as query!

**Inputs**:
**Input vectors**: **X** (Shape: $N_X$ x $D_X$)
**Key matrix**: **$W_K$** (Shape: $D_X$ x $D_Q$)
**Value matrix**: **$W_V$** (Shape: $D_X$ x $D_V$)
**Query matrix**: **$W_Q$** (Shape: $D_X$ x $D_Q$)

**Computation**:
**Query vectors**: **Q** = **X$W_Q$**
**Key vectors**: **K** = **X$W_K$** (Shape: $N_X$ x $D_Q$)
**Value vectors**: **V** = **X$W_V$** (Shape: $N_X$ x $D_V$)
**Similarities**: E = **QK$^T$** (Shape: $N_X$ x $N_X$) $E_{i,j}$ = **$Q_i$** · **$K_j$** / sqrt($D_Q$)
**Attention weights**: A = softmax(E, dim=1) (Shape: $N_X$ x $N_X$)
**Output vectors**: Y = A**V** (Shape: $N_X$ x $D_V$) $Y_i$ = $\sum_j A_{i,j}$**$V_j$**

# Self-Attention Layer

Sequence encode -> use each input element as query!

**Inputs**:
**Input vectors**: $X$ (Shape: $N_X \times D_X$)
**Key matrix**: $W_K$ (Shape: $D_X \times D_Q$)
**Value matrix**: $W_V$ (Shape: $D_X \times D_V$)
**Query matrix**: $W_Q$ (Shape: $D_X \times D_Q$)
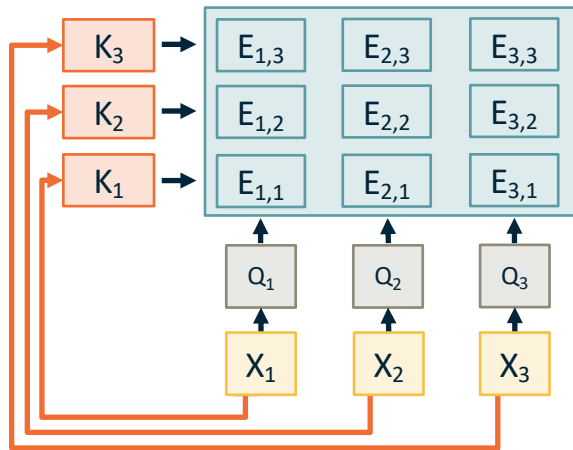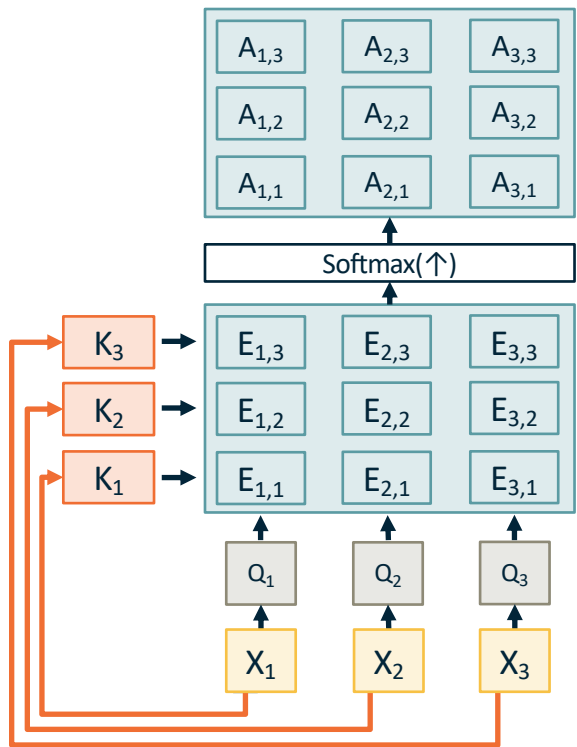
**Computation**:
**Query vectors**: $Q = XW_Q$
**Key vectors**: $K = XW_K$ (Shape: $N_X \times D_Q$)
**Value vectors**: $V = XW_V$ (Shape: $N_X \times D_V$)
**Similarities**: $E = QK^T$ (Shape: $N_X \times N_X$) $E_{i,j} = Q_i \cdot K_j / \sqrt{D_Q}$
**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_X \times N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} V_j$

# Self-Attention Layer

Sequence encode -> use each input element as query!

**Inputs**:
**Input vectors**: $X$ (Shape: $N_X \times D_X$)
**Key matrix**: $W_K$ (Shape: $D_X \times D_Q$)
**Value matrix**: $W_V$ (Shape: $D_X \times D_V$)
**Query matrix**: $W_Q$ (Shape: $D_X \times D_Q$)
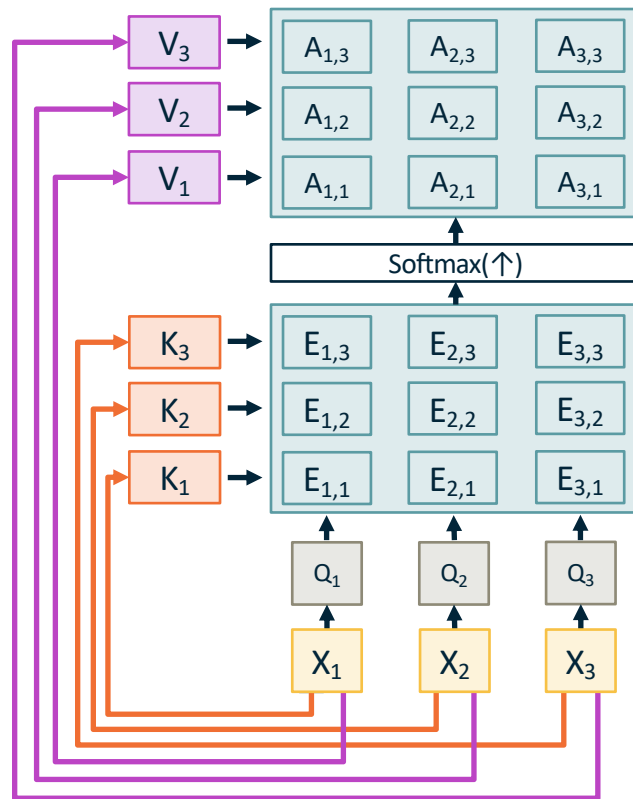
**Computation**:
**Query vectors**: $Q = XW_Q$
**Key vectors**: $K = XW_K$ (Shape: $N_X \times D_Q$)
**Value vectors**: $V = XW_V$ (Shape: $N_X \times D_V$)
**Similarities**: $E = QK^T$ (Shape: $N_X \times N_X$) $E_{i,j} = Q_i \cdot K_j / sqrt(D_Q)$
**Attention weights**: $A = softmax(E, dim=1)$ (Shape: $N_X \times N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} V_j$



Slide credit: Justin Johnson

# Self-Attention Layer

Sequence encode -> use each input element as query!

**Inputs**:
**Input vectors**: $X$ (Shape: $N_X \times D_X$)
**Key matrix**: $W_K$ (Shape: $D_X \times D_Q$)
**Value matrix**: $W_V$ (Shape: $D_X \times D_V$)
**Query matrix**: $W_Q$ (Shape: $D_X \times D_Q$)
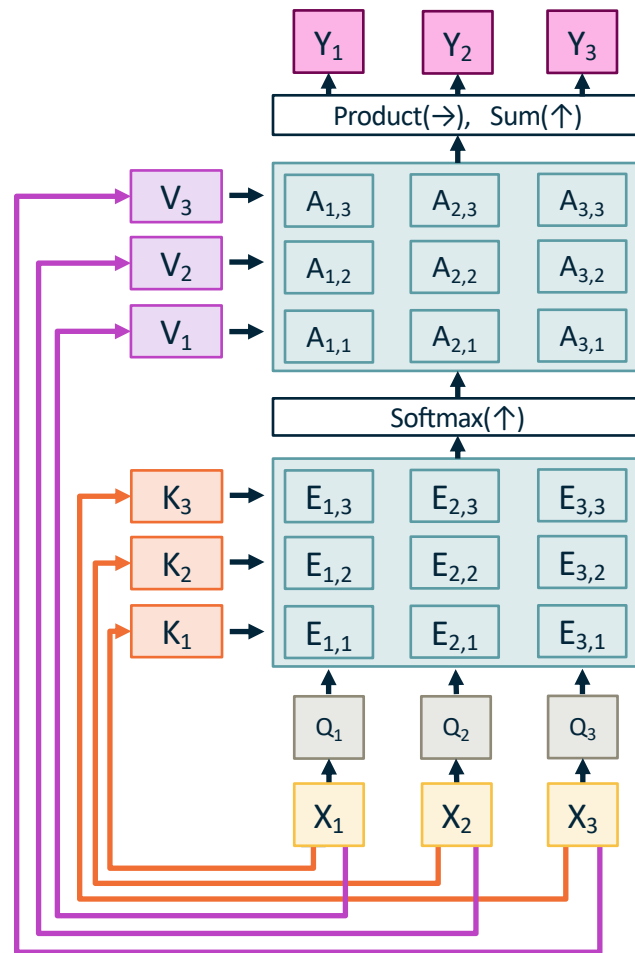
**Computation**:
**Query vectors**: $Q = XW_Q$
**Key vectors**: $K = XW_K$  (Shape: $N_X \times D_Q$)
**Value vectors**: $V = XW_V$ (Shape: $N_X \times D_V$)
**Similarities**: $E = QK^T$ (Shape: $N_X \times N_X$) $E_{i,j} = Q_i \cdot K_j / sqrt(D_Q)$
**Attention weights**: $A = softmax(E, dim=1)$  (Shape: $N_X \times N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} V_j$



Slide credit: Justin Johnson

# Self-Attention Layer

Sequence encode -> use each input element as query!

**Inputs**:
**Input vectors**: $X$ (Shape: $N_X \times D_X$)
**Key matrix**: $W_K$ (Shape: $D_X \times D_Q$)
**Value matrix**: $W_V$ (Shape: $D_X \times D_V$)
**Query matrix**: $W_Q$ (Shape: $D_X \times D_Q$)

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

**Computation**:
**Query vectors**: $Q = XW_Q$
**Key vectors**: $K = XW_K$ (Shape: $N_X \times D_Q$)
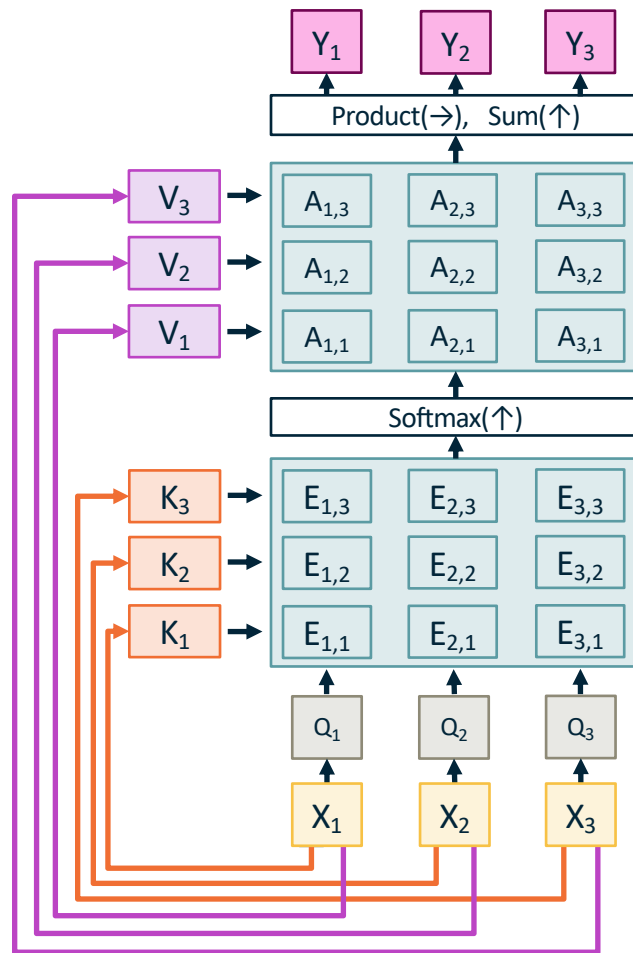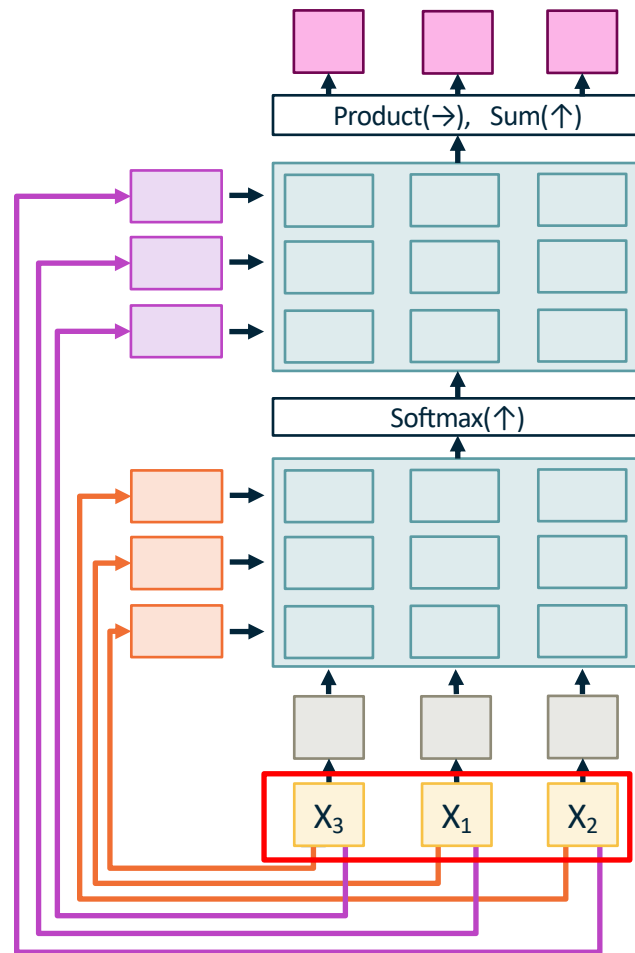**Value vectors**: $V = XW_V$ (Shape: $N_X \times D_V$)
**Similarities**: $E = QK^T$ (Shape: $N_X \times N_X$) $E_{i,j} = Q_i \cdot K_j / \sqrt{D_Q}$
**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_X \times N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} V_j$



Slide credit: Justin Johnson

# Self-Attention Layer

**Inputs**:
**Input vectors**: **X** (Shape: $N_X$ x $D_X$)
**Key matrix**: **$W_K$** (Shape: $D_X$ x $D_Q$)
**Value matrix**: **$W_V$** (Shape: $D_X$ x $D_V$)
**Query matrix**: **$W_Q$** (Shape: $D_X$ x $D_Q$)

**Computation**:
**Query vectors**: **Q** = **X$W_Q$**
**Key vectors**: **K** = **X$W_K$** (Shape: $N_X$ x $D_Q$)
**Value vectors**: **V** = **X$W_V$** (Shape: $N_X$ x $D_V$)
**Similarities**: E = **QK$^T$** (Shape: $N_X$ x $N_X$) $E_{i,j}$ = **$Q_i$** · **$K_j$** / sqrt($D_Q$)
**Attention weights**: A = softmax(E, dim=1)  (Shape: $N_X$ x $N_X$)
**Output vectors**: Y = A**V** (Shape: $N_X$ x $D_V$) $Y_i$ = $\sum_j A_{i,j}$**$V_j$**

Consider **permuting** the input vectors:



Product($\rightarrow$),  Sum($\uparrow$)

Softmax($\uparrow$)

$X_3$   $X_1$   $X_2$

# Self-Attention Layer

**Inputs**:
**Input vectors**: $X$ (Shape: $N_X$ x $D_X$)
**Key matrix**: $W_K$ (Shape: $D_X$ x $D_Q$)
**Value matrix**: $W_V$ (Shape: $D_X$ x $D_V$)
**Query matrix**: $W_Q$ (Shape: $D_X$ x $D_Q$)

**Computation**:
**Query vectors**: $Q = XW_Q$
**Key vectors**: $K = XW_K$ (Shape: $N_X$ x $D_Q$)
**Value Vectors**: $V = XW_V$ (Shape: $N_X$ x $D_V$)
**Similarities**: $E = QK^T$ (Shape: $N_X$ x $N_X$) $E_{i,j} = Q_i \cdot K_j$ / sqrt($D_Q$)
**Attention weights**: $A = $ softmax($E$, dim=1) (Shape: $N_X$ x $N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_X$ x $D_V$) $Y_i = \sum_j A_{i,j} V_j$

Consider **permuting** the input vectors:

Queries and Keys will be the same, but permuted

Product($\rightarrow$), Sum($\uparrow$)

Softmax($\uparrow$)

$K_2$
$K_1$
$K_3$

$Q_3$ $Q_1$ $Q_2$

$X_3$ $X_1$ $X_2$

# Self-Attention Layer

**Inputs**:
**Input vectors**: $X$ (Shape: $N_X \times D_X$)
**Key matrix**: $W_K$ (Shape: $D_X \times D_Q$)
**Value matrix**: $W_V$ (Shape: $D_X \times D_V$)
**Query matrix**: $W_Q$ (Shape: $D_X \times D_Q$)

**Computation**:
**Query vectors**: $Q = XW_Q$
**Key vectors**: $K = XW_K$ (Shape: $N_X \times D_Q$)
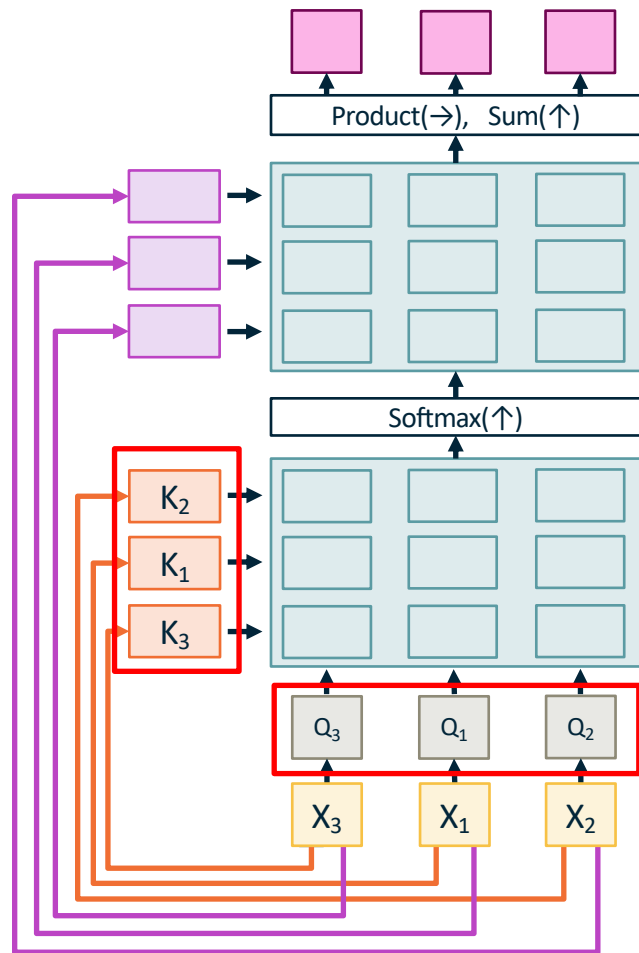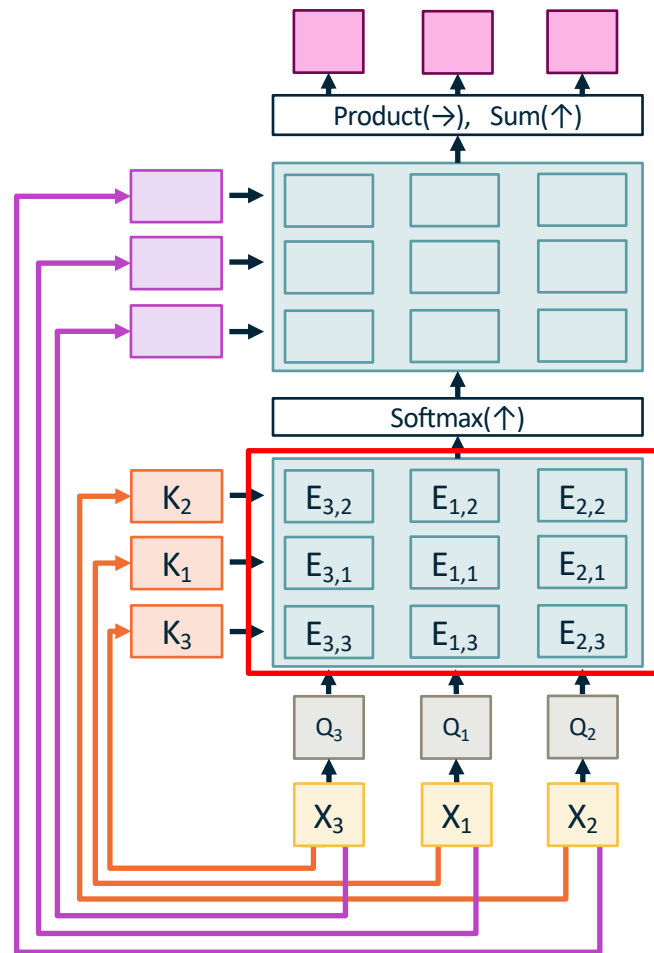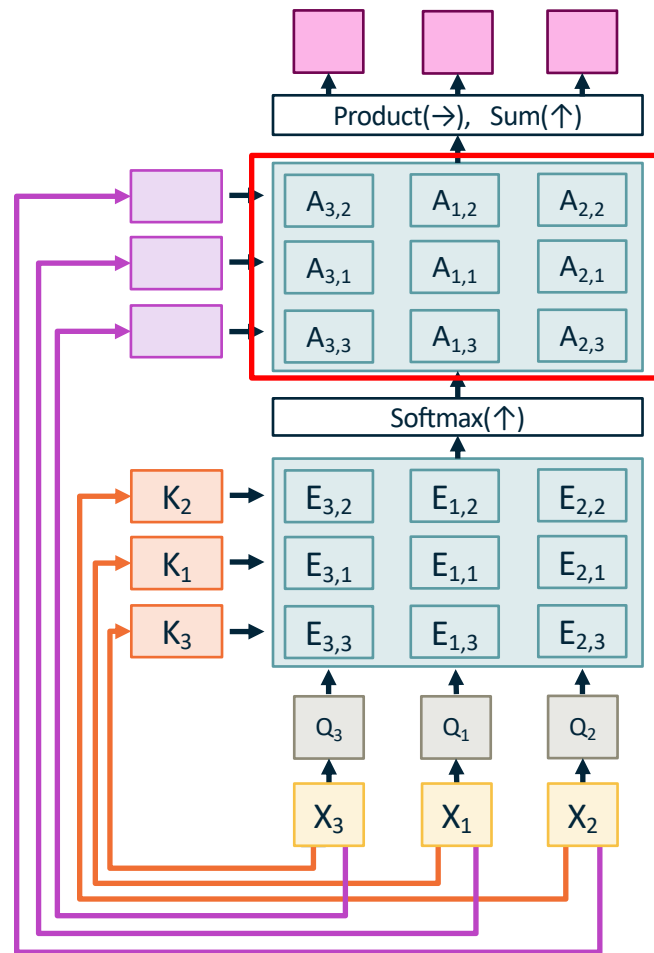**Value vectors**: $V = XW_V$ (Shape: $N_X \times D_V$)
**Similarities**: $E = QK^T$ (Shape: $N_X \times N_X$) $E_{i,j} = Q_i \cdot K_j / \text{sqrt}(D_Q)$
**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_X \times N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} V_j$

Consider **permuting** the input vectors:

Similarities will be the same, but permuted



Product($\rightarrow$), Sum($\uparrow$)

Softmax($\uparrow$)

| $E_{3,2}$ | $E_{1,2}$ | $E_{2,2}$ |
| $E_{3,1}$ | $E_{1,1}$ | $E_{2,1}$ |
| $E_{3,3}$ | $E_{1,3}$ | $E_{2,3}$ |

$K_2$  $K_1$  $K_3$

$Q_3$  $Q_1$  $Q_2$

$X_3$  $X_1$  $X_2$

# Self-Attention Layer

**Inputs**:

**Input vectors**: $X$ (Shape: $N_X$ x $D_X$)

**Key matrix**: $W_K$ (Shape: $D_X$ x $D_Q$)

**Value matrix**: $W_V$ (Shape: $D_X$ x $D_V$)

**Query matrix**: $W_Q$ (Shape: $D_X$ x $D_Q$)

**Computation**:

**Query vectors**: $Q = XW_Q$

**Key vectors**: $K = XW_K$ (Shape: $N_X$ x $D_Q$)
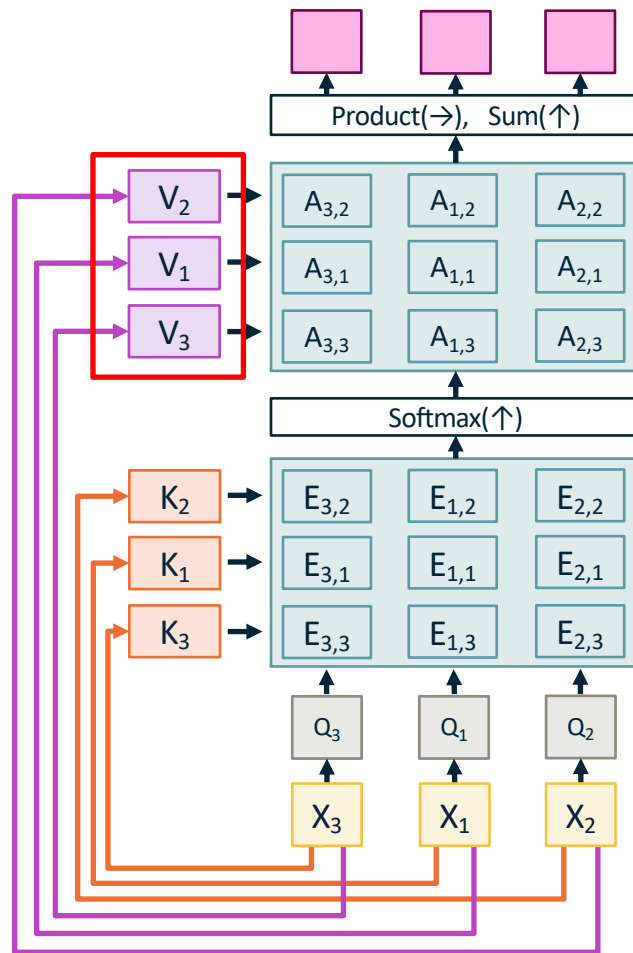
**Value vectors**: $V = XW_V$ (Shape: $N_X$ x $D_V$)

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

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

**Output vectors**: $Y = AV$ (Shape: $N_X$ x $D_V$) $Y_i = \sum_j A_{i,j} V_j$

Consider **permuting** the input vectors:

Attention weights will be the same, but permuted

# Self-Attention Layer

**Inputs**:
**Input vectors**: $X$ (Shape: $N_X$ x $D_X$)
**Key matrix**: $W_K$ (Shape: $D_X$ x $D_Q$)
**Value matrix**: $W_V$ (Shape: $D_X$ x $D_V$)
**Query matrix**: $W_Q$ (Shape: $D_X$ x $D_Q$)

**Computation**:
**Query vectors**: $Q = XW_Q$
**Key vectors**: $K = XW_K$  (Shape: $N_X$ x $D_Q$)
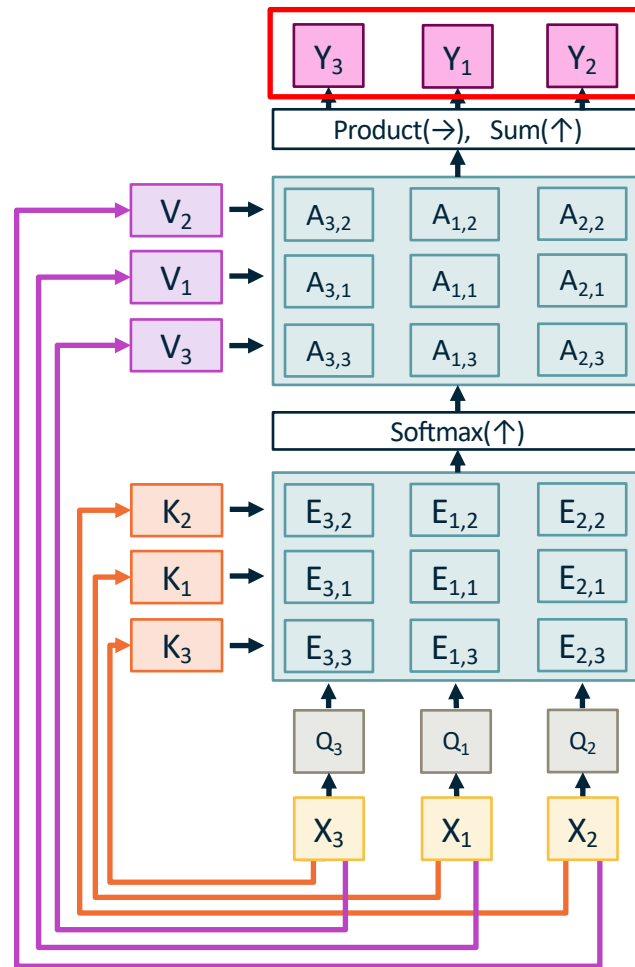**Value vectors**: $V = XW_V$ (Shape: $N_X$ x $D_V$)
**Similarities**: $E = QK^T$ (Shape: $N_X$ x $N_X$) $E_{i,j} = Q_i \cdot K_j / \sqrt{D_Q}$
**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$  (Shape: $N_X$ x $N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_X$ x $D_V$) $Y_i = \sum_j A_{i,j} V_j$

Consider **permuting** the input vectors:

Values will be the same, but permuted

# Self-Attention Layer

**Inputs**:
**Input vectors**: $X$ (Shape: $N_X \times D_X$)
**Key matrix**: $W_K$ (Shape: $D_X \times D_Q$)
**Value matrix**: $W_V$ (Shape: $D_X \times D_V$)
**Query matrix**: $W_Q$ (Shape: $D_X \times D_Q$)

Consider **permuting** the input vectors:

Outputs will be the same, but permuted

**Computation**:
**Query vectors**: $Q = XW_Q$
**Key vectors**: $K = XW_K$ (Shape: $N_X \times D_Q$)
**Value vectors**: $V = XW_V$ (Shape: $N_X \times D_V$)
**Similarities**: $E = QK^T$ (Shape: $N_X \times N_X$) $E_{i,j} = Q_i \cdot K_j / \sqrt{D_Q}$
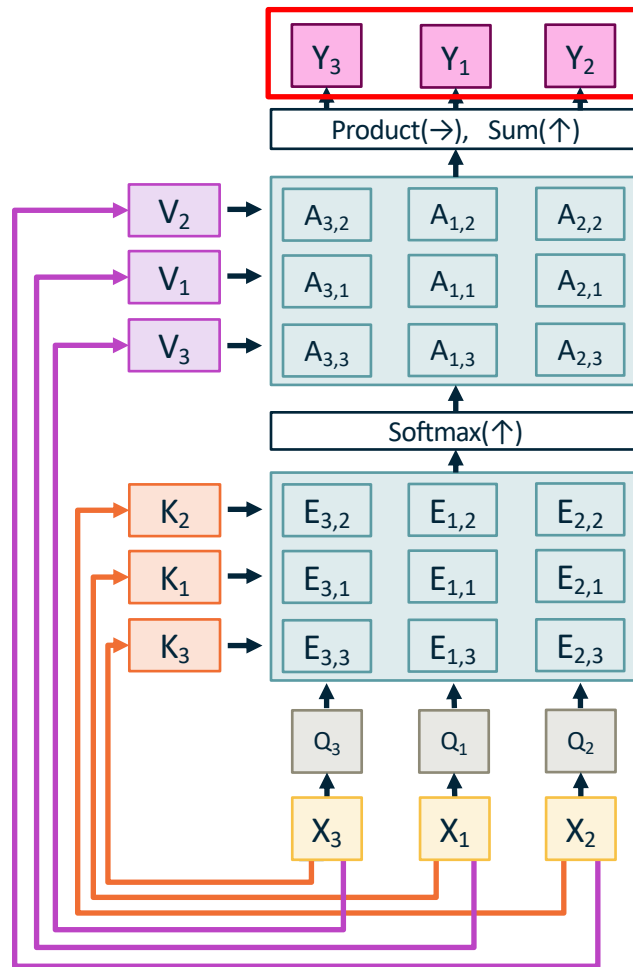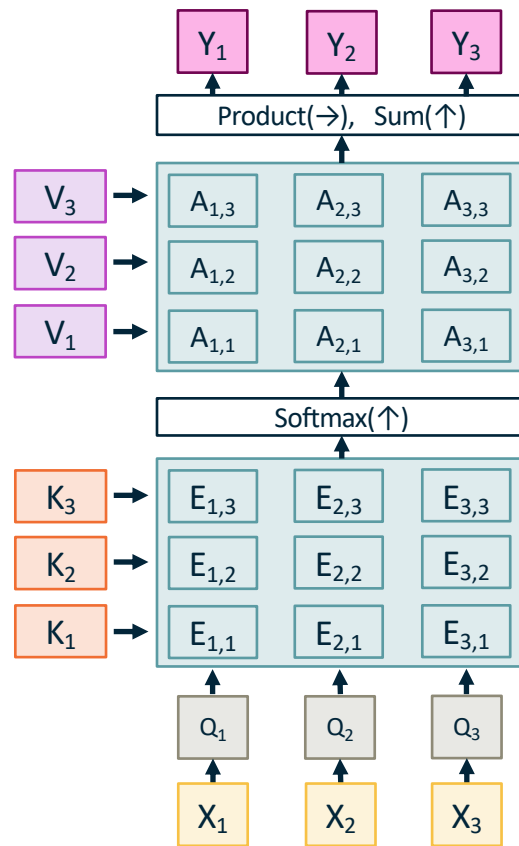**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_X \times N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} V_j$



Slide credit: Justin Johnson

# Self-Attention Layer

**Inputs**:
**Input vectors**: **X** (Shape: $N_X$ x $D_X$)
**Key matrix**: **W$_K$** (Shape: $D_X$ x $D_Q$)
**Value matrix**: **W$_V$** (Shape: $D_X$ x $D_V$)
**Query matrix**: **W$_Q$** (Shape: $D_X$ x $D_Q$)

**Computation**:
**Query vectors**: **Q** = **XW$_Q$**
**Key vectors**: **K** = **XW$_K$**  (Shape: $N_X$ x $D_Q$)
**Value vectors**: **V** = **XW$_V$** (Shape: $N_X$ x $D_V$)
**Similarities**: E = **QK$^T$** (Shape: $N_X$ x $N_X$) $E_{i,j}$ = **Q$_i$** · **K$_j$** / sqrt($D_Q$)
**Attention weights**: A = softmax(E, dim=1)  (Shape: $N_X$ x $N_X$)
**Output vectors**: Y = A**V** (Shape: $N_X$ x $D_V$) $Y_i$ = $\sum_j A_{i,j}$**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**: $X$ (Shape: $N_X$ x $D_X$)
**Key matrix**: $W_K$ (Shape: $D_X$ x $D_Q$)
**Value matrix**: $W_V$ (Shape: $D_X$ x $D_V$)
**Query matrix**: $W_Q$ (Shape: $D_X$ x $D_Q$)

**Computation**:
**Query vectors**: $Q = XW_Q$
**Key vectors**: $K = XW_K$ (Shape: $N_X$ x $D_Q$)
**Value vectors**: $V = XW_V$ (Shape: $N_X$ x $D_V$)
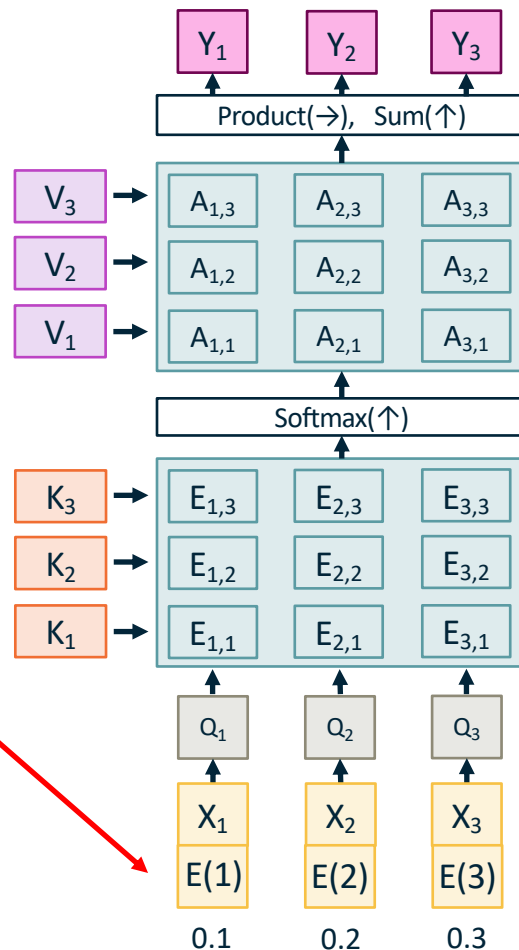**Similarities**: $E = QK^T$ (Shape: $N_X$ x $N_X$) $E_{i,j} = Q_i \cdot K_j / sqrt(D_Q)$
**Attention weights**: $A = softmax(E, dim=1)$ (Shape: $N_X$ x $N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_X$ x $D_V$) $Y_i = \sum_j A_{i,j} V_j$

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

# Self-Attention Layer

**Inputs**:
**Input vectors**: $X$ (Shape: $N_X \times D_X$)
**Key matrix**: $W_K$ (Shape: $D_X \times D_Q$)
**Value matrix**: $W_V$ (Shape: $D_X \times D_V$)
**Query matrix**: $W_Q$ (Shape: $D_X \times D_Q$)

**Computation**:
**Query vectors**: $Q = XW_Q$
**Key vectors**: $K = XW_K$  (Shape: $N_X \times D_Q$)
**Value vectors**: $V = XW_V$ (Shape: $N_X \times D_V$)
**Similarities**: $E = QK^T$ (Shape: $N_X \times N_X$) $E_{i,j} = Q_i \cdot K_j / \sqrt{D_Q}$
**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$  (Shape: $N_X \times N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} V_j$

In order to make processing position-aware, concatenate input with **positional encoding 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.

**Learned Positional Encoding**: $PE_\theta(pos, i)$.
Learn the most suitable position embedding for the training set.

# **Masked** Self-Attention Layer

**Inputs**:
**Input vectors**: $X$ (Shape: $N_X \times D_X$)
**Key matrix**: $W_K$ (Shape: $D_X \times D_Q$)
**Value matrix**: $W_V$ (Shape: $D_X \times D_V$)
**Query matrix**: $W_Q$ (Shape: $D_X \times D_Q$)

**Computation**:
**Query vectors**: $Q = XW_Q$
**Key vectors**: $K = XW_K$  (Shape: $N_X \times D_Q$)
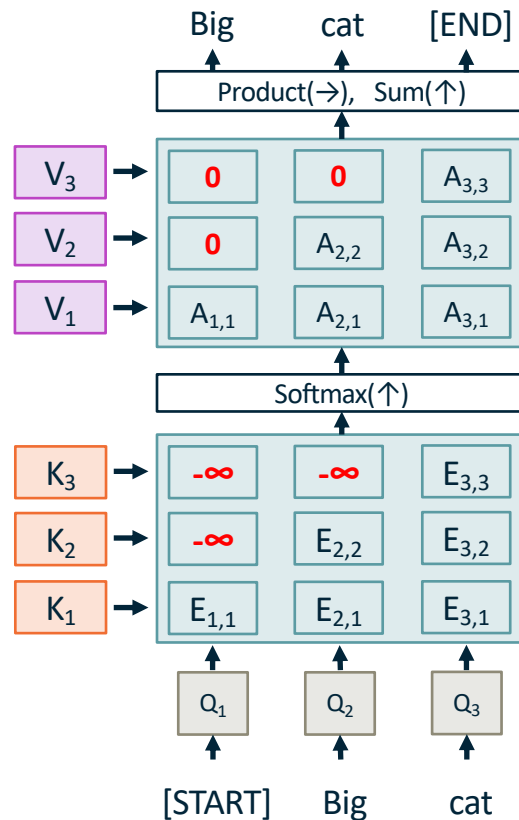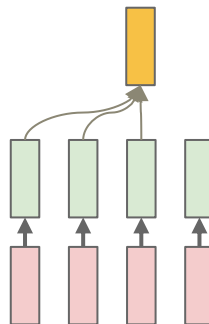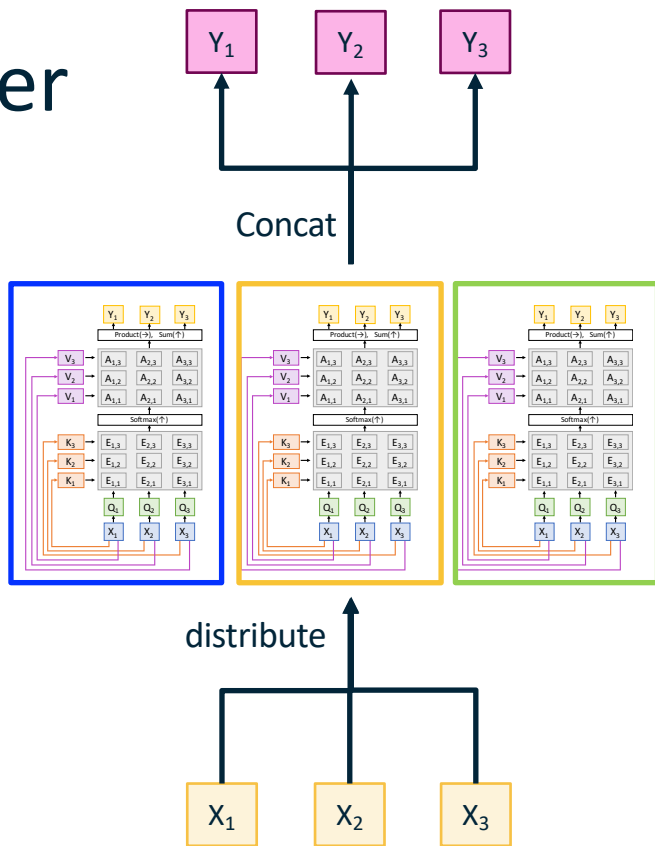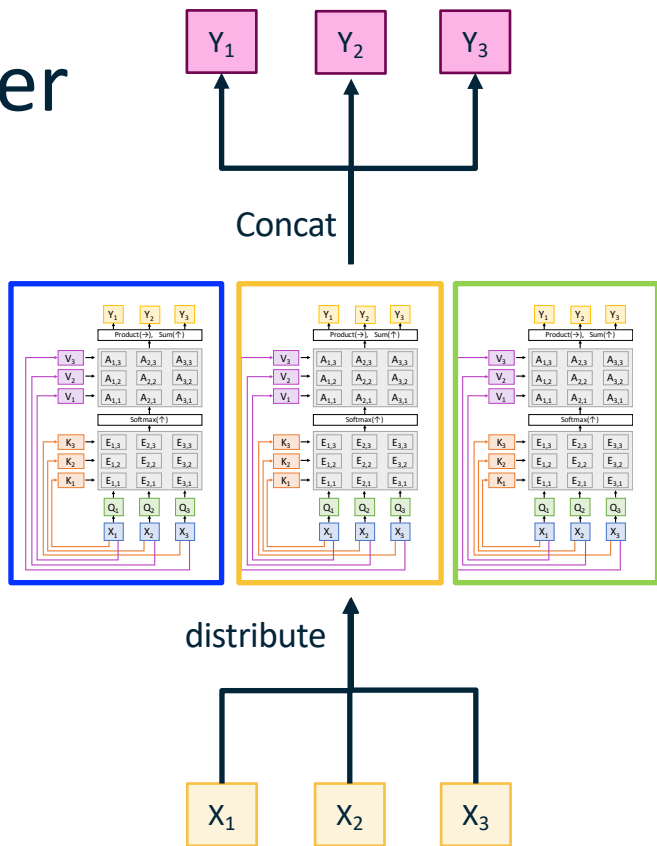**Value vectors**: $V = XW_V$ (Shape: $N_X \times D_V$)
**Similarities**: $E = QK^T$ (Shape: $N_X \times N_X$) $E_{i,j} = Q_i \cdot K_j / sqrt(D_Q)$
**Attention weights**: $A = softmax(E, dim=1)$  (Shape: $N_X \times N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} V_j$

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

Used for sequence decoding (predict next word)

Big    cat    [END]

Product($\rightarrow$),  Sum($\uparrow$)

| | | | |
|---|---|---|---|
| $V_3$ | **0** | **0** | $A_{3,3}$ |
| $V_2$ | **0** | $A_{2,2}$ | $A_{3,2}$ |
| $V_1$ | $A_{1,1}$ | $A_{2,1}$ | $A_{3,1}$ |

Softmax($\uparrow$)

| | | | |
|---|---|---|---|
| $K_3$ | **-∞** | **-∞** | $E_{3,3}$ |
| $K_2$ | **-∞** | $E_{2,2}$ | $E_{3,2}$ |
| $K_1$ | $E_{1,1}$ | $E_{2,1}$ | $E_{3,1}$ |

$Q_1$    $Q_2$    $Q_3$

[START]    Big    cat

# **Multi-headed** Self-Attention Layer

**Inputs**:
**Input vectors**: **X** (Shape: $N_X$ x $D_X$)
**Key matrix**: $W_K$ (Shape: $D_X$ x $D_Q$)
**Value matrix**: $W_V$ (Shape: $D_X$ x $D_V$)
**Query matrix**: $W_Q$ (Shape: $D_X$ x $D_Q$)

Use H independent "Attention Heads" in parallel

**Computation**:
**Query vectors**: **Q** = $XW_Q$
**Key vectors**: **K** = $XW_K$  (Shape: $N_X$ x $D_Q$)
**Value vectors**: **V** = $XW_V$ (Shape: $N_X$ x $D_V$)
**Similarities**: E = $QK^T$ (Shape: $N_X$ x $N_X$) $E_{i,j}$ = $Q_i \cdot K_j$ / sqrt($D_Q$)
**Attention weights**: A = softmax(E, dim=1)  (Shape: $N_X$ x $N_X$)
**Output vectors**: Y = A**V** (Shape: $N_X$ x $D_V$) $Y_i = \sum_j A_{i,j} V_j$



Concat

distribute

# **Multi-headed** Self-Attention Layer

**Inputs**:
**Input vectors**: **X** (Shape: $N_X$ x $D_X$)
**Key matrix**: **$W_K$** (Shape: $D_X$ x $D_Q$)
**Value matrix**: **$W_V$** (Shape: $D_X$ x $D_V$)
**Query matrix**: **$W_Q$** (Shape: $D_X$ x $D_Q$)


**Computation**:
**Query vectors**: **Q** = **X$W_Q$**
**Key vectors**: **K** = **X$W_K$** (Shape: $N_X$ x $D_Q$)
**Value vectors**: **V** = **X$W_V$** (Shape: $N_X$ x $D_V$)
**Similarities**: E = **QK$^T$** (Shape: $N_X$ x $N_X$) $E_{i,j}$ = **$Q_i$** · **$K_j$** / sqrt($D_Q$)
**Attention weights**: A = softmax(E, dim=1)  (Shape: $N_X$ x $N_X$)
**Output vectors**: Y = A**V** (Shape: $N_X$ x $D_V$) $Y_i$ = $\sum_j A_{i,j}$**$V_j$**

Use H independent "Attention Heads" in parallel

Concat

distribute

Highly parallelizable: Can compute attentions for all input element from all head 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**

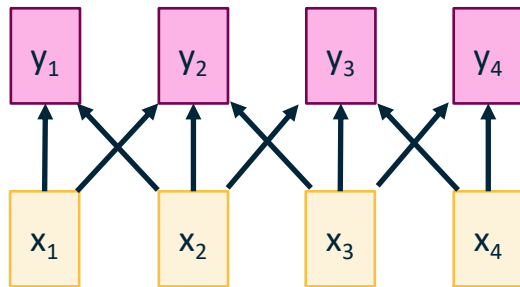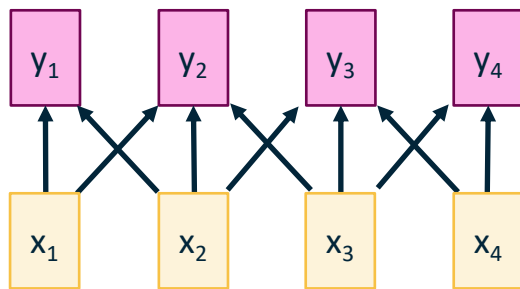# Three Ways of Processing Sequences

## 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**
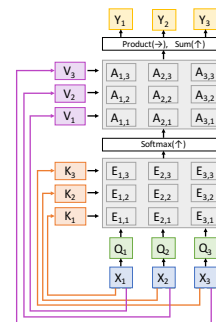
## 1D Convolution



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**
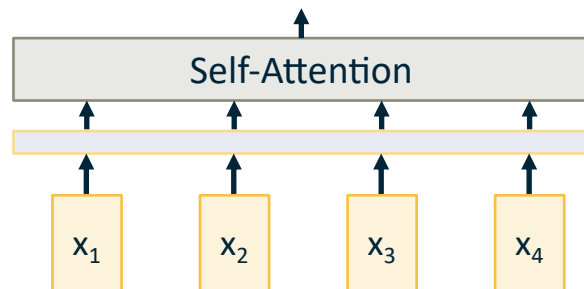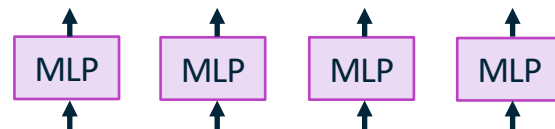(-) **Very memory intensive**
(-) **Requires positional encoding**

Slide credit: Justin Johnson
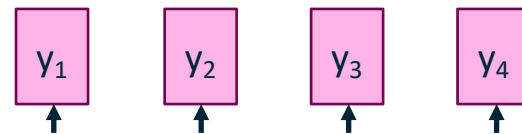
# Three Ways of Processing Sequences

Recurrent Neural Network          1D Convolution                    Self-Attention

## Attention is all you need

Vaswani et al, NeurIPS 2017

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

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

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

# The Transformer Block

$x_1$    $x_2$    $x_3$    $x_4$

Vaswani et al, "Attention is all you need", NeurIPS 2017

# The Transformer Block

All vectors interact
with each other



Vaswani et al, "Attention is all you need", NeurIPS 2017

# The Transformer Block

$y_1$ $y_2$ $y_3$ $y_4$

MLP independently on each vector

MLP   MLP   MLP   MLP

All vectors interact with each other

Self-Attention

$x_1$ $x_2$ $x_3$ $x_4$

Vaswani et al, "Attention is all you need", NeurIPS 2017

# The Transformer Block

$y_1$  $y_2$  $y_3$  $y_4$

MLP independently on
each vector

MLP  MLP  MLP  MLP

Residual connection

⊕

Self-Attention

All vectors interact
with each other

$x_1$  $x_2$  $x_3$  $x_4$

Vaswani et al, "Attention is all you need", NeurIPS 2017

# The Transformer Block

Recall **Layer Normalization**:

Given $h_1, ..., h_N$   (shape: D)

scale: $\gamma$          (shape: D)

shift: $\beta$           (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 element**, not across the sequence

MLP independently on each vector

Residual connection

All vectors interact with each other



Vaswani et al, "Attention is all you need", NeurIPS 2017

# The Transformer Block



Residual connection

MLP independently on each vector

Residual connection

All vectors interact with each other

Vaswani et al, "Attention is all you need", NeurIPS 2017
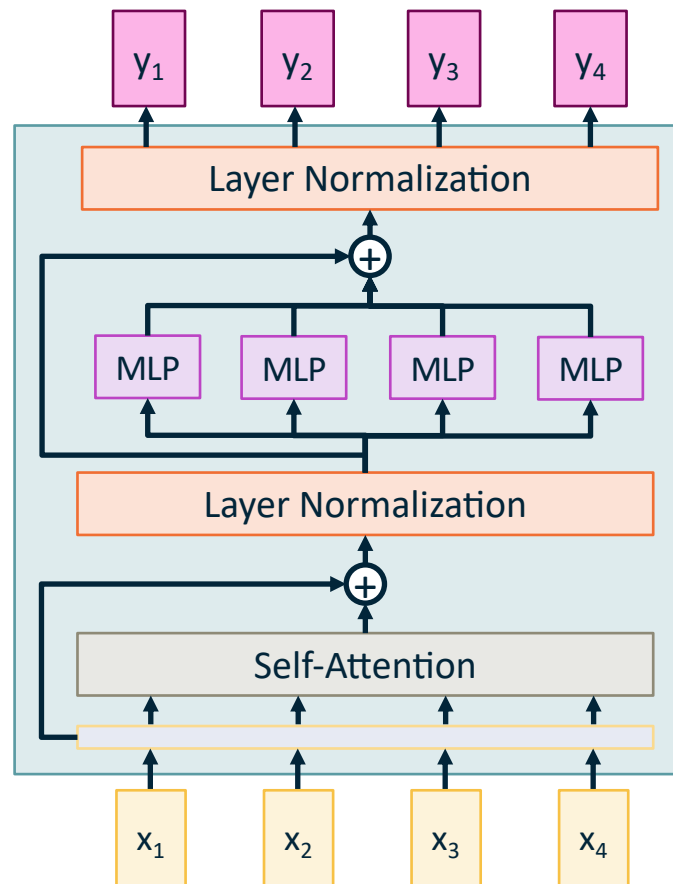
# The Transformer Block
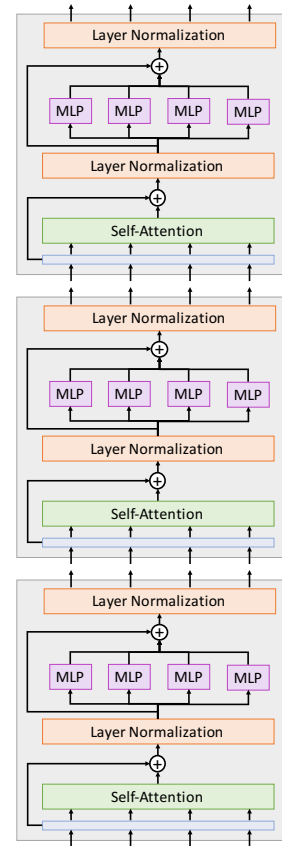
**Transformer Block:**
**Input**: Set of vectors x
**Output**: Set of vectors y

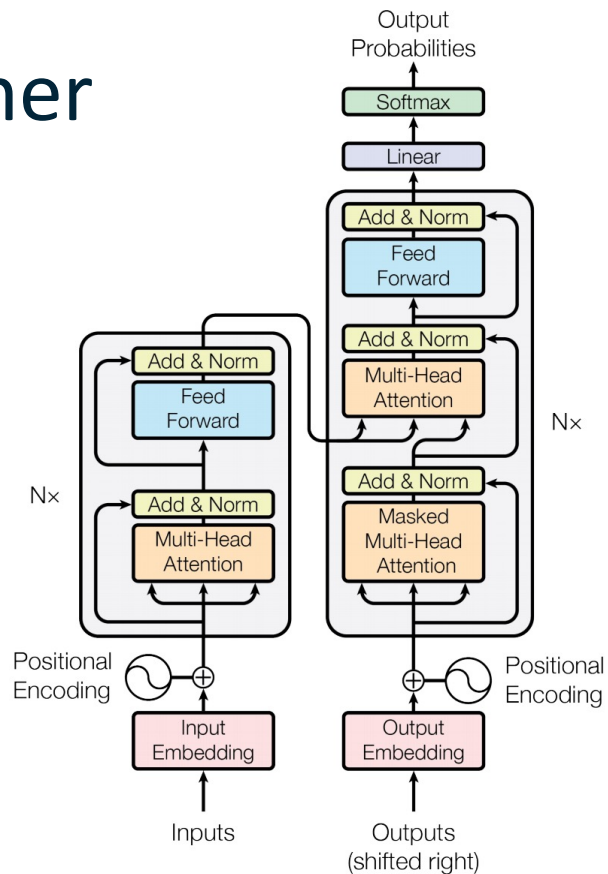Self-attention is the only interaction among vectors!

Layer norm and MLP work independently per vector

Highly scalable, highly parallelizable

Vaswani et al, "Attention is all you need", NeurIPS 2017

# 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

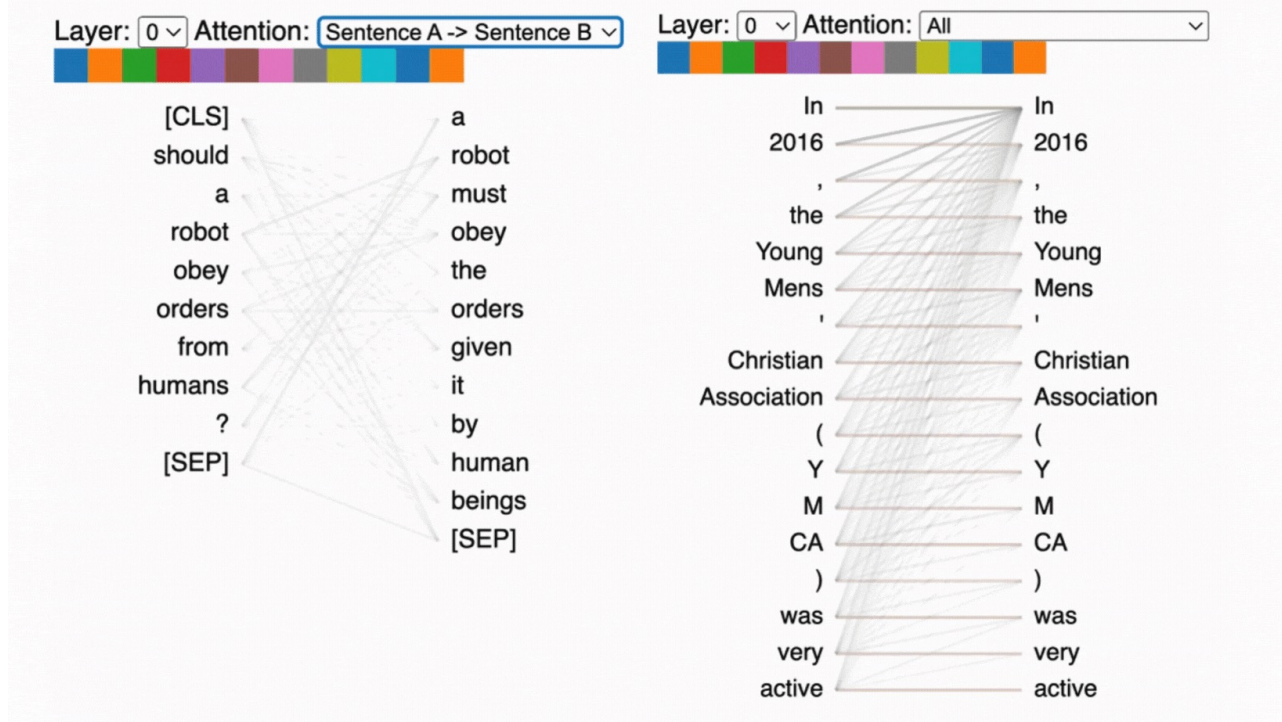A **Transformer** is a sequence
of transformer blocks



Vaswani et al, "Attention is all you need", NeurIPS 2017

Slide credit: Justin Johnson

# The Transformer



Encoder-Decoder

Vaswani et al, "Attention is all you need", NeurIPS 2017

# Visualizing Transformer Attentions

In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

```
MODEL COMPLETION (MACHINE-WRITTEN, 10 TRIES)
The scientist named the population, after their distinctive horn, Ovid's
Unicorn. These four-horned, silver-white unicorns were previously unknown to
science.

Now, after almost two centuries, the mystery of what sparked this odd
phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and
several companions, were exploring the Andes Mountains when they found a small
valley, with no other animals or humans. Pérez noticed that the valley had
what appeared to be a natural fountain, surrounded by two peaks of rock and
silver snow.

Pérez and the others then ventured further into the valley. "By the time we
reached the top of one peak, the water looked blue, with some crystals on
top," said Pérez.
```
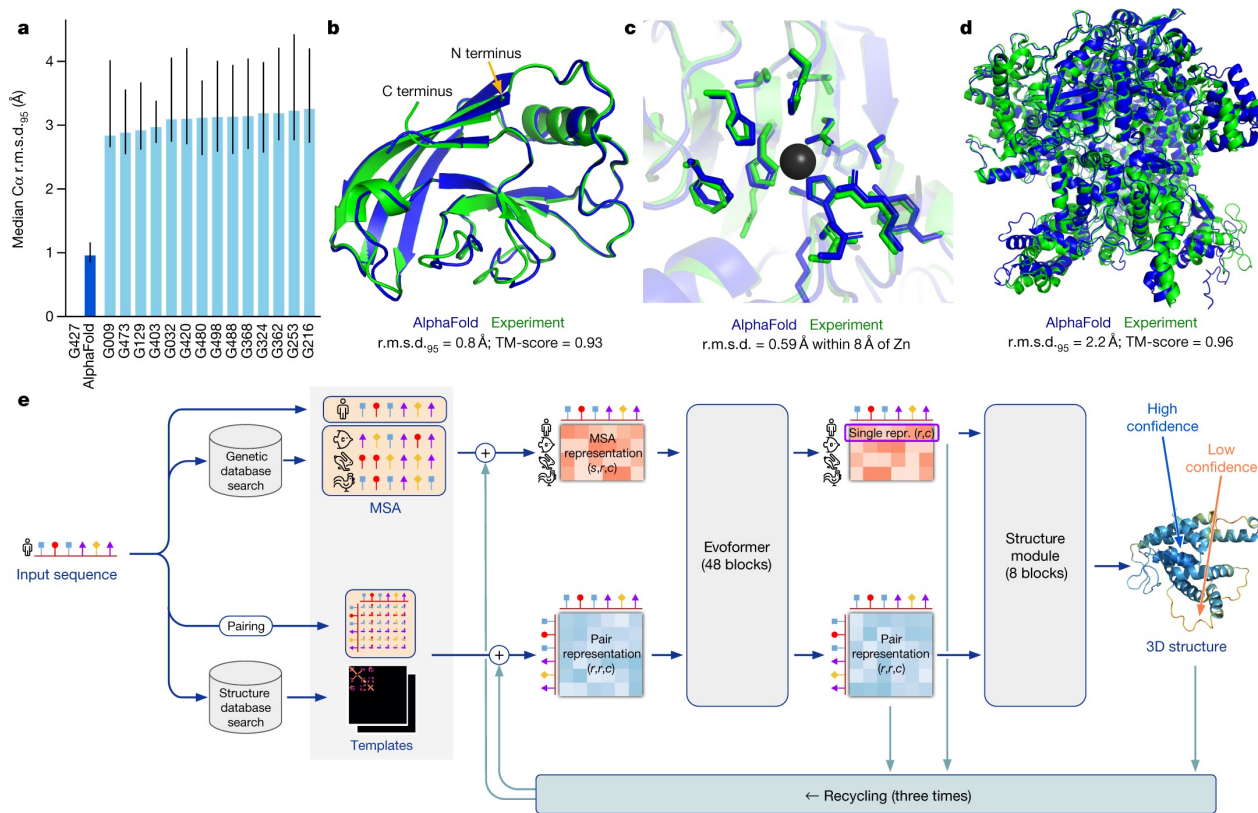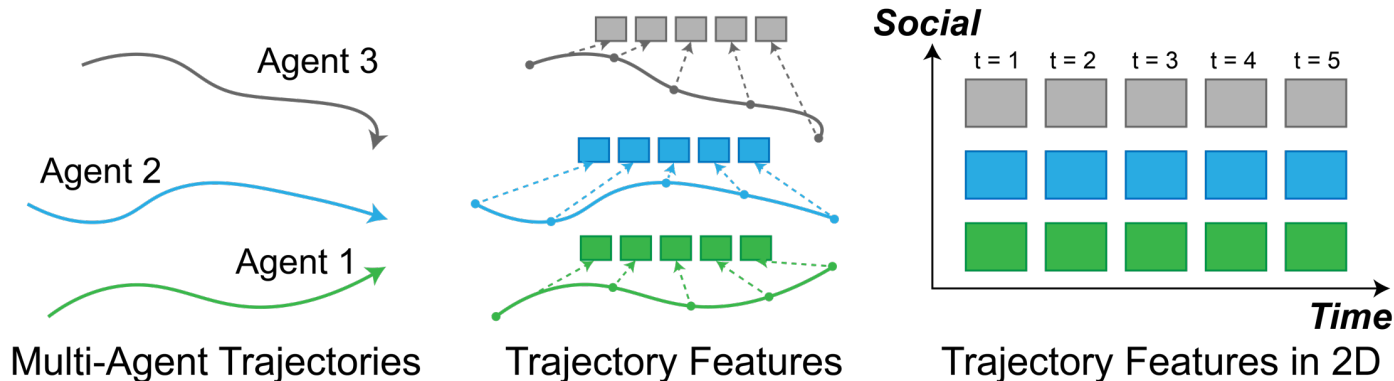
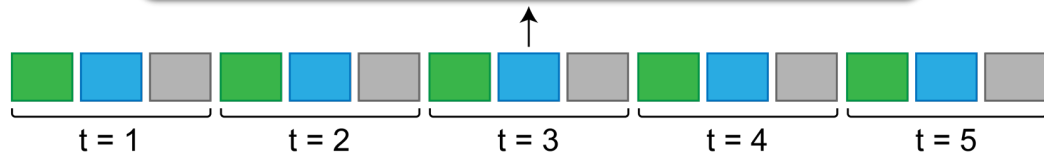# Can Attention/Transformers be used from more than text processing?

# Encoding/Decoding Protein Structures (AlphaFold)



https://www.nature.com/articles/s41586-021-03819-2
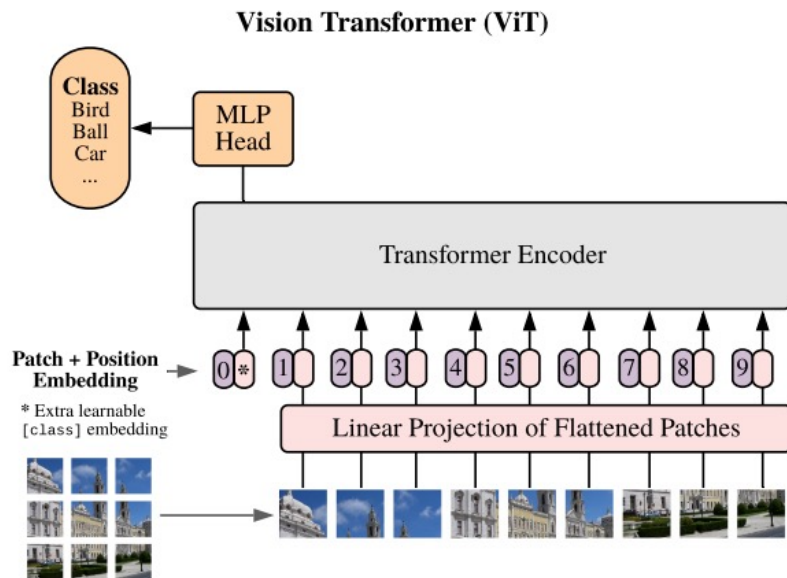
# Predicting Multi-agent Behaviors



Multi-Agent Trajectories

Trajectory Features

Trajectory Features in 2D

**Agent-Aware Transformer**

(**Joint** Social & Temporal Modeling +
Preserve **Time** & **Agent** Information)

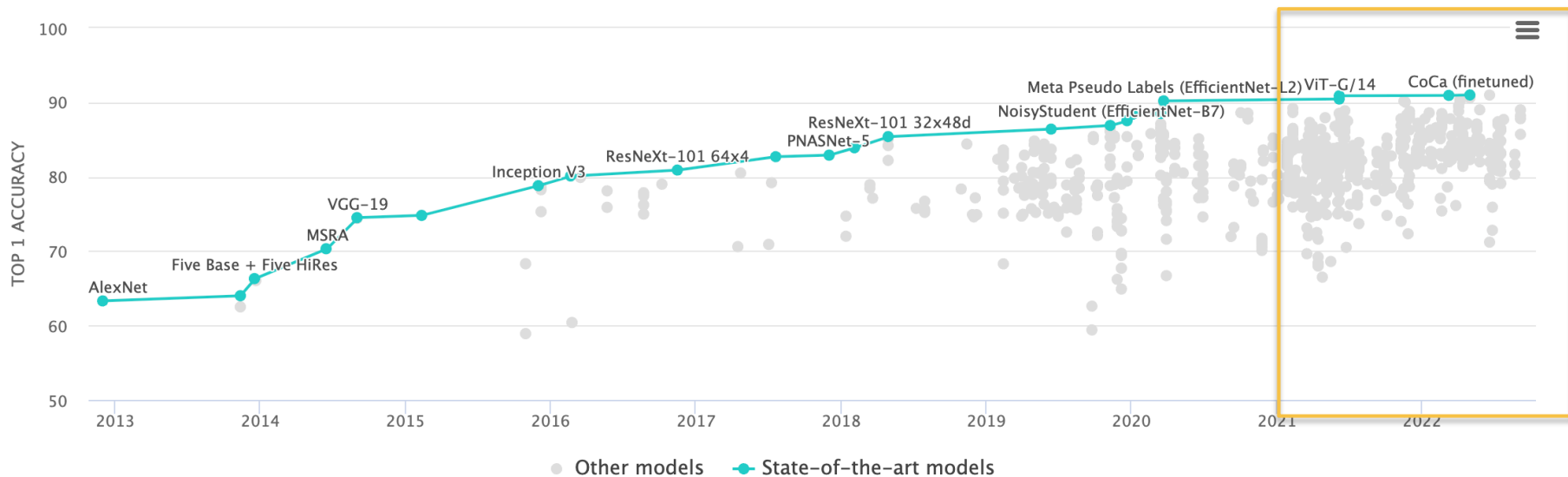Yuan et al., 2021 AgentFormer: Agent-Aware Transformers for Socio-Temporal Multi-Agent Forecasting

# ViT: Vision Transformer



An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale (Dosovitskiy *et al.*, 2021)

# ViT: Vision Transformer



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

# Formal Algorithms for Transformers

**Mary Phuong**[1] **and Marcus Hutter**[1]

[1]DeepMind

**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.**
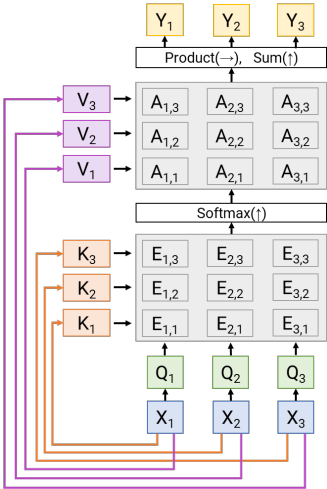
## Contents

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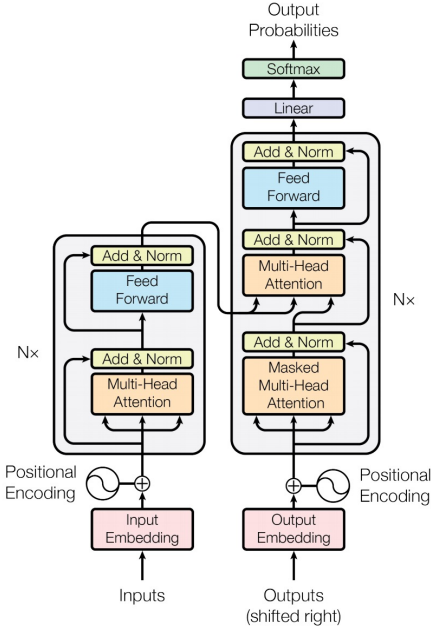
The essentially complete pseudocode is about 50 lines, compared to thousands of lines of actual real source code. We believe these formal algorithms will be useful for theoreticians who require compact, complete, and precise formulations, experimental researchers interested in implementing a Transformer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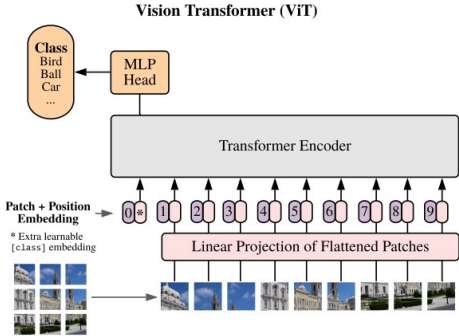
# Summary

## Self-Attention



## Transformer Model



## Beyond Language

Next time - Training Large Language Models
Instructor: Will Held