

Lecture 13: Attention

Midterm

Grades will be out in ~1 week

Please do not discuss midterm questions on Piazza

Someone left a waterbottle in exam room – Post on Piazza if it is yours

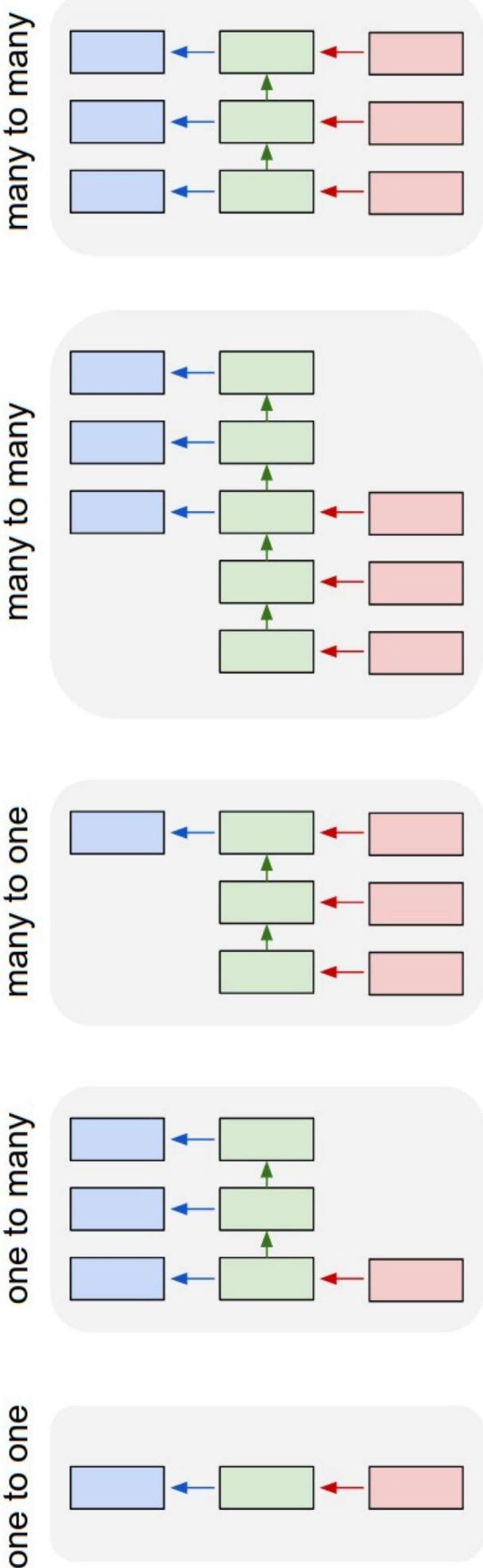
Assignment 4

A4 will be released today or tomorrow
Due 2 weeks from the time it is released

Will cover:

- PyTorch autograd
- Residual networks
- Recurrent neural networks
- Attention
- Feature visualization
- Style transfer
- Adversarial examples

Last Time: Recurrent Neural Networks

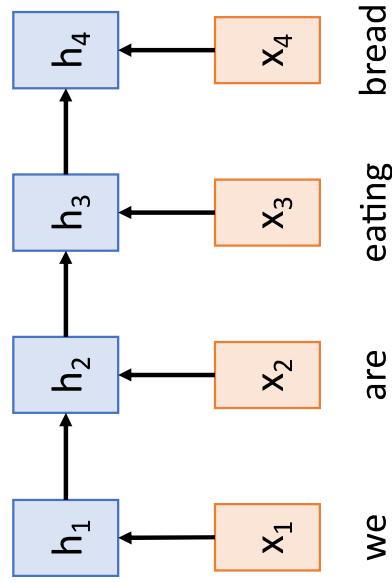


Sequence-to-Sequence with RNNs

Input: Sequence x_1, \dots, x_T

Output: Sequence $y_1, \dots, y_{T'}$

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



Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

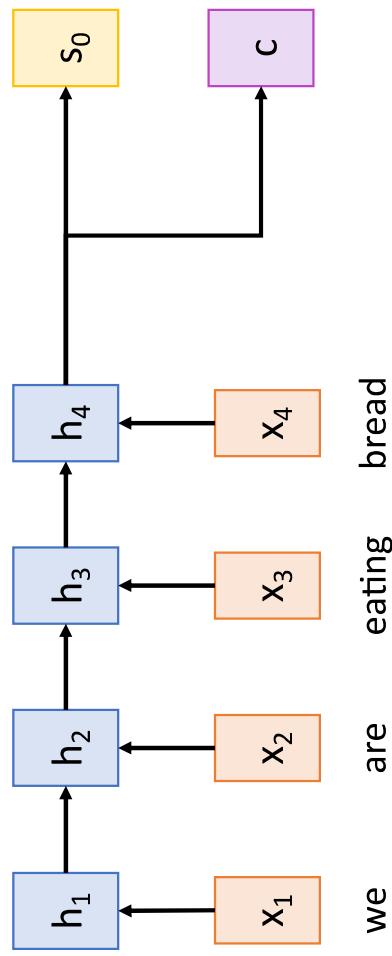
Sequence-to-Sequence with RNNs

Input: Sequence x_1, \dots, x_T

Output: Sequence $y_1, \dots, y_{T'}$

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

From final hidden state predict:
Initial decoder state s_0
Context vector c (often $c=h_T$)



Sutskever et al., "Sequence to sequence learning with neural networks", NeurIPS 2014

Sequence-to-Sequence with RNNs

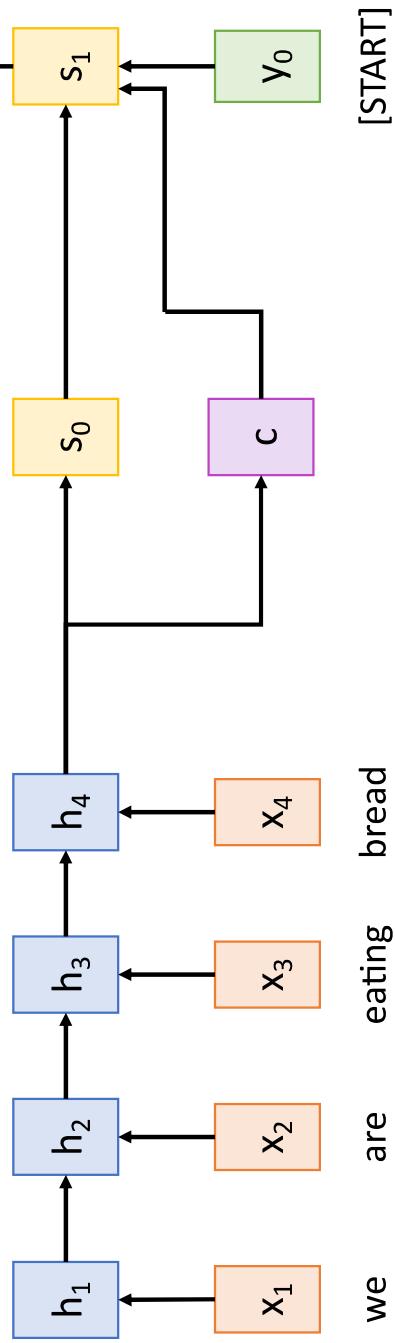
Input: Sequence x_1, \dots, x_T

Output: Sequence $y_1, \dots, y_{T'}$

Decoder: $s_t = g_U(y_{t-1}, h_{t-1}, c)$

estamos

Encoder: $h_t = f_W(x_t, h_{t-1})$
Decoder: s_t predict:
Initial decoder state s_0
Context vector c (often $c=h_T$)



Sutskever et al., "Sequence to sequence learning with neural networks", NeurIPS 2014

Sequence-to-Sequence with RNNs

Input: Sequence x_1, \dots, x_T

Output: Sequence $y_1, \dots, y_{T'}$

Decoder: $s_t = g_U(y_{t-1}, h_{t-1}, c)$

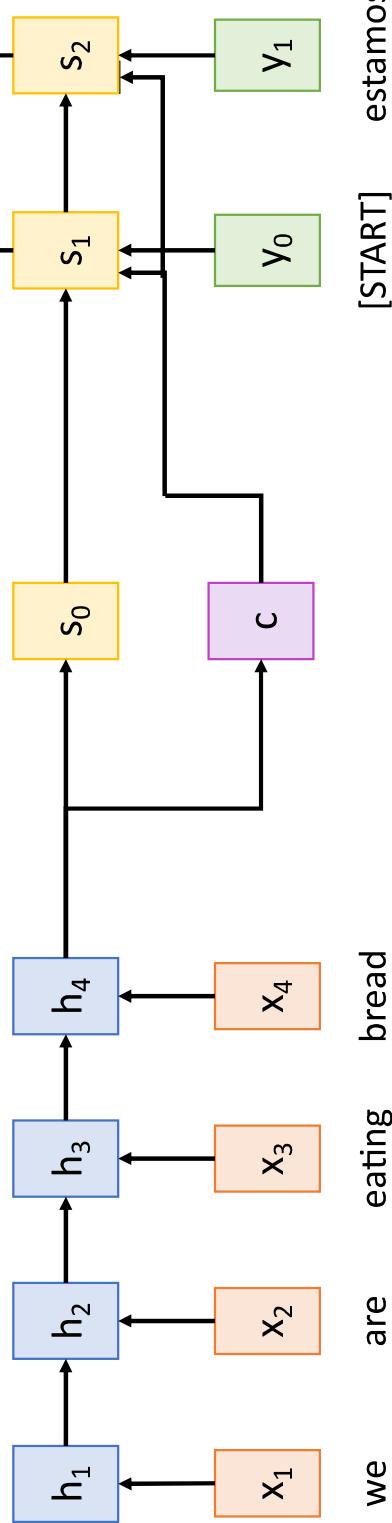
estamos comiendo

Encoder: $h_t = f_W(x_t, h_{t-1})$
Output: Sequence $y_1, \dots, y_{T'}$

From final hidden state predict:

Initial decoder state s_0

Context vector c (often $c=h_T$)



Sutskever et al., "Sequence to sequence learning with neural networks", NeurIPS 2014

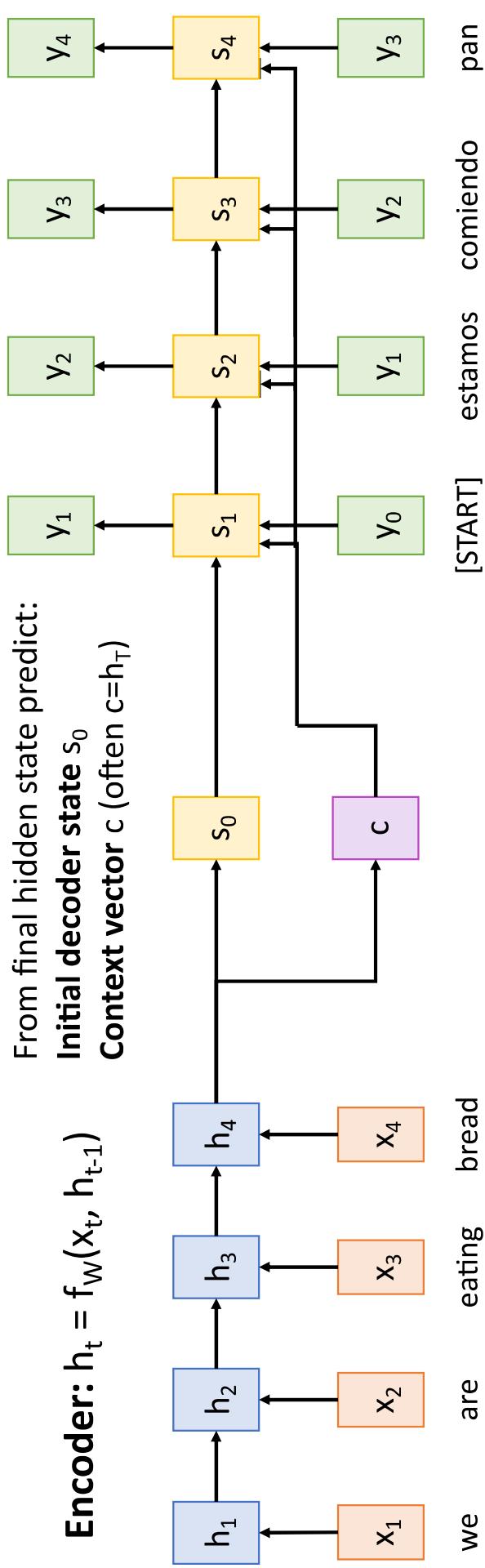
Sequence-to-Sequence with RNNs

Input: Sequence x_1, \dots, x_T

Output: Sequence $y_1, \dots, y_{T'}$

Decoder: $s_t = g_U(y_{t-1}, h_{t-1}, c)$

From final hidden state predict:
Initial decoder state s_0
Context vector c (often $c=h_T$)



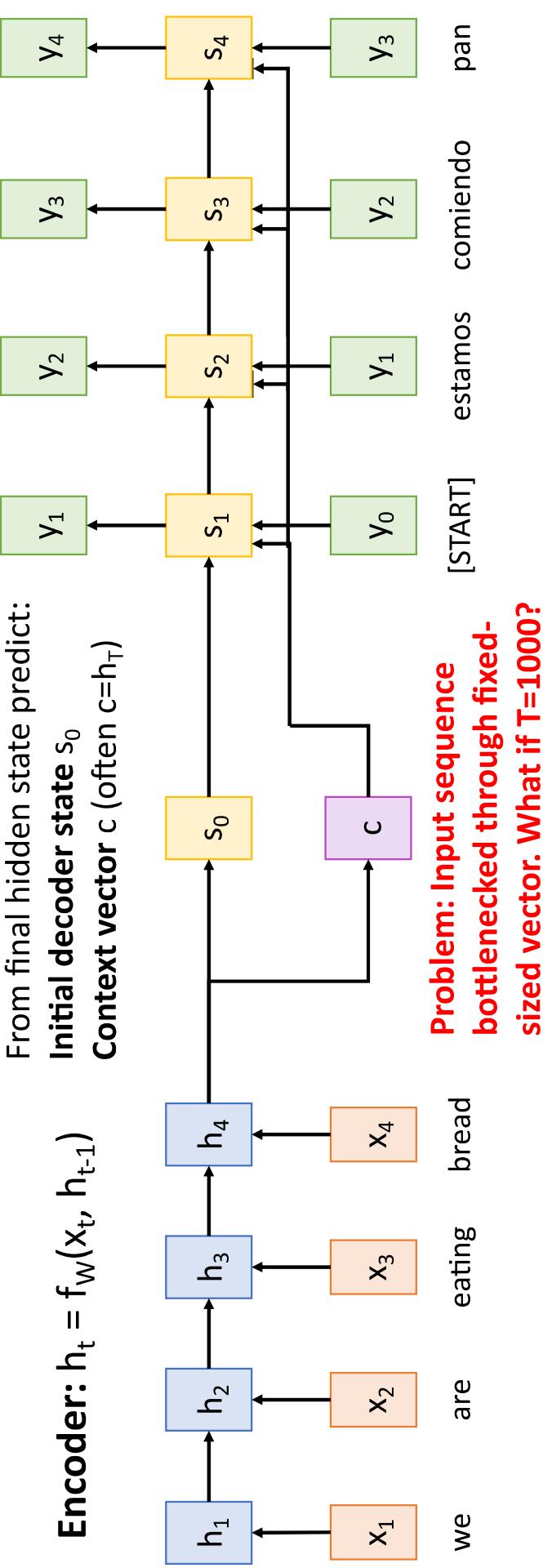
Sequence-to-Sequence with RNNs

Input: Sequence x_1, \dots, x_T

Output: Sequence $y_1, \dots, y_{T'}$

Decoder: $s_t = g_U(y_{t-1}, h_{t-1}, c)$

From final hidden state predict:
Initial decoder state s_0
Context vector c (often $c=h_T$)

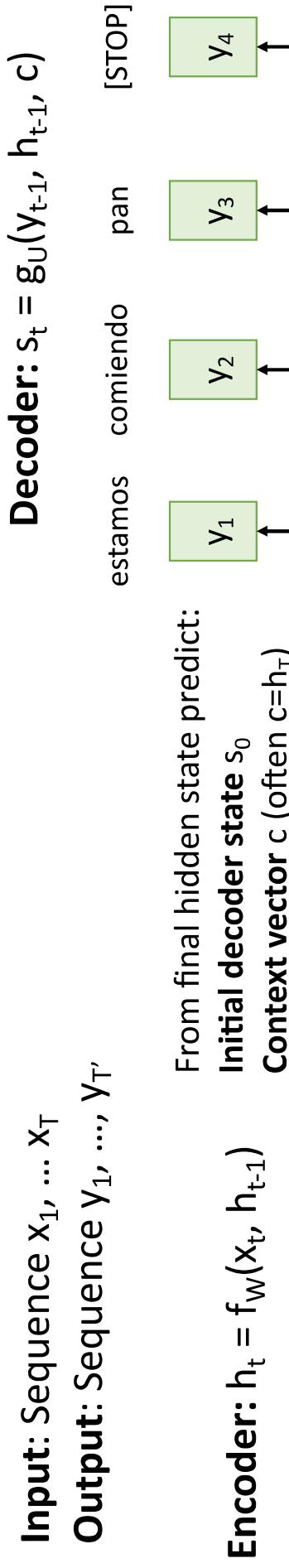


Sutskever et al., "Sequence to sequence learning with neural networks", NeurIPS 2014

Sequence-to-Sequence with RNNs

Input: Sequence x_1, \dots, x_T

Output: Sequence $y_1, \dots, y_{T'}$



we **are** **eating** **bread**

[START] **estamos** **comiendo** **pan**

Problem: Input sequence bottlenecked through fixed-sized vector. What if $T=1000$?

Idea: use new context vector at each step of decoder!

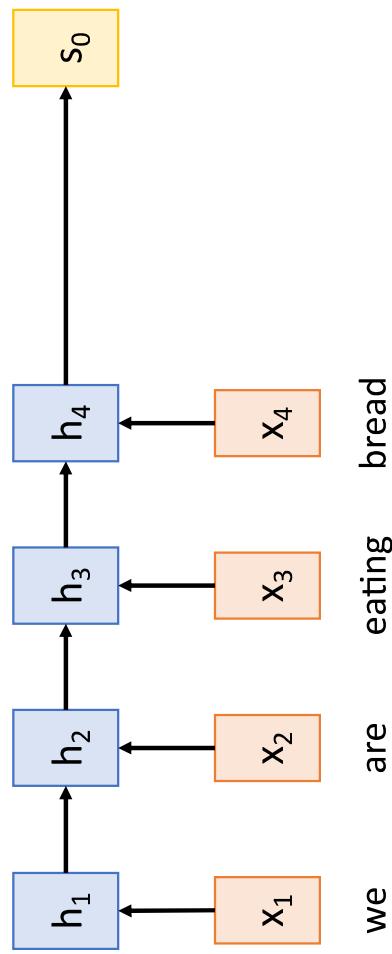
Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

Sequence-to-Sequence with RNNs and Attention

Input: Sequence x_1, \dots, x_T

Output: Sequence $y_1, \dots, y_{T'}$

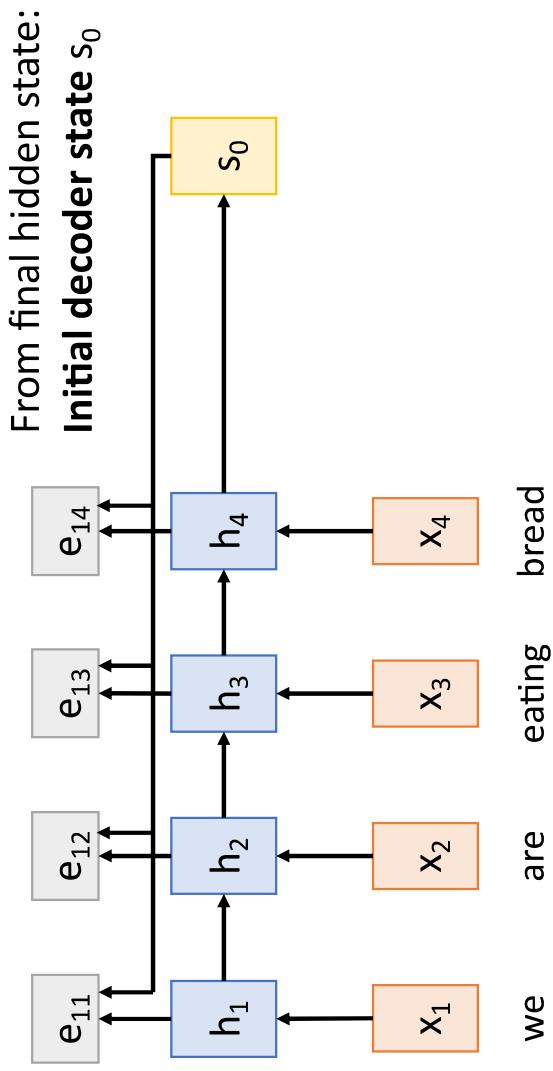
Encoder: $h_t = f_W(x_t, h_{t-1})$ From final hidden state:
Initial decoder state s_0



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

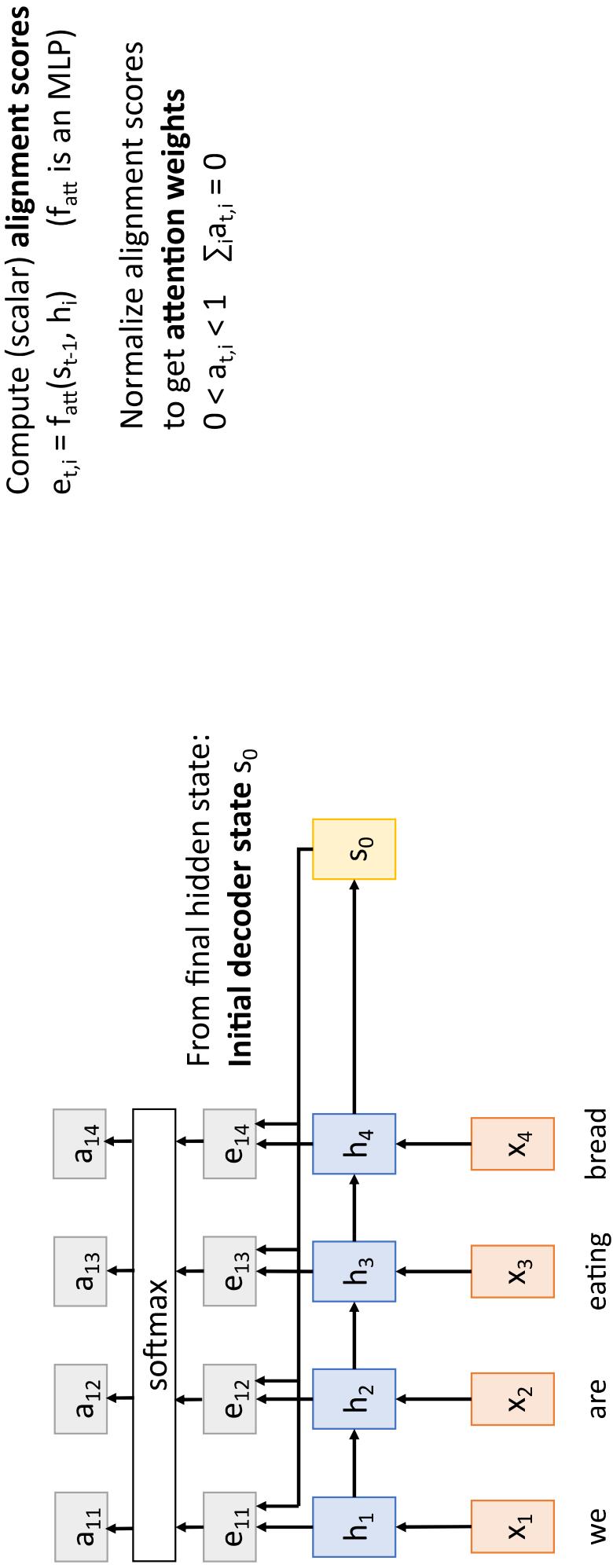
Sequence-to-Sequence with RNNs and Attention

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



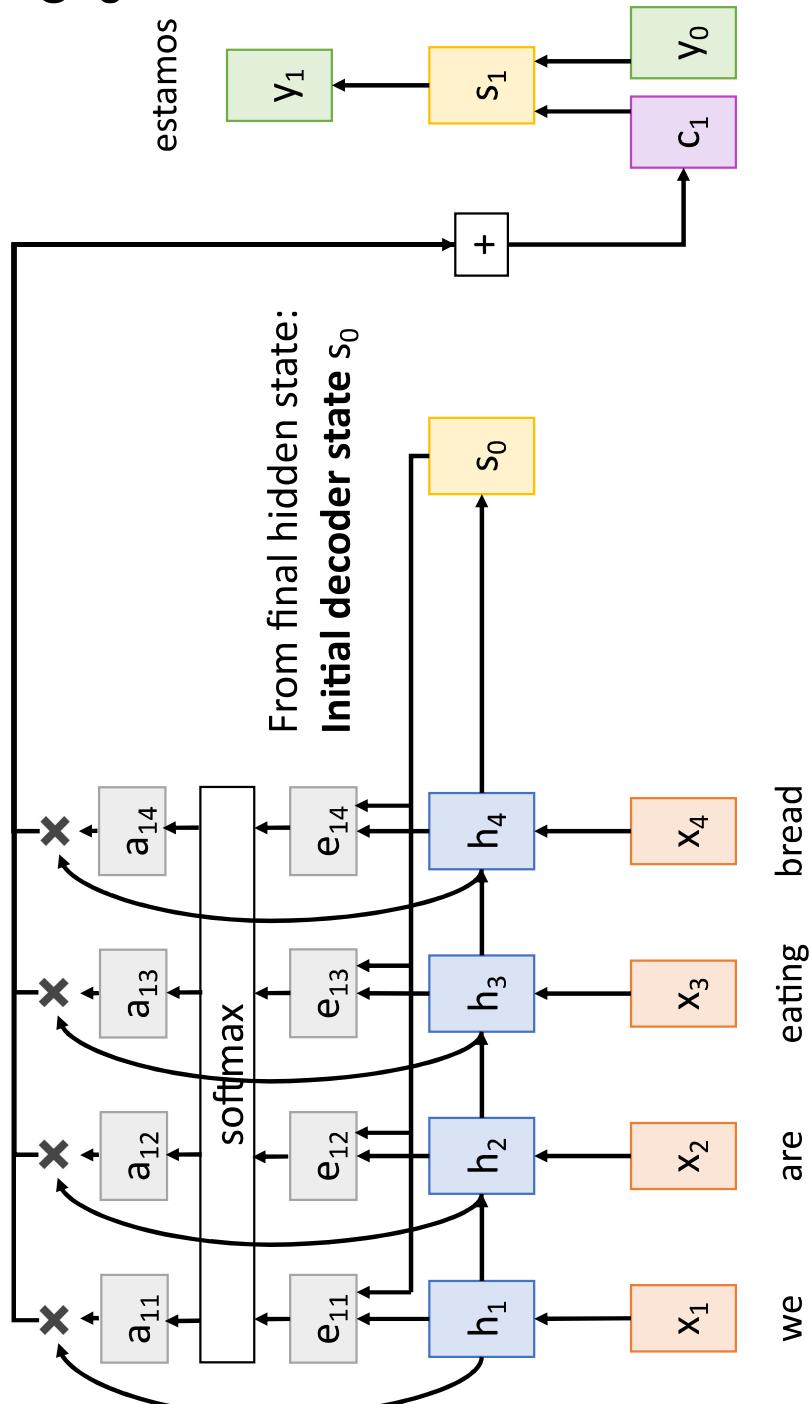
Bahdanau et al., "Neural machine translation by jointly learning to align and translate", ICLR 2015

Sequence-to-Sequence with RNNs and Attention



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

Sequence-to-Sequence with RNNs and Attention

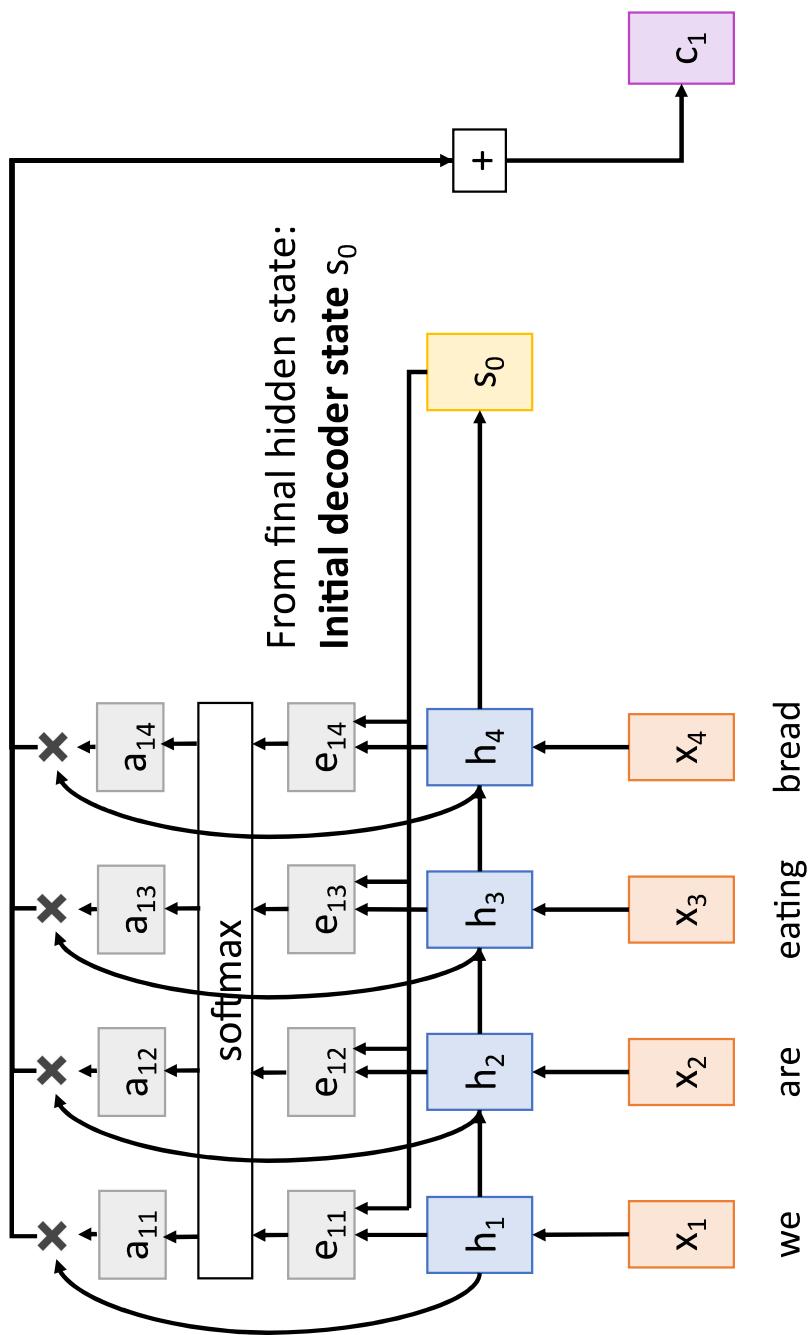


Bahdanau et al., "Neural machine translation by jointly learning to align and translate", ICLR 2015

Justin Johnson
Lecture 13 - 15

October 23, 2019

Sequence-to-Sequence with RNNs and Attention

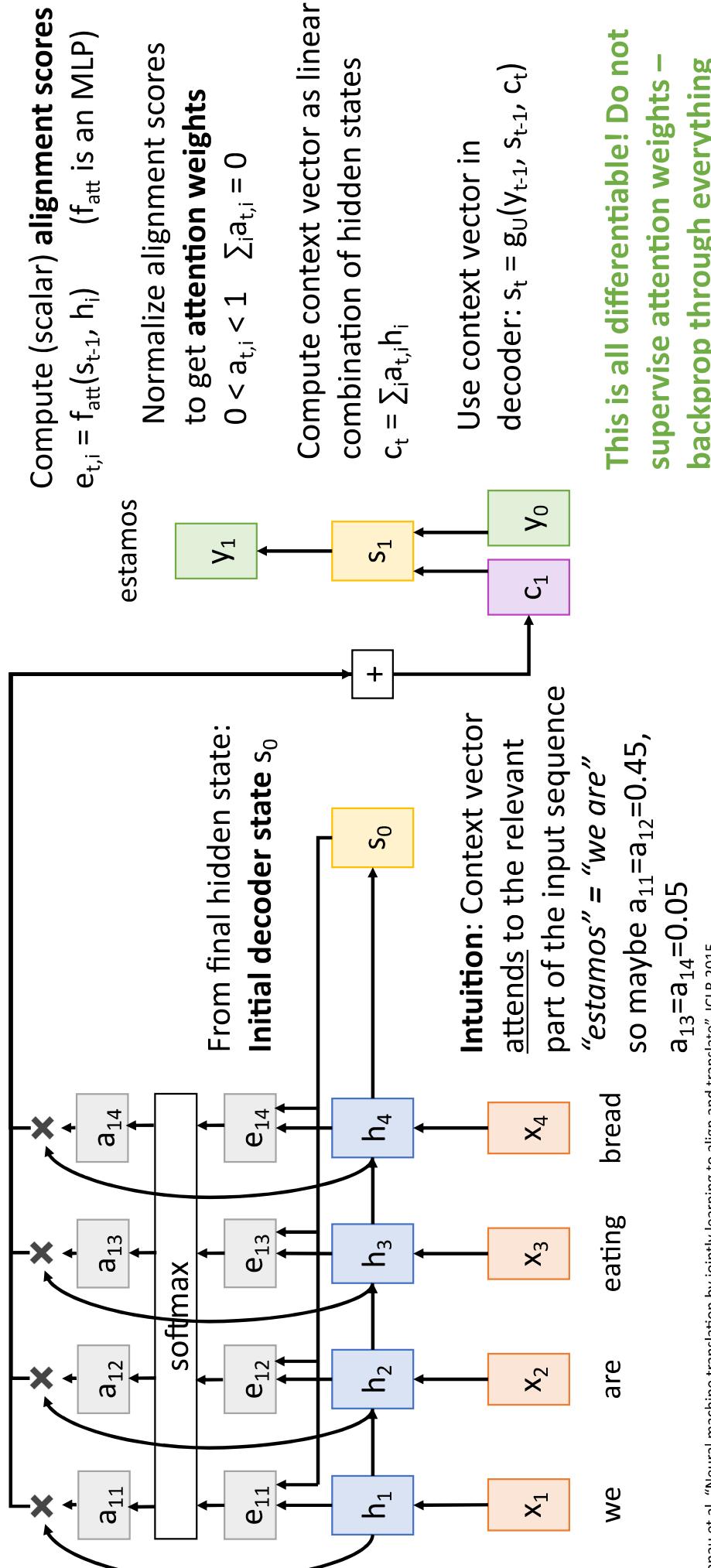


Compute (scalar) alignment scores
 $e_{t,i} = f_{\text{att}}(s_{t-1}, h_i)$ (f_{att} is an MLP)

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

Compute context vector as linear
 combination of hidden states
 $c_t = \sum_i a_{t,i} h_i$

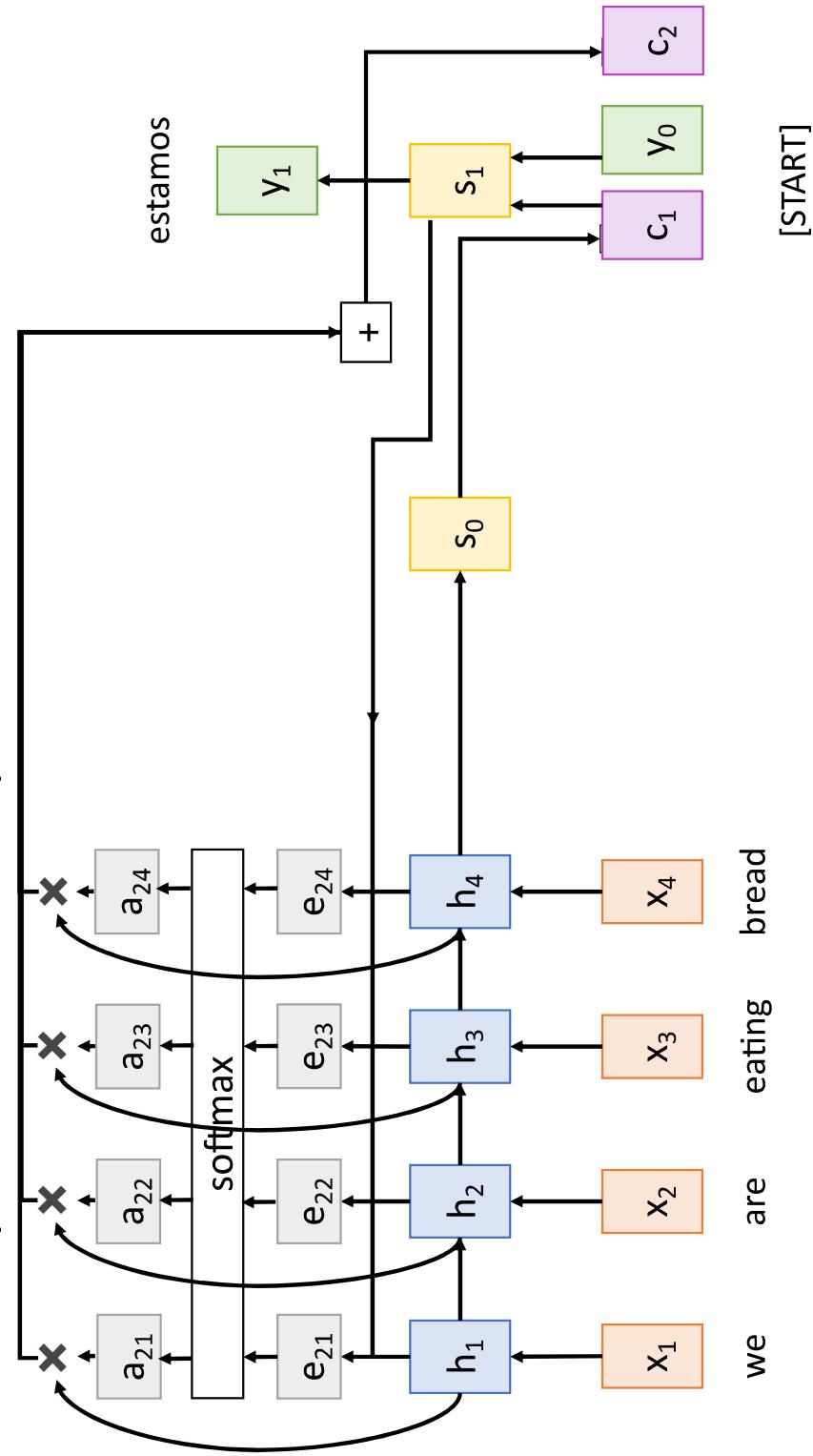
Sequence-to-Sequence with RNNs and Attention



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

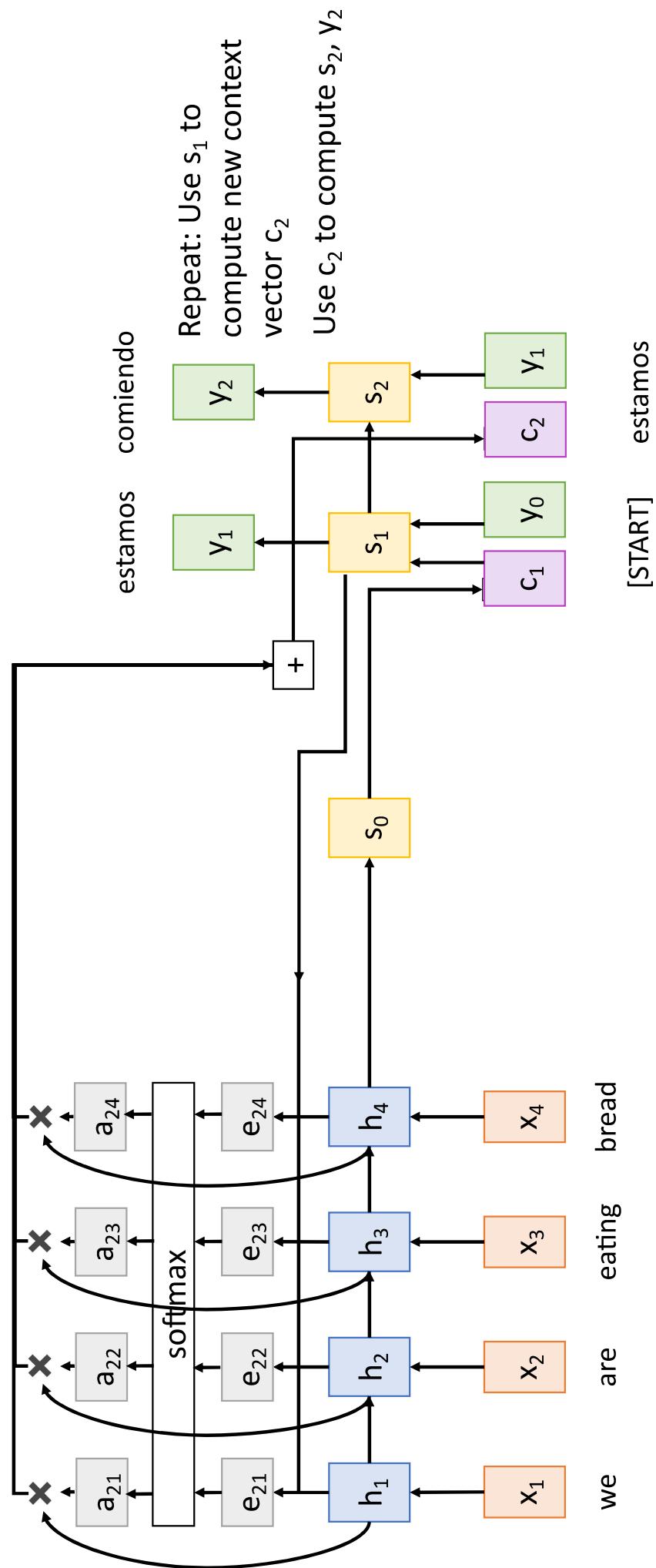
Sequence-to-Sequence with RNNs

Repeat: Use s_1 to compute
new context vector c_2



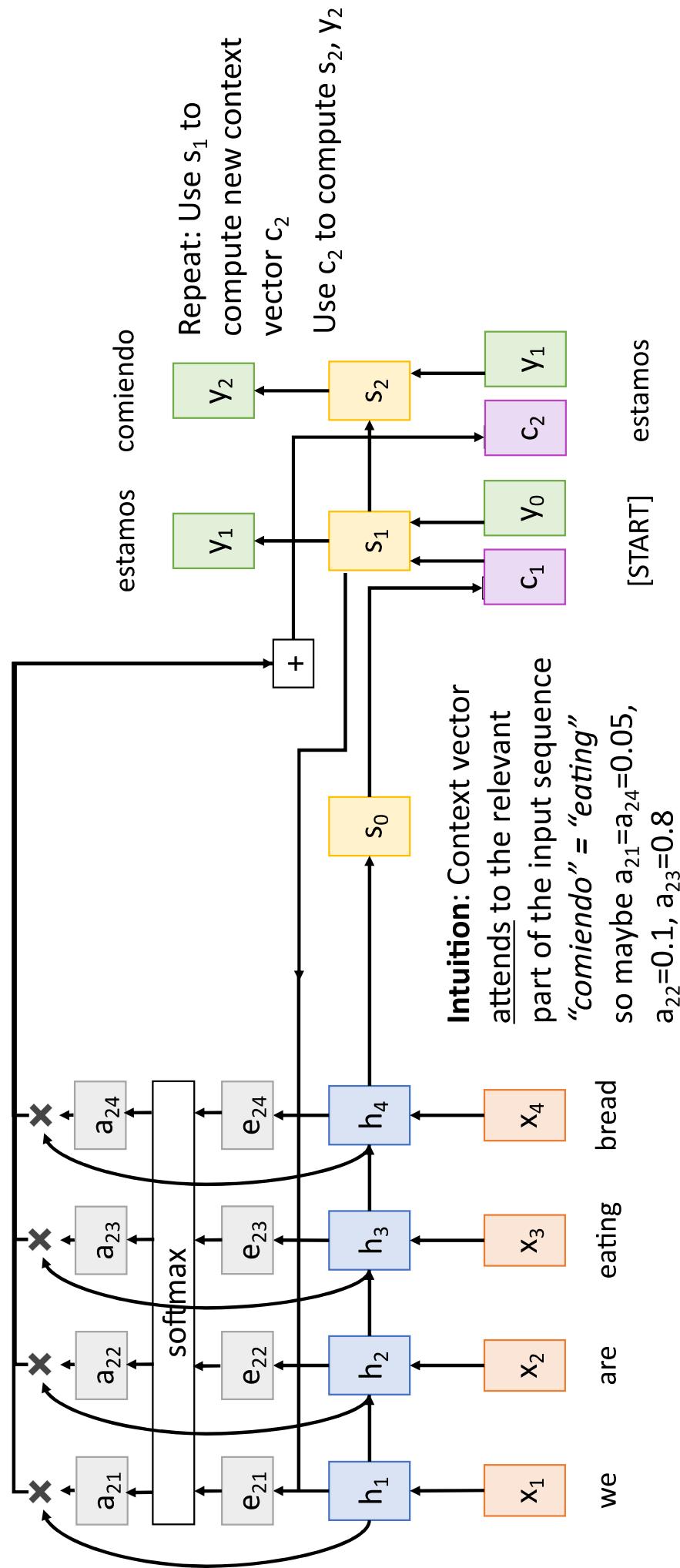
Bahdanau et al., "Neural machine translation by jointly learning to align and translate", ICLR 2015

Sequence-to-Sequence with RNNs and Attention



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

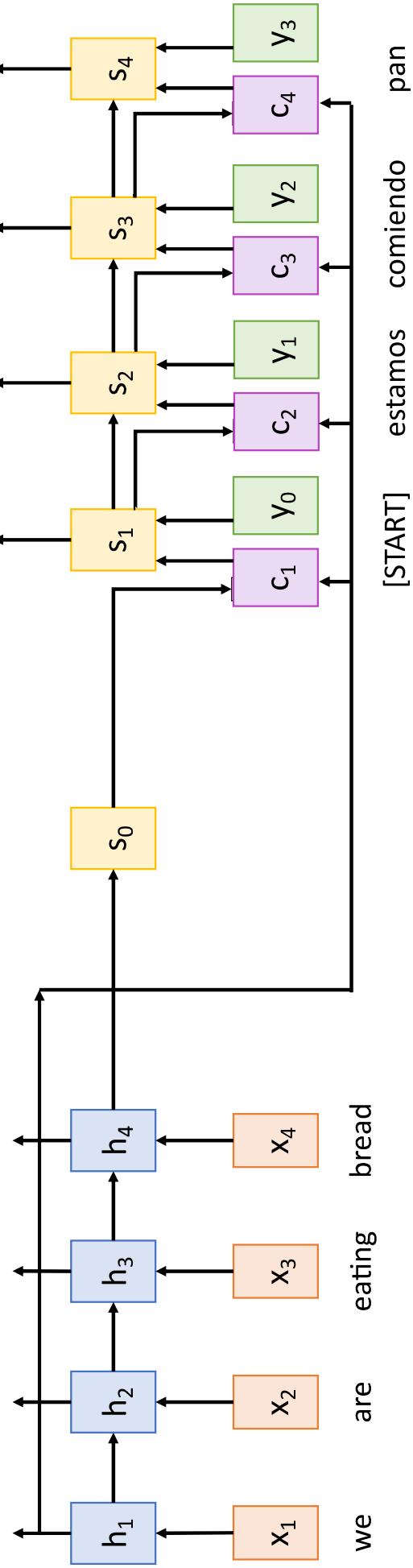
Sequence-to-Sequence with RNNs and Attention



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

Sequence-to-Sequence 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



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

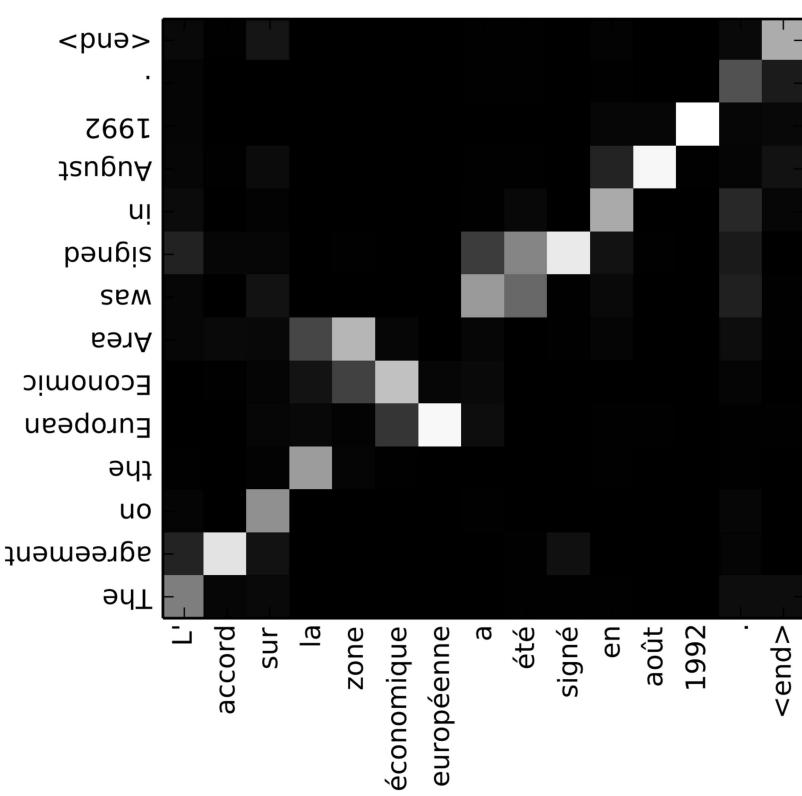
Sequence-to-Sequence with RNNs and Attention

Visualize attention weights $a_{t,i}$

Example: English to French translation

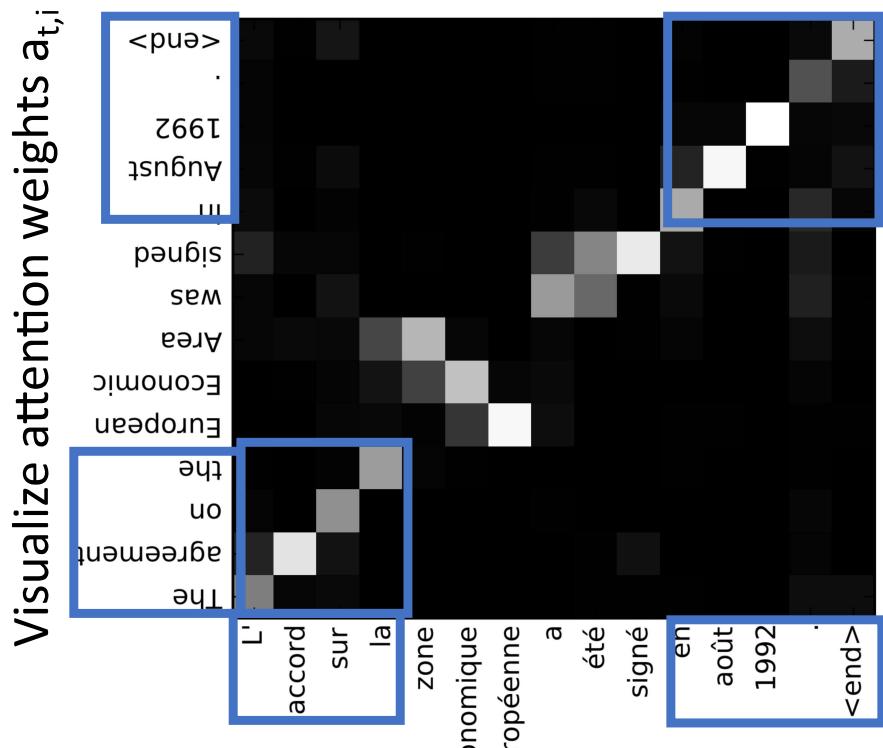
Input: “The agreement on the European Economic Area was signed in August 1992.”

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

Sequence-to-Sequence with RNNs and Attention



Example: English to French translation

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

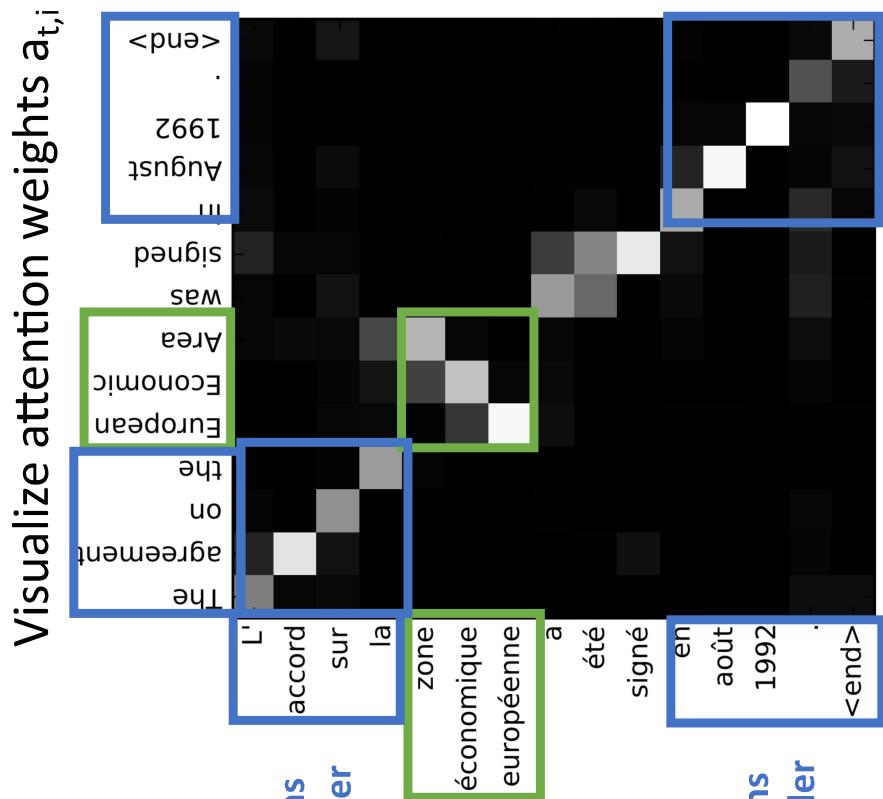
Diagonal attention means words correspond in order

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

Diagonal attention means words correspond in order

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

Sequence-to-Sequence with RNNs and Attention



Example: English to French translation

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

Diagonal attention means words correspond in order

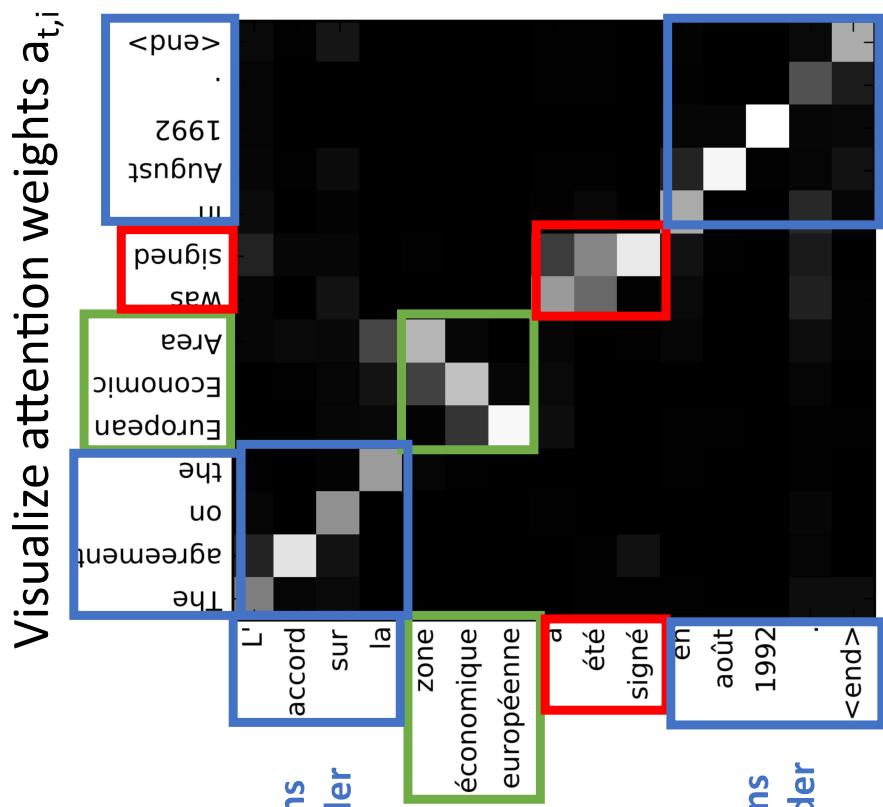
Attention figures out different word orders

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

Diagonal attention means words correspond in order

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

Sequence-to-Sequence with RNNs and Attention



Example: English to French translation

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

Diagonal attention means words correspond in order

Attention figures out different word orders

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

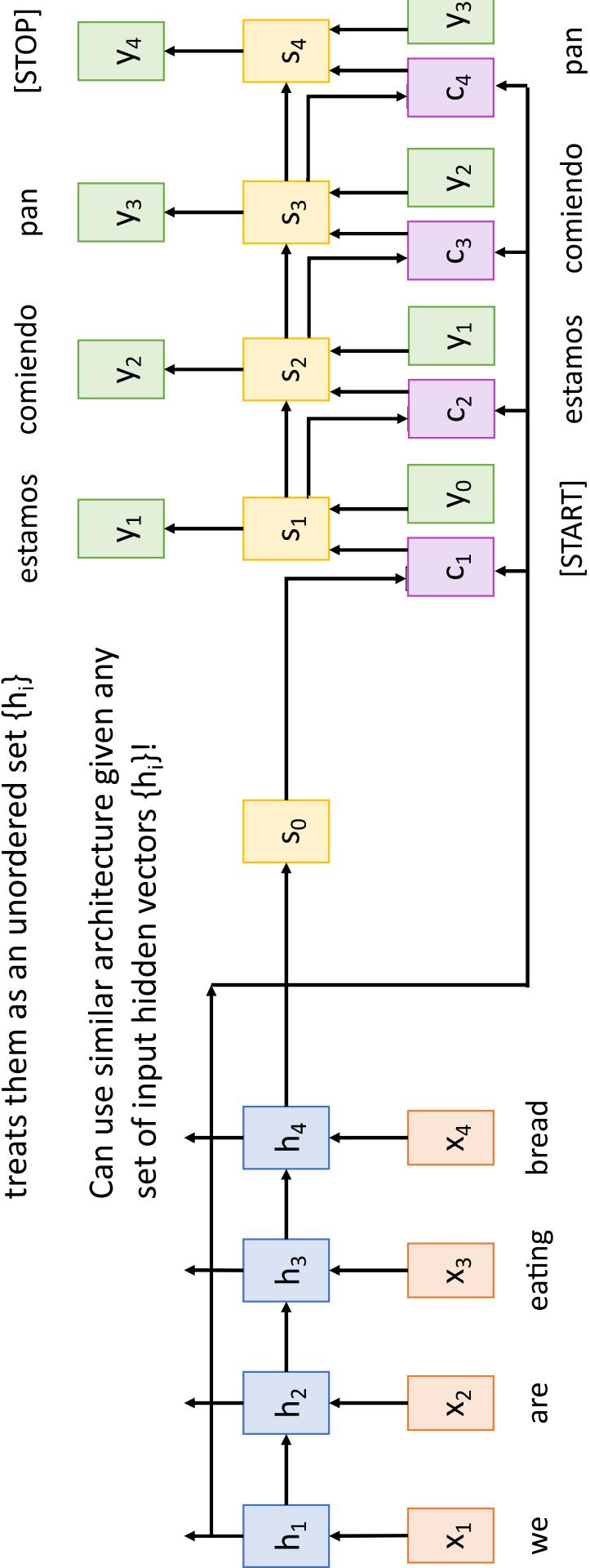
Verb conjugation

Diagonal attention means words correspond in order

Sequence-to-Sequence with RNNs and Attention

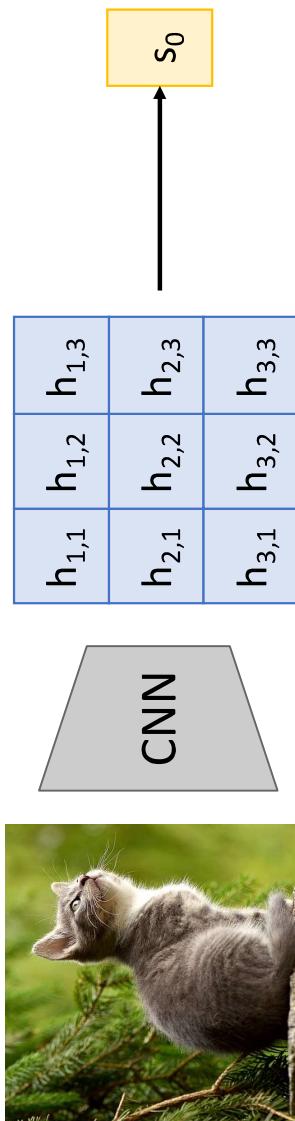
The decoder doesn't use the fact that h_i form an ordered sequence – it just treats them as an unordered set $\{h_i\}$

Can use similar architecture given any set of input hidden vectors $\{h_i\}$!



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

Image Captioning with RNNs and Attention



Use a CNN to compute a
grid of features for an image

Cat image is free to use under the Pixabay license.

Xu et al., "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Justin Johnson

Lecture 13 - 27

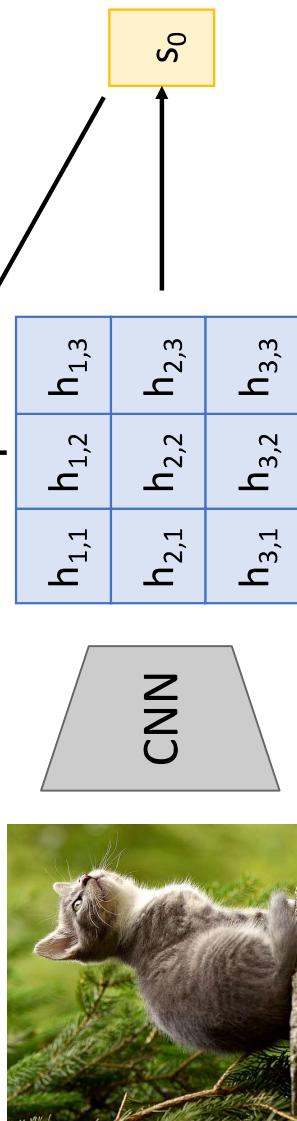
October 23, 2019

Image Captioning with RNNs and Attention

Alignment scores

$e_{1,1,1}$	$e_{1,1,2}$	$e_{1,1,3}$
$e_{1,2,1}$	$e_{1,2,2}$	$e_{1,2,3}$
$e_{1,3,1}$	$e_{1,3,2}$	$e_{1,3,3}$

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



Use a CNN to compute a grid of features for an image

Xu et al., "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Justin Johnson

Lecture 13 - 28

October 23, 2019

Image Captioning with RNNs and Attention

Alignment scores

$e_{1,1,1}$	$e_{1,1,2}$	$e_{1,1,3}$
$e_{1,2,1}$	$e_{1,2,2}$	$e_{1,2,3}$
$e_{1,3,1}$	$e_{1,3,2}$	$e_{1,3,3}$

$$e_{t,i,j} = f_{\text{att}}(s_{t-1}, h_{i,j})$$
$$a_{t,:,:} = \text{softmax}(e_{t,:,:})$$

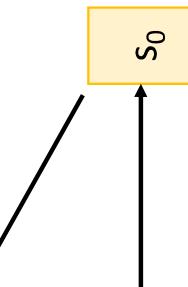
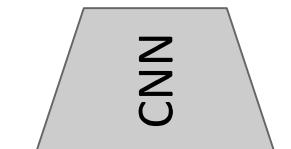
softmax



Attention weights

$a_{1,1,1}$	$a_{1,1,2}$	$a_{1,1,3}$
$a_{1,2,1}$	$a_{1,2,2}$	$a_{1,2,3}$
$a_{1,3,1}$	$a_{1,3,2}$	$a_{1,3,3}$

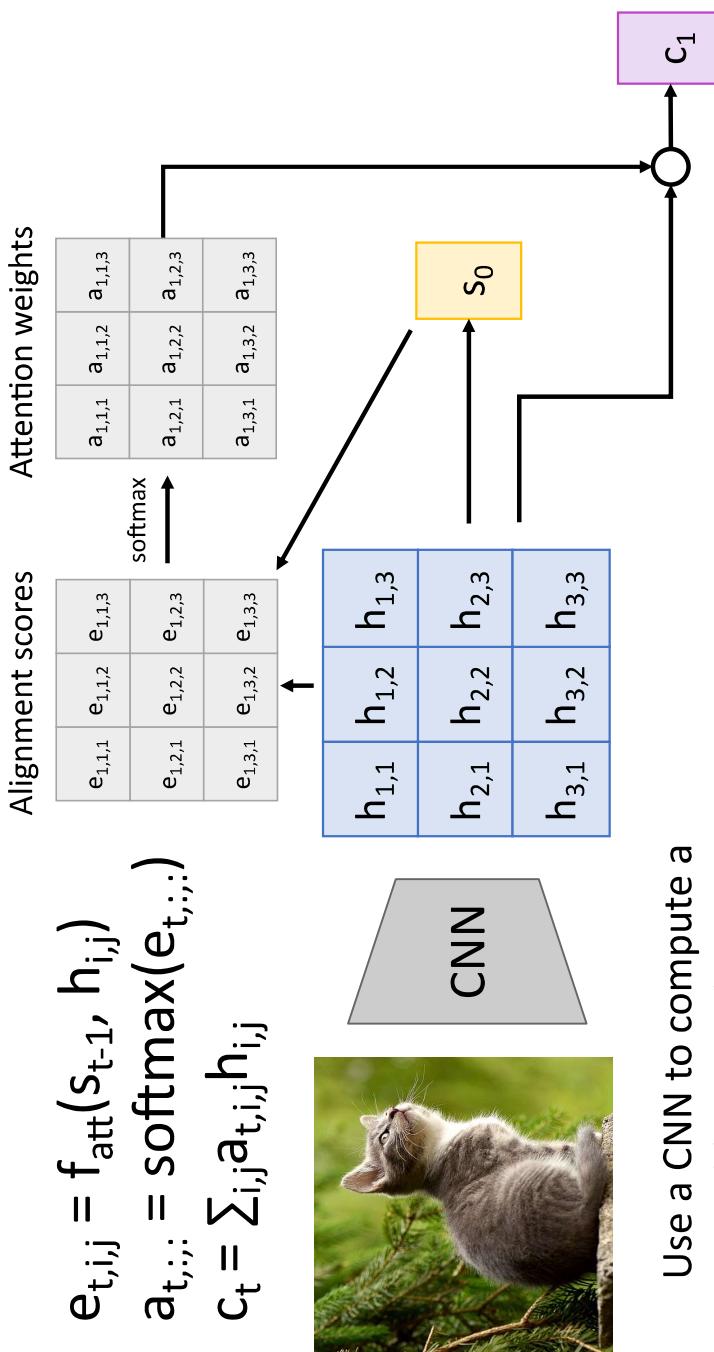
$h_{1,1}$	$h_{1,2}$	$h_{1,3}$
$h_{2,1}$	$h_{2,2}$	$h_{2,3}$
$h_{3,1}$	$h_{3,2}$	$h_{3,3}$



Use a CNN to compute a grid of features for an image

Xu et al., "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

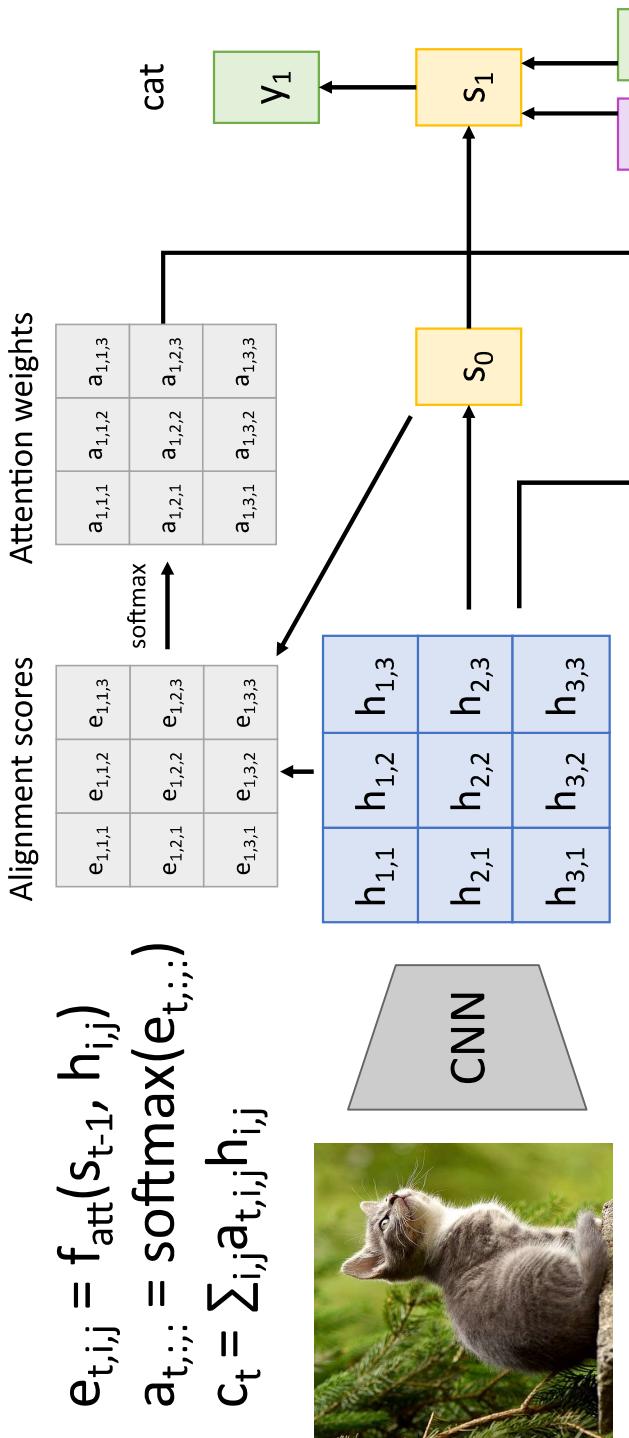
Image Captioning with RNNs and Attention



Use a CNN to compute a grid of features for an image

Xu et al., "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Image Captioning with RNNs and Attention



Use a CNN to compute a grid of features for an image

Xu et al., "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Justin Johnson

Lecture 13 - 31

October 23, 2019

Image Captioning with RNNs and Attention

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

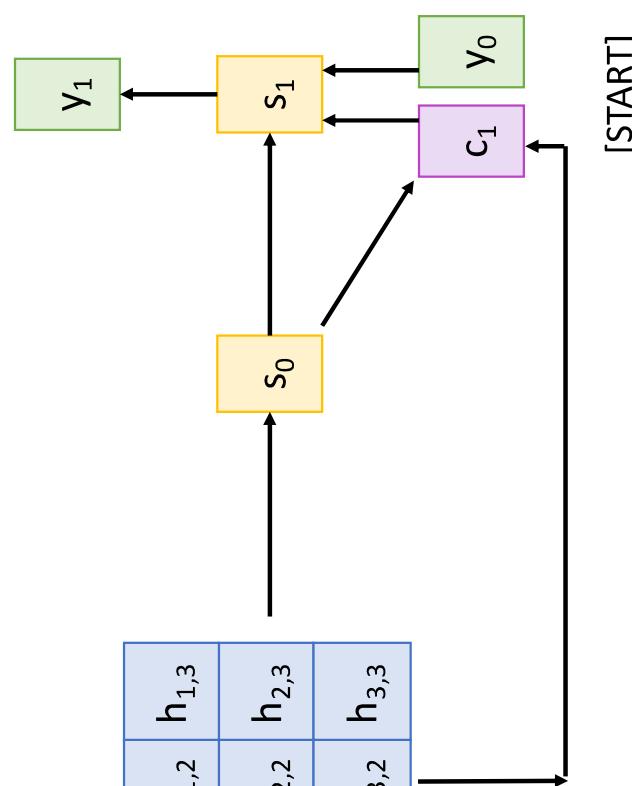
$$a_{t,:,:} = \text{softmax}(e_{t,:,:})$$

$$c_t = \sum_{i,j} a_{t,i,j} h_{i,j}$$



$h_{1,1}$	$h_{1,2}$	$h_{1,3}$
$h_{2,1}$	$h_{2,2}$	$h_{2,3}$
$h_{3,1}$	$h_{3,2}$	$h_{3,3}$

cat



Use a CNN to compute a
grid of features for an image

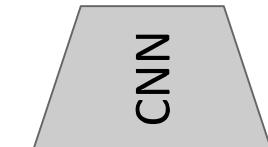
Xu et al., "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Image Captioning with RNNs and Attention

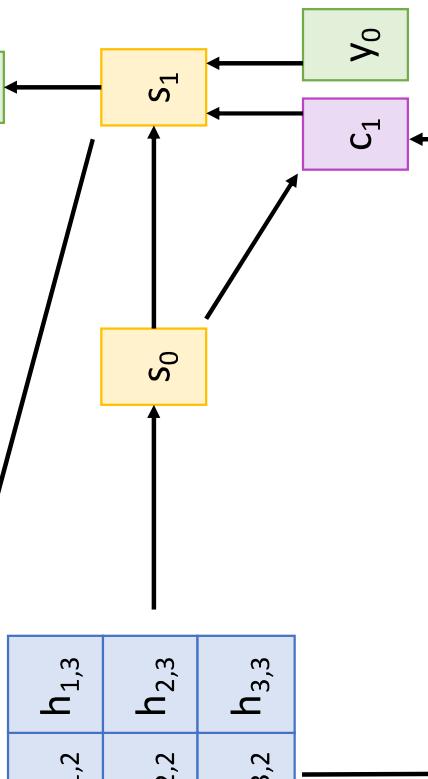
Alignment scores

$e_{2,1,1}$	$e_{2,1,2}$	$e_{2,1,3}$
$e_{2,2,1}$	$e_{2,2,2}$	$e_{2,2,3}$
$e_{2,3,1}$	$e_{2,3,2}$	$e_{2,3,3}$

$$\begin{aligned} e_{t,i,j} &= f_{\text{att}}(s_{t-1}, h_{i,j}) \\ a_{t,:,:} &= \text{softmax}(e_{t,:,:}) \\ c_t &= \sum_{i,j} a_{t,i,j} h_{i,j} \end{aligned}$$



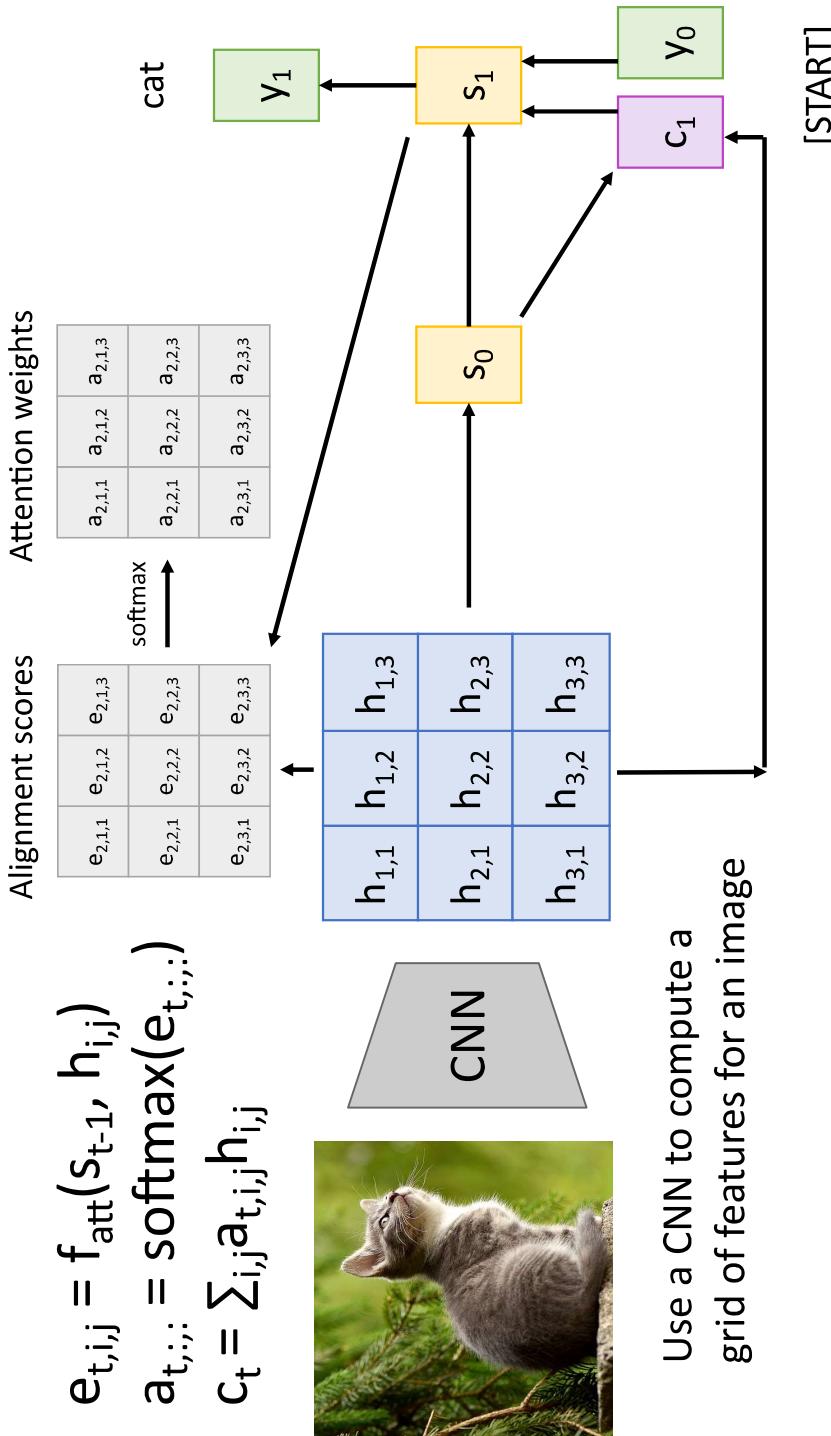
$h_{1,1}$	$h_{1,2}$	$h_{1,3}$
$h_{2,1}$	$h_{2,2}$	$h_{2,3}$
$h_{3,1}$	$h_{3,2}$	$h_{3,3}$



Use a CNN to compute a grid of features for an image

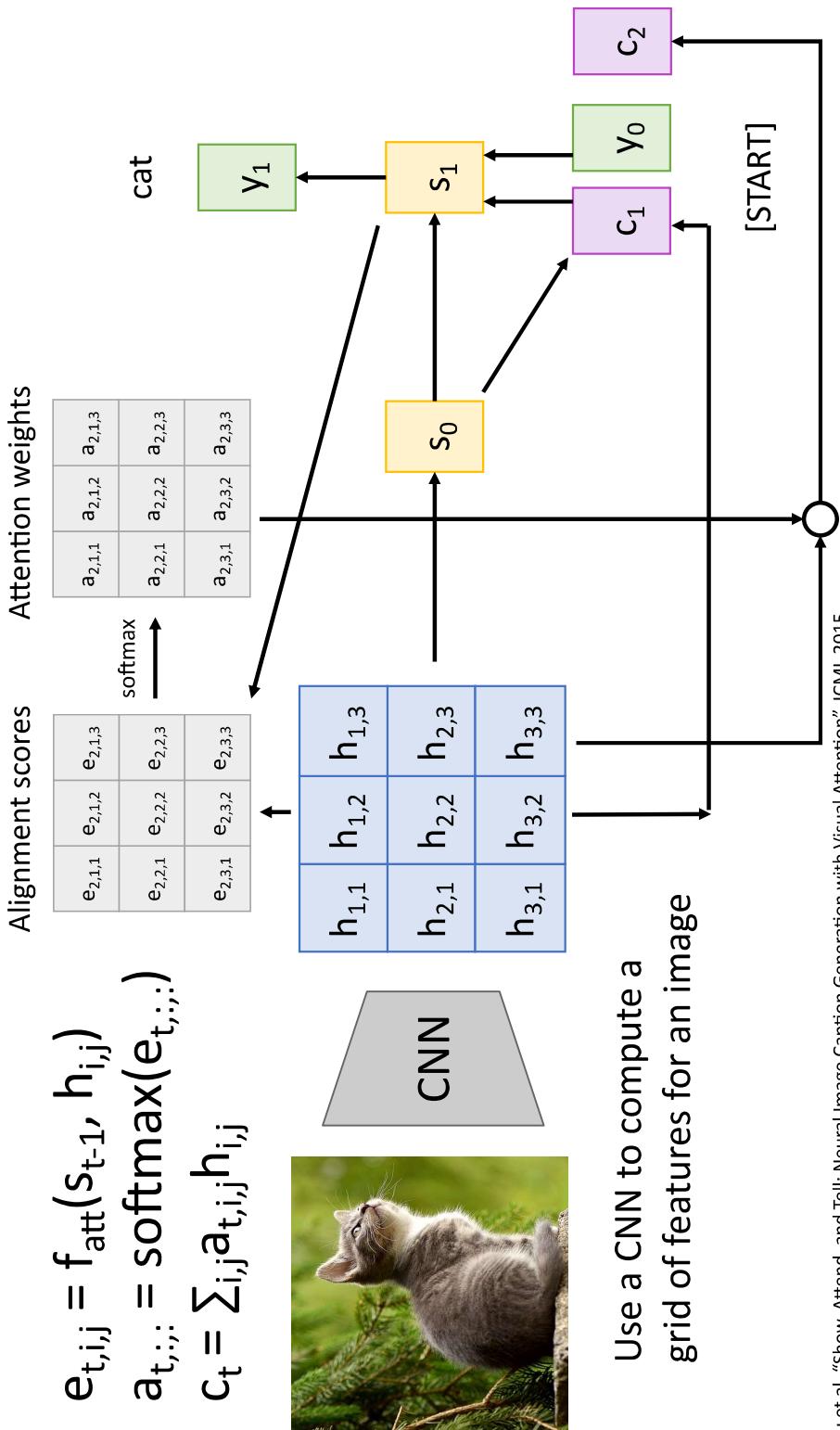
Xu et al., "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Image Captioning with RNNs and Attention



Xu et al., "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Image Captioning with RNNs and Attention



Xu et al., "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Image Captioning with RNNs and Attention

Alignment scores

$e_{2,1,1}$	$e_{2,1,2}$	$e_{2,1,3}$
$e_{2,2,1}$	$e_{2,2,2}$	$e_{2,2,3}$
$e_{2,3,1}$	$e_{2,3,2}$	$e_{2,3,3}$

softmax

Attention weights

sitting

cat

$a_{2,1,1} \quad a_{2,1,2} \quad a_{2,1,3}$

$a_{2,2,1} \quad a_{2,2,2} \quad a_{2,2,3}$

$a_{2,3,1} \quad a_{2,3,2} \quad a_{2,3,3}$

$h_{1,1}$	$h_{1,2}$	$h_{1,3}$
$h_{2,1}$	$h_{2,2}$	$h_{2,3}$
$h_{3,1}$	$h_{3,2}$	$h_{3,3}$

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

$$a_{t,:,:} = \text{softmax}(e_{t,:,:})$$

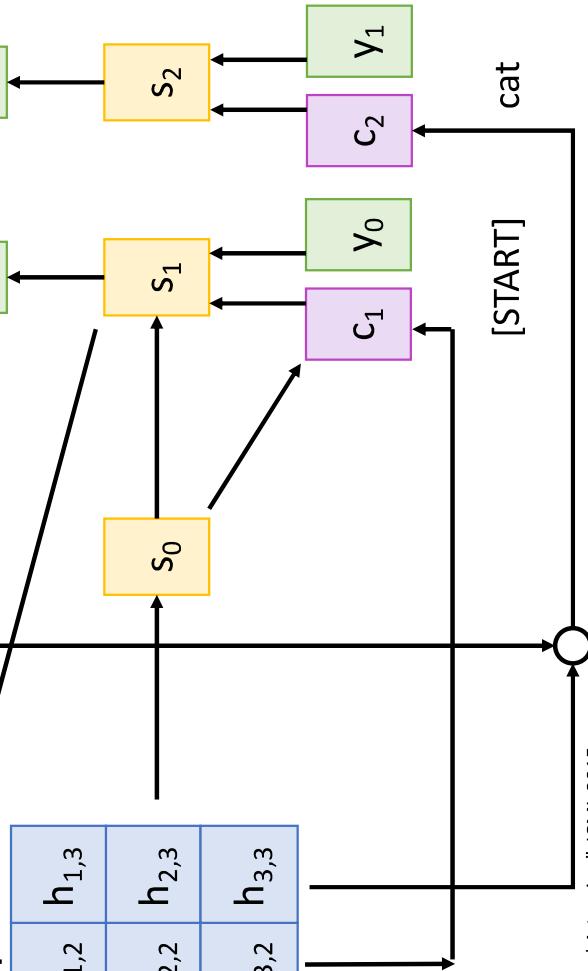
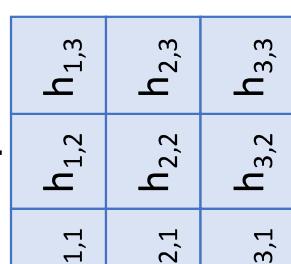
$$c_t = \sum_{i,j} a_{t,i,j} h_{i,j}$$



CNN



Use a CNN to compute a grid of features for an image



Xu et al., "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Image Captioning with RNNs and Attention

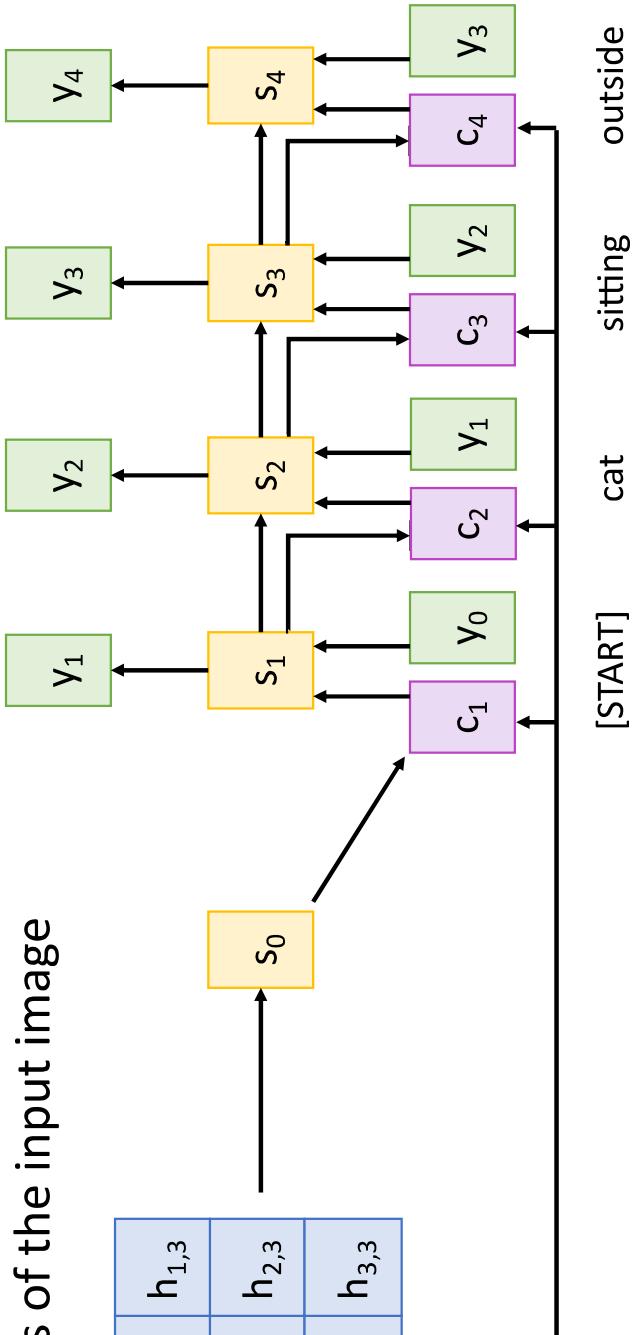
$$\begin{aligned} e_{t,i,j} &= f_{\text{att}}(s_{t-1}, h_{i,j}) \\ a_{t,:,:} &= \text{softmax}(e_{t,:,:}) \\ c_t &= \sum_{i,j} a_{t,i,j} h_{i,j} \end{aligned}$$

Each timestep of decoder
uses a different context
vector that looks at different
parts of the input image



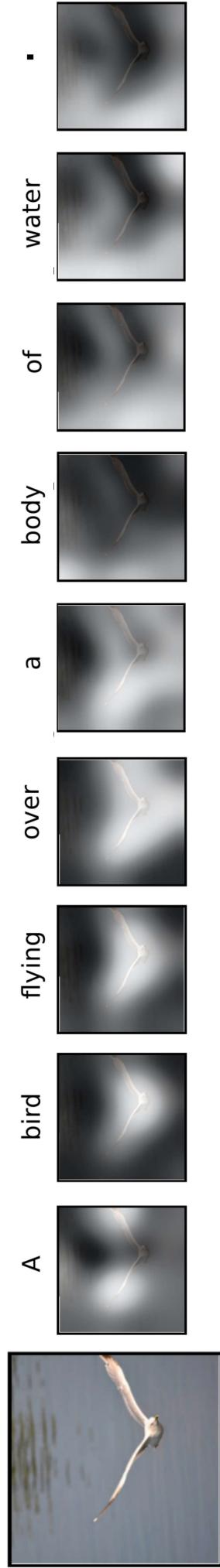
$h_{1,1}$	$h_{1,2}$	$h_{1,3}$
$h_{2,1}$	$h_{2,2}$	$h_{2,3}$
$h_{3,1}$	$h_{3,2}$	$h_{3,3}$

Use a CNN to compute a
grid of features for an image



Xu et al., "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Image Captioning with RNNs and Attention



Xu et al., "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Justin Johnson

Lecture 13 - 38

October 23, 2019

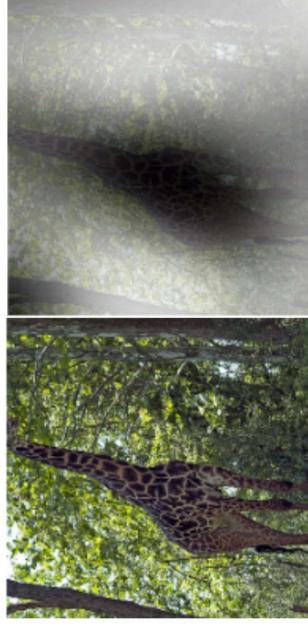
Image Captioning with RNNs and Attention



A dog is standing on a hardwood floor.



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

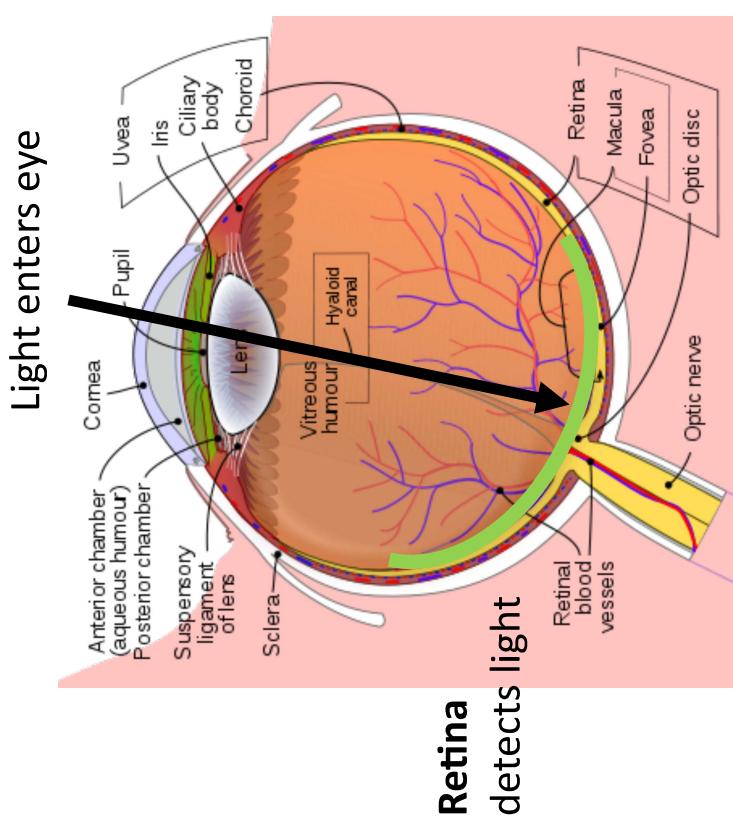
Xu et al., "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Justin Johnson

Lecture 13 - 39

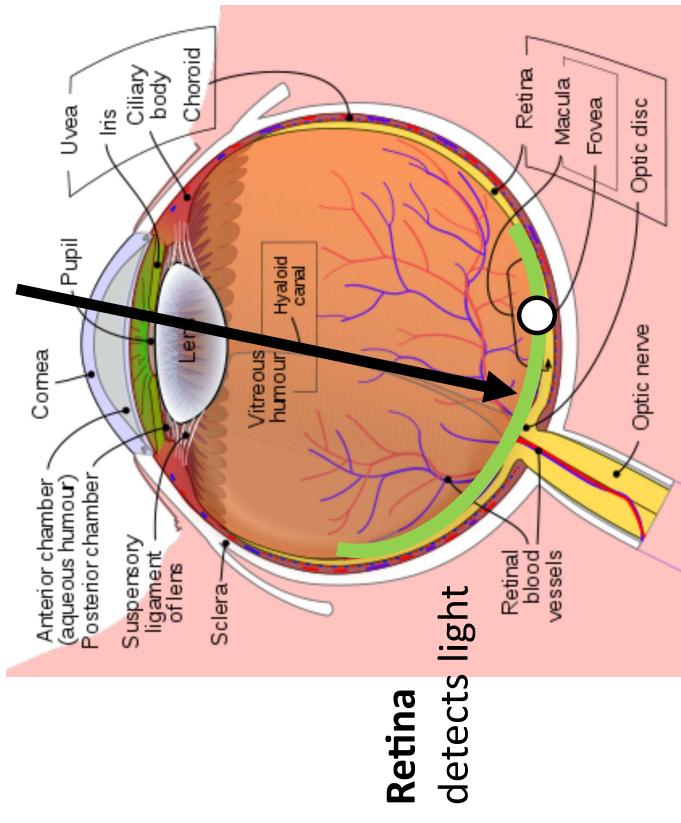
October 23, 2019

Human Vision: Fovea

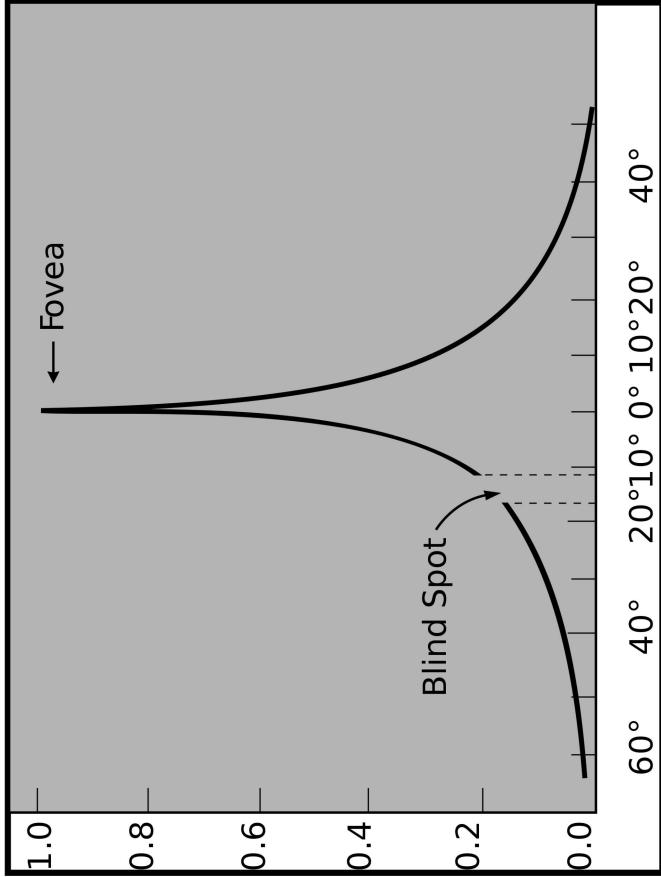


Human Vision: Fovea

Light enters eye



The **fovea** is a tiny region of the retina that can see with high acuity



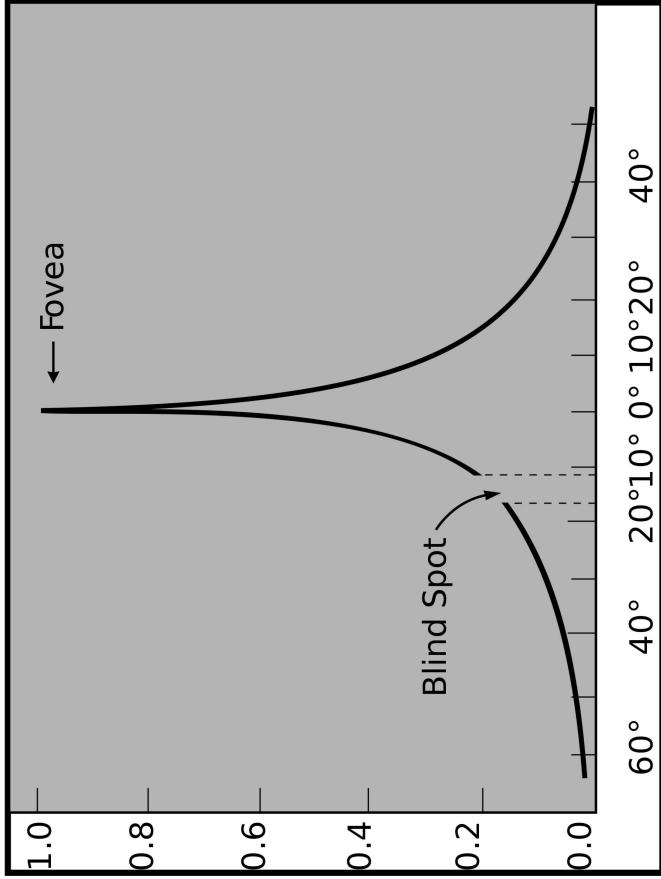
Eye image is licensed under [CC ASA 3.0 Unported](#) (added black arrow, green arc, and white circle)

Acuity graph is licensed under [CC ASA 3.0 Unported](#) (No changes made)

Human Vision: Saccades

Human eyes are constantly moving so we don't notice

The **fovea** is a tiny region of the retina that can see with high acuity



[Saccade video](#) is licensed under [CC A-SA 4.0 International](#) (no changes made)

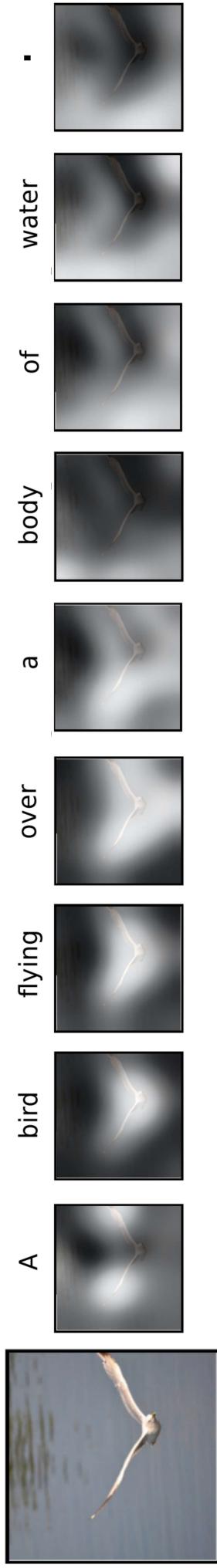
[Acuity graph](#) is licensed under [CC A-SA 3.0 Unported](#) (No changes made)

Justin Johnson

Lecture 13 - 42

October 23, 2019

Image Captioning with RNNs and Attention



Attention weights at each
timestep kind of like
saccades of human eye

Xu et al., "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Saccade video is licensed under CC ASA 4.0 International (no changes made)

X , Attend, and Y

“Show, attend, and tell” (*Xu et al, ICML 2015*)

Look at image, attend to image regions, produce question

“Ask, attend, and answer” (*Xu and Saenko, ECCV 2016*)

“Show, ask, attend, and answer” (*Kazemi and Elqursh, 2017*)

Read text of question, attend to image regions, produce answer

“Listen, attend, and spell” (*Chan et al, ICASSP 2016*)

Process raw audio, attend to audio regions while producing text

“Listen, attend, and walk” (*Mei et al, AAAI 2016*)

Process text, attend to text regions, output navigation commands

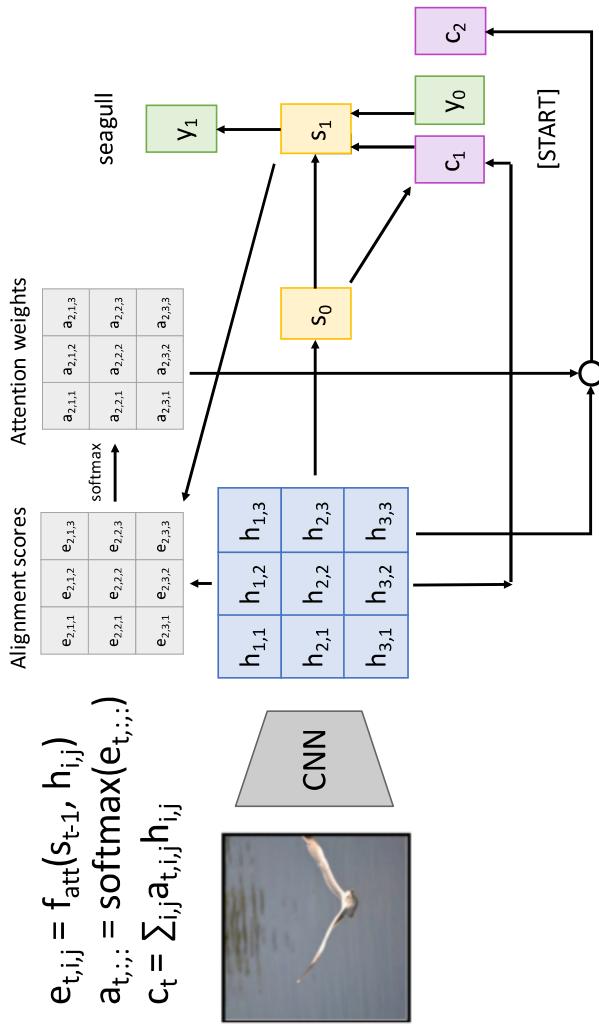
“Show, attend, and interact” (*Qureshi et al, ICRA 2017*)

Process image, attend to image regions, output robot control commands

“Show, attend, and read” (*Li et al, AAAI 2019*)

Process image, attend to image regions, output text

Attention Layer



Computation:

Similarities: e (Shape: N_x) $e_i = f_{att}(\mathbf{q}, \mathbf{X}_i)$

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

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

Attention Layer

$$\begin{aligned}
 e_{t,j} &= f_{\text{att}}(s_{t-1}, h_{i,j}) \\
 a_{t,:} &= \text{softmax}(e_{t,:}) \\
 c_t &= \sum_{i,j} a_{t,i,j} h_{i,j}
 \end{aligned}$$

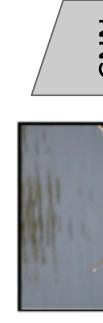
Inputs:

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

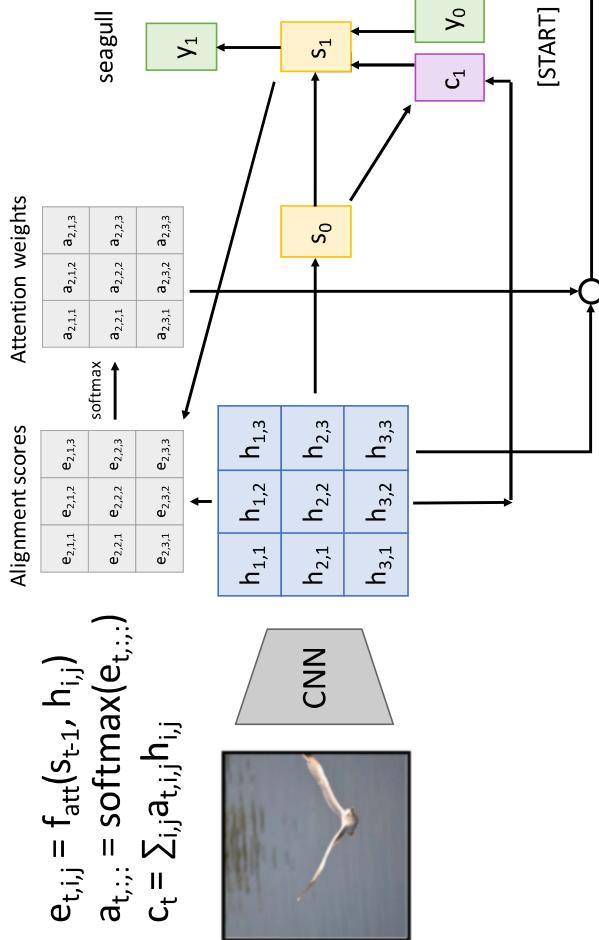
Input vectors: \mathbf{X} (Shape: $N_x \times D_Q$)

Similarity function: **dot product**

Outputs:



CNN



Computation:

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

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

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

Changes:

- Use dot product for similarity

Attention Layer

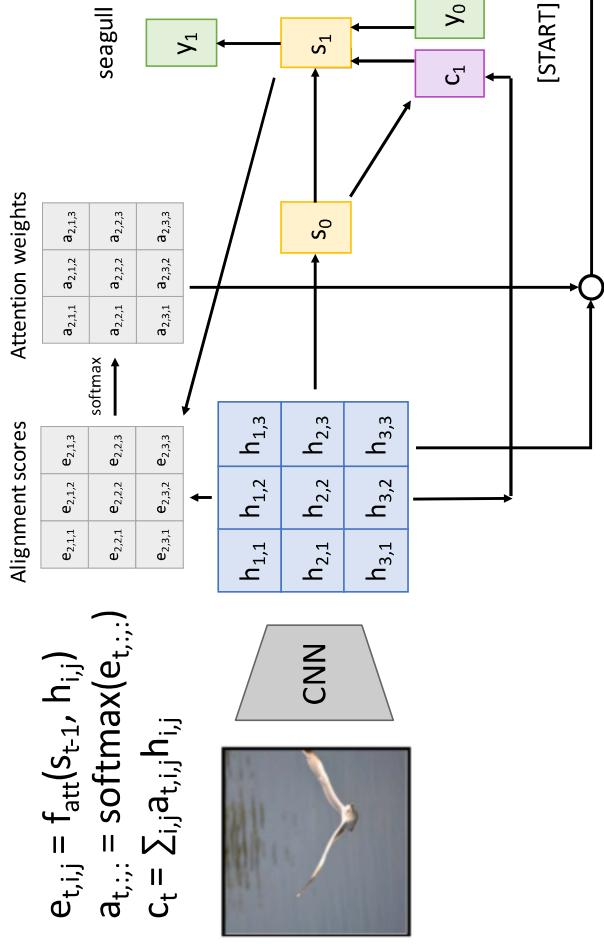
$$\begin{aligned}
 e_{t,i,j} &= f_{\text{att}}(s_{t-1}, h_{i,j}) \\
 a_{t,:,:} &= \text{softmax}(e_{t,:,:}) \\
 c_t &= \sum_{i,j} a_{t,i,j} h_{i,j}
 \end{aligned}$$

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: e (Shape: N_x) $e_i = \mathbf{q} \cdot \mathbf{x}_i / \sqrt{D_Q}$

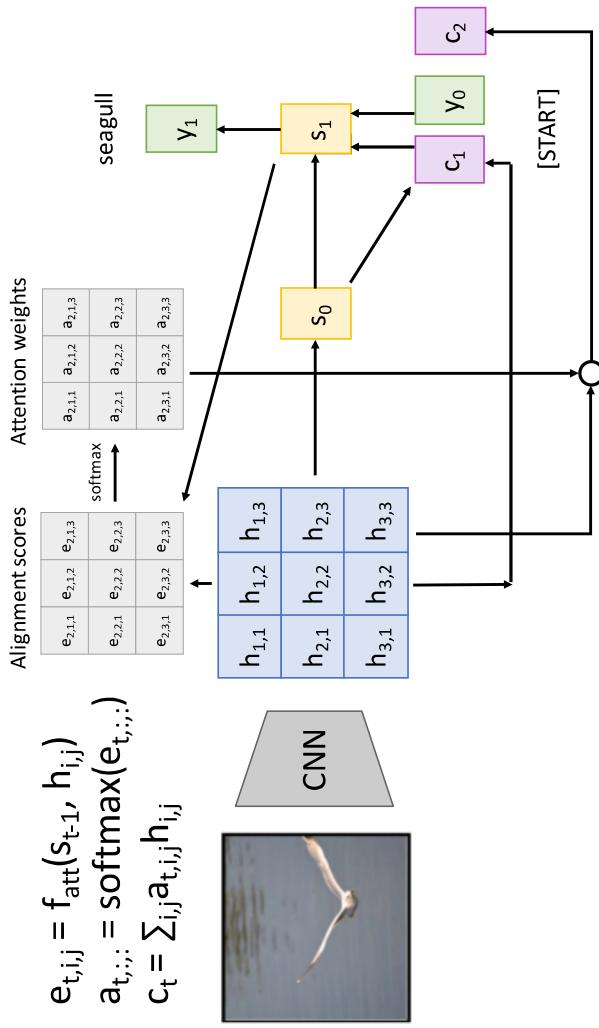
Attention weights: $a = \text{softmax}(e)$ (Shape: N_x)

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

Changes:

- Use **scaled** dot product for similarity

Attention Layer



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

October 23, 2019

Lecture 13 - 48

Justin Johnson

Attention Layer

$$\begin{aligned} e_{t,j} &= f_{\text{att}}(s_{t-1}, h_{i,j}) \\ a_{t,:} &= \text{softmax}(e_{t,:}) \\ c_t &= \sum_{i,j} a_{t,i,j} h_{i,j} \end{aligned}$$

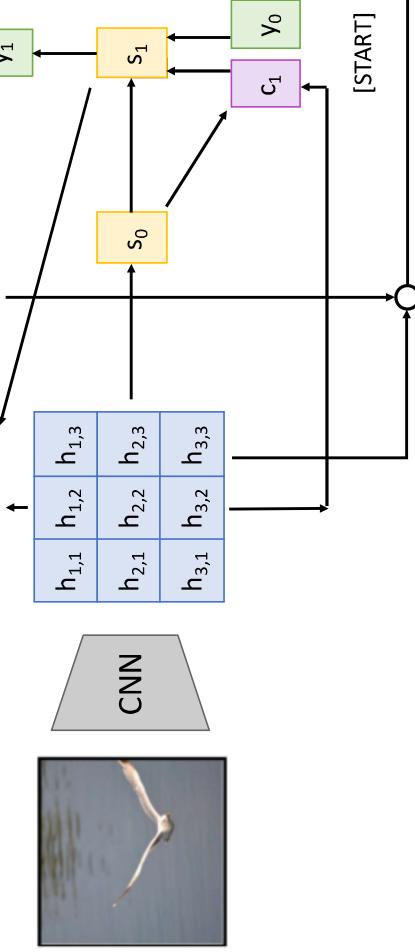
Inputs:

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

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

Alignment scores

Attention weights		
$e_{2,1,1}$	$e_{2,1,2}$	$e_{2,1,3}$
$e_{2,2,1}$	$e_{2,2,2}$	$e_{2,2,3}$
$e_{2,3,1}$	$e_{2,3,2}$	$e_{2,3,3}$



Computation:

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

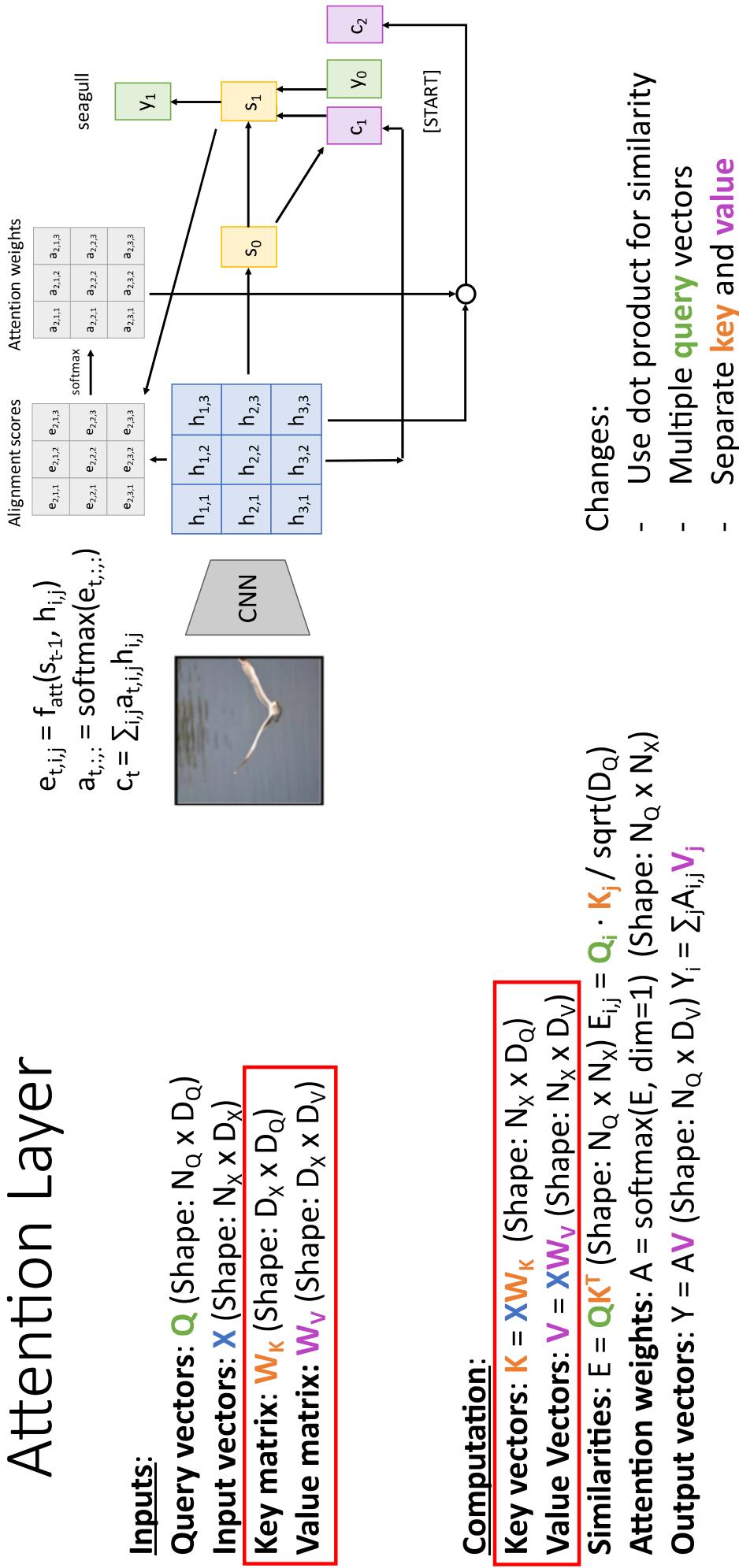
Attention weights: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_Q \times N_X$)

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

Changes:

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

Attention Layer



Changes:

- Use dot product for similarity
- 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:

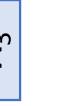
Key vectors: $\mathbf{K} = \mathbf{X} \mathbf{W_K}$ (Shape: $N_x \times D_Q$)



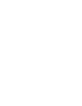
Value Vectors: $\mathbf{V} = \mathbf{X} \mathbf{W_V}$ (Shape: $N_x \times D_V$)



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



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 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:

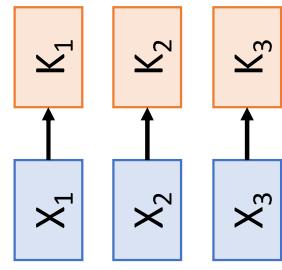
Key vectors: $\mathbf{K} = \mathbf{X} \mathbf{W_K}$ (Shape: $N_x \times D_Q$)

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

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

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

Output vectors: $\mathbf{Y} = A \mathbf{V}$ (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:

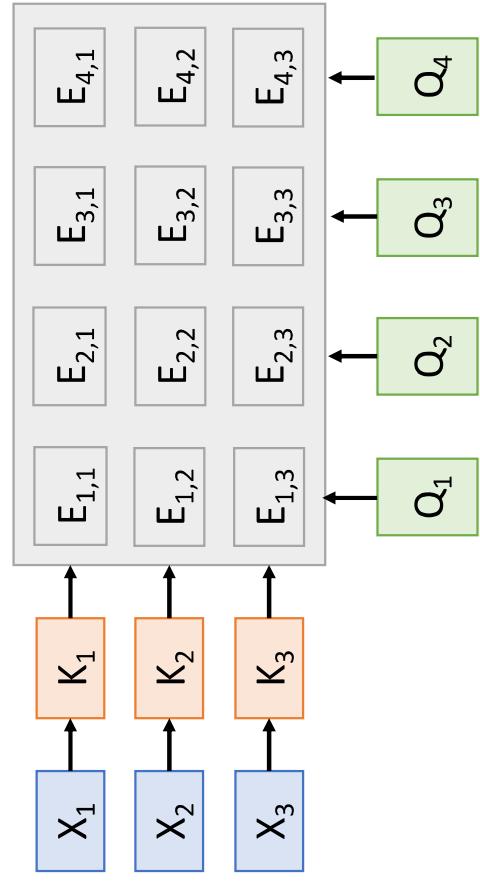
Key vectors: $\mathbf{K} = \mathbf{X} \mathbf{W}_K$ (Shape: $N_x \times D_Q$)

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

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

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

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



Attention Layer

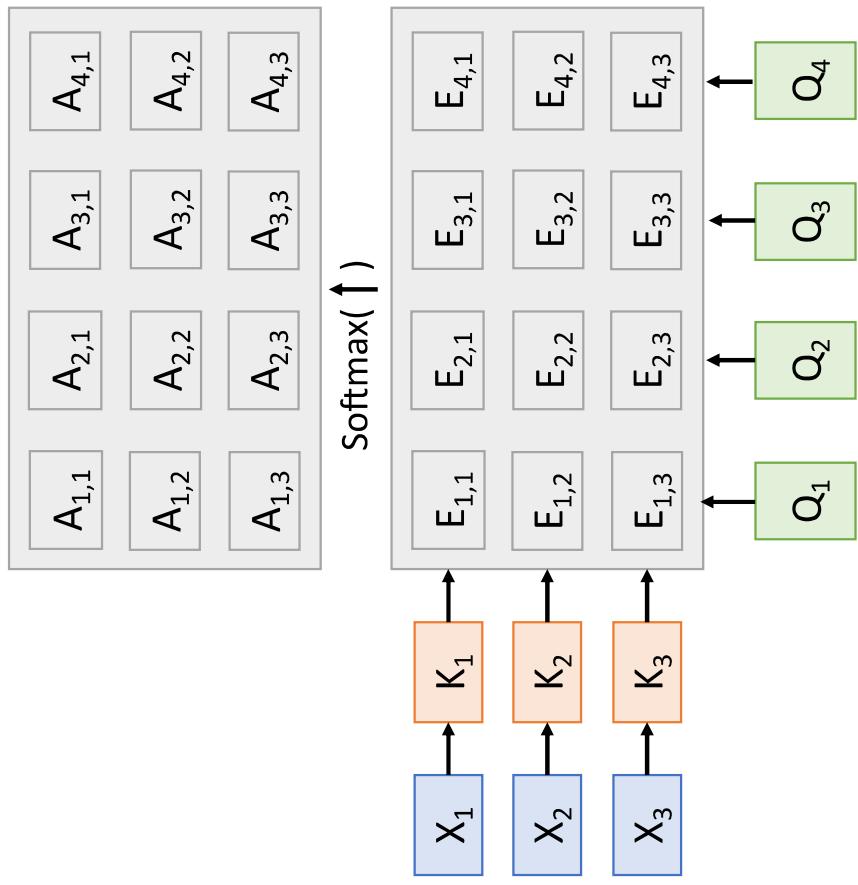
Inputs:

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: $\mathbf{K} = \mathbf{X}\mathbf{W}_K$ (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: $N_Q \times N_x$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

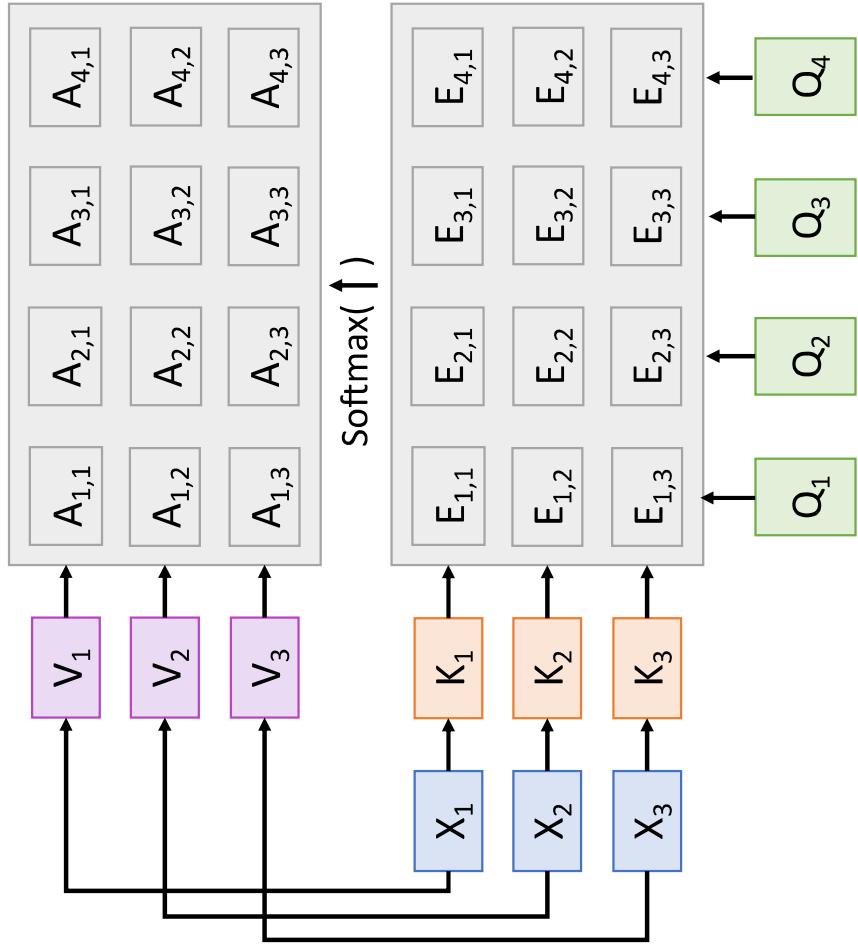
Attention weights: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_Q \times N_x$)

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

Attention Layer

Inputs:

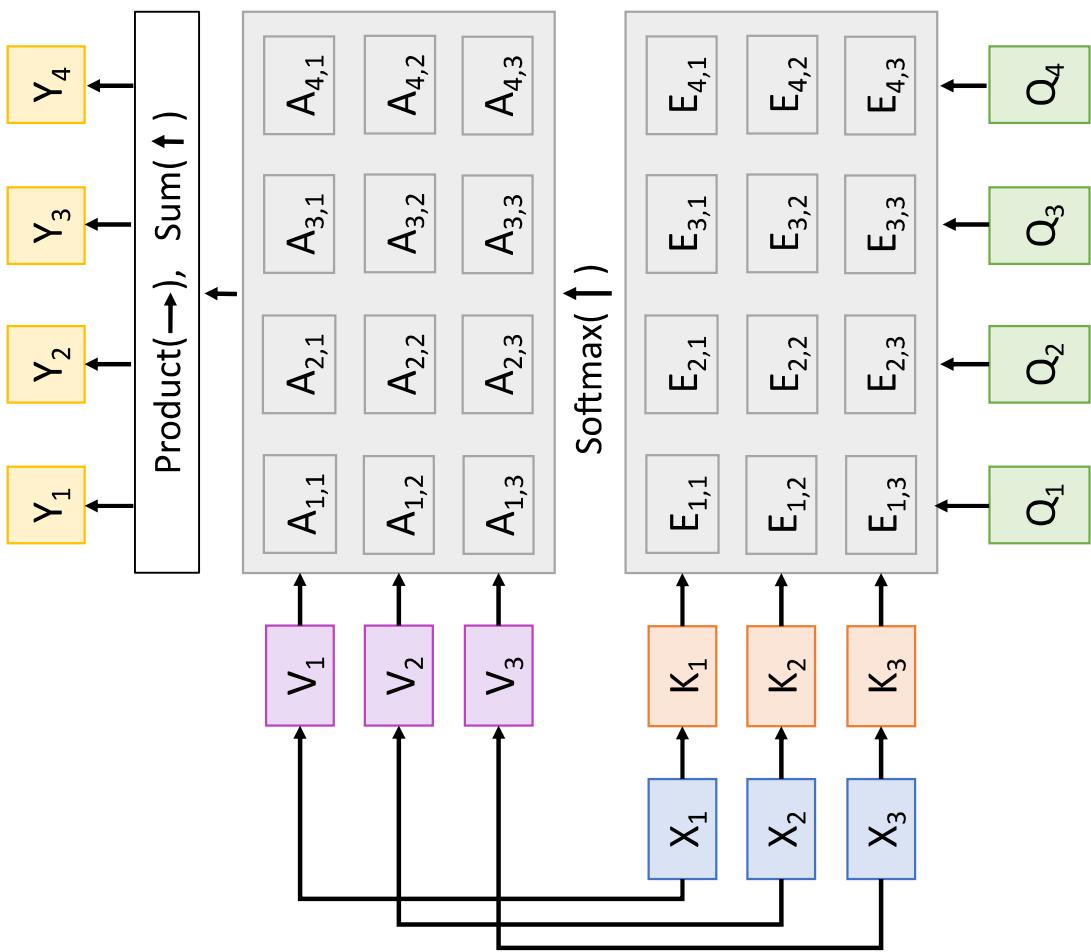
- Query vectors: \mathbf{Q} (Shape: $N_Q \times D_Q$)
- Input vectors: \mathbf{X} (Shape: $N_x \times D_x$)
- Key matrix: \mathbf{W}_K (Shape: $D_x \times D_Q$)
- Value matrix: \mathbf{W}_V (Shape: $D_x \times D_V$)



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: $N_Q \times N_x$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$
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 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: $N_Q \times D_V$) $Y_i = \sum_j A_{i,j} V_j$

Self-Attention Layer

One **query** per **input vector**

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: $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q$

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

Attention weights: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_x \times N_x$)
Output vectors: $\mathbf{Y} = A\mathbf{V}$ (Shape: $N_x \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$



Self-Attention Layer

One **query** per **input vector**

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: $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q$

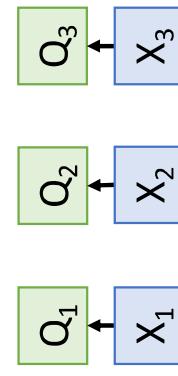
Key vectors: $\mathbf{K} = \mathbf{X}\mathbf{W}_K$ (Shape: $N_x \times D_Q$)

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

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

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

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



Self-Attention Layer

One **query** per **input vector**

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: $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q$

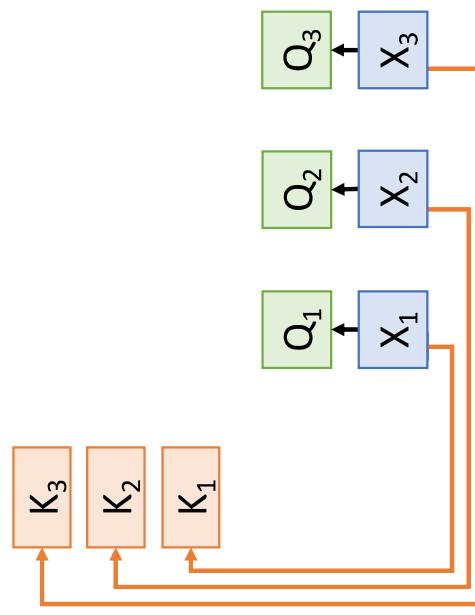
Key vectors: $\mathbf{K} = \mathbf{X}\mathbf{W}_K$ (Shape: $N_x \times D_Q$)

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

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

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

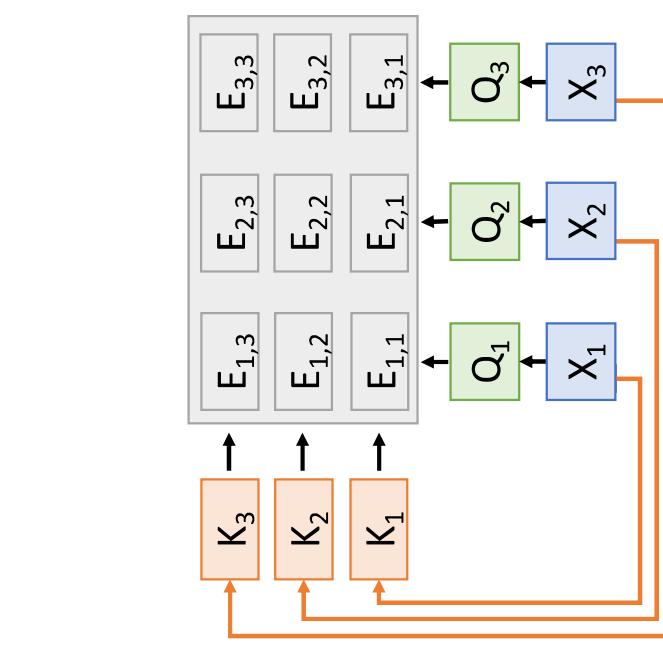
Output vectors: $\mathbf{Y} = A\mathbf{V}$ (Shape: $N_x \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$



Self-Attention Layer

One **query** per **input vector**

Inputs:



Computation:

- Query vectors: $Q = XW_Q$ (Shape: $N_x \times D_Q$)
- 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: $N_x \times D_V$) $Y_i = \sum_j A_{i,j} V_j$

Self-Attention Layer

One **query** per **input vector**

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: $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q$

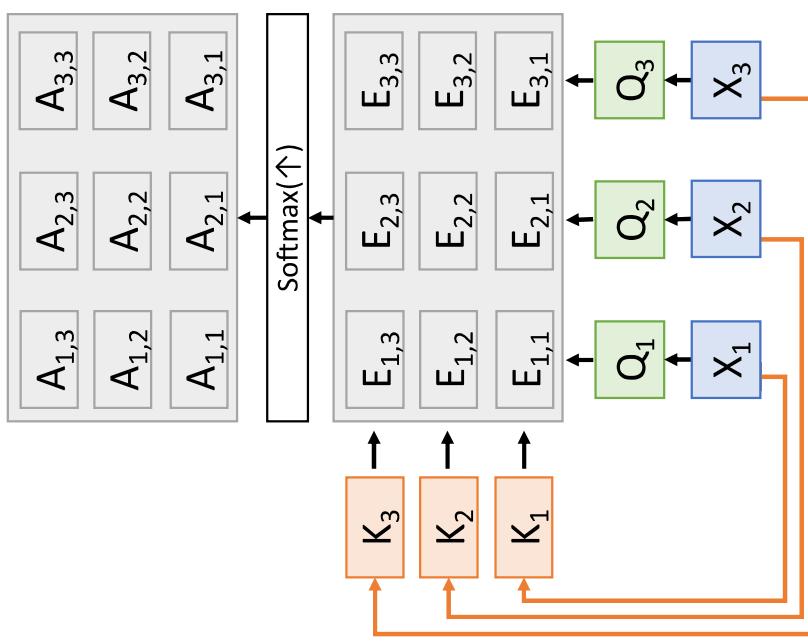
Key vectors: $\mathbf{K} = \mathbf{X}\mathbf{W}_K$ (Shape: $N_x \times D_Q$)

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

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

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

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



Self-Attention Layer

One **query** per **input vector**

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: $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q$

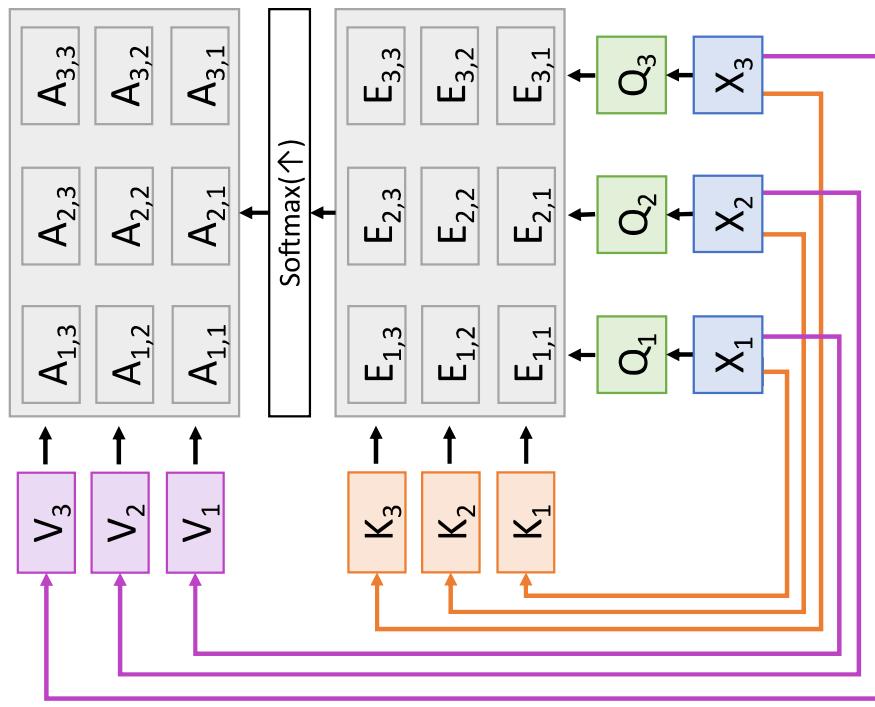
Key vectors: $\mathbf{K} = \mathbf{X}\mathbf{W}_K$ (Shape: $N_x \times D_Q$)

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

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

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

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



Self-Attention Layer

One **query** per **input vector**

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: $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q$

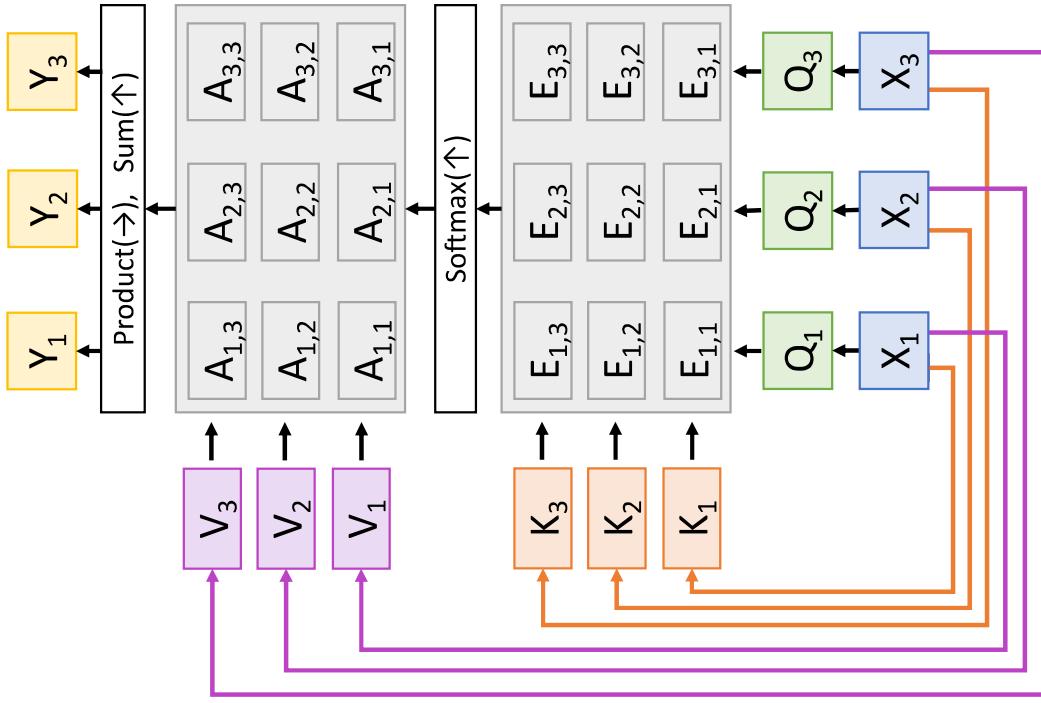
Key vectors: $\mathbf{K} = \mathbf{X}\mathbf{W}_K$ (Shape: $N_x \times D_Q$)

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

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

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

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



Self-Attention Layer

Consider permuting
the input vectors:

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: $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q$

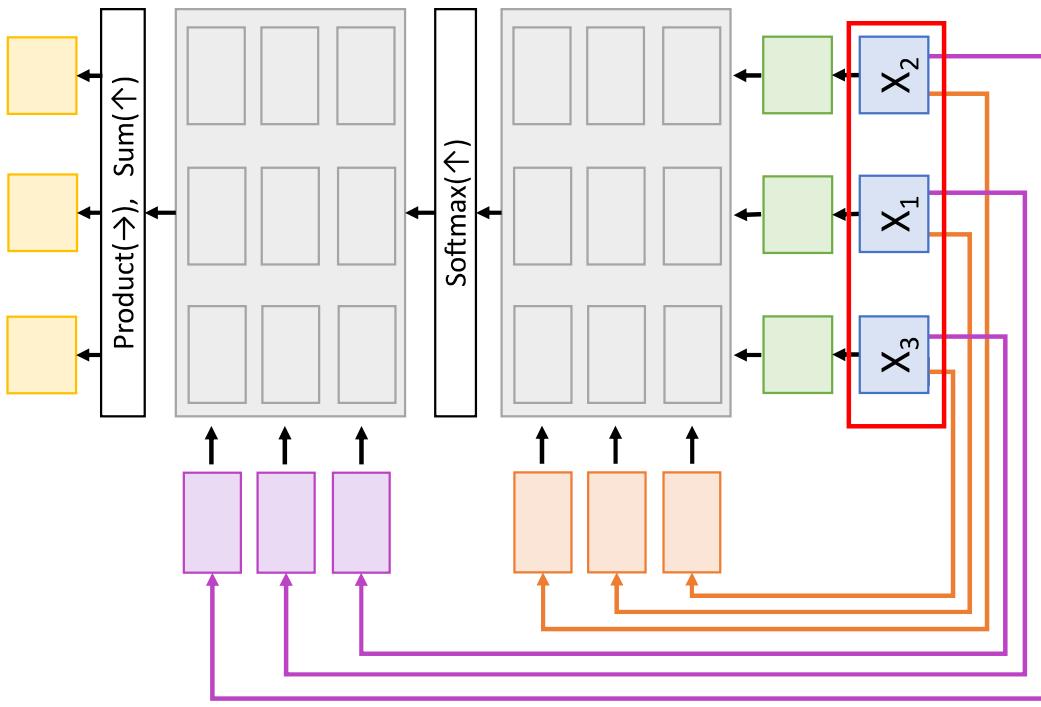
Key vectors: $\mathbf{K} = \mathbf{X}\mathbf{W}_K$ (Shape: $N_x \times D_Q$)

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

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

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

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



Self-Attention Layer

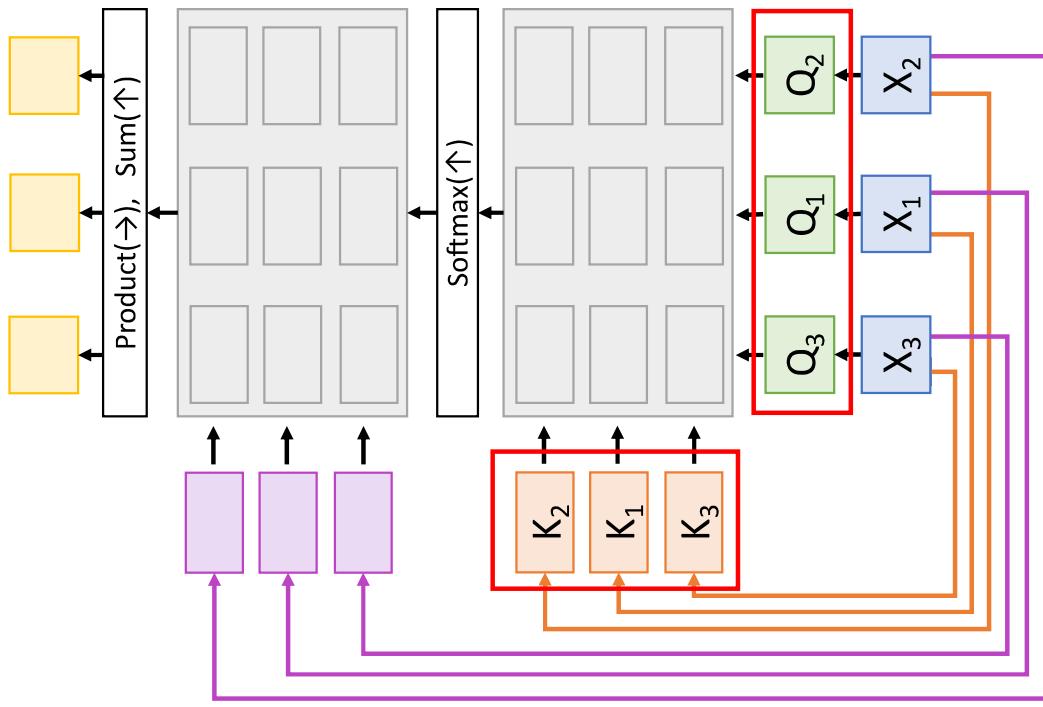
Inputs:
Input vectors: \mathbf{X} (Shape: $N_x \times D_x$)
Key matrix: \mathbf{W}_K (Shape: $D_x \times D_Q$)
Value matrix: \mathbf{W}_V (Shape: $D_x \times D_V$)
Query matrix: \mathbf{W}_Q (Shape: $D_x \times D_Q$)

Computation:

Query vectors: $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q$
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: $N_x \times N_x$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \sqrt{D_Q}$
Attention weights: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_x \times N_x$)
Output vectors: $\mathbf{Y} = A\mathbf{V}$ (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

Consider **permuting**
the input vectors:

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: $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q$

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

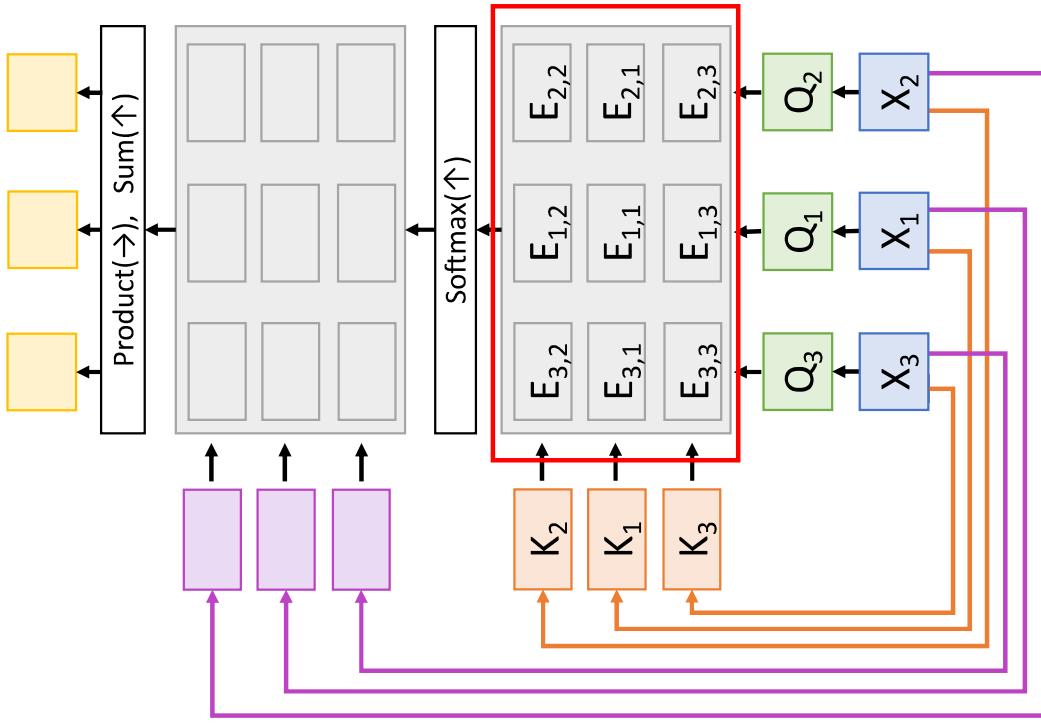
Value vectors: $\mathbf{V} = \mathbf{X}\mathbf{W}_V$ (Shape: $N_x \times D_V$)

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

Similarities: $E = softmax(E, dim=1)$ (Shape: $N_x \times N_x$)

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

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



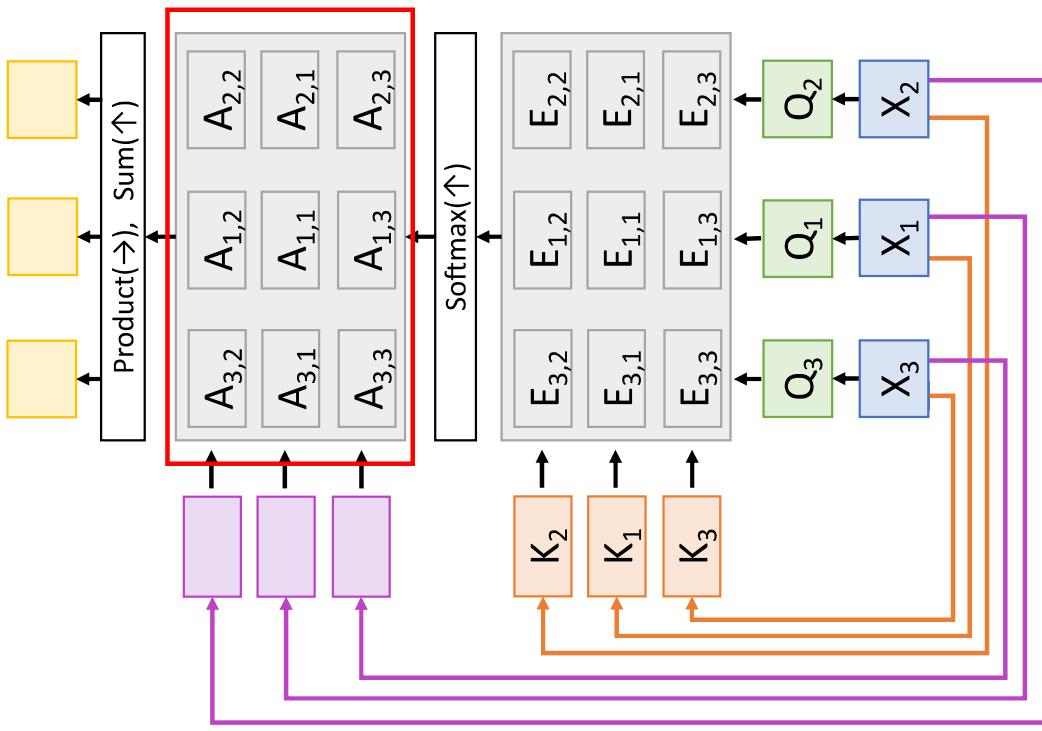
Self-Attention Layer

Consider permuting
the input vectors:

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: $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q$ (Shape: $N_x \times D_Q$)
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: $N_x \times N_x$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \sqrt{D_Q}$
Attention weights: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_x \times N_x$)
Output vectors: $\mathbf{Y} = A\mathbf{V}$ (Shape: $N_x \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$



Self-Attention Layer

Consider permuting
the input vectors:

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: $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q$

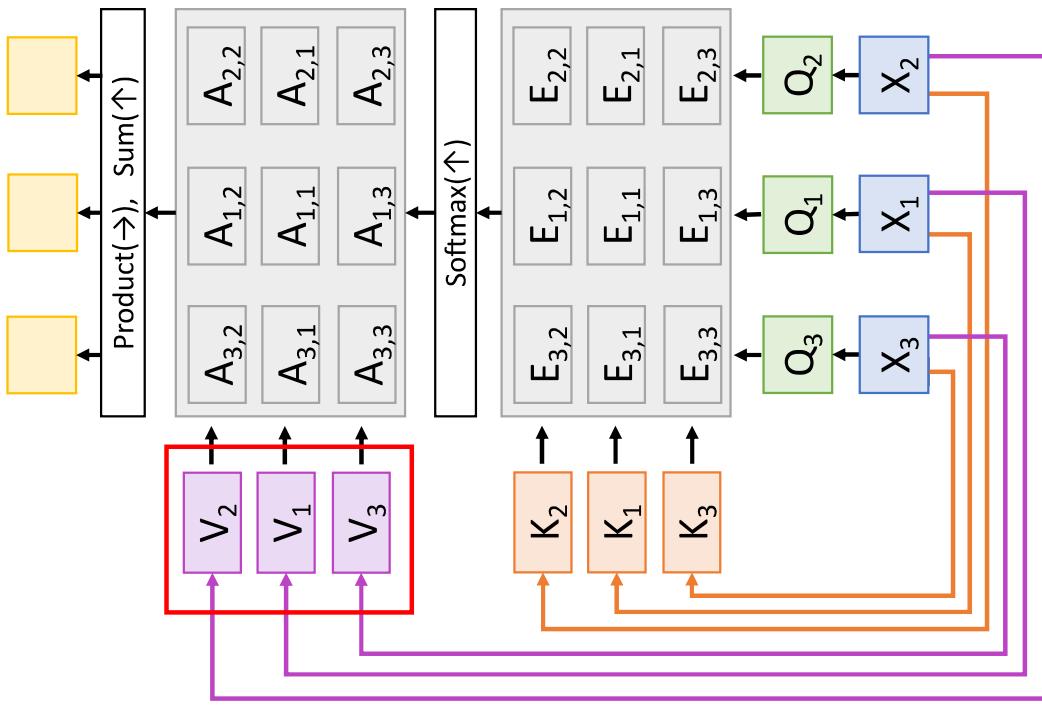
Key vectors: $\mathbf{K} = \mathbf{X}\mathbf{W}_K$ (Shape: $N_x \times D_Q$)

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

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

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

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



Self-Attention Layer

Consider **permuting**
the input vectors:

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:

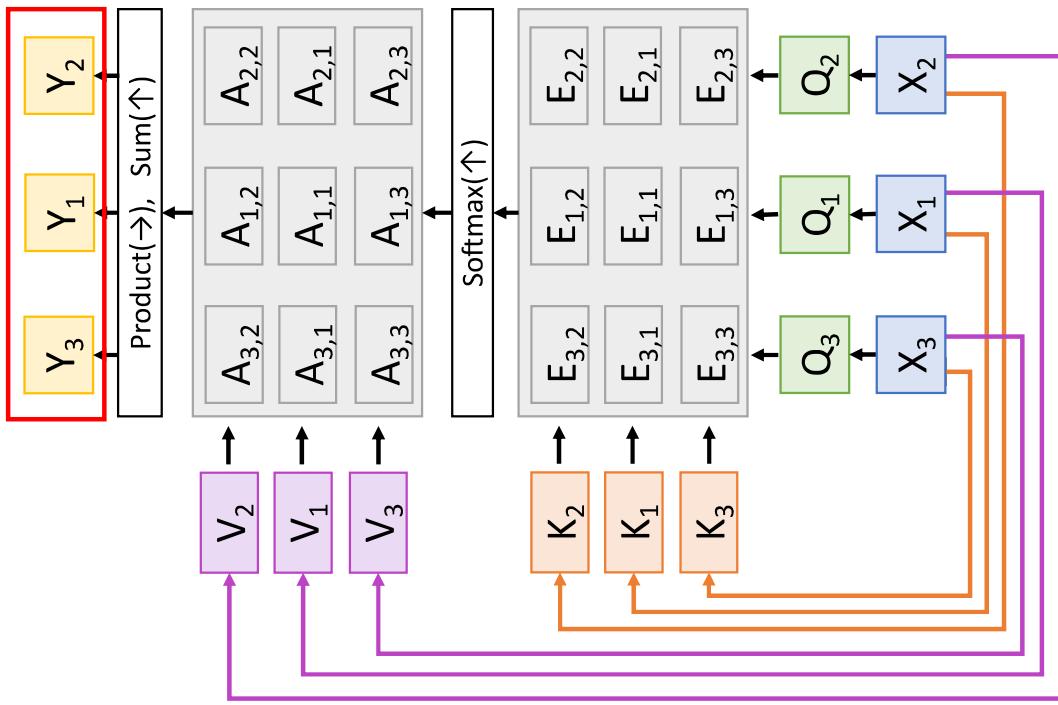
Query vectors: $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q$

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: $N_x \times N_x$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \sqrt{D_Q}$

Attention weights: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_x \times N_x$)
Output vectors: $\mathbf{Y} = A\mathbf{V}$ (Shape: $N_x \times D_V$) $\mathbf{Y}_i = \sum_j A_{i,j} \mathbf{V}_j$



Self-Attention Layer

Consider **permuting**
the input vectors:

Inputs:

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

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

**Outputs will be the
same, but permuted**

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

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

**Self-attention layer is
Permutation Equivariant**
 $f(s(\mathbf{x})) = s(f(\mathbf{x}))$

Computation:

Query vectors: $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q$ (Shape: $N_x \times D_Q$)

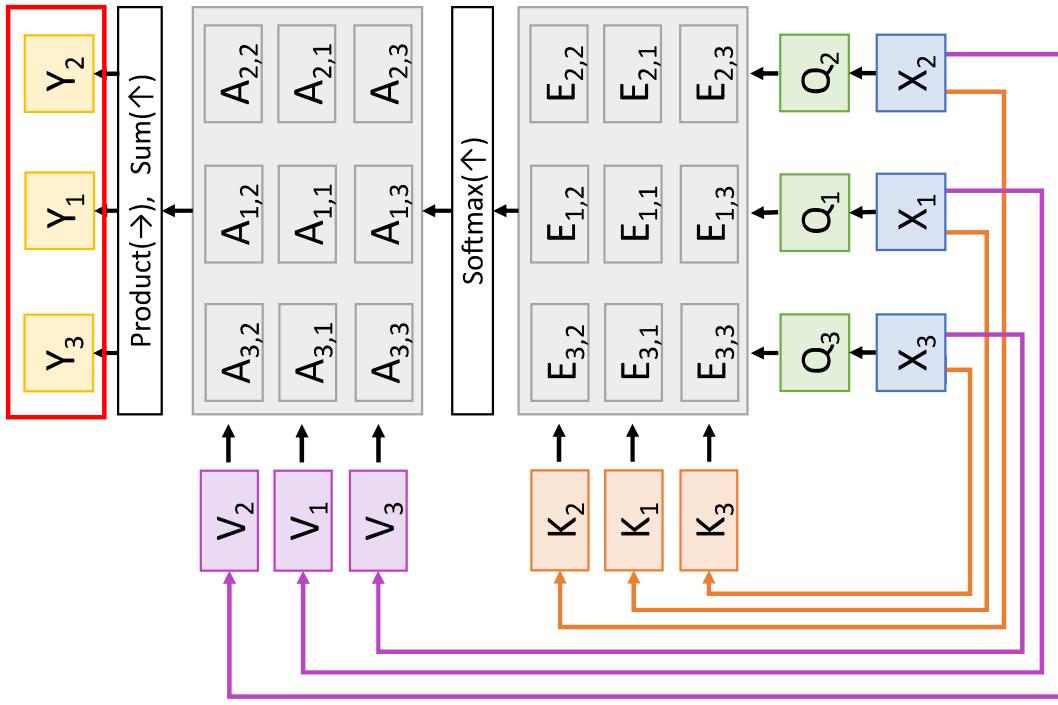
Key vectors: $\mathbf{K} = \mathbf{X}\mathbf{W}_K$ (Shape: $N_x \times D_Q$)

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

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

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

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



Self-Attention Layer

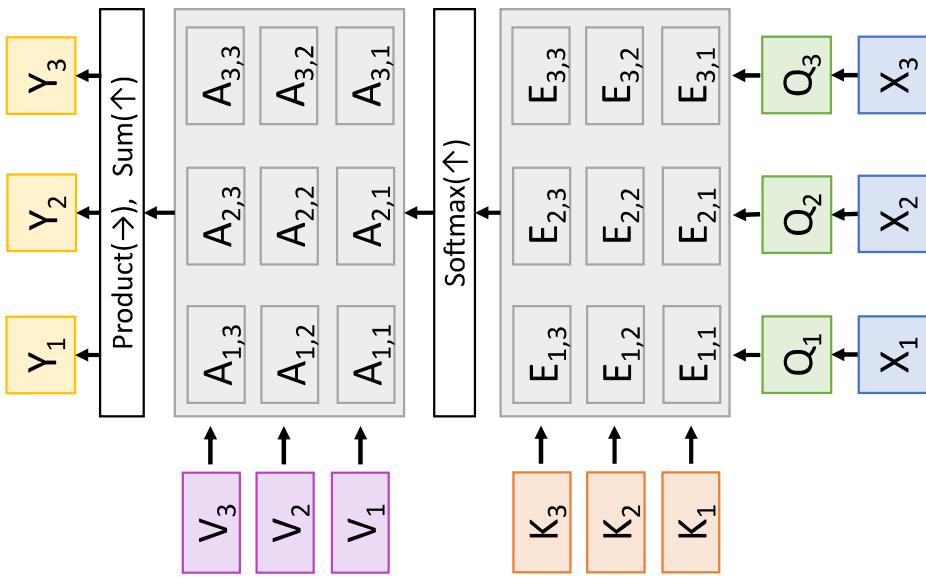
Inputs:
Input vectors: \mathbf{X} (Shape: $N_x \times D_x$)
Key matrix: \mathbf{W}_K (Shape: $D_x \times D_Q$)
Value matrix: \mathbf{W}_V (Shape: $D_x \times D_V$)
Query matrix: \mathbf{W}_Q (Shape: $D_x \times D_Q$)

Computation:

Query vectors: $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q$
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: $N_x \times N_x$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$

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



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: $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q$ (Shape: $N_x \times D_Q$)
 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: $N_x \times N_x$) $E_{i,j} = \mathbf{Q}_i \cdot \mathbf{K}_j / \text{sqrt}(D_Q)$
 Attention weights: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_x \times N_x$)
 Output vectors: $\mathbf{Y} = A\mathbf{V}$ (Shape: $N_x \times D_V$) $\mathbf{Y}_i = \sum_j A_{i,j} \mathbf{V}_j$

Self attention doesn't "know" the order of the vectors it is processing!
 In order to make processing position-aware, concatenate input with **positional encoding**

E can be learned lookup table, or fixed function

Product(\rightarrow), Sum(\uparrow)
 Softmax(\uparrow)
 $\mathbf{Y}_1 \leftarrow A_{1,1} \mathbf{V}_1$
 $\mathbf{Y}_2 \leftarrow A_{2,2} \mathbf{V}_2$
 $\mathbf{Y}_3 \leftarrow A_{3,3} \mathbf{V}_3$

Masked Self-Attention Layer

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

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: $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q$

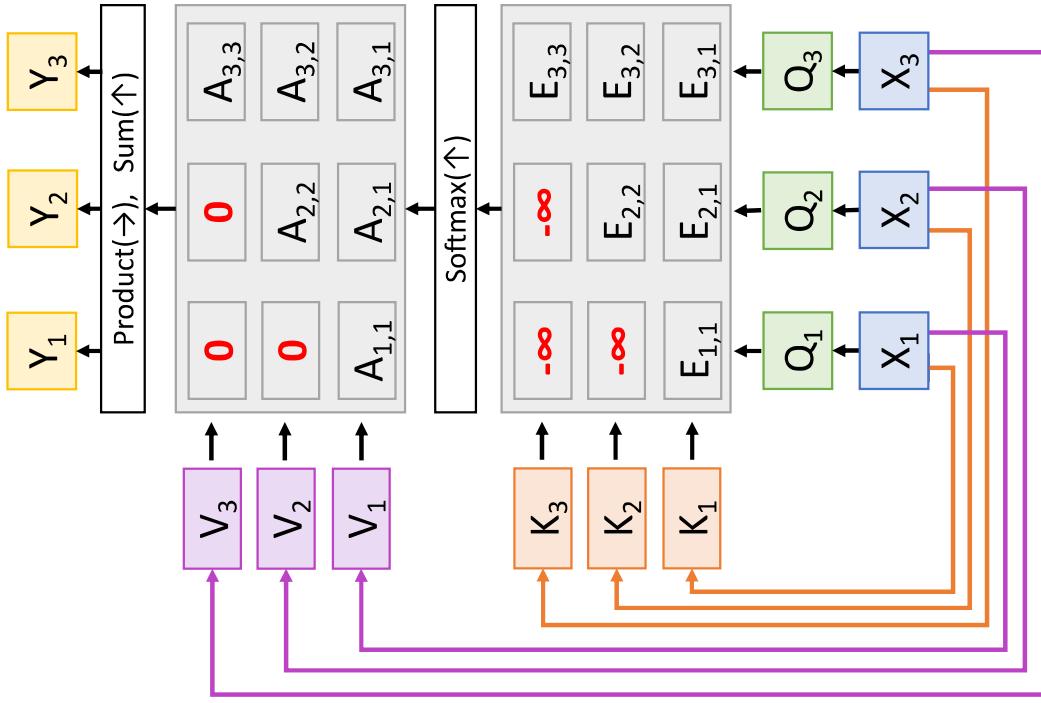
Key vectors: $\mathbf{K} = \mathbf{X}\mathbf{W}_K$ (Shape: $N_x \times D_Q$)

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

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

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

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



Masked Self-Attention Layer

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

Used for language modeling (predict next word)

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: $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q$

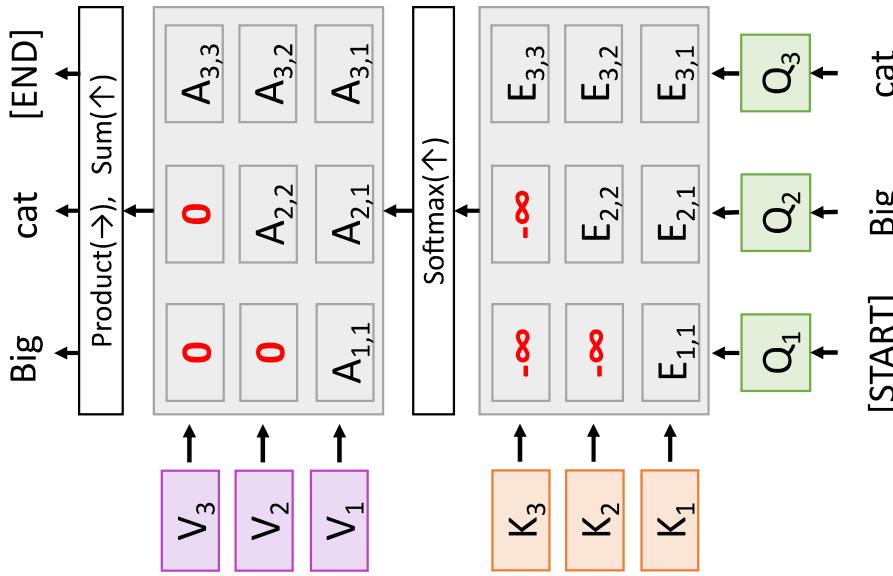
Key vectors: $\mathbf{K} = \mathbf{X}\mathbf{W}_K$ (Shape: $N_x \times D_Q$)

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

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

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

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



Multihead Self-Attention Layer

Use H independent
“Attention Heads” in parallel

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:

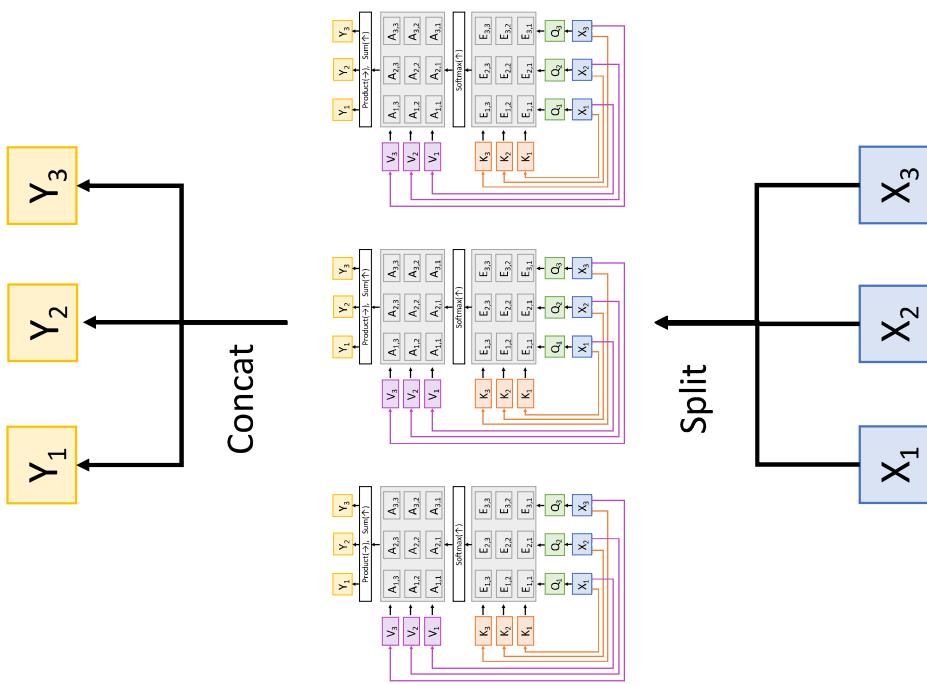
Query vectors: $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q$

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

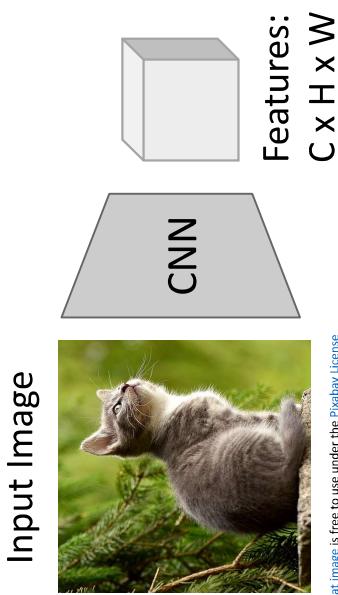
Value Vectors: $\mathbf{V} = \mathbf{X}\mathbf{W}_V$ (Shape: $N_x \times D_V$)

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

Attention weights: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_x \times N_x$)
Output vectors: $\mathbf{Y} = A\mathbf{V}$ (Shape: $N_x \times D_V$) $Y_i = \sum_j A_{i,j} \mathbf{V}_j$



Example: CNN with Self-Attention



Cat image is free to use under the [Pikabay License](#).

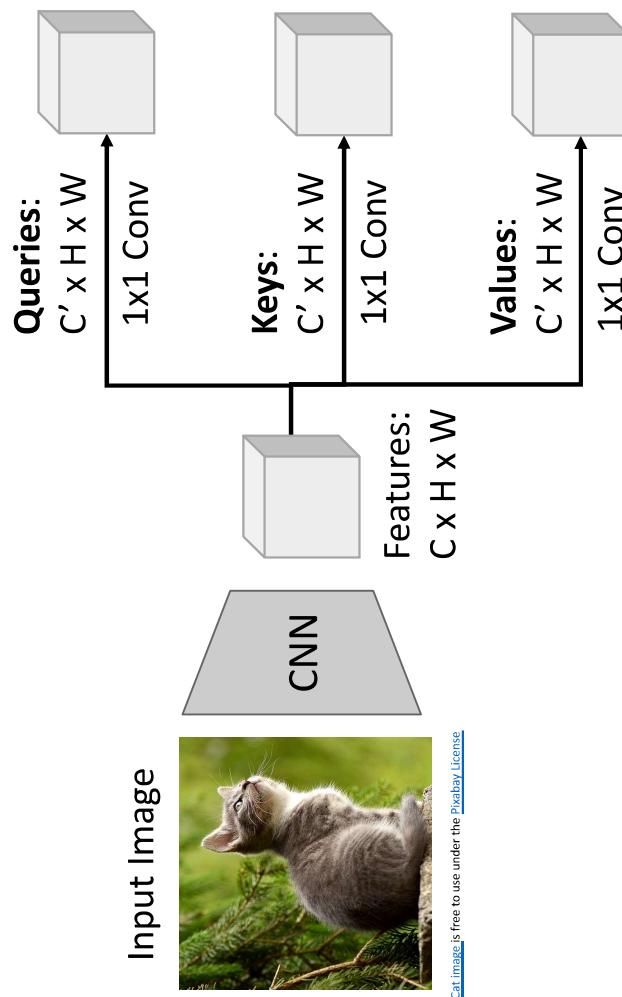
Zhang et al, "Self-Attention Generative Adversarial Networks", ICML 2018

Justin Johnson

Lecture 13 - 76

October 23, 2019

Example: CNN with Self-Attention



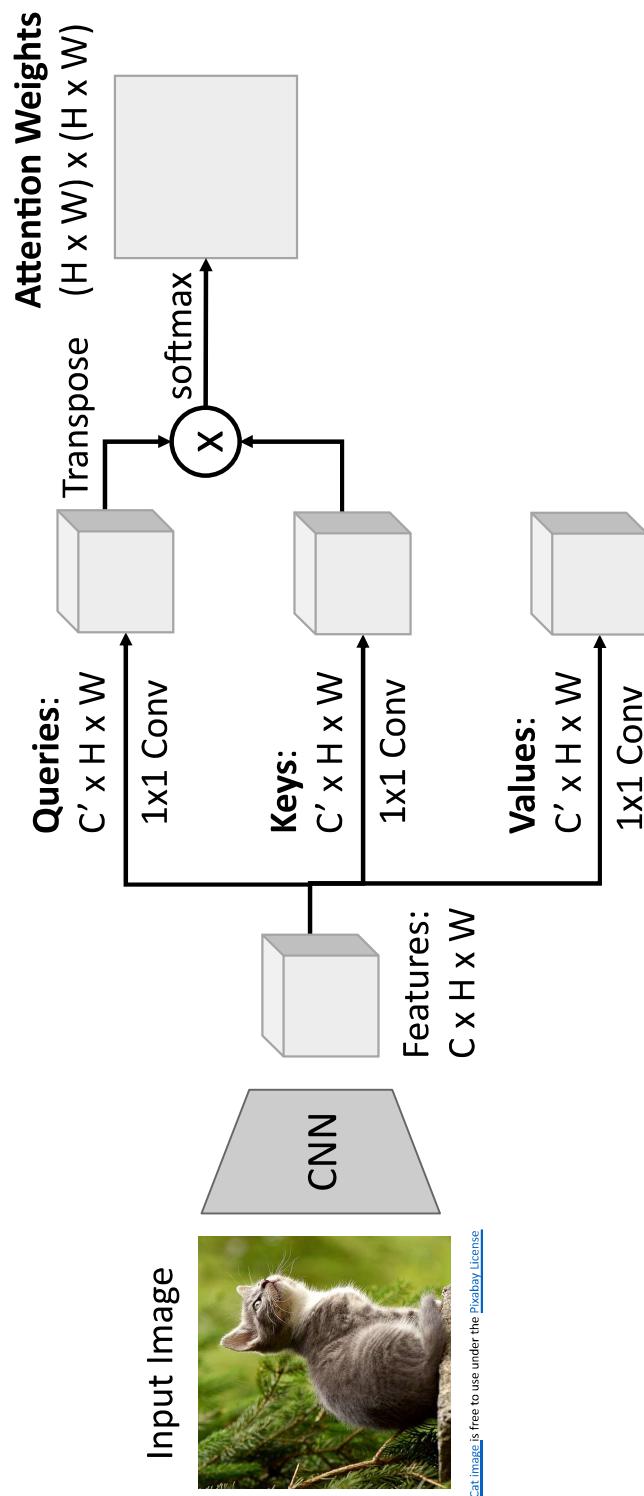
Zhang et al, "Self-Attention Generative Adversarial Networks", ICML 2018

Justin Johnson

Lecture 13 - 77

October 23, 2019

Example: CNN with Self-Attention

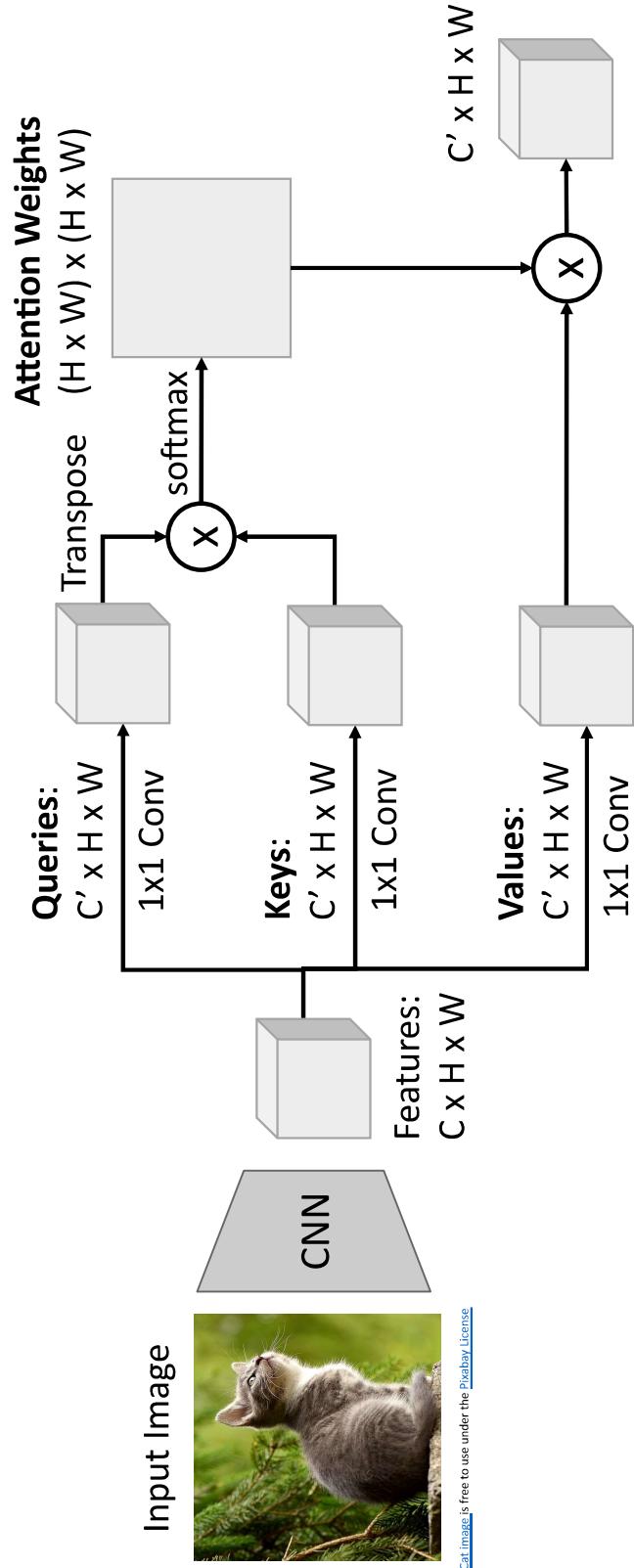


Zhang et al, "Self-Attention Generative Adversarial Networks", ICML 2018

Justin Johnson
Lecture 13 - 78

October 23, 2019

Example: CNN with Self-Attention



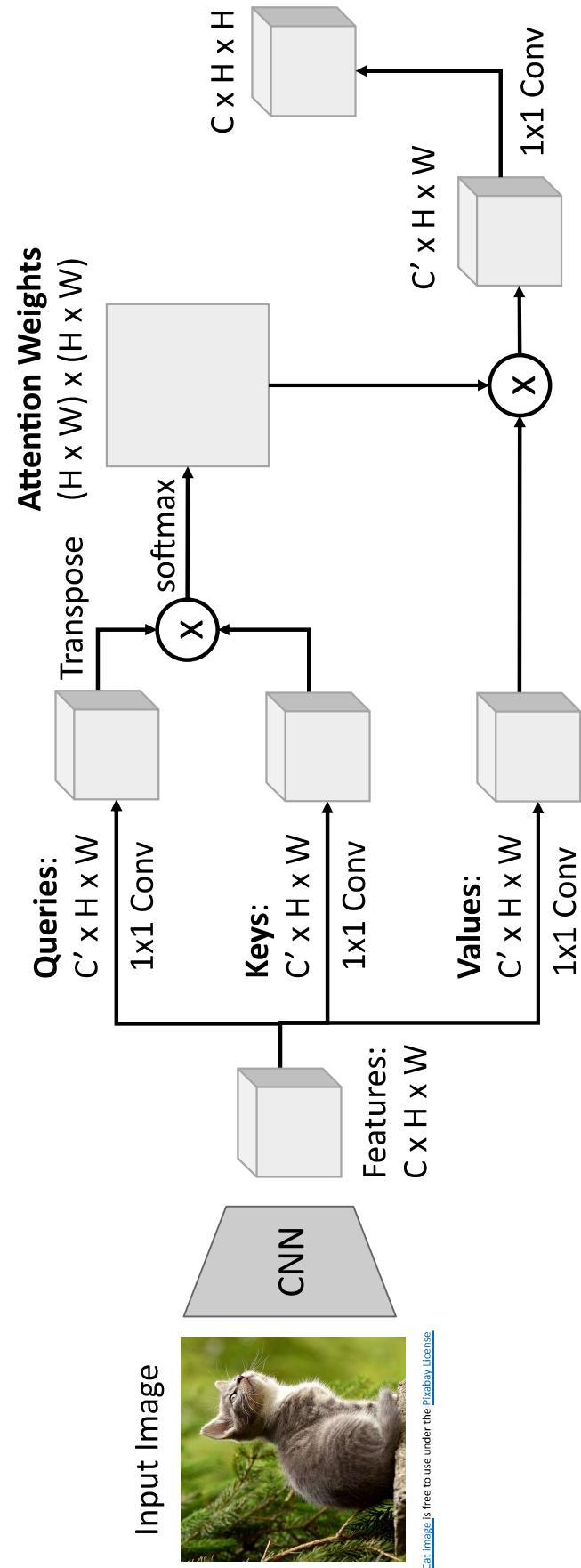
Zhang et al, "Self-Attention Generative Adversarial Networks", ICML 2018

Justin Johnson

Lecture 13 - 79

October 23, 2019

Example: CNN with Self-Attention



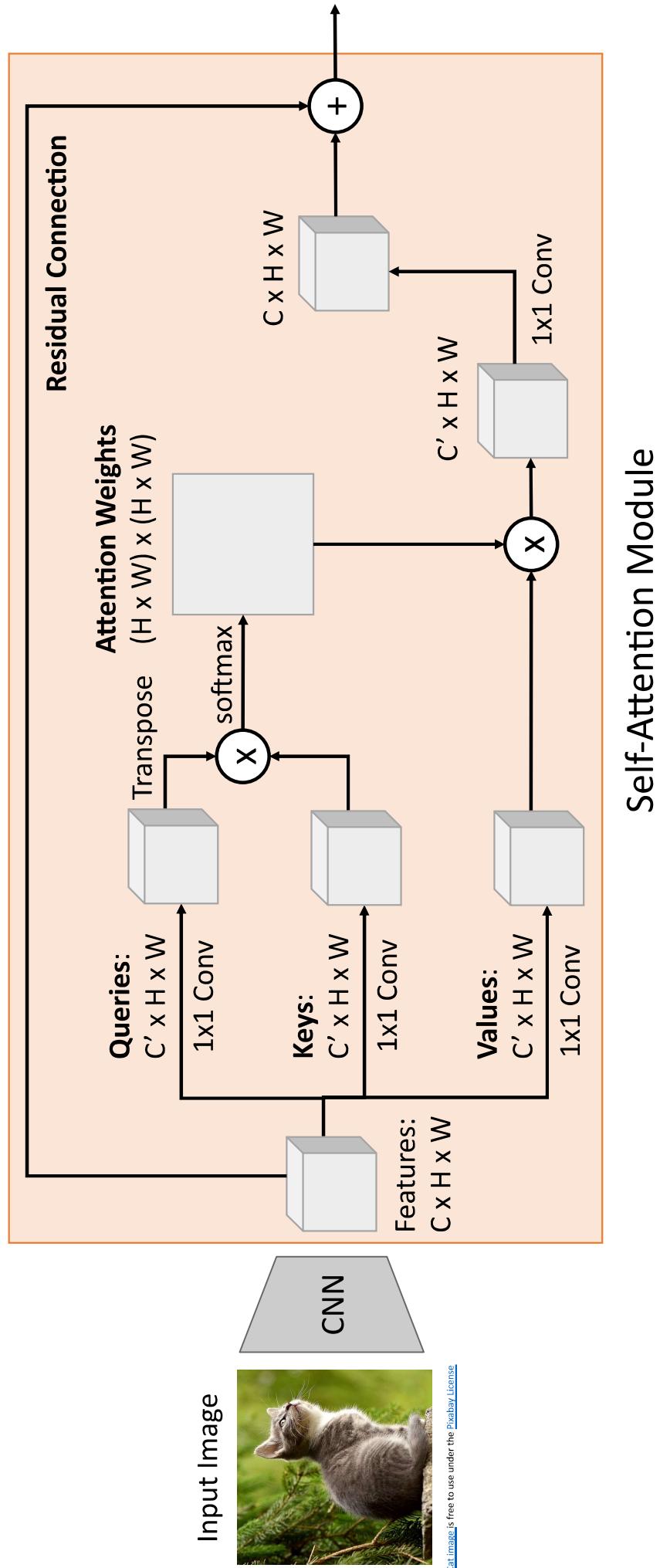
Zhang et al., "Self-Attention Generative Adversarial Networks", ICML 2018

Justin Johnson

Lecture 13 - 80

October 23, 2019

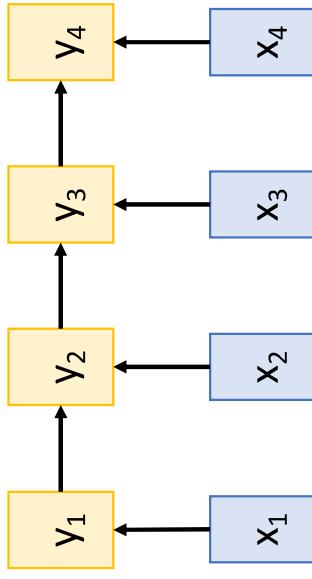
Example: CNN with Self-Attention



Zhang et al, "Self-Attention Generative Adversarial Networks", ICML 2018

Three Ways of Processing Sequences

Recurrent Neural Network

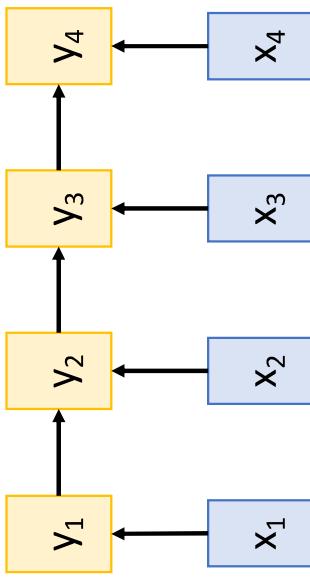


Works on **Ordered Sequences**

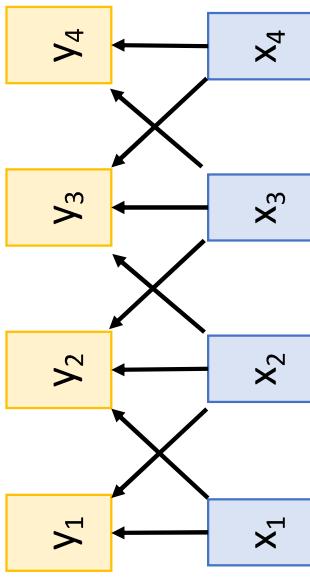
- (+) **Good at long sequences:** After one RNN layer, h_T "sees" the whole sequence
- (-) **Not parallelizable:** need to **compute hidden states sequentially**

Three Ways of Processing Sequences

Recurrent Neural Network



1D Convolution



Works on Ordered Sequences

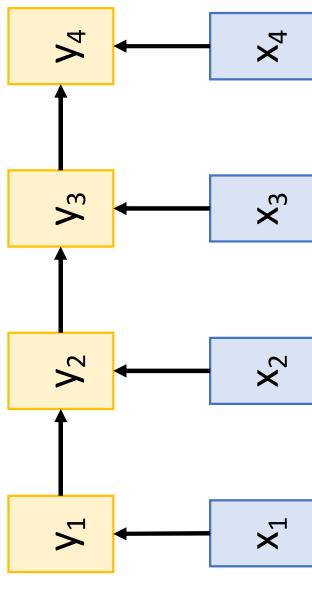
- (+) Good at long sequences: After one RNN layer, h_T "sees" the whole sequence
- (-) Not parallelizable: need to compute hidden states sequentially

Works on Multidimensional Grids

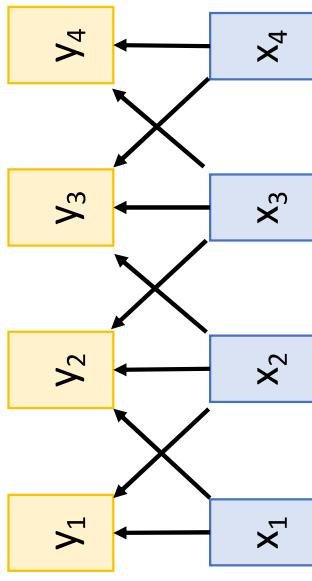
- (-) Bad at long sequences: Need to stack many conv layers for outputs to "see" the whole sequence
- (+) Highly parallel: Each output can be computed in parallel

Three Ways of Processing Sequences

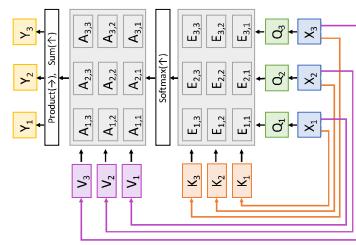
Recurrent Neural Network



1D Convolution



Self-Attention



Works on Ordered Sequences

- (+) Good at long sequences: After one RNN layer, h_T "sees" the whole sequence
- (-) 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

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**

- (+) Good at long sequences: After one RNN layer, h_T "sees" the whole sequence
- (-) 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**

The Transformer



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

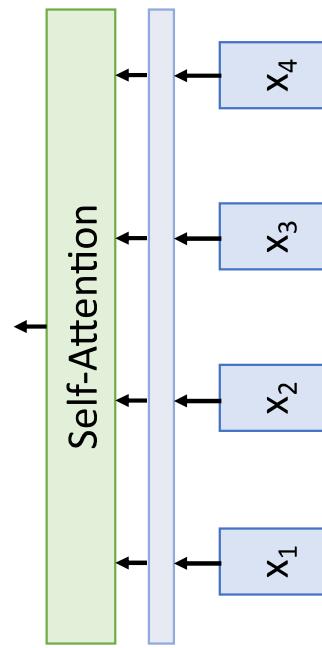
Justin Johnson

Lecture 13 - 86

October 23, 2019

The Transformer

All vectors interact
with each other



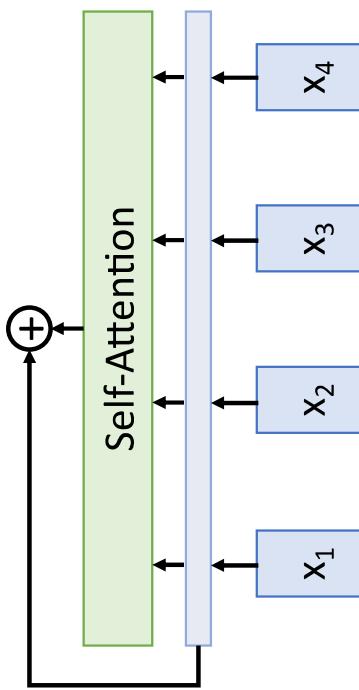
Vaswani et al, "Attention is all you need", NeurIPS 2017

Justin Johnson

Lecture 13 - 87

October 23, 2019

The Transformer



Residual connection
All vectors interact
with each other

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

Justin Johnson

Lecture 13 - 88

October 23, 2019

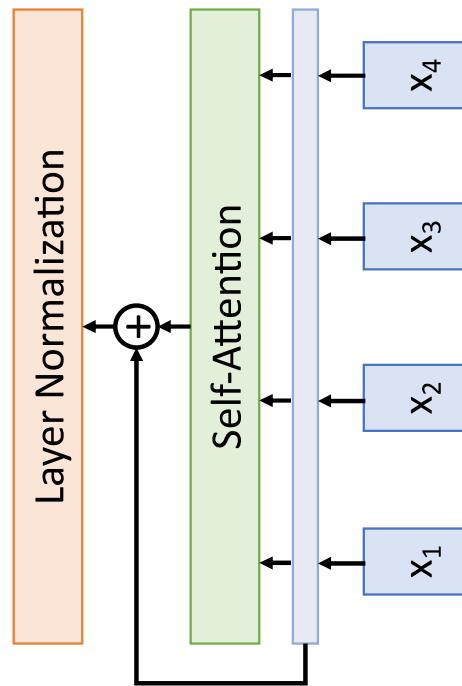
The Transformer

Recall Layer Normalization:

Given h_1, \dots, h_N (Shape: D)
scale: γ (Shape: D)
shift: β (Shape: D)

$$\mu_i = (1/D) \sum_j h_{i,j} \quad (\text{scalar})$$
$$\sigma_i = (\sum_j (h_{i,j} - \mu_i)^2)^{1/2} \quad (\text{scalar})$$
$$z_i = (h_i - \mu_i) / \sigma_i$$
$$y_i = \gamma * z_i + \beta$$

Residual connection
All vectors interact
with each other



Ba et al, 2016

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

Justin Johnson

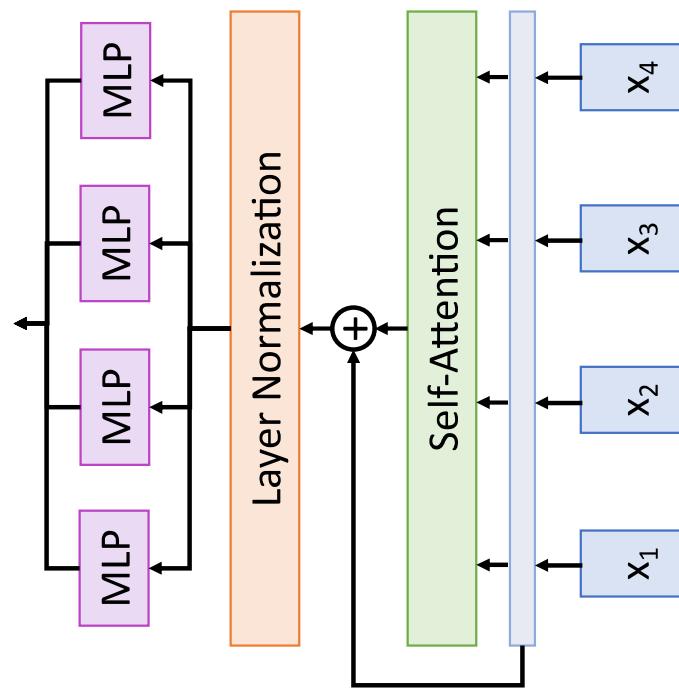
Lecture 13 - 89

October 23, 2019

The Transformer

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)
 $z_i = (h_i - \mu_i) / \sigma_i$
 $y_i = \gamma * z_i + \beta$



Ba et al, 2016

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

Justin Johnson

Lecture 13 - 90

October 23, 2019

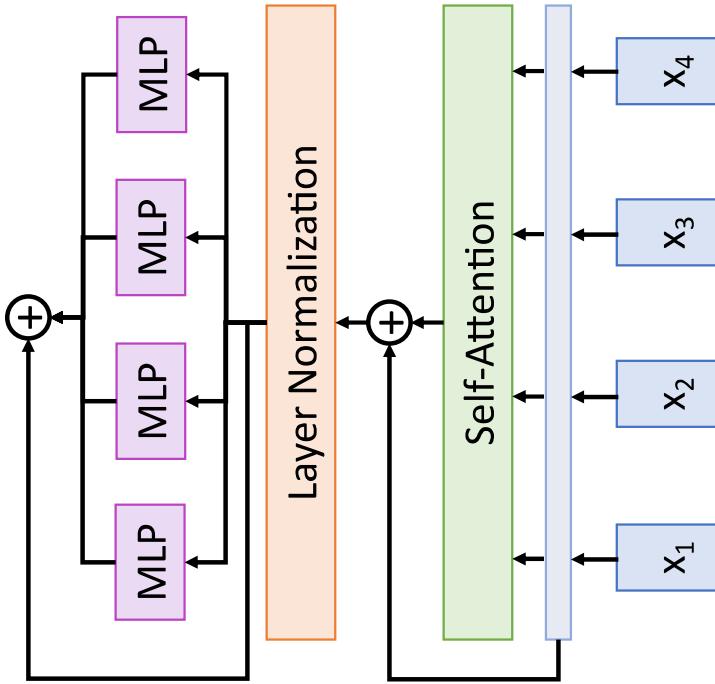
The Transformer

Recall Layer Normalization:

Given h_1, \dots, h_N (Shape: D)
scale: γ
shift: β
 $\mu_i = (1/D)\sum_j h_{ij}$ (scalar)
 $\sigma_i = (\sum_j (h_{ij} - \mu_i)^2)^{1/2}$ (scalar)
 $z_i = (h_i - \mu_i) / \sigma_i$
 $y_i = \gamma * z_i + \beta$

Residual connection

MLP independently
on each vector



Residual connection
All vectors interact
with each other

Ba et al, 2016

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

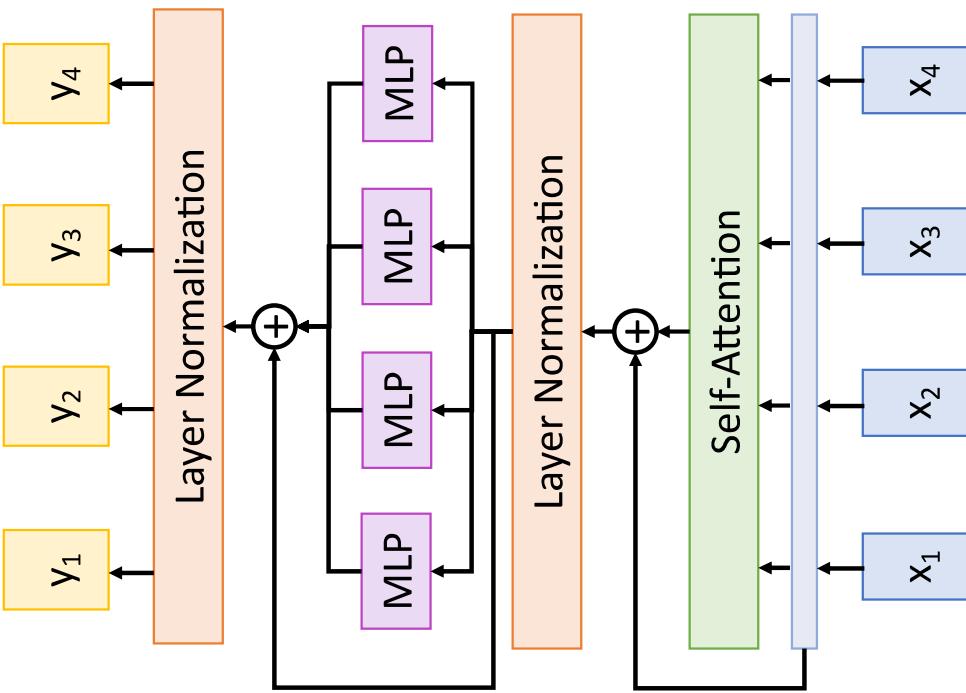
Justin Johnson

Lecture 13 - 91

October 23, 2019

The Transformer

Recall Layer Normalization:
Given h_1, \dots, h_N (Shape: D)
scale: γ (Shape: D)
shift: β (Shape: D)
 $\mu_i = (1/D)\sum_j h_{ij}$ (scalar)
 $\sigma_i = (\sum_j (h_{ij} - \mu_i)^2)^{1/2}$ (scalar)
 $z_i = (h_i - \mu_i) / \sigma_i$
 $y_i = \gamma * z_i + \beta$



Ba et al, 2016

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

Justin Johnson

Lecture 13 - 92

October 23, 2019

The Transformer

Transformer Block:

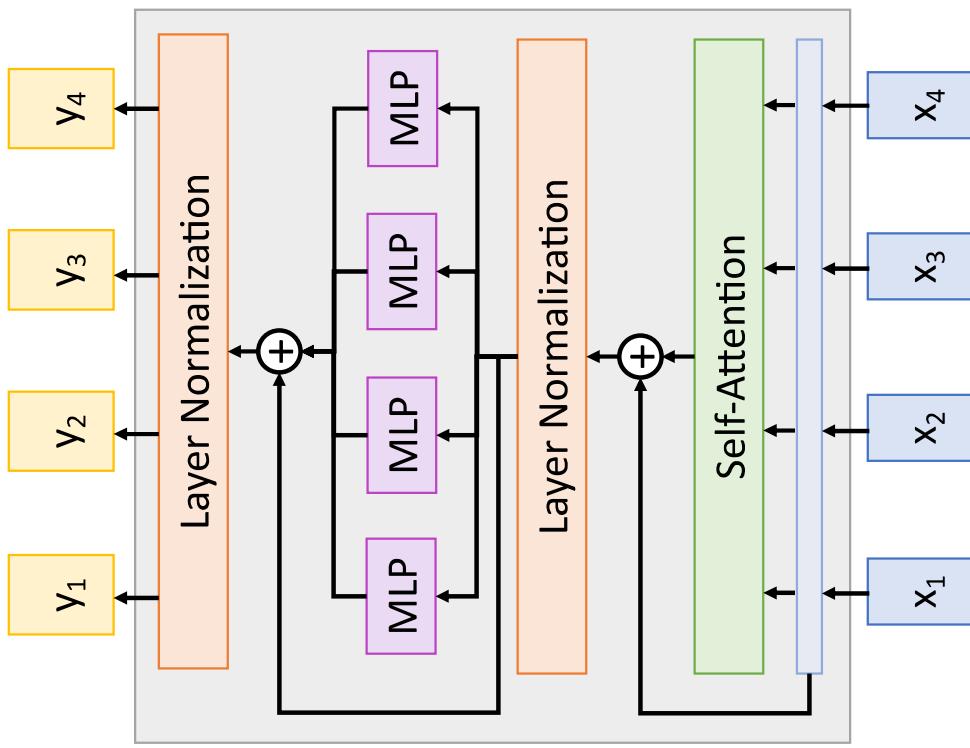
Input: Set of vectors x

Output: Set of vectors y

Self-attention is the only
interaction between 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

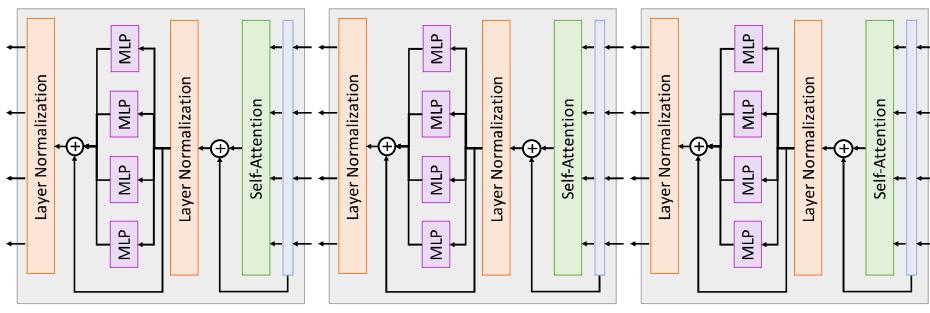
Output: Set of vectors y

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

Self-attention is the only
interaction between vectors!

Layer norm and MLP work
independently per vector

Highly scalable, highly
parallelizable



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

Justin Johnson

Lecture 13 - 94

October 23, 2019

The Transformer: Transfer Learning

“ImageNet Moment for Natural Language Processing”

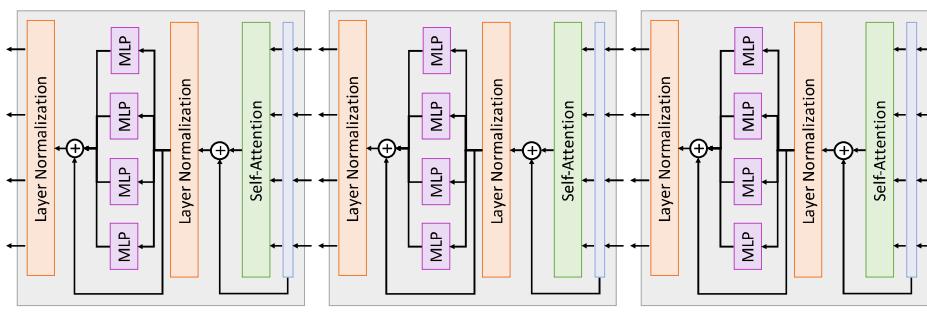
Pretraining:

Download a lot of text from the internet

Train a giant Transformer model for language modeling

Finetuning:

Fine-tune the Transformer on your own NLP task



Scaling up Transformers

Model	Layers	Width	Heads	Params	Data	Training
Transformer-Base	12	512	8	65M		8x P100 (12 hours)
Transformer-Large	12	1024	16	213M		8x P100 (3.5 days)

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

Justin Johnson

Lecture 13 - 96

October 23, 2019

Scaling up Transformers

Model	Layers	Width	Heads	Params	Data	Training
Transformer-Base	12	512	8	65M		8x P100 (12 hours)
Transformer-Large	12	1024	16	213M		8x P100 (3.5 days)
BERT-Base	12	768	12	110M	13 GB	
BERT-Large	24	1024	16	340M	13 GB	

Devlin et al, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding", EMNLP 2018

Scaling up Transformers

Model	Layers	Width	Heads	Params	Data	Training
Transformer-Base	12	512	8	65M		8x P100 (12 hours)
Transformer-Large	12	1024	16	213M		8x P100 (3.5 days)
BERT-Base	12	768	12	110M		13 GB
BERT-Large	24	1024	16	340M		13 GB
XLNet-Large	24	1024	16	~340M	126 GB	512x TPU-v3 (2.5 days)
RoBERTa	24	1024	16	355M	160 GB	1024x V100 GPU (1 day)

Yang et al, "XLNet: Generalized Autoregressive Pretraining for Language Understanding", 2019
Liu et al, "RoBERTa: A Robustly Optimized BERT Pretraining Approach", 2019

Scaling up Transformers

Model	Layers	Width	Heads	Params	Data	Training
Transformer-Base	12	512	8	65M		8x P100 (12 hours)
Transformer-Large	12	1024	16	213M		8x P100 (3.5 days)
BERT-Base	12	768	12	110M		
BERT-Large	24	1024	16	340M		13 GB
XLNet-Large	24	1024	16	~340M		126 GB
RoBERTa	24	1024	16	355M		160 GB
GPT-2	12	768	?	117M		40 GB
GPT-2	24	1024	?	345M		40 GB
GPT-2	36	1280	?	762M		40 GB
GPT-2	48	1600	?	1.5B		40 GB

Scaling up Transformers

Model	Layers	Width	Heads	Params	Data	Training
Transformer-Base	12	512	8	65M		8x P100 (12 hours)
Transformer-Large	12	1024	16	213M		8x P100 (3.5 days)
BERT-Base	12	768	12	110M	13 GB	
BERT-Large	24	1024	16	340M	13 GB	
XLNet-Large	24	1024	16	~340M	126 GB	512x TPU-v3 (2.5 days)
RoBERTa	24	1024	16	355M	160 GB	1024x V100 GPU (1 day)
GPT-2	12	768	?	117M	40 GB	
GPT-2	24	1024	?	345M	40 GB	
GPT-2	36	1280	?	762M	40 GB	
GPT-2	48	1600	?	1.5B	40 GB	
Megatron-LM	40	1536	16	1.2B	174 GB	64x V100 GPU
Megatron-LM	54	1920	20	2.5B	174 GB	128x V100 GPU
Megatron-LM	64	2304	24	4.2B	174 GB	256x V100 GPU (10 days)
Megatron-LM	72	3072	32	8.3B	174 GB	512x V100 GPU (9 days)

Shoeybi et al, "Megatron-LM: Training Multi-Billion Parameter Language Models using Model Parallelism", 2019

Scaling up Transformers

~\$430,000 on Amazon AWS!

Model	Layers	Width	Heads	Params	Data	Training
Transformer-Base	12	512	8	65M		8x P100 (12 hours)
Transformer-Large	12	1024	16	213M		8x P100 (3.5 days)
BERT-Base	12	768	12	110M		
BERT-Large	24	1024	16	340M		13 GB
XLNet-Large	24	1024	16	~340M		126 GB
RoBERTa	24	1024	16	355M		160 GB
GPT-2	12	768	?	117M		40 GB
GPT-2	24	1024	?	345M		40 GB
GPT-2	36	1280	?	762M		40 GB
GPT-2	48	1600	?	1.5B		40 GB
Megatron-LM	40	1536	16	1.2B		64x V100 GPU
Megatron-LM	54	1920	20	2.5B		174 GB
Megatron-LM	64	2304	24	4.2B		128x V100 GPU (10 days)
Megatron-LM	72	3072	32	8.3B		256x V100 GPU (10 days)
					174 GB	512x V100 GPU (9 days)

Shoeybi et al, "Megatron-LM: Training Multi-Billion Parameter Language Models using Model Parallelism", 2019

PROMPT (Human-written): *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.*

PROMPT (Human-written): *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.*

COMPLETION (Transformer-written): 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.

Pérez and his friends were astonished to see the unicorn herd. These creatures could be seen from the air without having to move too much to see them – they were so close they could touch their horns.

While examining these bizarre creatures the scientists discovered that the creatures also spoke some fairly regular English. Pérez stated, “We can see, for example, that they have a common ‘language,’ something like a dialect or dialectic.” Dr. Pérez believes that the unicorns may have originated in Argentina, where the animals were believed to be descendants of a lost race of people who lived there before the arrival of humans in those parts of South America.

While their origins are still unclear, some believe that perhaps the creatures were created when a human and a unicorn met each other in a time before human civilization. According to Pérez, “In South America, such incidents seem to be quite common.” However, Pérez also pointed out that it is likely that the only way of knowing for sure if unicorns are indeed the descendants of a lost alien race is through DNA. “But they seem to be able to communicate in English quite well, which I believe is a sign of evolution, or at least a change in social organization,” said the scientist.

OpenAI, “Better Language Models and their Implications”, 2019, <https://openai.com/blog/better-language-models/>

PROMPT (Human-written): 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.

COMPLETION (Transformer-written): 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

Try it yourself:

<https://talktotransformer.com>

Dr. Pérez believes that the unicorns may have originated in Argentina, where the animals were believed to be descendants of a lost race of people who lived there before the arrival of humans in those parts of South America. While their origins are still unclear, some believe that perhaps the creatures were created when a human and a unicorn met each other in a time before human civilization. According to Pérez, "In South America, such incidents seem to be quite common." However, Pérez also pointed out that it is likely that the only way of knowing for sure if unicorns are indeed the descendants of a lost alien race is through DNA. "But they seem to be able to communicate in English quite well, which I believe is a sign of evolution, or at least a change in social organization," said the scientist.

OpenAI, "Better Language Models and their Implications", 2019, <https://openai.com/blog/better-language-models/>

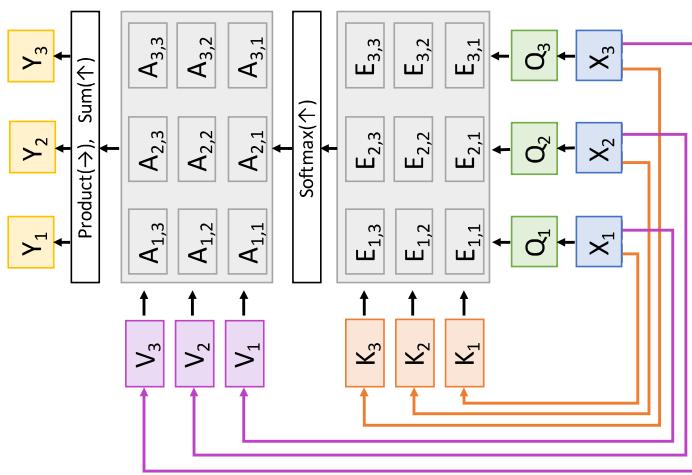
Summary

Generalized **Self-Attention**
is new, powerful neural
network primitive

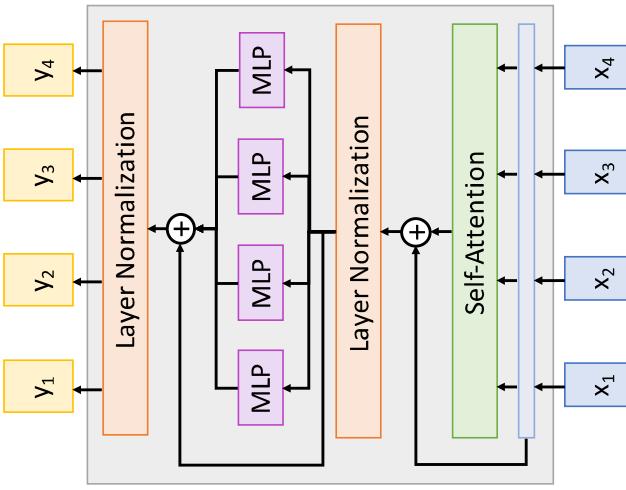
Adding **Attention** to RNN
models lets them look at
different parts of the
input at each timestep



A dog is standing on a hardwood floor.



Transformers are a new
neural network model
that only uses attention



Next Week: Guest Lectures



Monday 10/28
Luowei Zhou
Vision and Language



Wednesday 10/30
Prof. Atul Prakash
Adversarial Machine Learning

October 23, 2019

Lecture 13 - 106

Justin Johnson