

# CS163: Deep Learning for Computer Vision

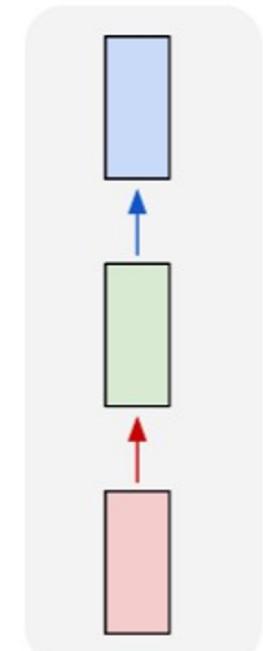
## Lecture 11: Recursive Neural Network (RNN) and Attention

# Announcement

- Assignment 3 is released at  
<https://github.com/UCLAdeepvision/CS163-Assessments-2024Fall/tree/main/Assignment3>
  - Due: Nov. 24, 2024 (three weeks to go)
  - Build a Vision Transformer (ViT) from scratch.
  - Build a Semantic Segmentation model using a ViT encoder.
- Midway presentation: Please follow the 20min time limits when preparing your midway presentations

# So far: “Feedforward” Neural Networks

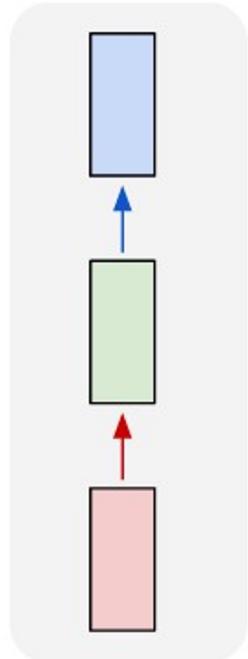
one to one



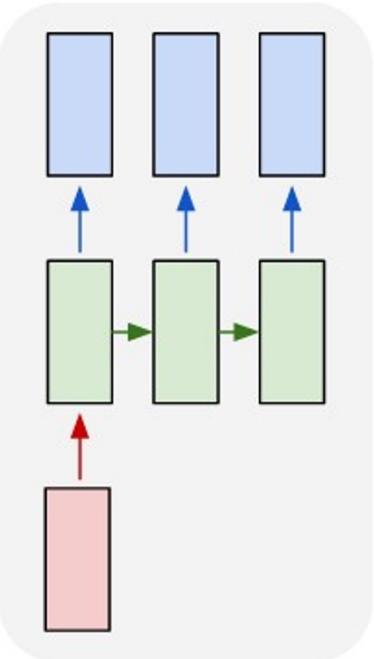
e.g. **Image classification**  
Image -> Label

# Recurrent Neural Networks: Process Sequences

one to one



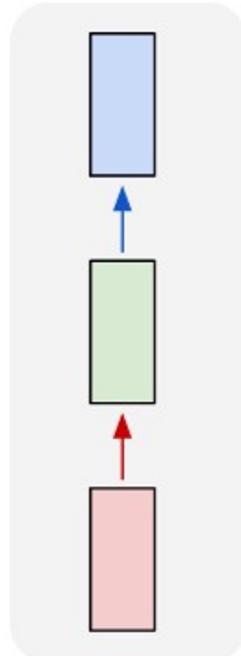
one to many



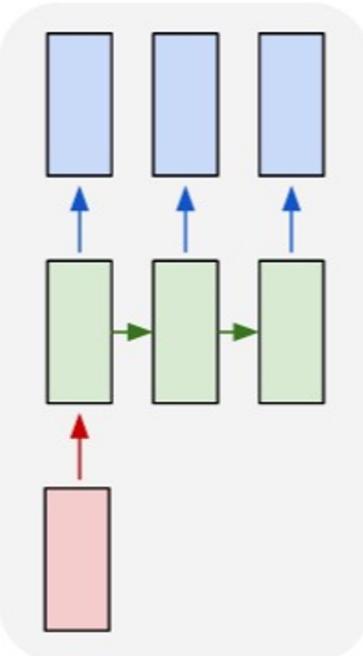
e.g. **Image Captioning:**  
Image -> sequence of words

# Recurrent Neural Networks: Process Sequences

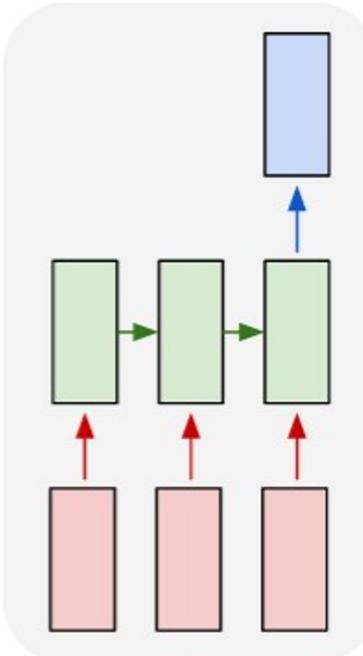
one to one



one to many



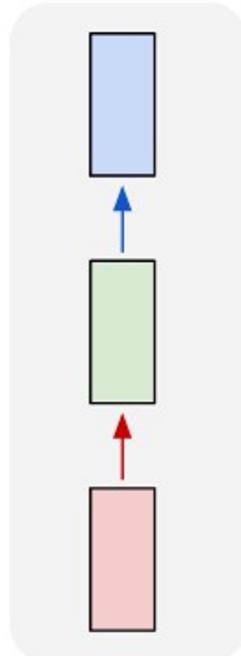
many to one



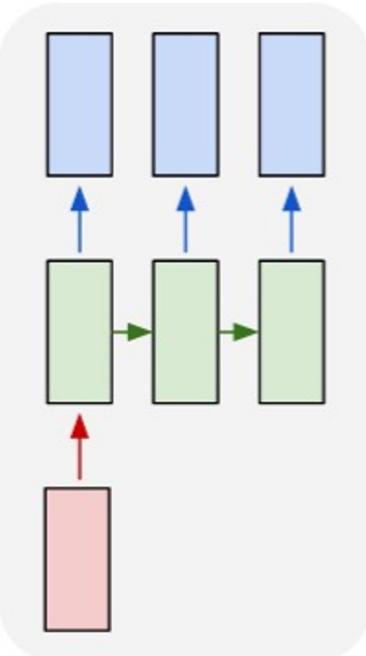
e.g. **Video classification:**  
Sequence of images -> label

# Recurrent Neural Networks: Process Sequences

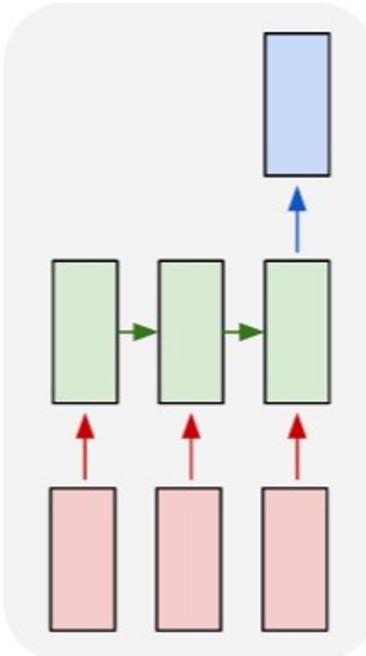
one to one



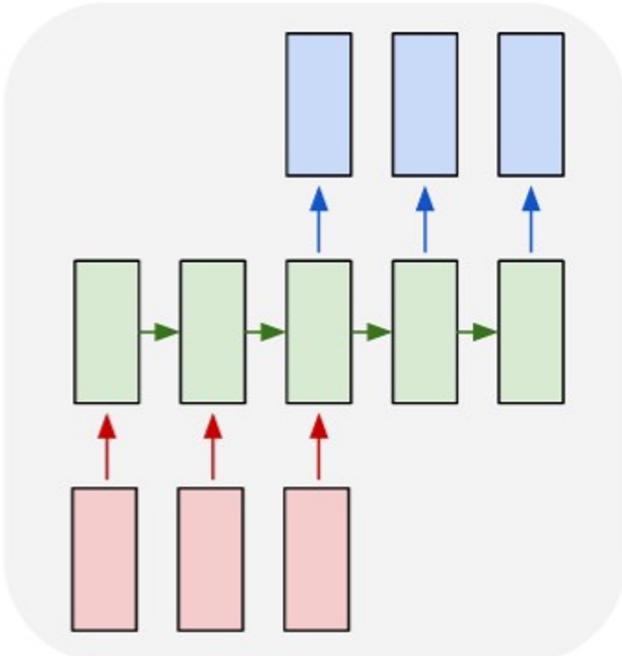
one to many



many to one



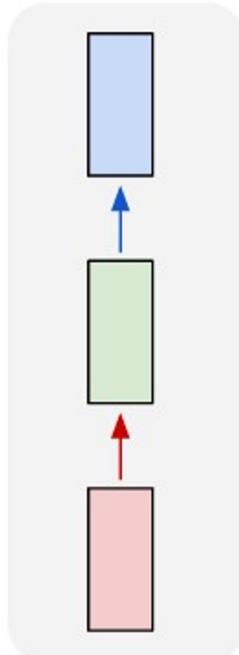
many to many



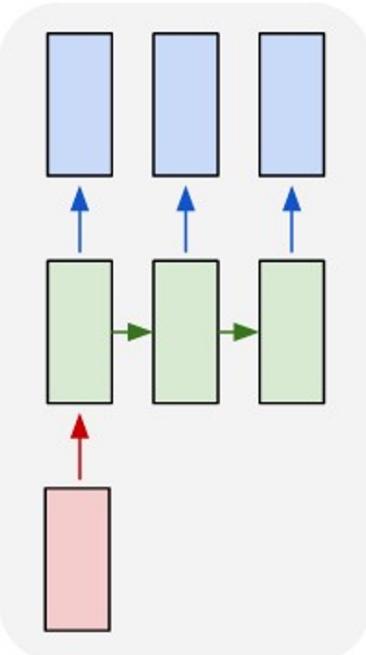
e.g. **Machine Translation:**  
Sequence of words -> Sequence of words

# Recurrent Neural Networks: Process Sequences

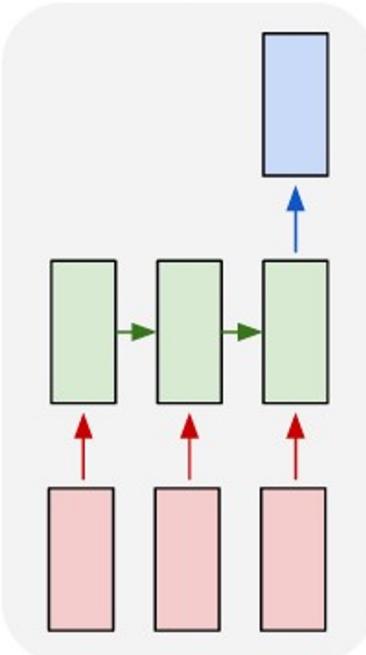
one to one



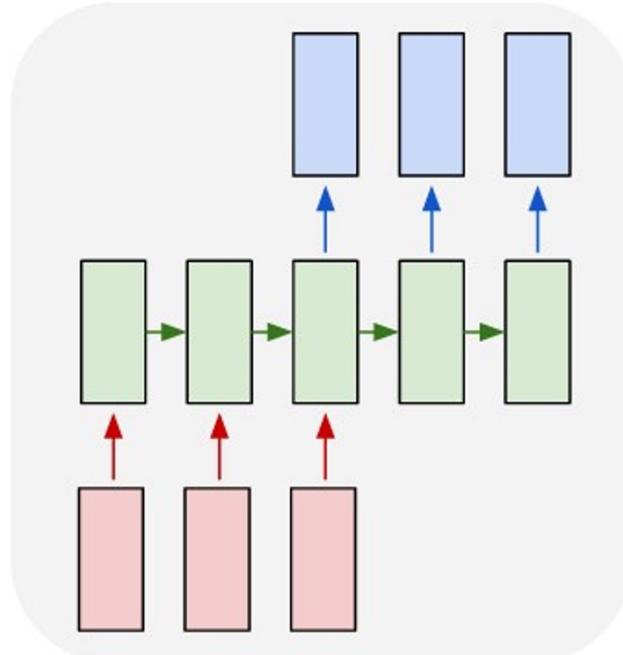
one to many



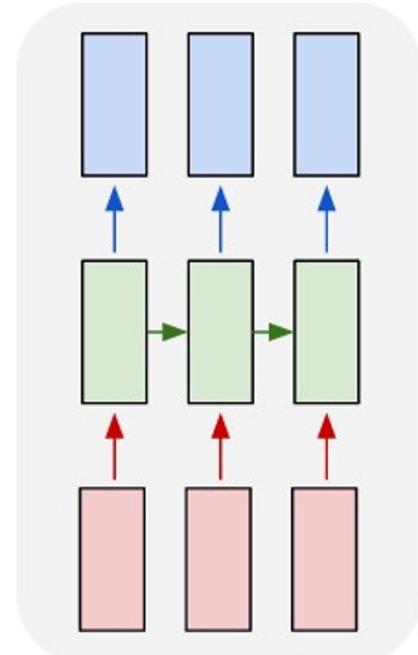
many to one



many to many

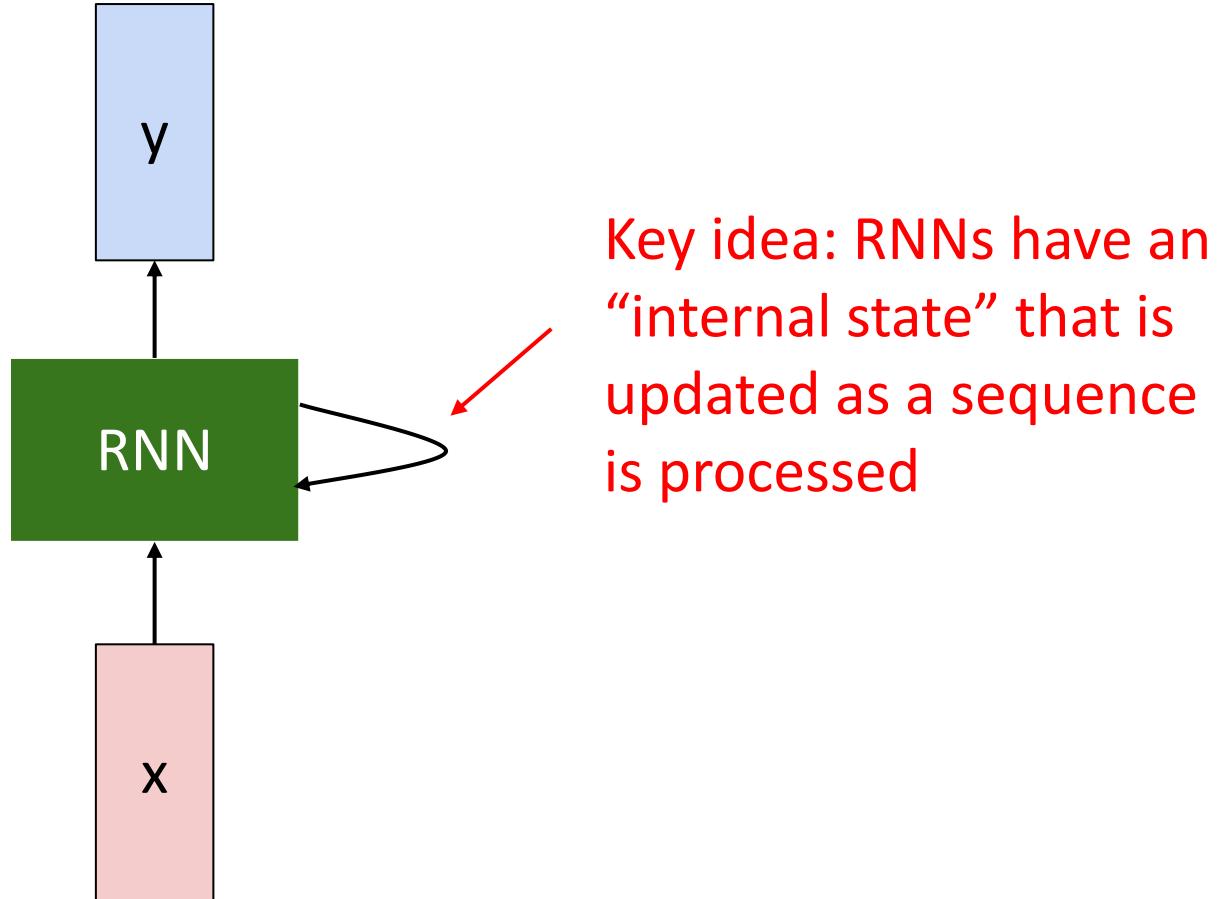


many to many



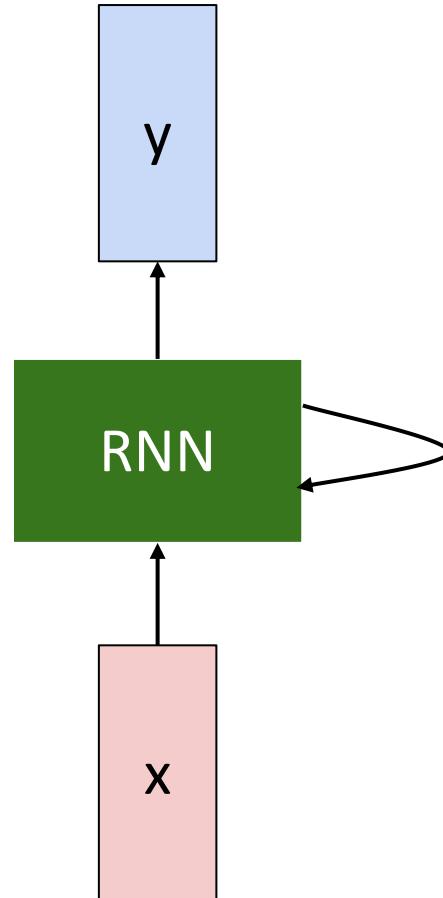
e.g. **Per-frame video classification:**  
Sequence of images -> Sequence of labels

# Recurrent Neural Networks



Key idea: RNNs have an  
“internal state” that is  
updated as a sequence  
is processed

# Recurrent Neural Networks



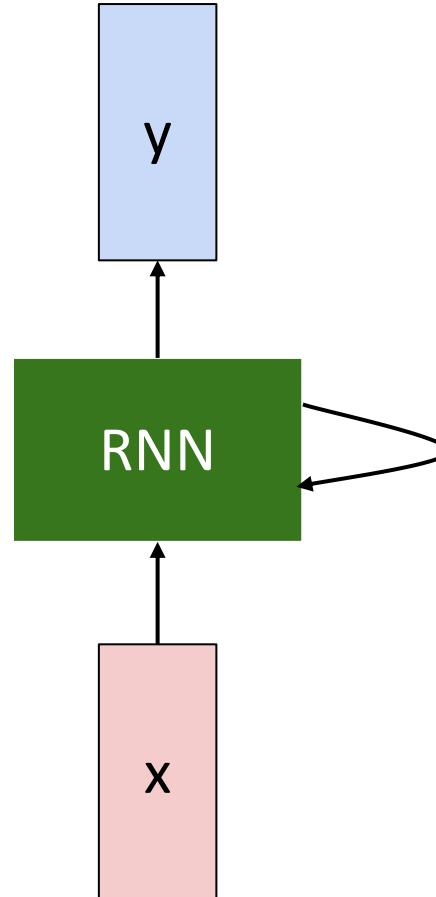
We can process a sequence of vectors  $\mathbf{x}$  by applying a **recurrence formula** at every time step:

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

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

# Recurrent Neural Networks

Notice: the same function and the same set of parameters are used at every time step.



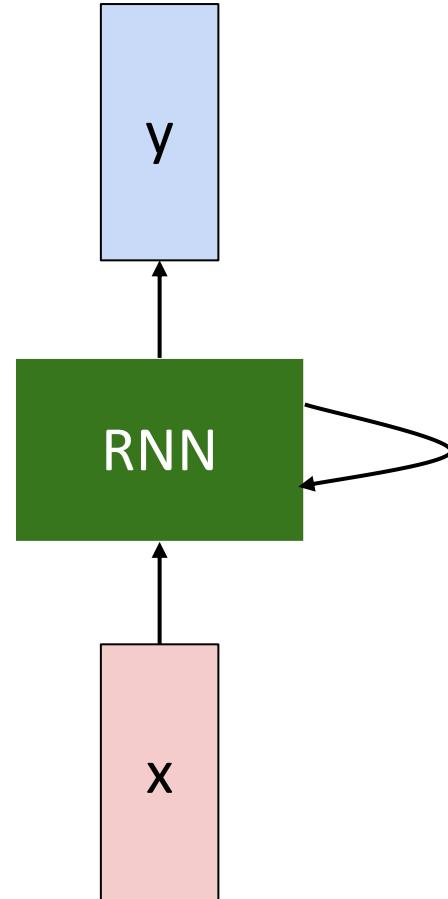
We can process a sequence of vectors  $x$  by applying a **recurrence formula** at every time step:

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

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

# (Vanilla) Recurrent Neural Networks

The state consists of a single “*hidden*” vector  $\mathbf{h}$ :



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

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

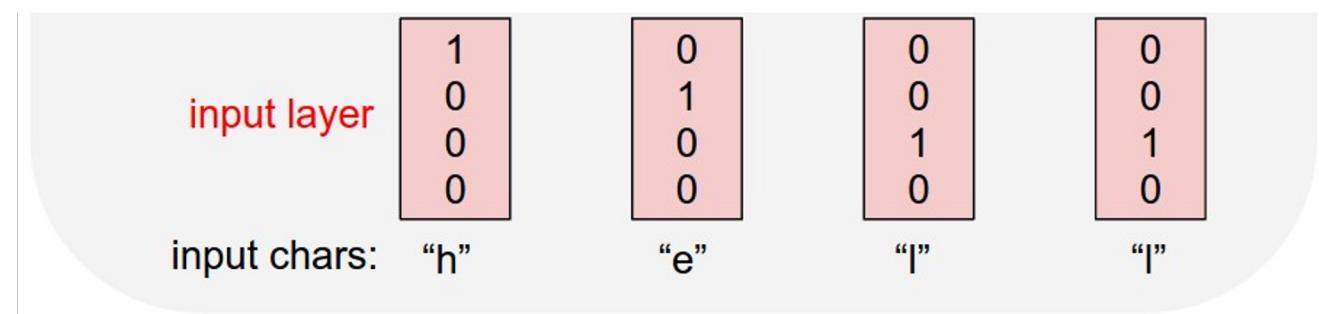
$$y_t = W_{hy}h_t + b_y$$

# Example: Language Modeling

Given characters 1, 2, ..., t-1,  
model predicts character t

Training sequence: "hello"

Vocabulary: [h, e, l, o]



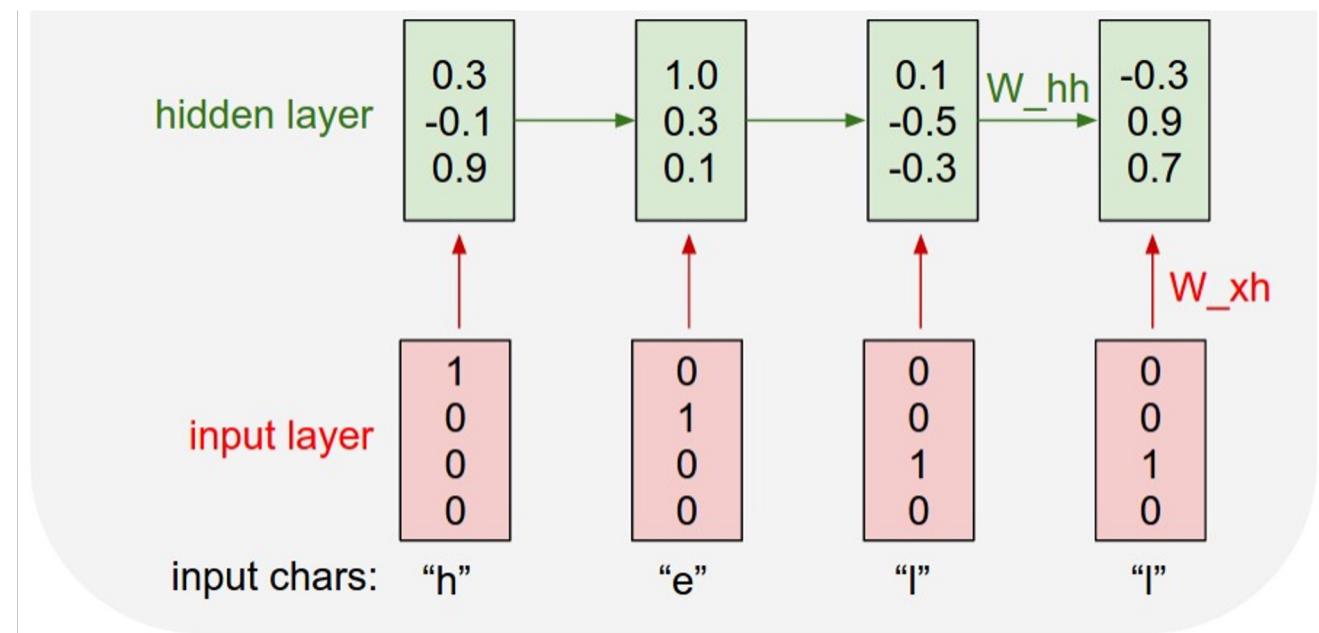
# Example: Language Modeling

Given characters 1, 2, ..., t-1,  
model predicts character t

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

Training sequence: "hello"

Vocabulary: [h, e, l, o]



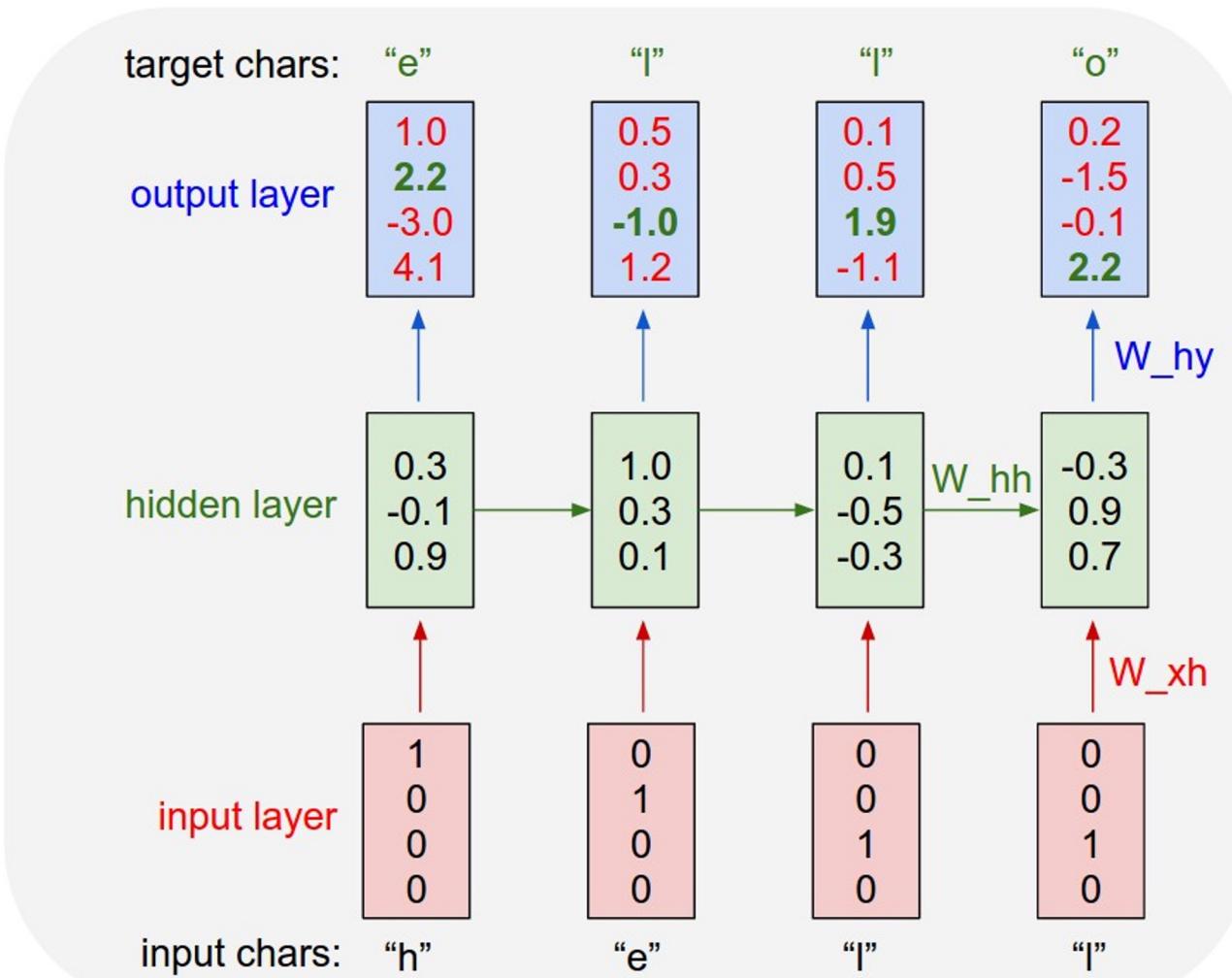
# Example: Language Modeling

Given characters 1, 2, ..., t-1,  
model predicts character t

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

Training sequence: "hello"

Vocabulary: [h, e, l, o]



# Example: Language Modeling

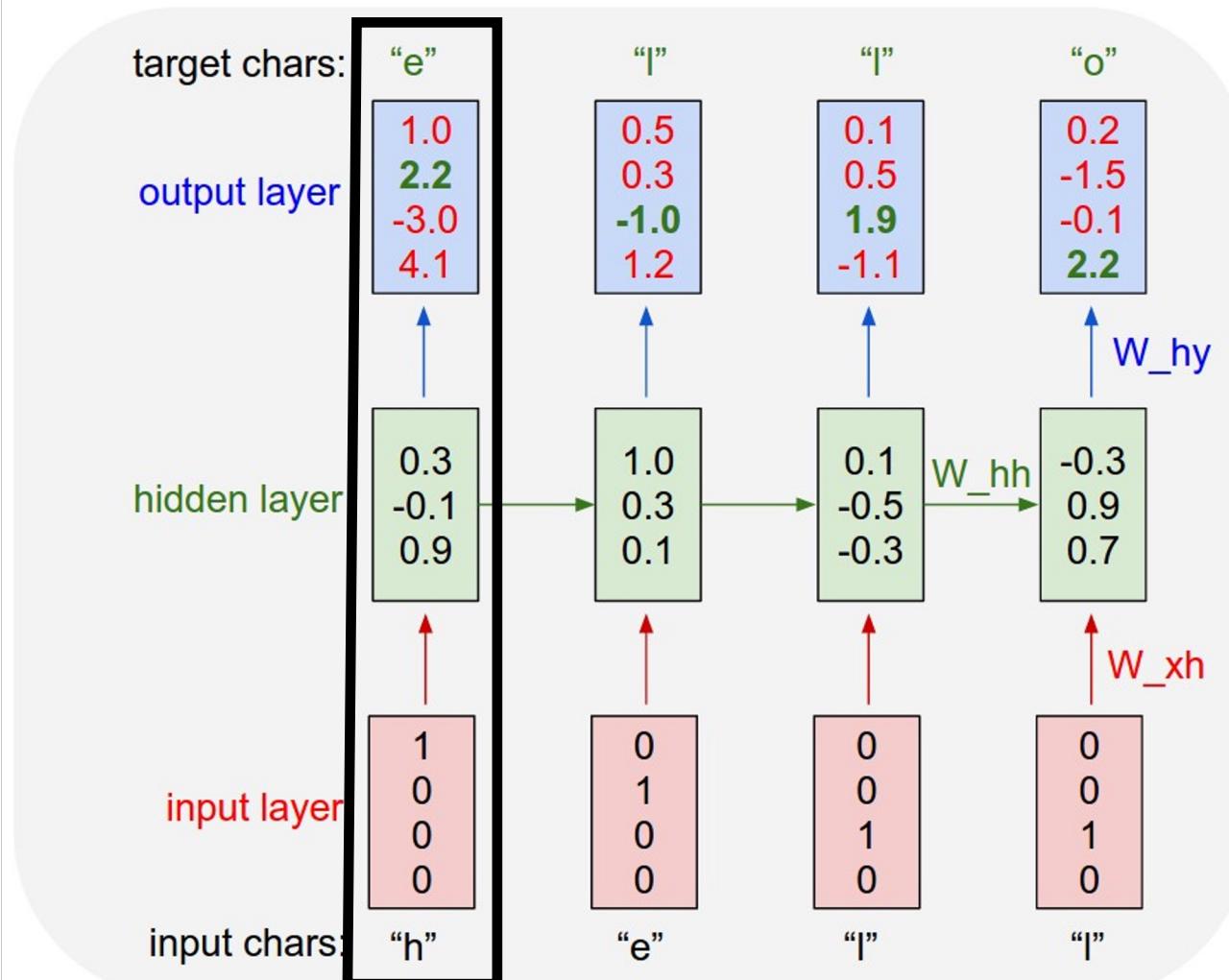
Given “h”, predict “e”

Given characters 1, 2, ..., t-1,  
model predicts character t

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

Training sequence: “hello”

Vocabulary: [h, e, l, o]



# Example: Language Modeling

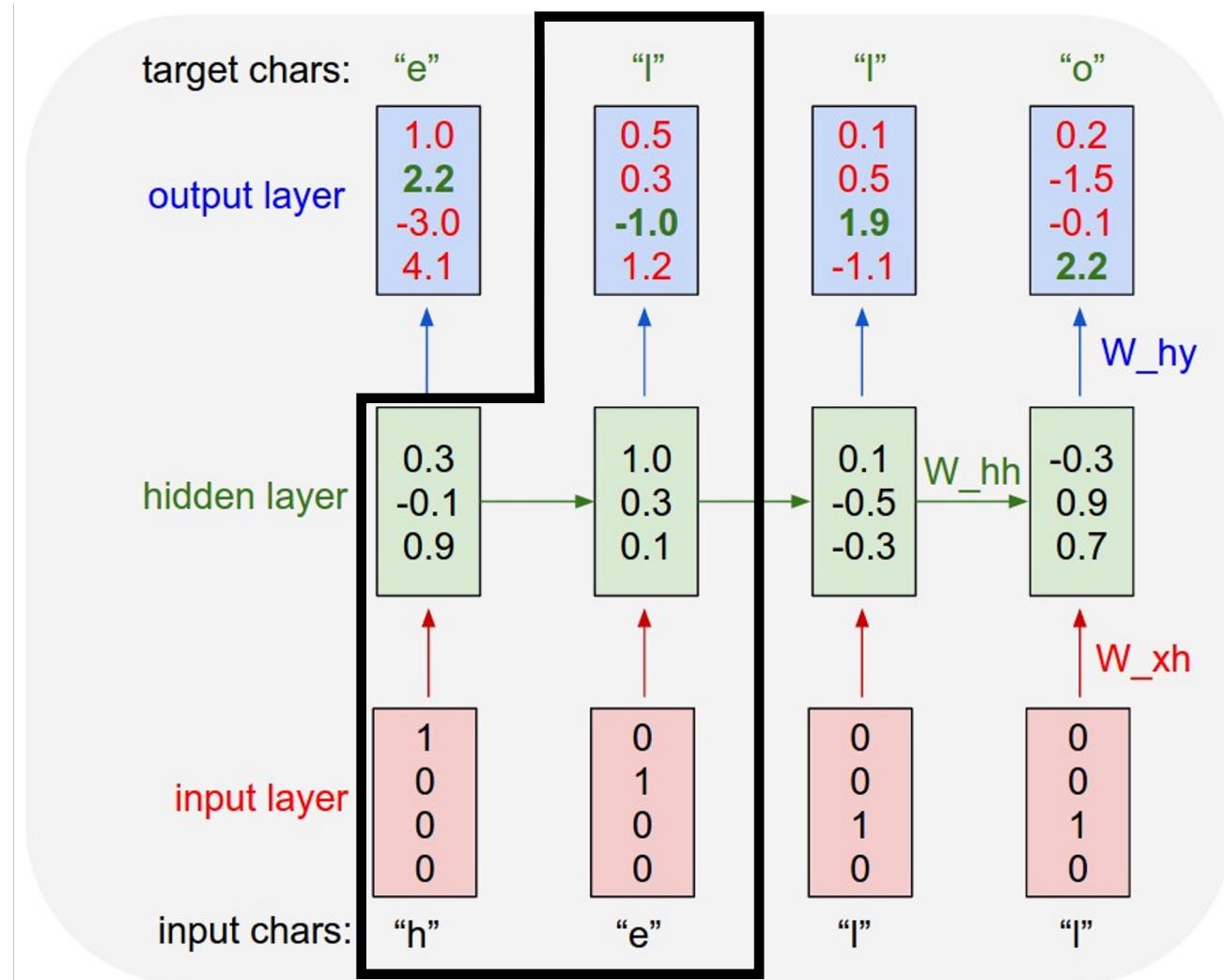
Given “he”, predict “l”

Given characters 1, 2, ..., t-1,  
model predicts character t

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

Training sequence: “hello”

Vocabulary: [h, e, l, o]



# Example: Language Modeling

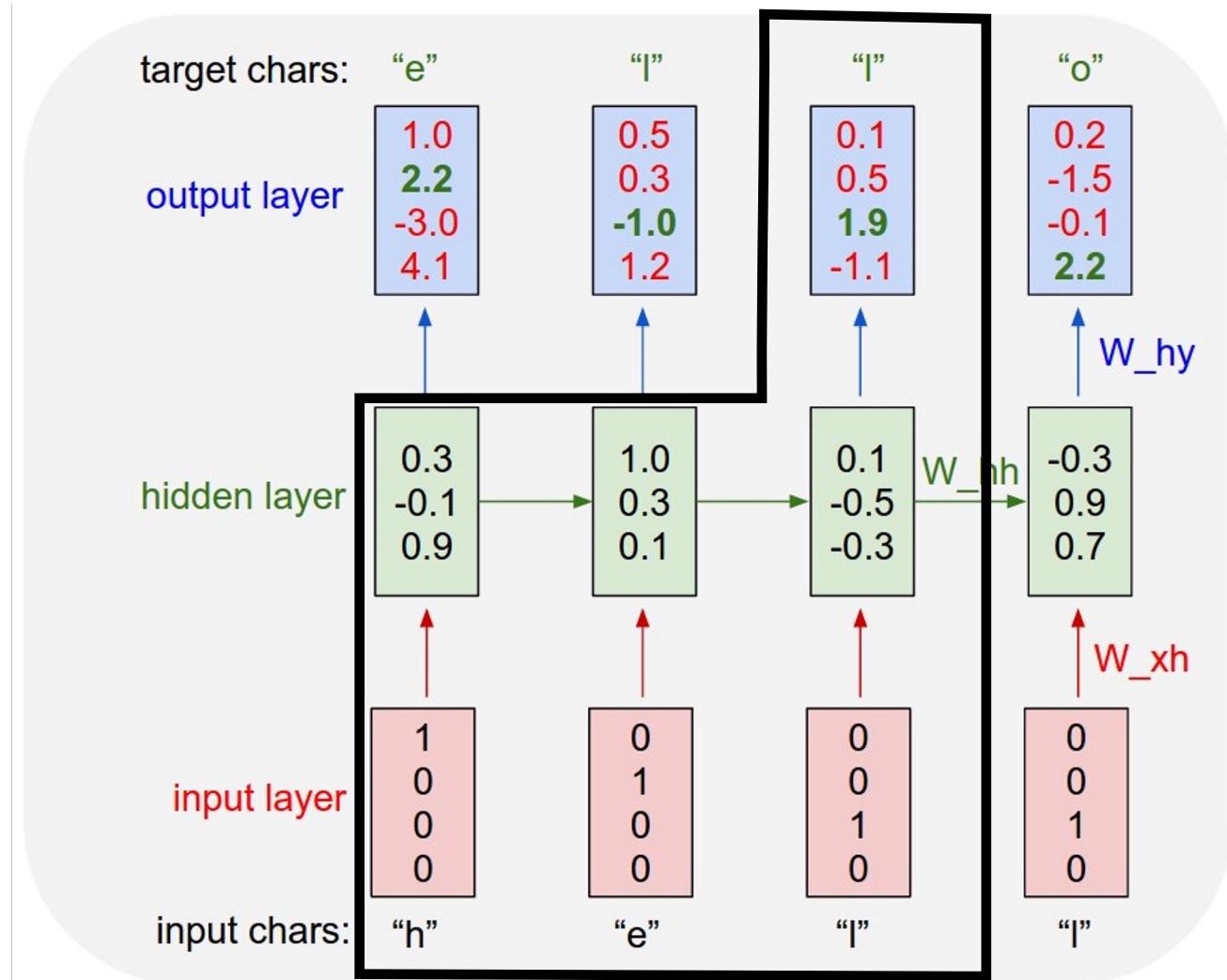
Given “hel”, predict “l”

Given characters 1, 2, ..., t-1,  
model predicts character t

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

Training sequence: “hello”

Vocabulary: [h, e, l, o]



# Example: Language Modeling

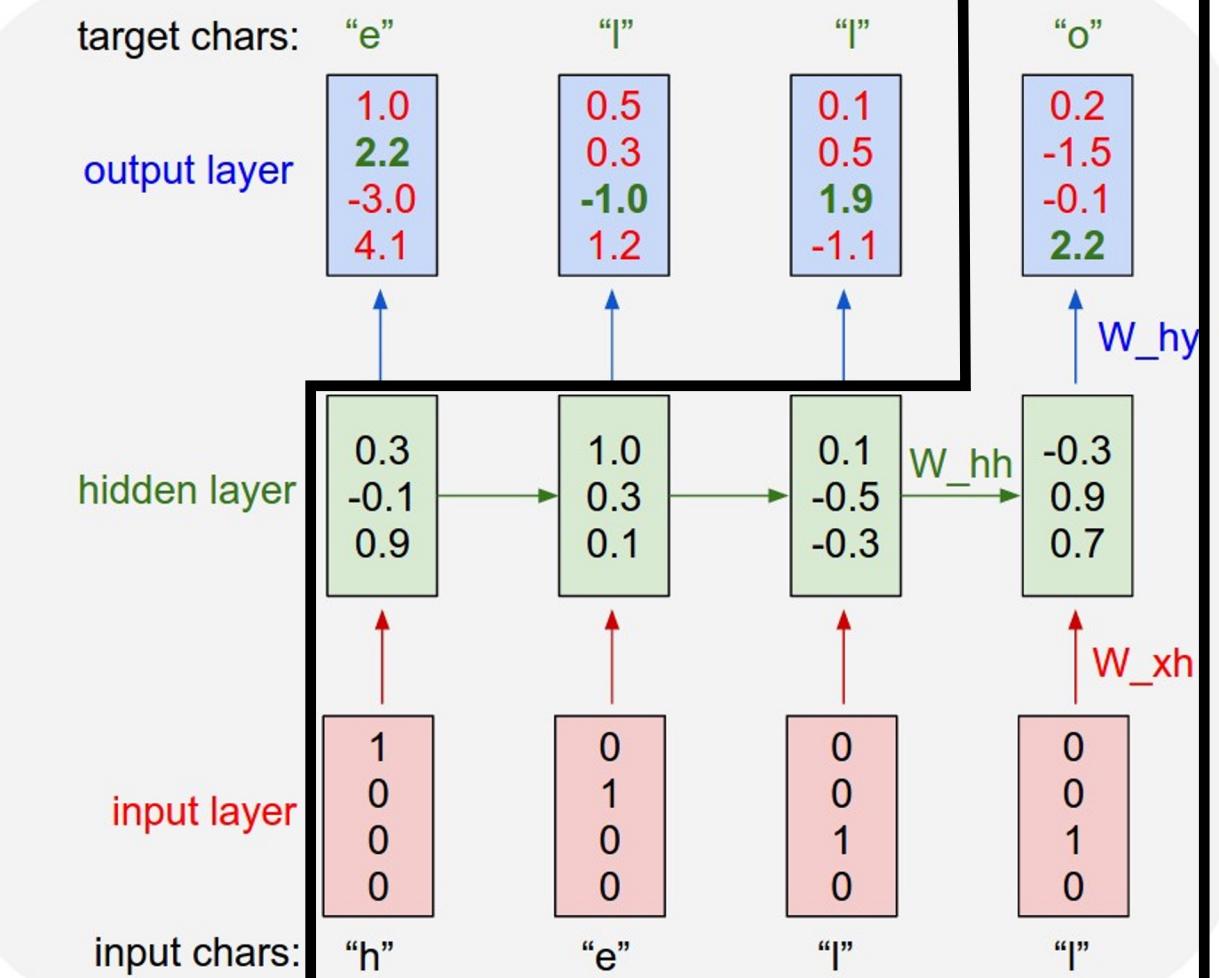
Given “hell”, predict “o”

Given characters 1, 2, ..., t-1,  
model predicts character t

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

Training sequence: “hello”

Vocabulary: [h, e, l, o]

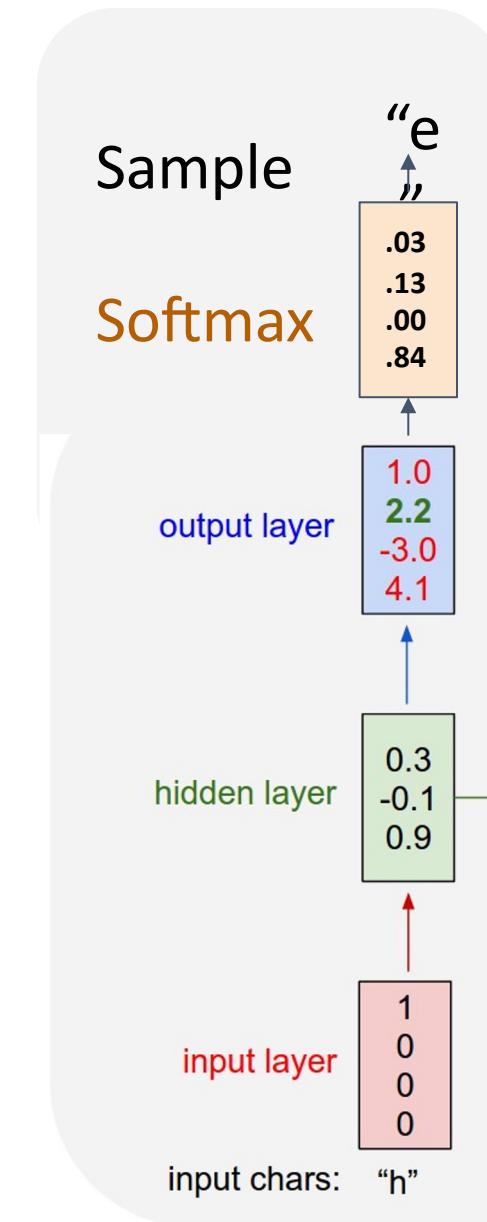


# Example: Language Modeling

At test-time, **generate** new text: sample characters one at a time, feed back to model

Training sequence: "hello"

Vocabulary: [h, e, l, o]

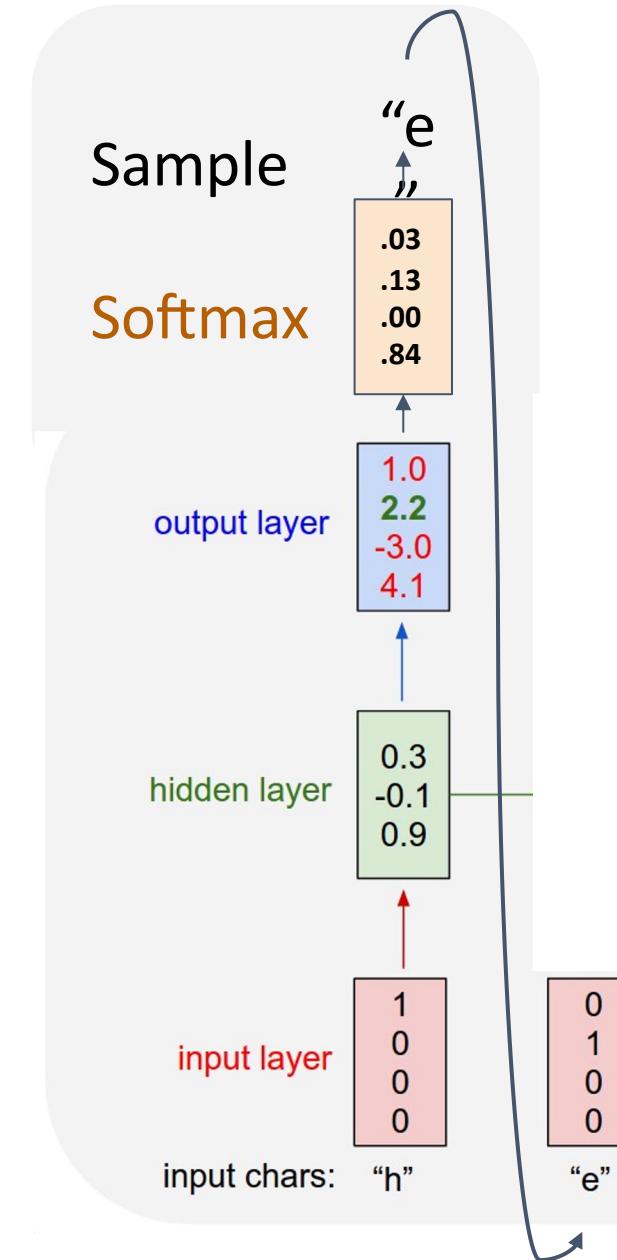


# Example: Language Modeling

At test-time, **generate** new text: sample characters one at a time, feed back to model

Training sequence: "hello"

Vocabulary: [h, e, l, o]

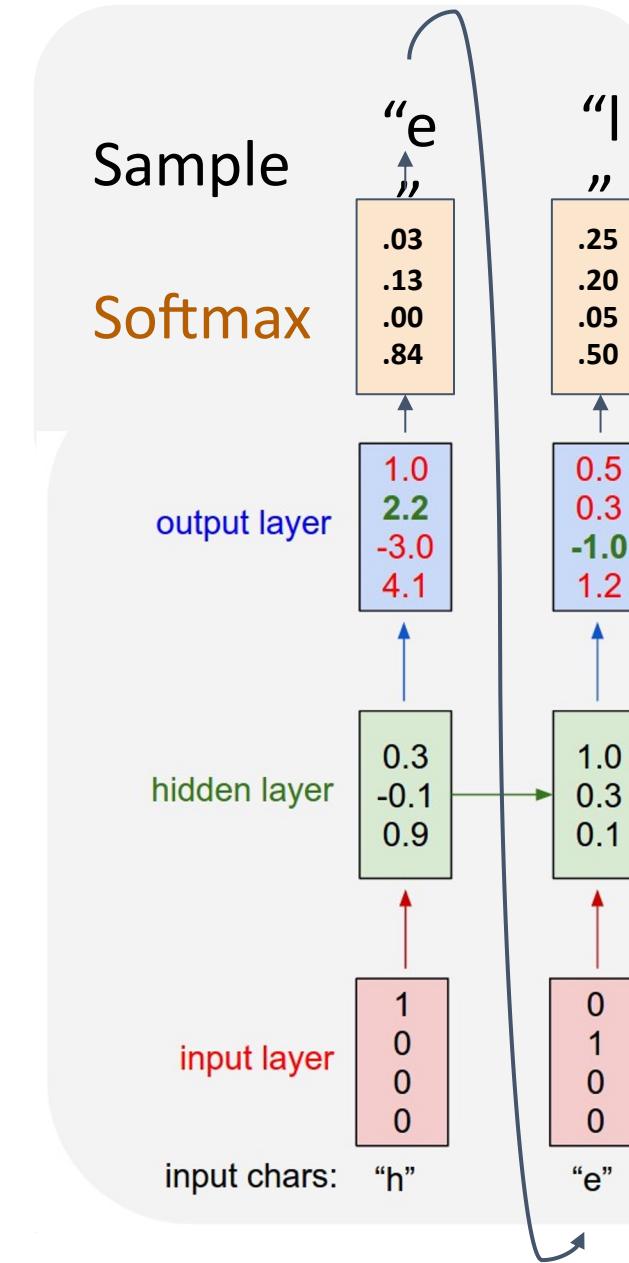


# Example: Language Modeling

At test-time, **generate** new text: sample characters one at a time, feed back to model

Training sequence: "hello"

Vocabulary: [h, e, l, o]

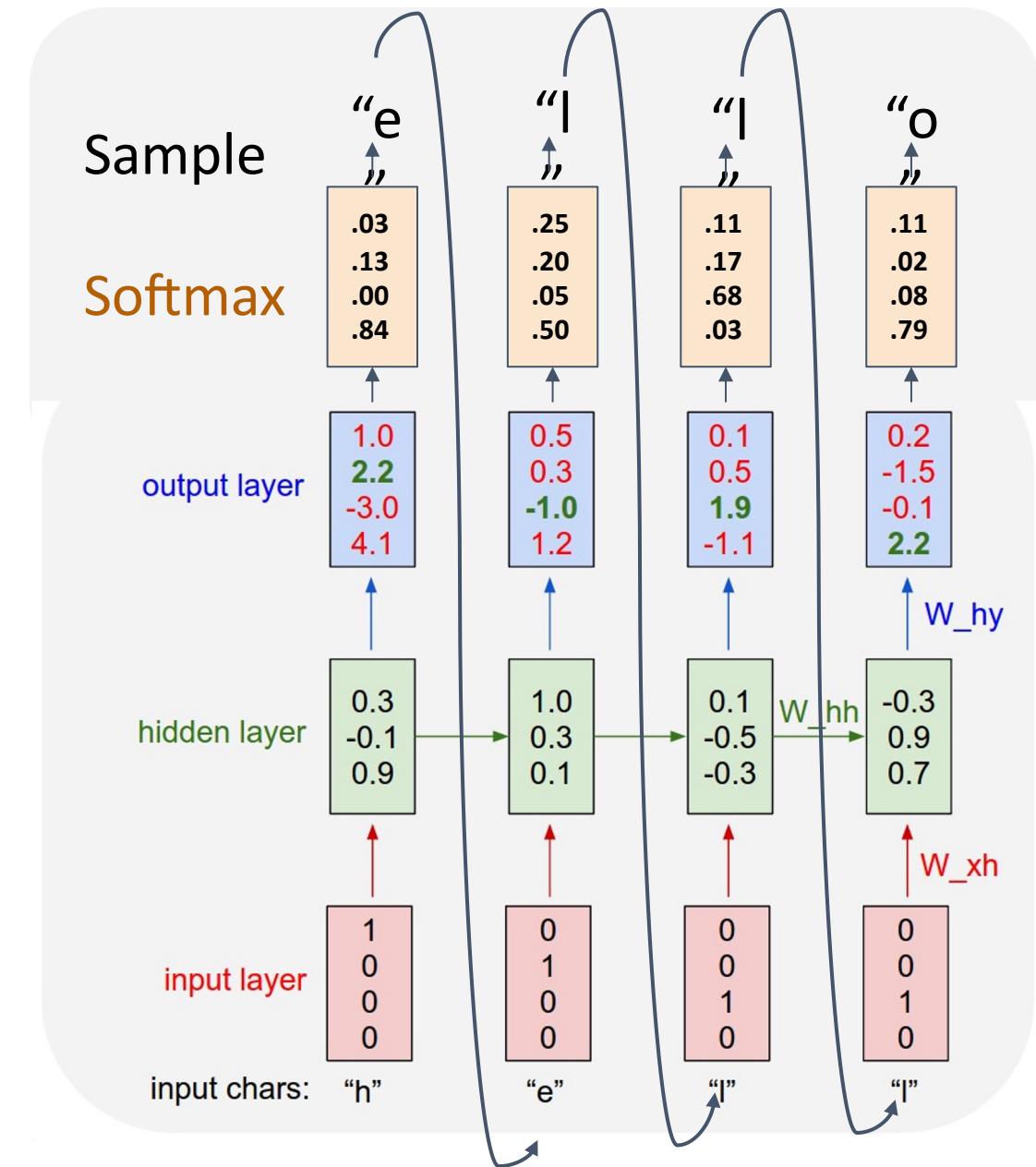


# Example: Language Modeling

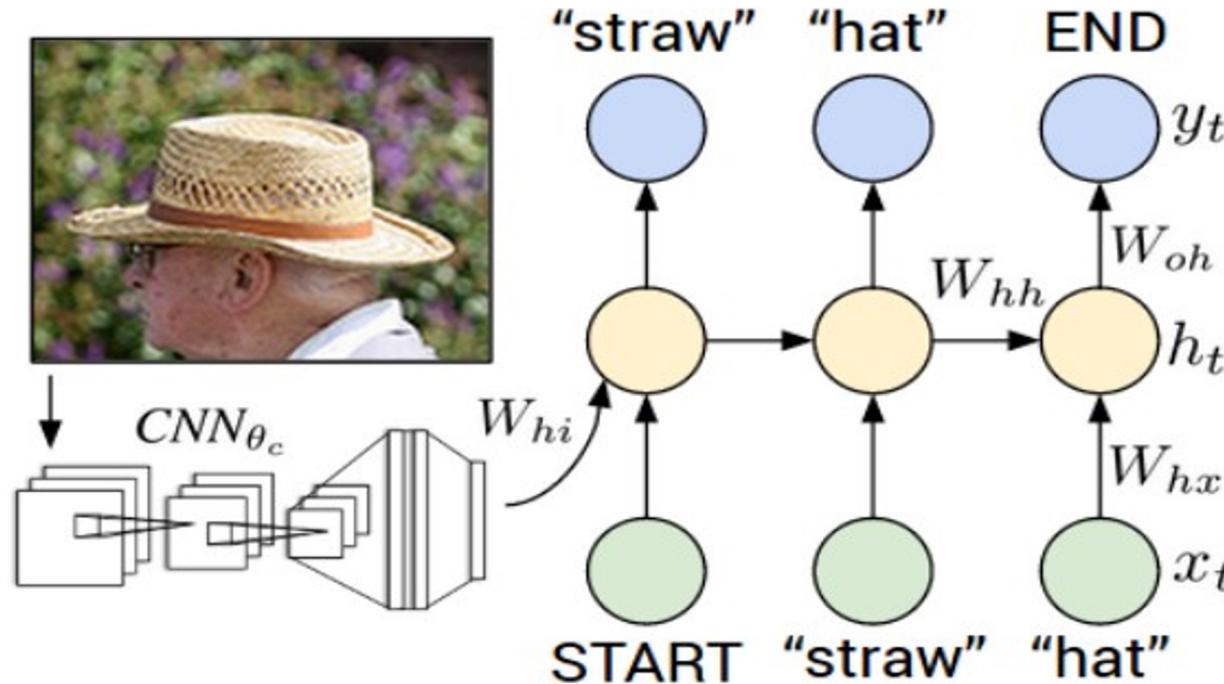
At test-time, **generate** new text: sample characters one at a time, feed back to model

Training sequence: "hello"

Vocabulary: [h, e, l, o]



# Example: Image Captioning



Mao et al, "Explain Images with Multimodal Recurrent Neural Networks", NeurIPS 2014 Deep Learning and Representation Workshop

Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015

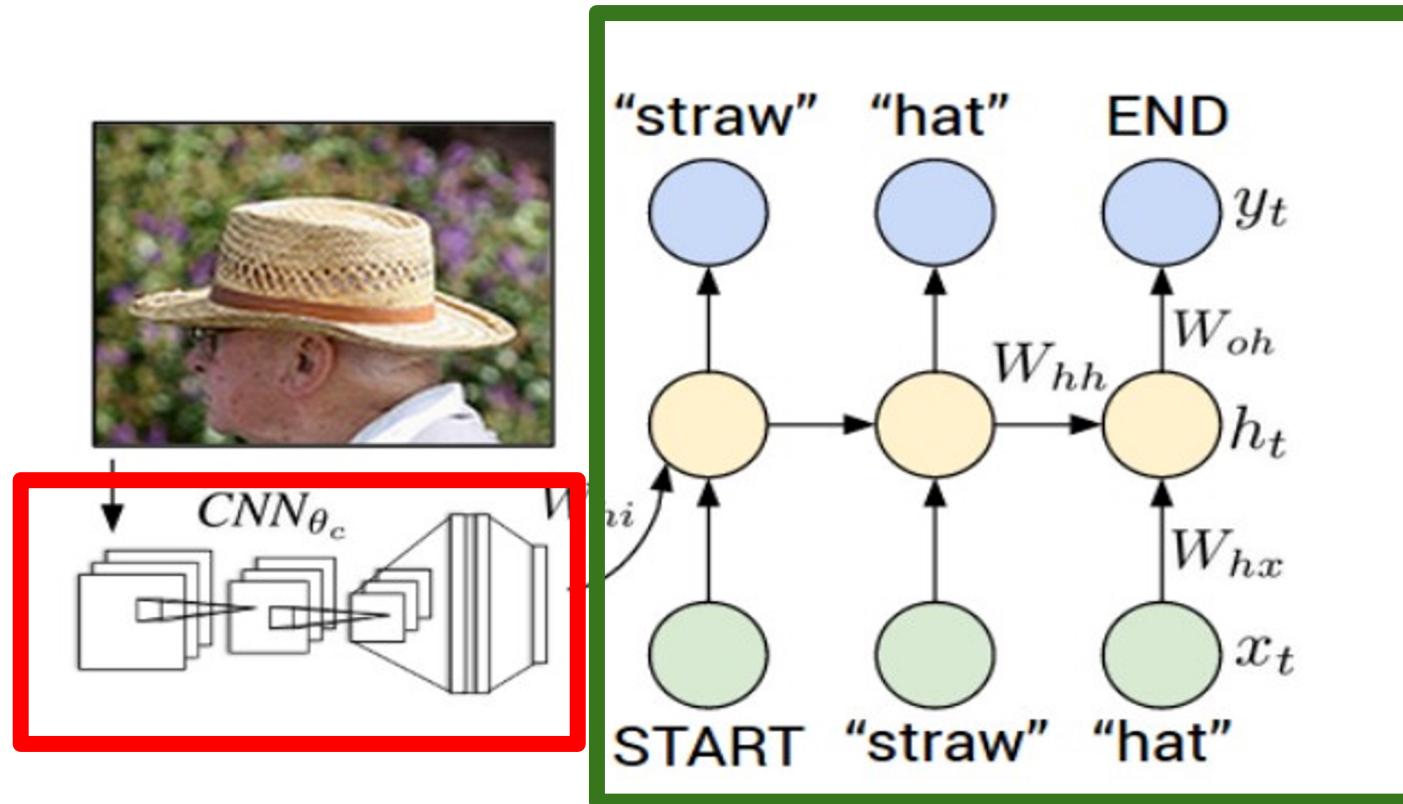
Vinyals et al, "Show and Tell: A Neural Image Caption Generator", CVPR 2015

Donahue et al, "Long-term Recurrent Convolutional Networks for Visual Recognition and Description", CVPR 2015

Chen and Zitnick, "Learning a Recurrent Visual Representation for Image Caption Generation", CVPR 2015

Figure from Karpathy et al, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015

# Example: Image Captioning



## Convolutional Neural Network

Recurrent  
Neural  
Network

Figure from Karpathy et al, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015

image

conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

FC-4096

FC-1000

softmax



**Transfer learning:** Take  
CNN trained on ImageNet,  
chop off last layer



<START>

image
conv-64
conv-64
maxpool
conv-128
conv-128
maxpool
conv-256
conv-256
maxpool
conv-512
conv-512
maxpool
conv-512
conv-512
maxpool
FC-4096
FC-4096



**Before:**

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

**Now:**

$$\tanh(W_{hh}h_{t-1} + W_{xh}x_t + W_{ih}v + b_h)$$

$W_{ih}$

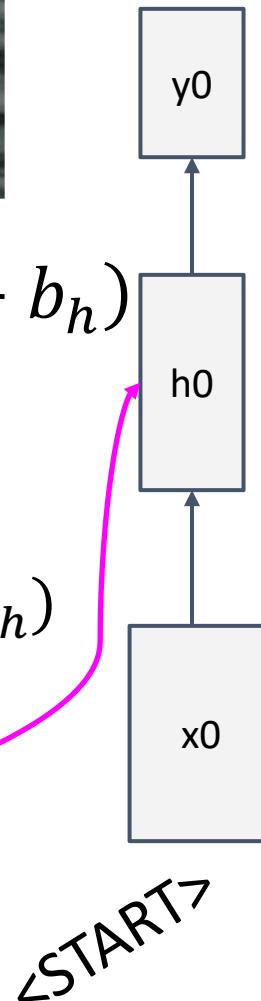


image
conv-64
conv-64
maxpool
conv-128
conv-128
maxpool
conv-256
conv-256
maxpool
conv-512
conv-512
maxpool
conv-512
conv-512
maxpool
FC-4096
FC-4096



**Before:**

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

**Now:**

$$\tanh(W_{hh}h_{t-1} + W_{xh}x_t + W_{ih}v + b_h)$$

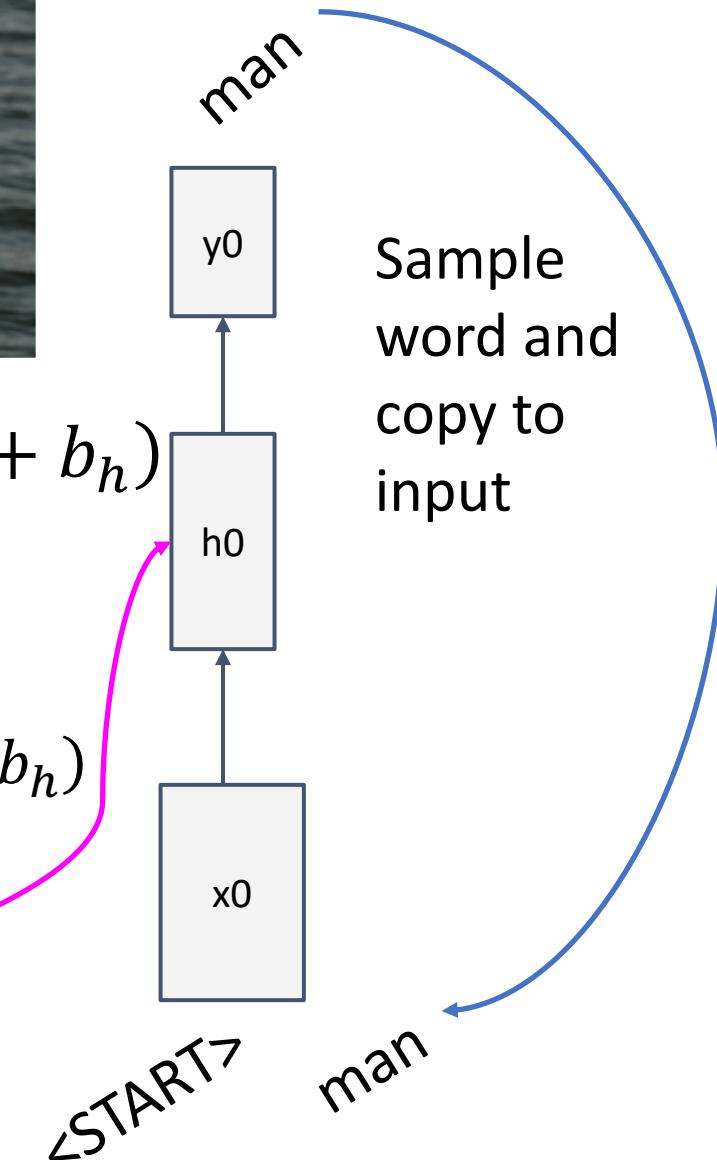


image
conv-64
conv-64
maxpool
conv-128
conv-128
maxpool
conv-256
conv-256
maxpool
conv-512
conv-512
maxpool
conv-512
conv-512
maxpool
FC-4096
FC-4096



**Before:**

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

**Now:**

$$\tanh(W_{hh}h_{t-1} + W_{xh}x_t + W_{ih}v + b_h)$$

$W_{ih}$

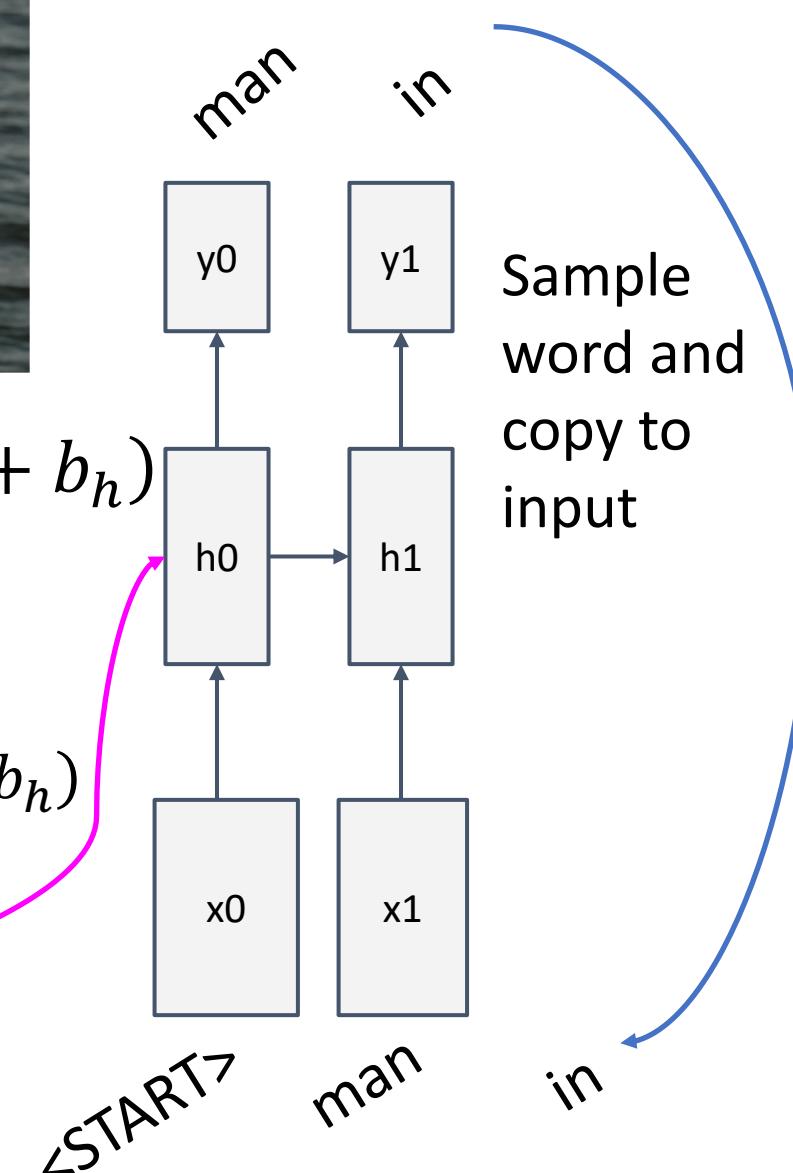


image
conv-64
conv-64
maxpool
conv-128
conv-128
maxpool
conv-256
conv-256
maxpool
conv-512
conv-512
maxpool
conv-512
conv-512
maxpool
FC-4096
FC-4096



Before:

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

Now:

$$\tanh(W_{hh}h_{t-1} + W_{xh}x_t + W_{ih}v + b_h)$$

$W_{ih}$

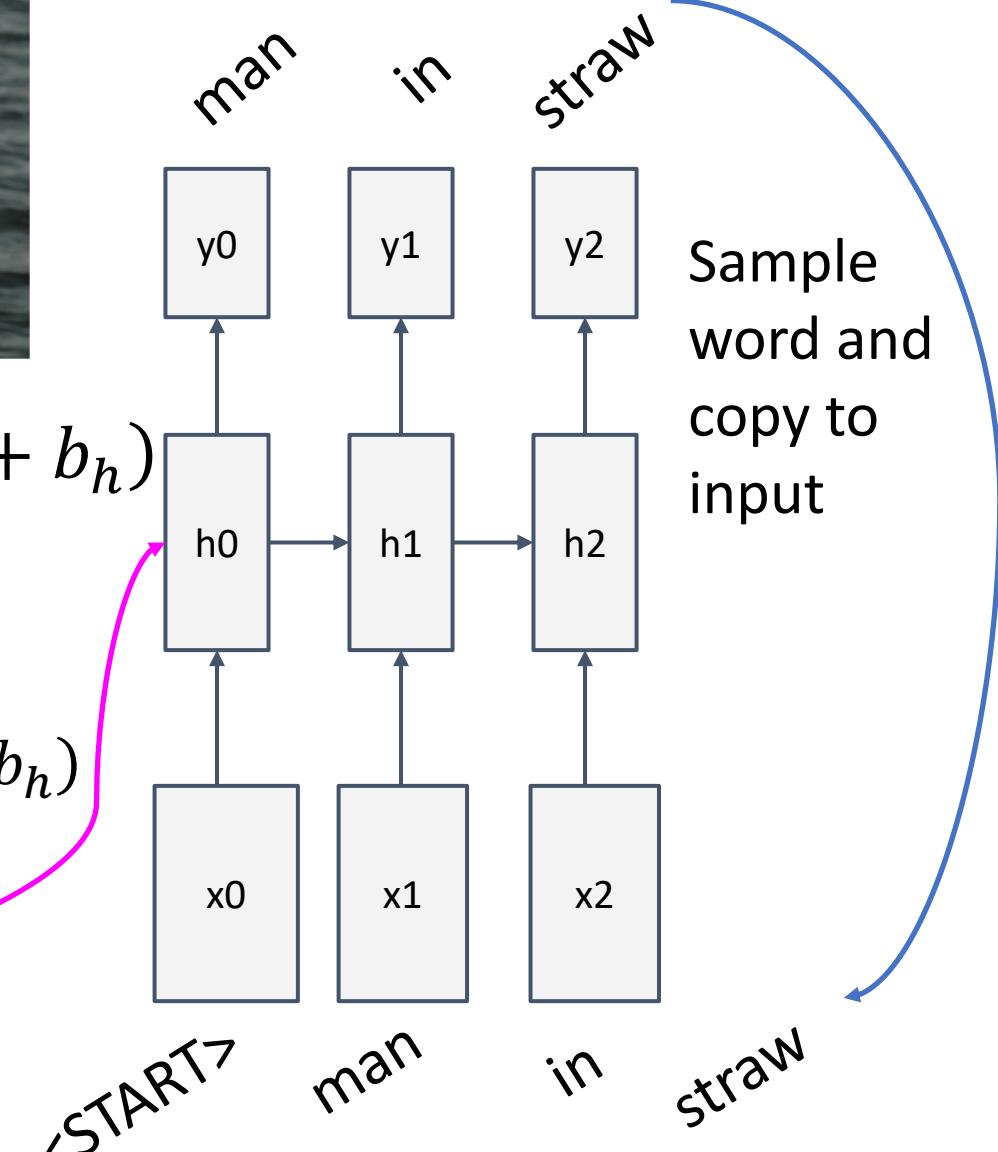


image
conv-64
conv-64
maxpool
conv-128
conv-128
maxpool
conv-256
conv-256
maxpool
conv-512
conv-512
maxpool
conv-512
conv-512
maxpool
FC-4096
FC-4096



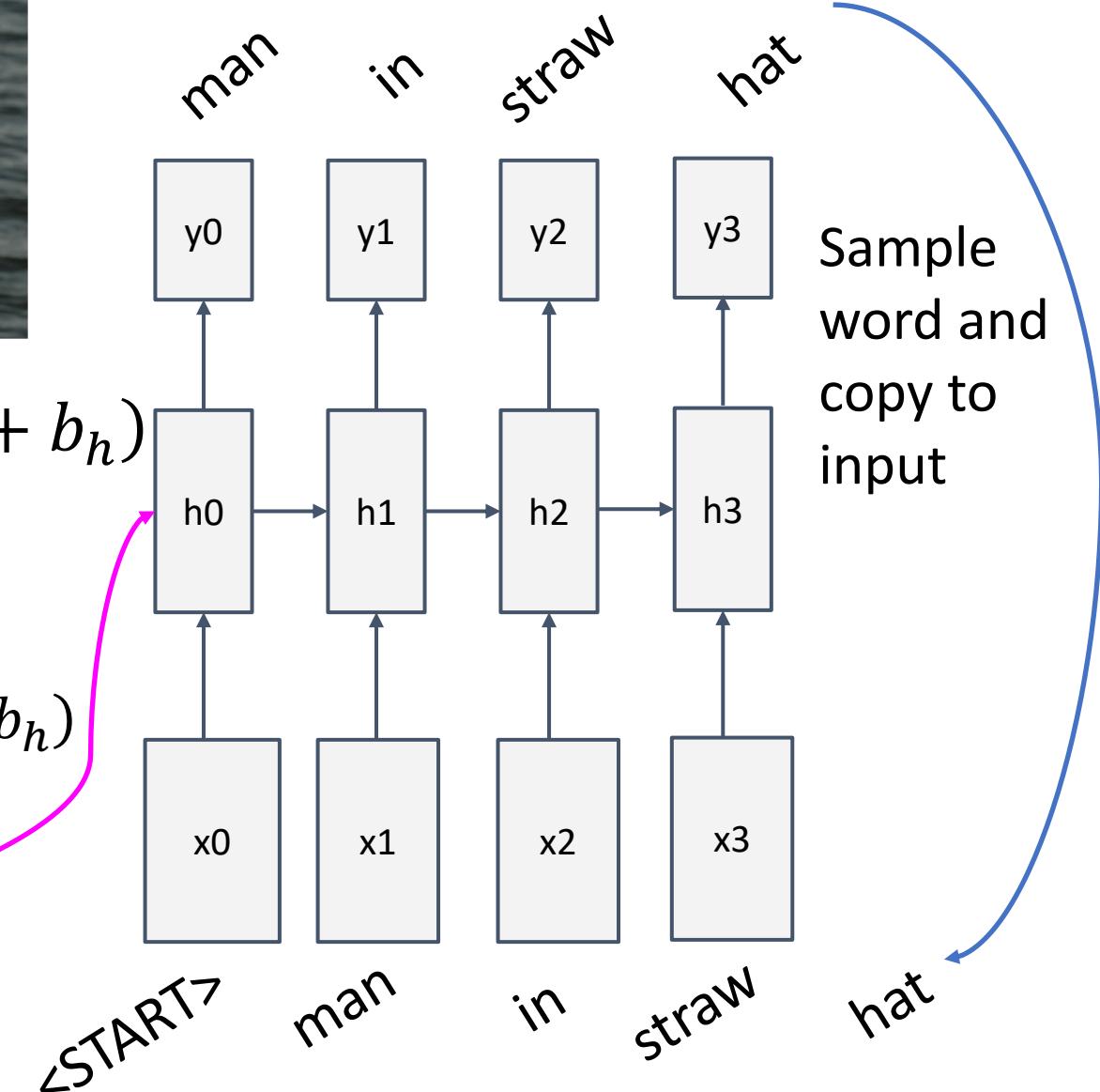
**Before:**

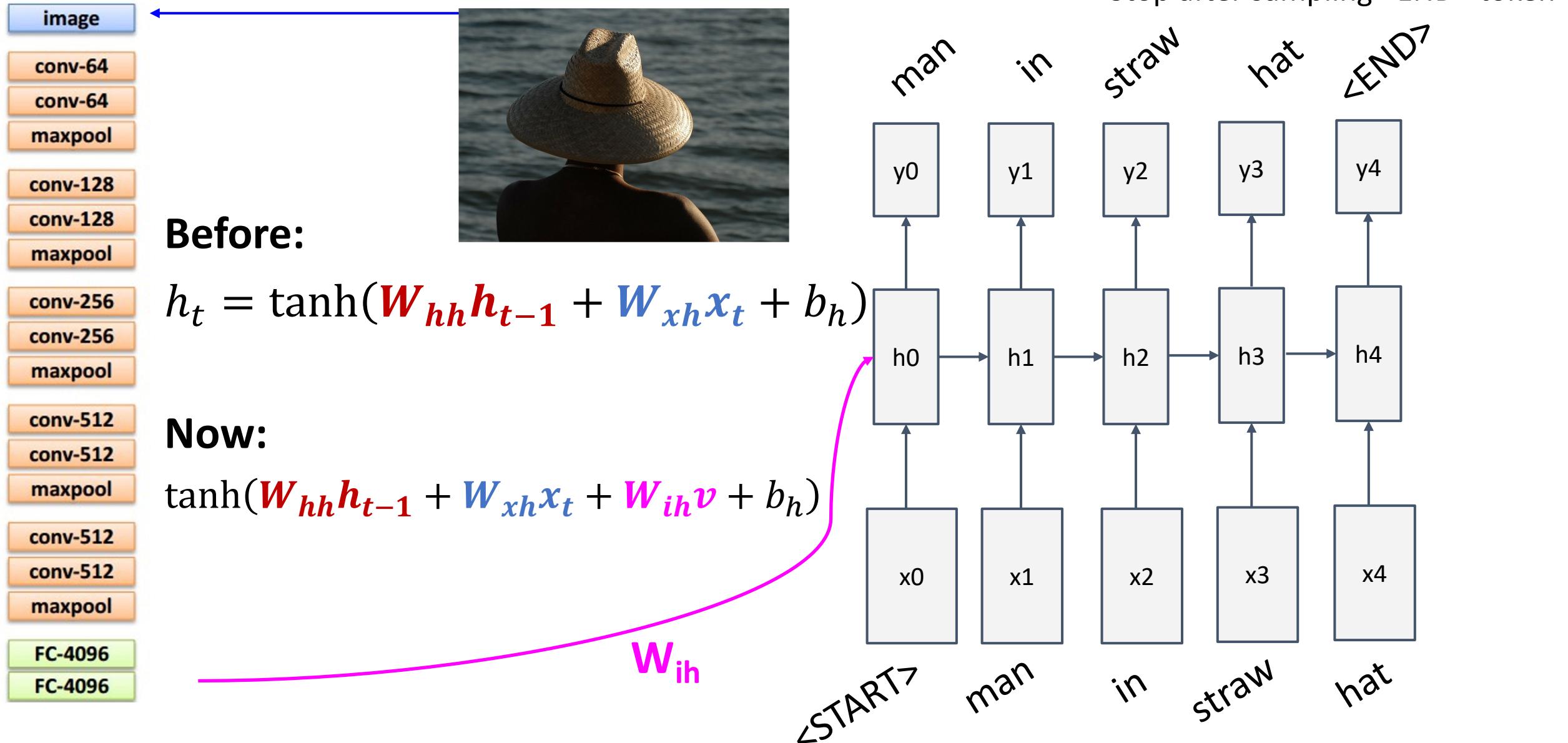
$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$$

**Now:**

$$\tanh(W_{hh}h_{t-1} + W_{xh}x_t + W_{ih}v + b_h)$$

$W_{ih}$





# Image Captioning: Example Results



*A cat sitting on a suitcase on the floor*



*A cat is sitting on a tree branch*



*A dog is running in the grass with a frisbee*



*A white teddy bear sitting in the grass*



*Two people walking on the beach with surfboards*



*A tennis player in action on the court*



*Two giraffes standing in a grassy field*



*A man riding a dirt bike on a dirt track*

# Image Captioning: Failure Cases



*A woman is holding a cat  
in her hand*



*A person holding a computer  
mouse on a desk*



*A woman standing on a beach  
holding a surfboard*



*A bird is perched on a  
tree branch*



*A man in a  
baseball uniform  
throwing a ball*

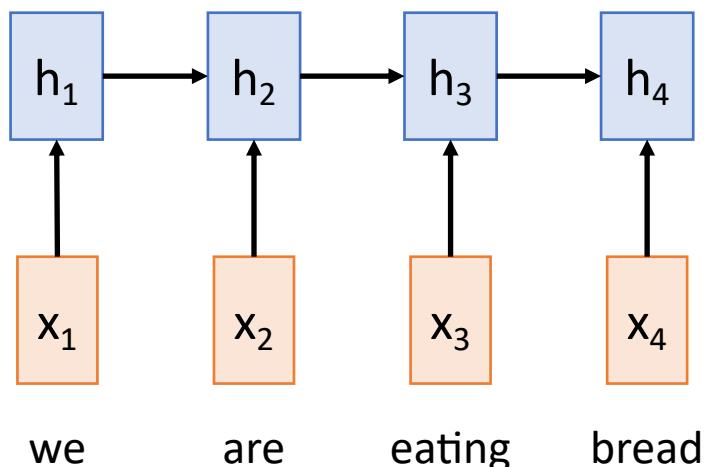
# RNN with Attention

# Sequence-to-Sequence with RNNs

**Input:** Sequence  $x_1, \dots, x_T$

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

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



# Sequence-to-Sequence with RNNs

**Input:** Sequence  $x_1, \dots, x_T$

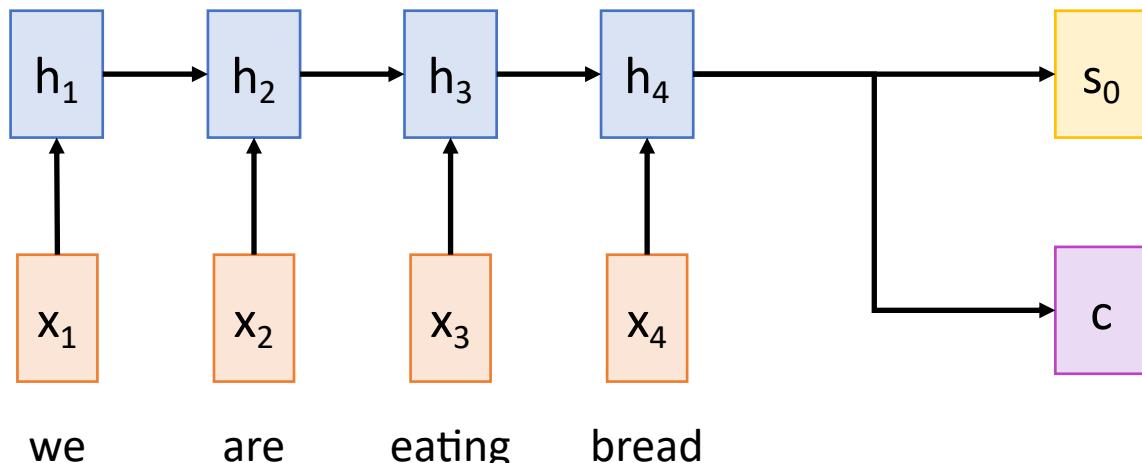
**Output:** Sequence  $y_1, \dots, y_{T'}$

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

From final hidden state predict:

**Initial decoder state**  $s_0$

**Context vector**  $c$  (often  $c=h_T$ )



# Sequence-to-Sequence with RNNs

**Input:** Sequence  $x_1, \dots, x_T$

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

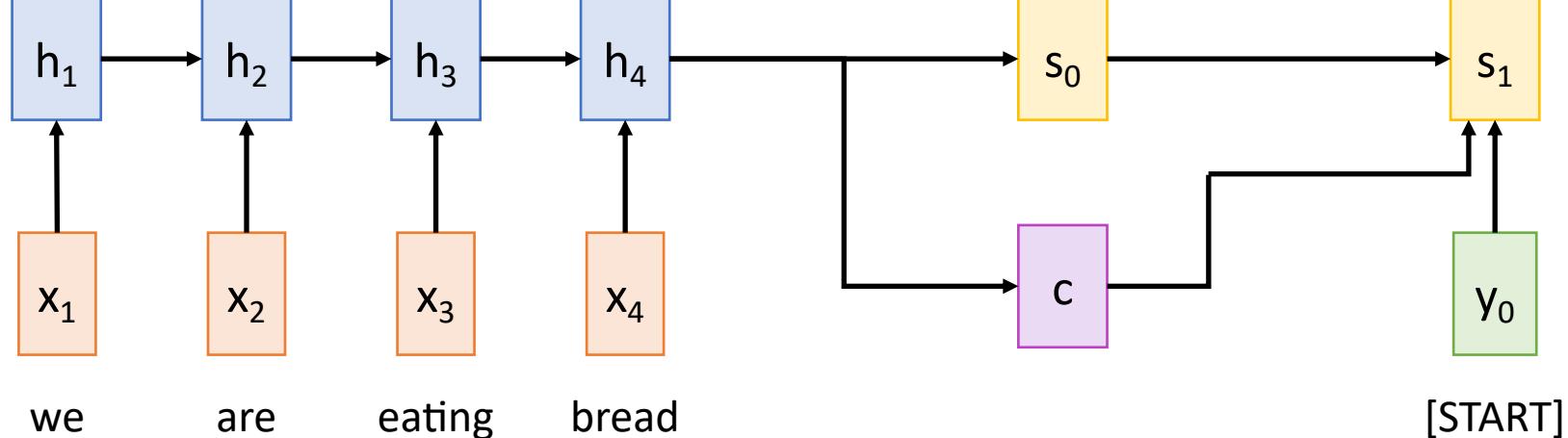
**Decoder:**  $s_t = g_U(y_{t-1}, s_{t-1}, c)$

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

From final hidden state predict:

**Initial decoder state**  $s_0$

**Context vector**  $c$  (often  $c=h_T$ )



# Sequence-to-Sequence with RNNs

**Input:** Sequence  $x_1, \dots, x_T$

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

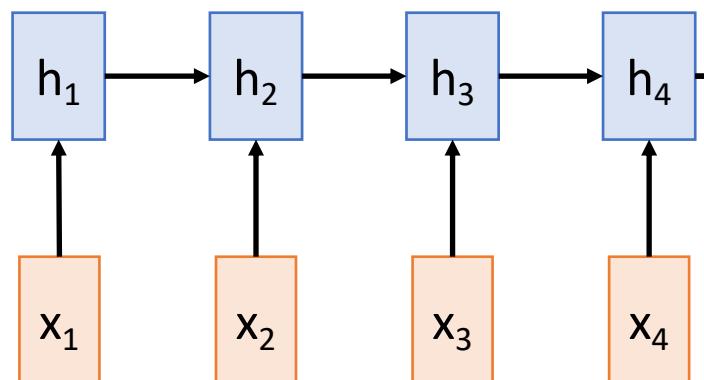
**Decoder:**  $s_t = g_U(y_{t-1}, s_{t-1}, c)$

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

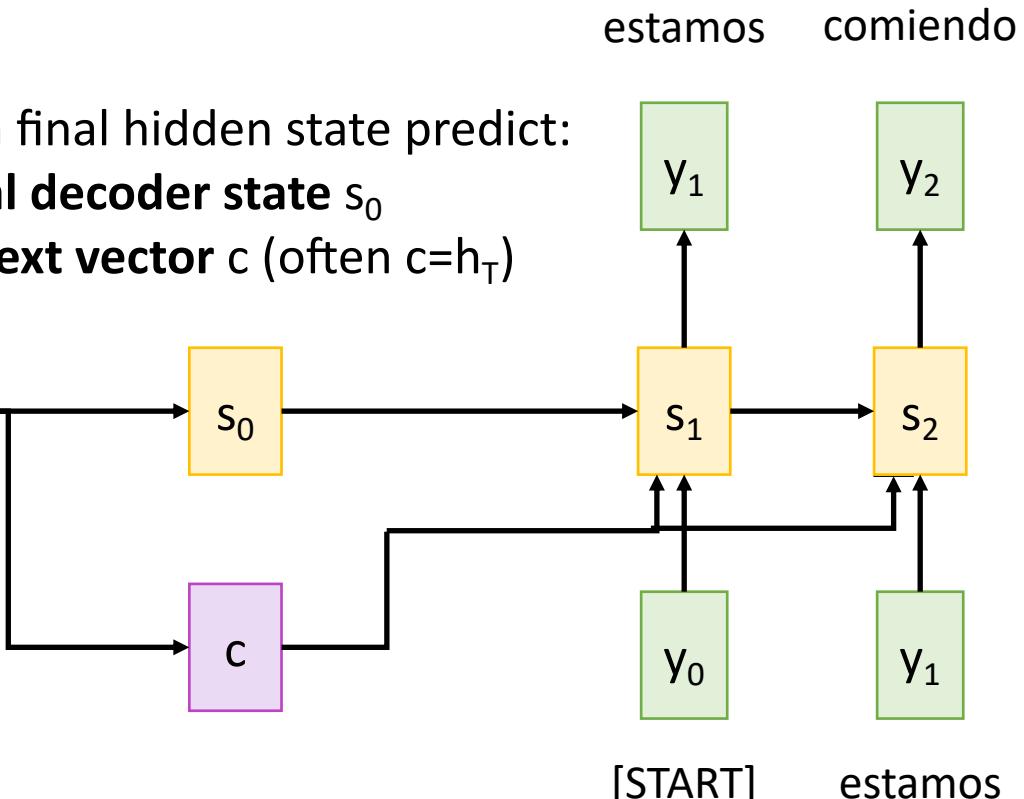
From final hidden state predict:

**Initial decoder state**  $s_0$

**Context vector**  $c$  (often  $c=h_T$ )



we      are      eating      bread



[START]      estamos

# Sequence-to-Sequence with RNNs

**Input:** Sequence  $x_1, \dots, x_T$

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

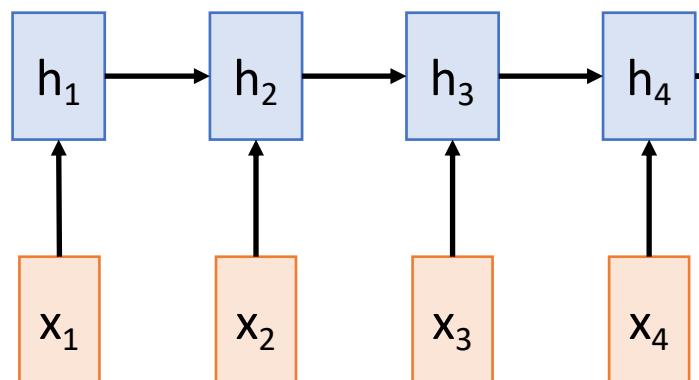
**Decoder:**  $s_t = g_U(y_{t-1}, s_{t-1}, c)$

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

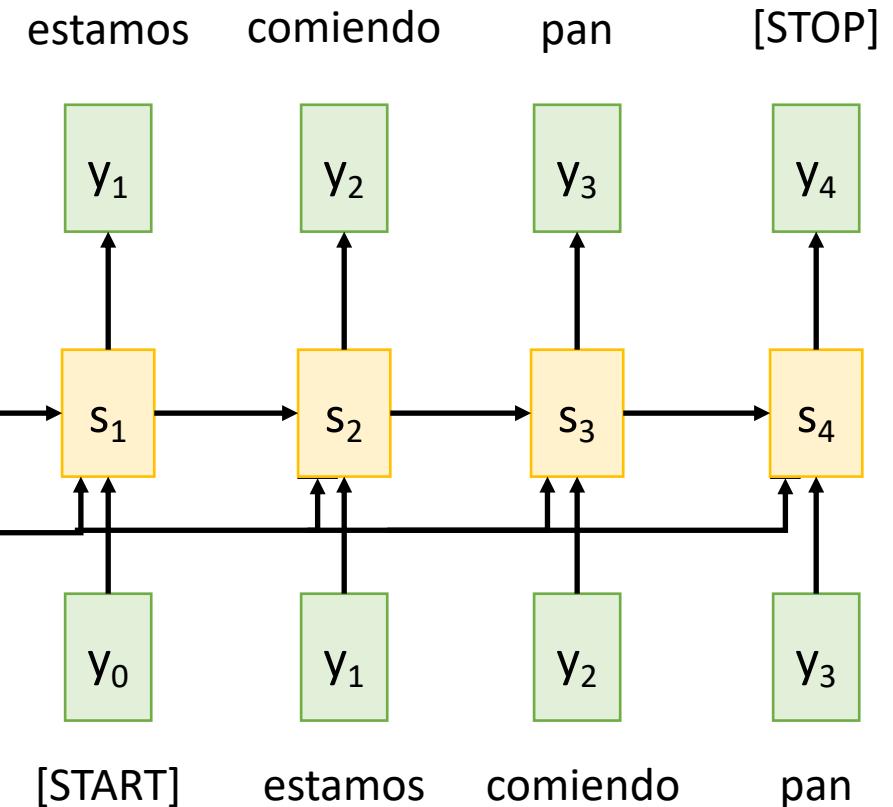
From final hidden state predict:

**Initial decoder state**  $s_0$

**Context vector**  $c$  (often  $c=h_T$ )



we      are      eating      bread



# Sequence-to-Sequence with RNNs

**Input:** Sequence  $x_1, \dots, x_T$

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

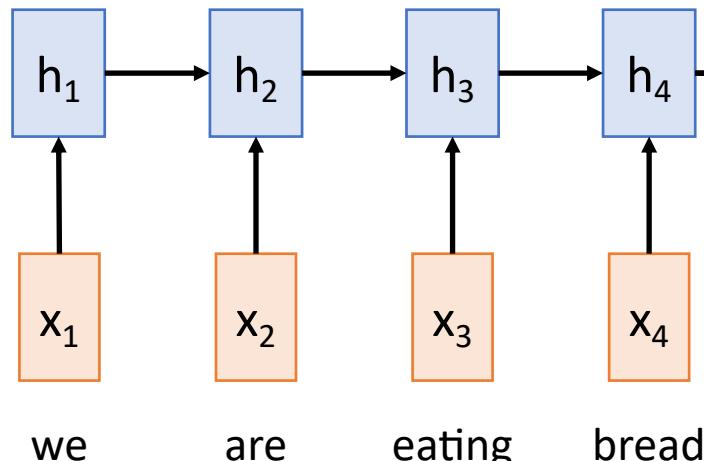
**Decoder:**  $s_t = g_U(y_{t-1}, s_{t-1}, c)$

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

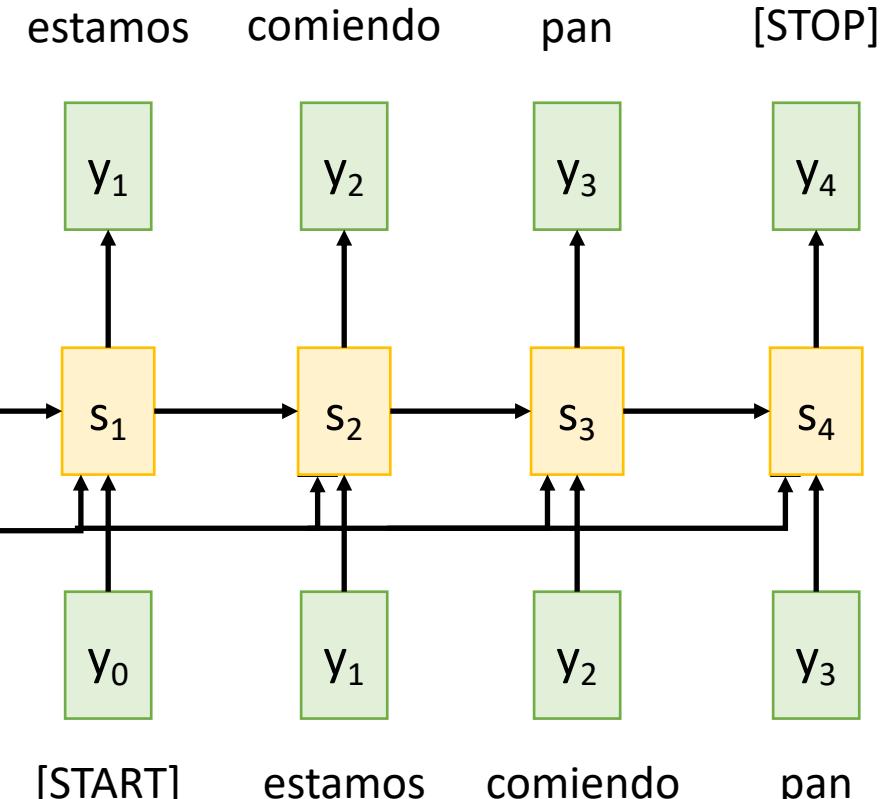
From final hidden state predict:

**Initial decoder state**  $s_0$

**Context vector**  $c$  (often  $c=h_T$ )



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

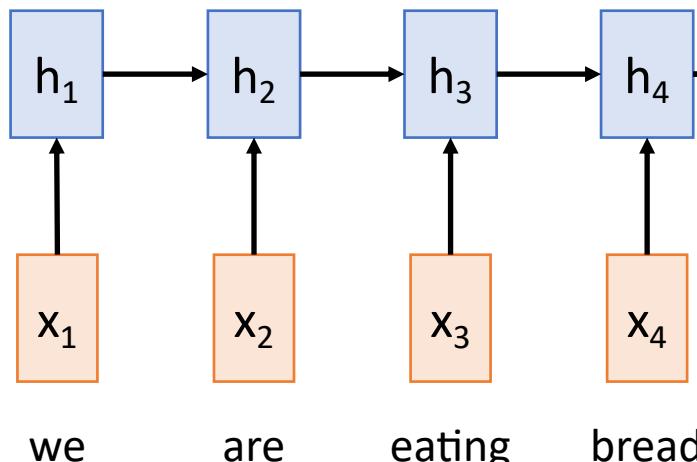


# Sequence-to-Sequence with RNNs

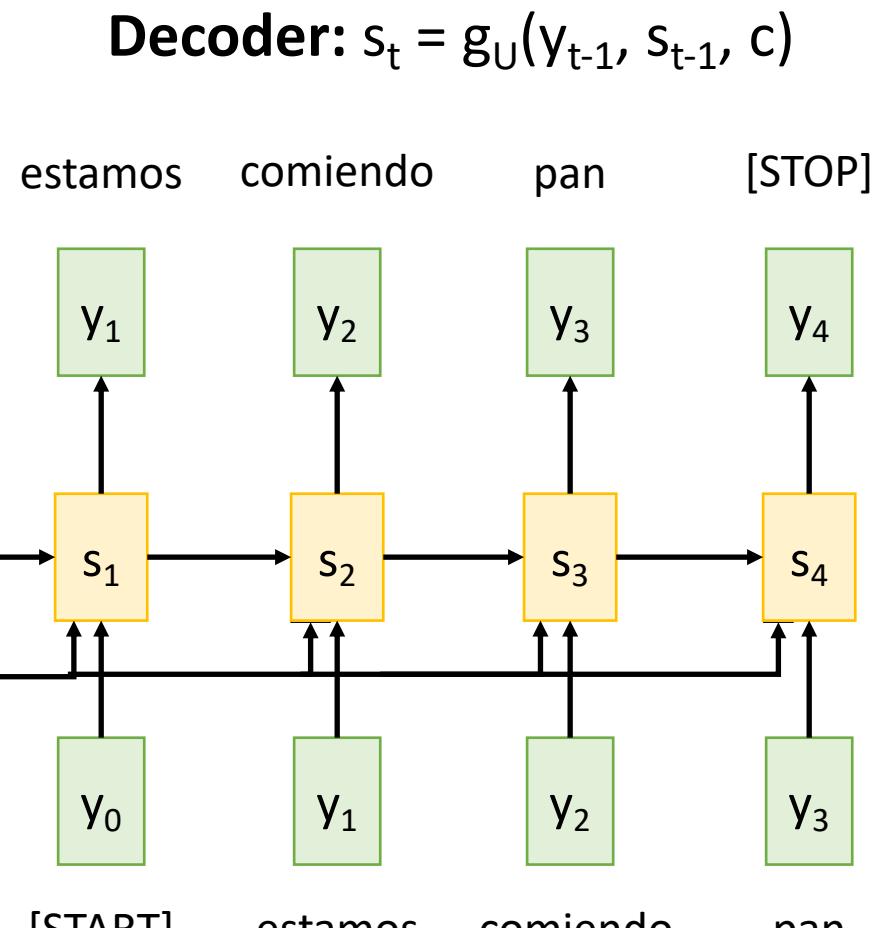
**Input:** Sequence  $x_1, \dots, x_T$

**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  $c=h_T$ )



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

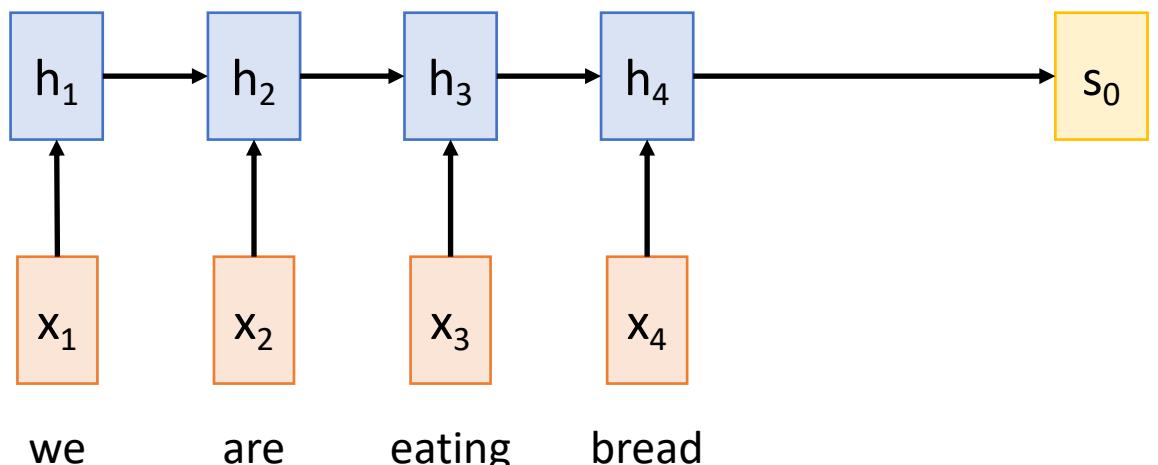
**Idea: use new context vector  
at each step of decoder!**

# Sequence-to-Sequence with RNNs and Attention

**Input:** Sequence  $x_1, \dots, x_T$

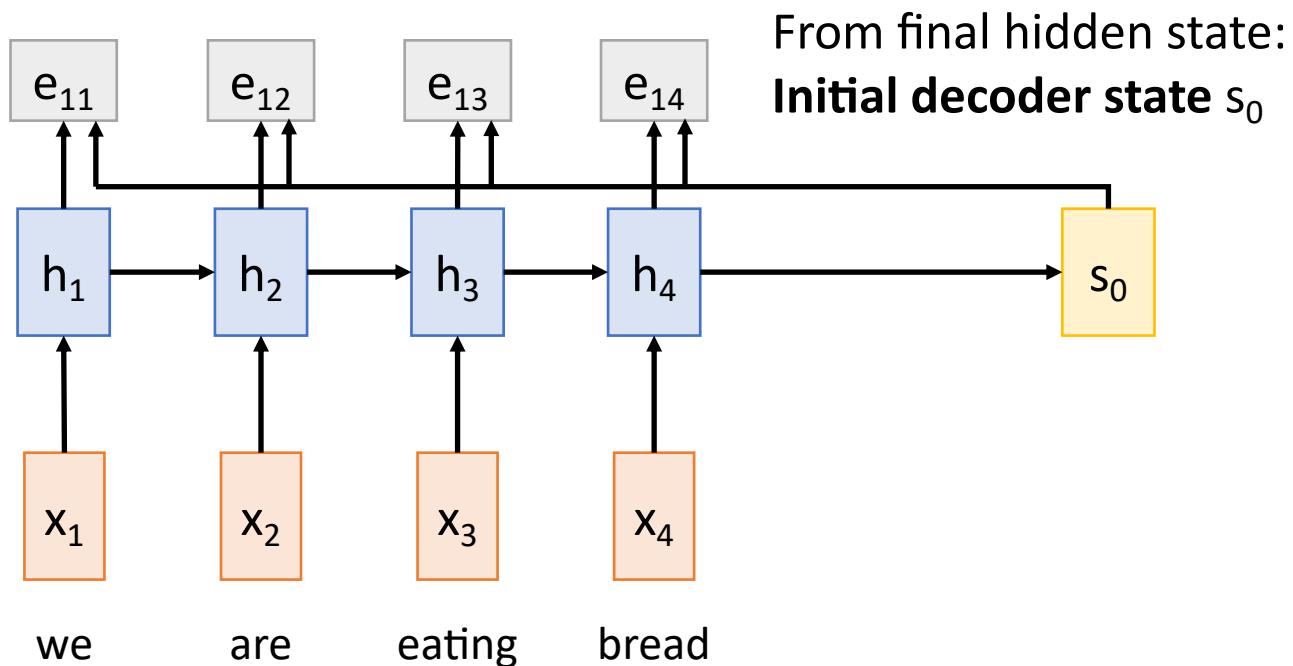
**Output:** Sequence  $y_1, \dots, y_{T'}$

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

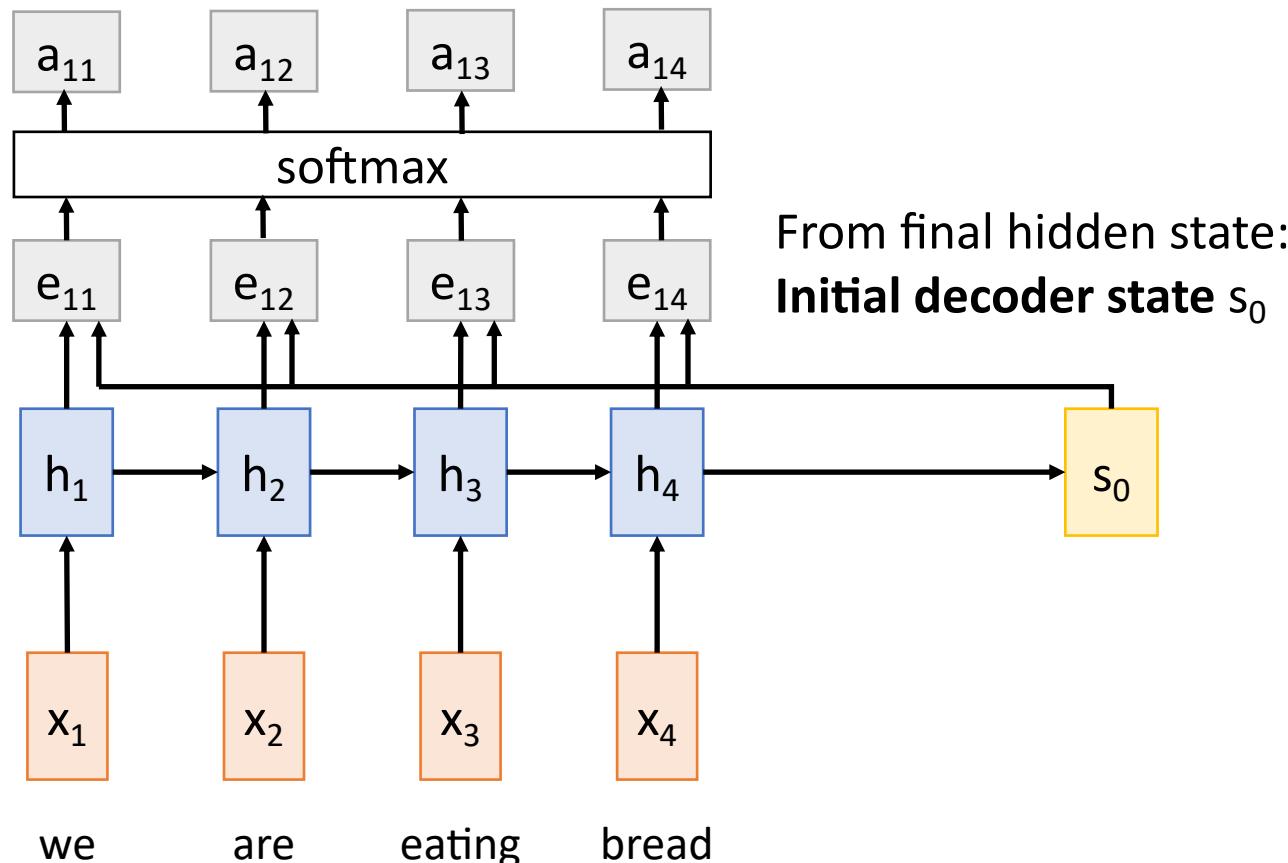


# Sequence-to-Sequence with RNNs and Attention

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



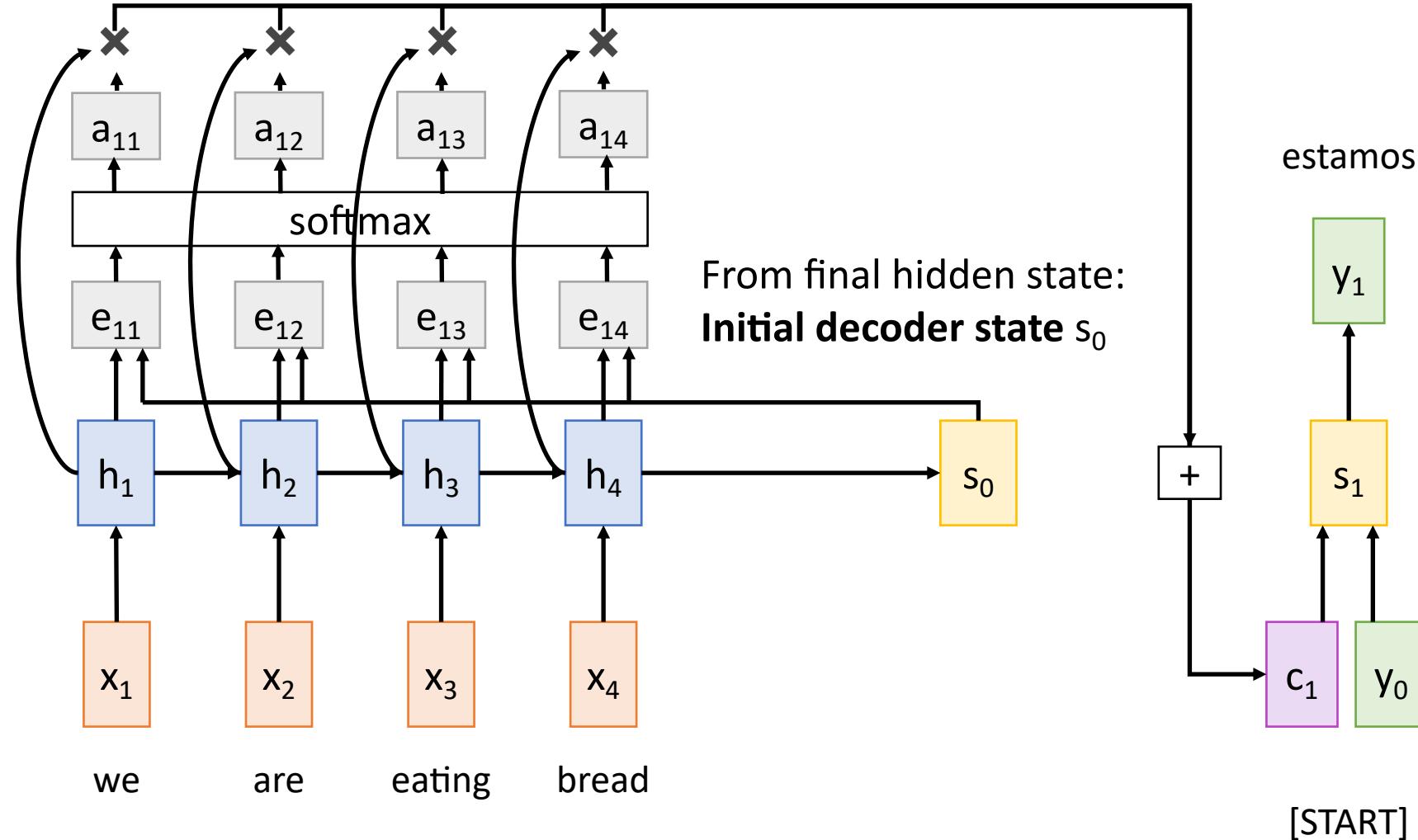
# Sequence-to-Sequence with RNNs and Attention



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

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

# Sequence-to-Sequence with RNNs and Attention



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

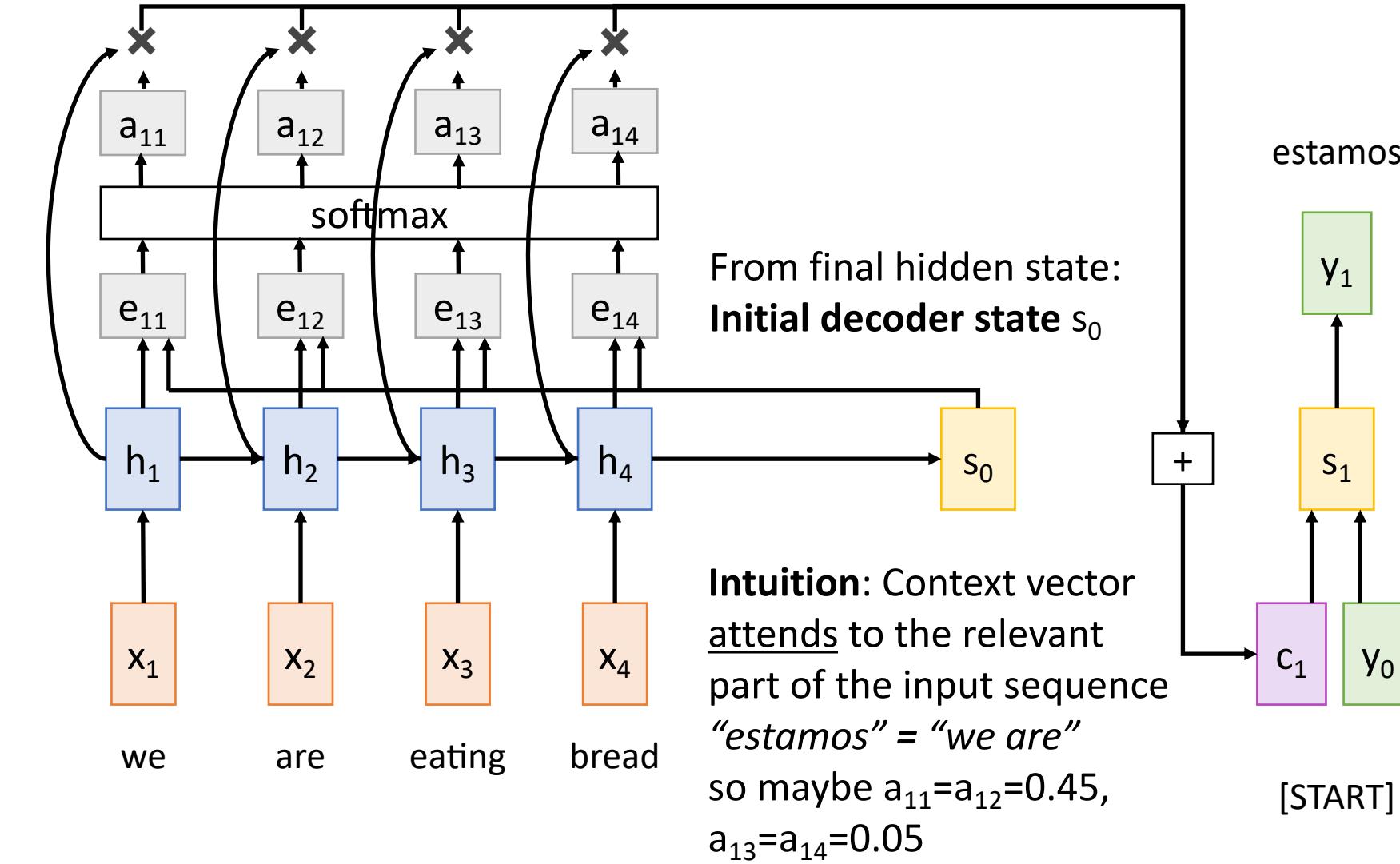
Normalize alignment scores  
to get **attention weights**  
 $0 < a_{t,i} < 1 \quad \sum_i a_{t,i} = 1$

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

Use context vector in  
decoder:  $s_t = g_U(y_{t-1}, s_{t-1}, c_t)$

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

# Sequence-to-Sequence with RNNs and Attention



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

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

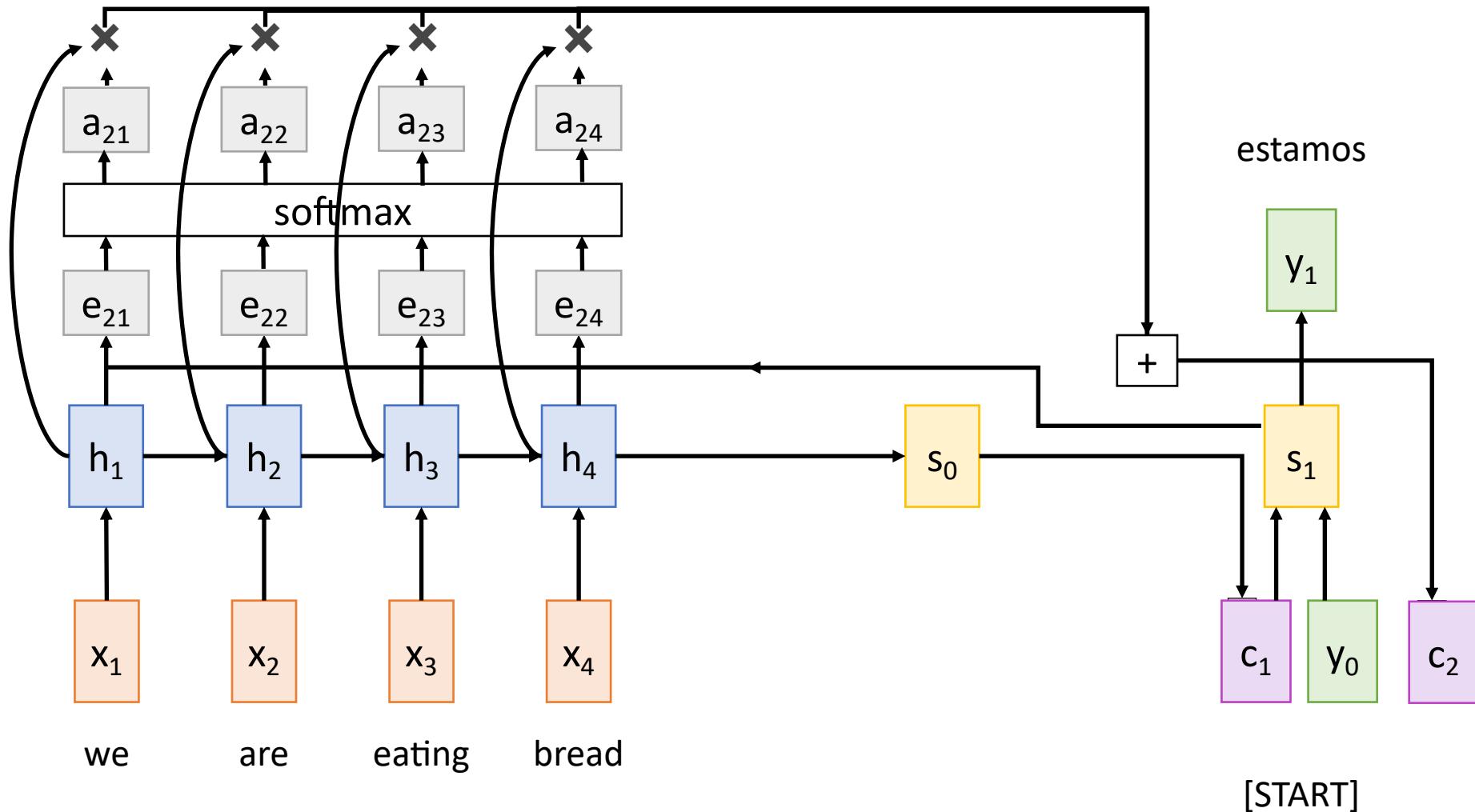
Compute context vector as linear  
combination of hidden states  
 $c_t = \sum_i a_{t,i} h_i$

Use context vector in  
decoder:  $s_t = g_U(y_{t-1}, s_{t-1}, c_t)$

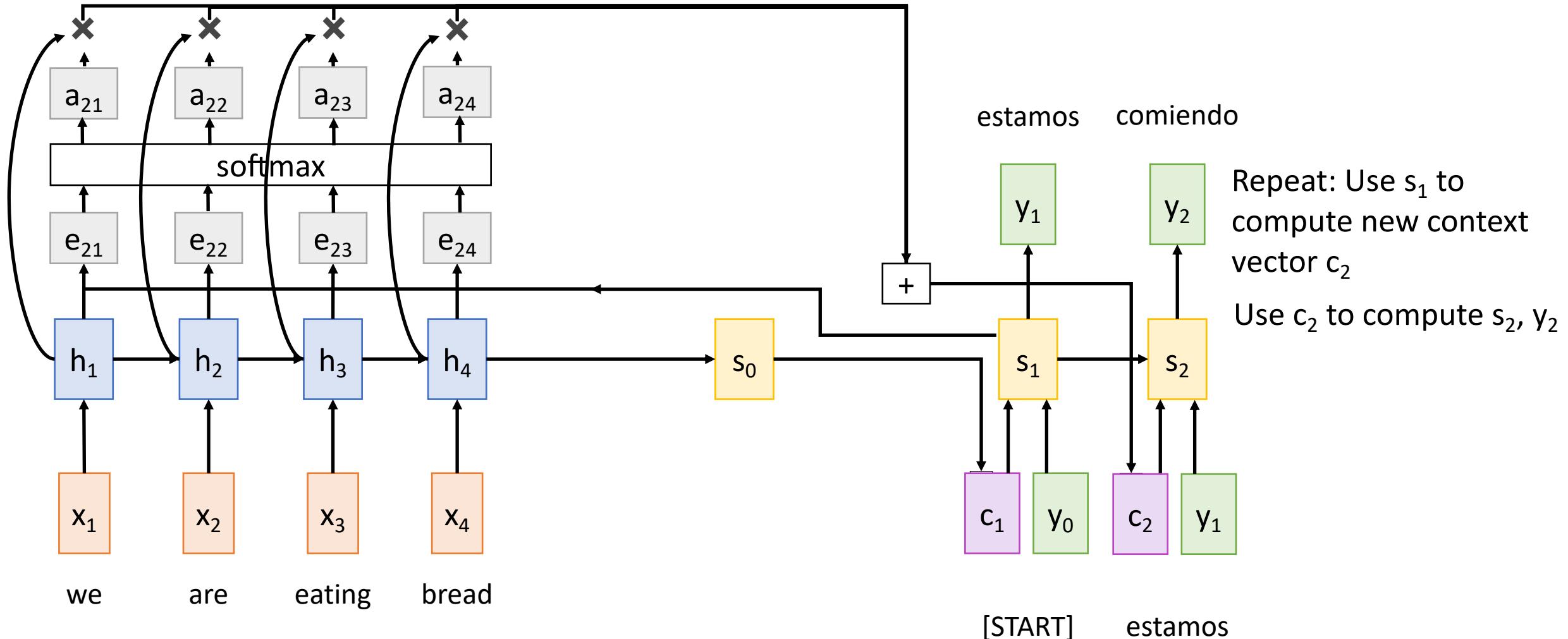
This is all differentiable! Do not  
supervise attention weights –  
backprop through everything

# Sequence-to-Sequence with RNNs

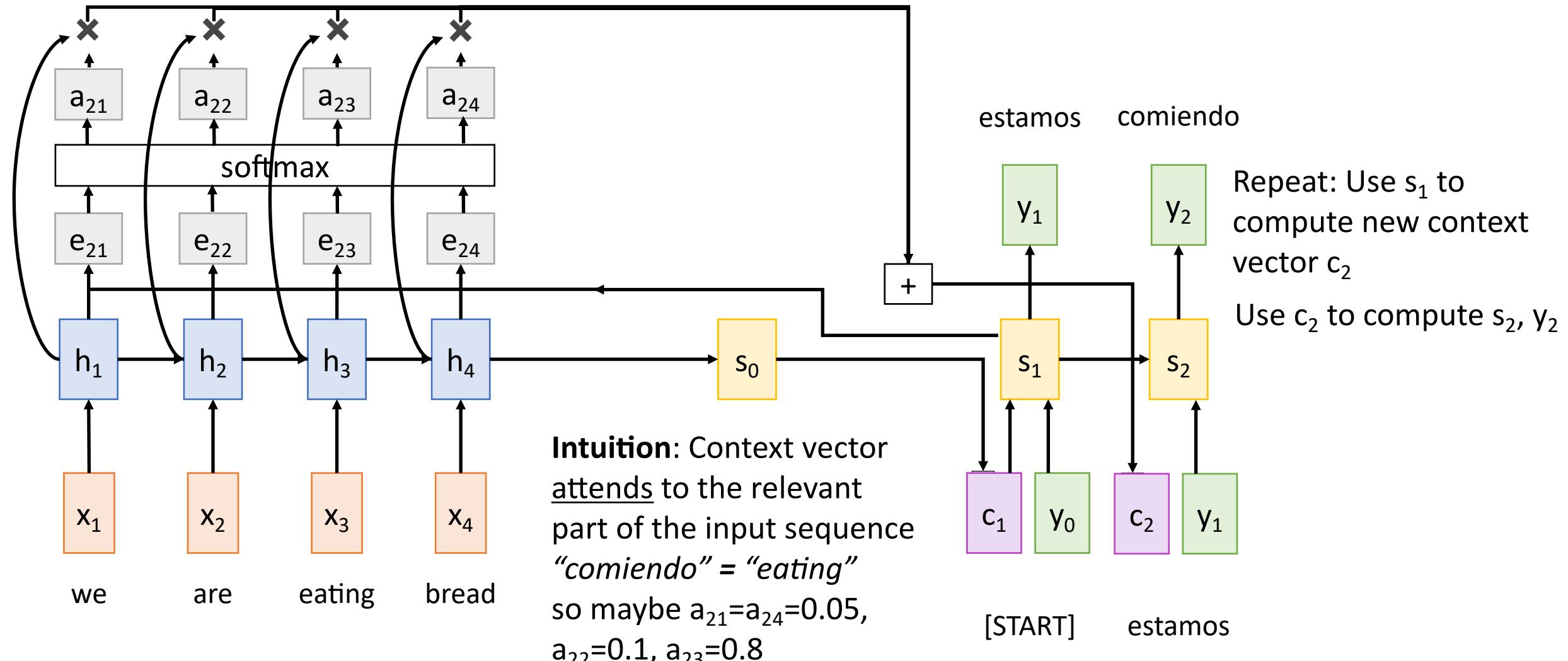
Repeat: Use  $s_1$  to compute new context vector  $c_2$



# Sequence-to-Sequence with RNNs and Attention



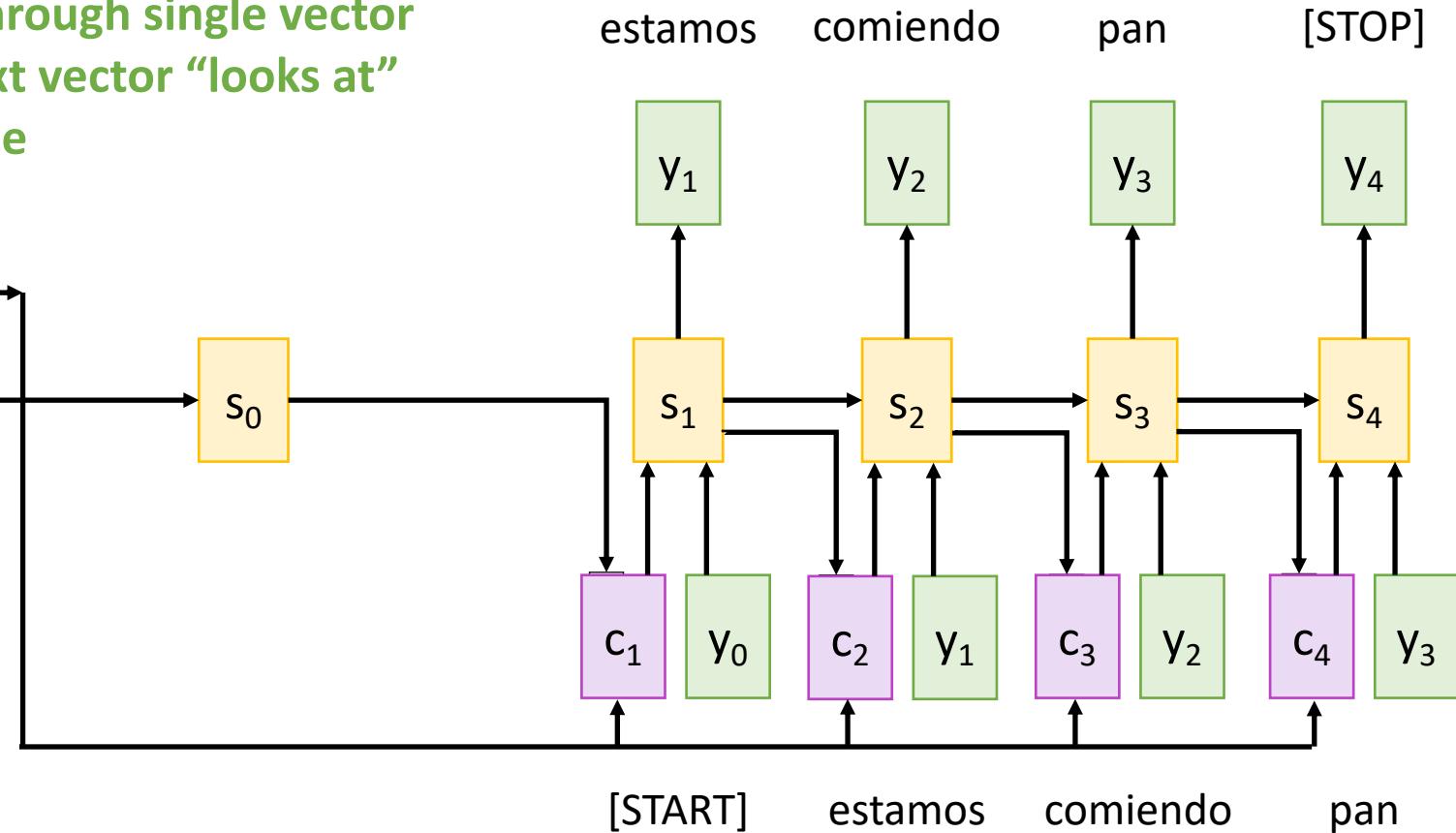
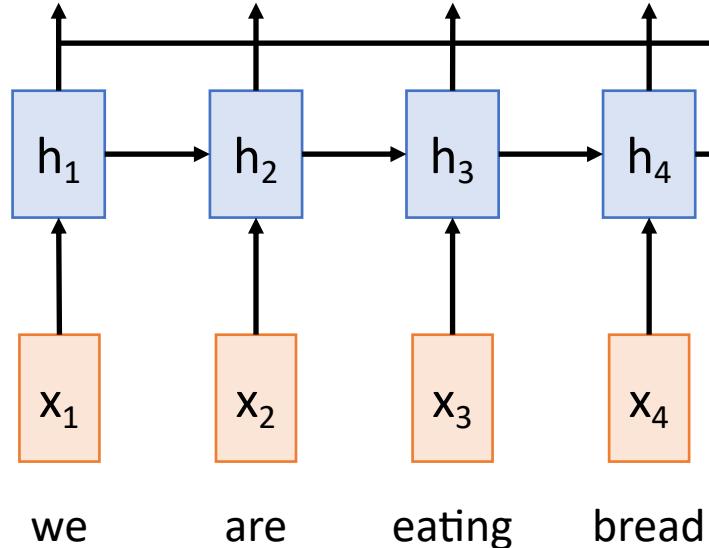
# Sequence-to-Sequence with RNNs and Attention



# 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



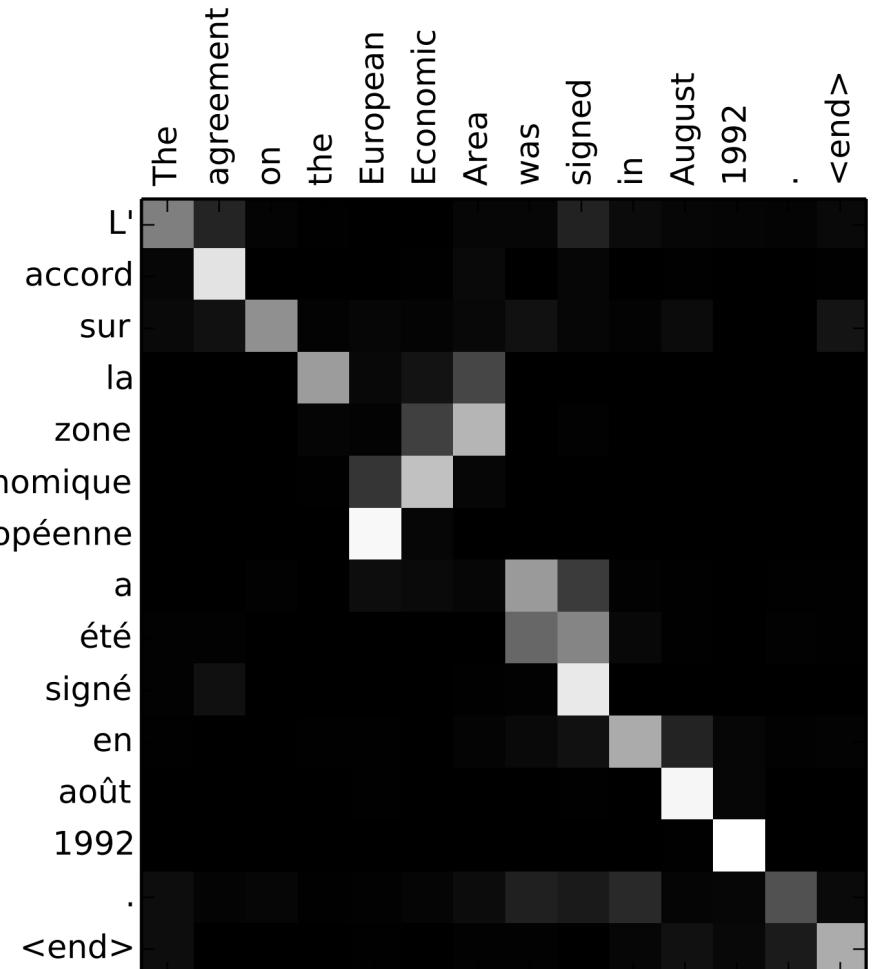
# Sequence-to-Sequence with RNNs and Attention

**Example:** English to French translation

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

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

Visualize attention weights  $a_{t,i}$



# Sequence-to-Sequence with RNNs and Attention

**Example:** English to French translation

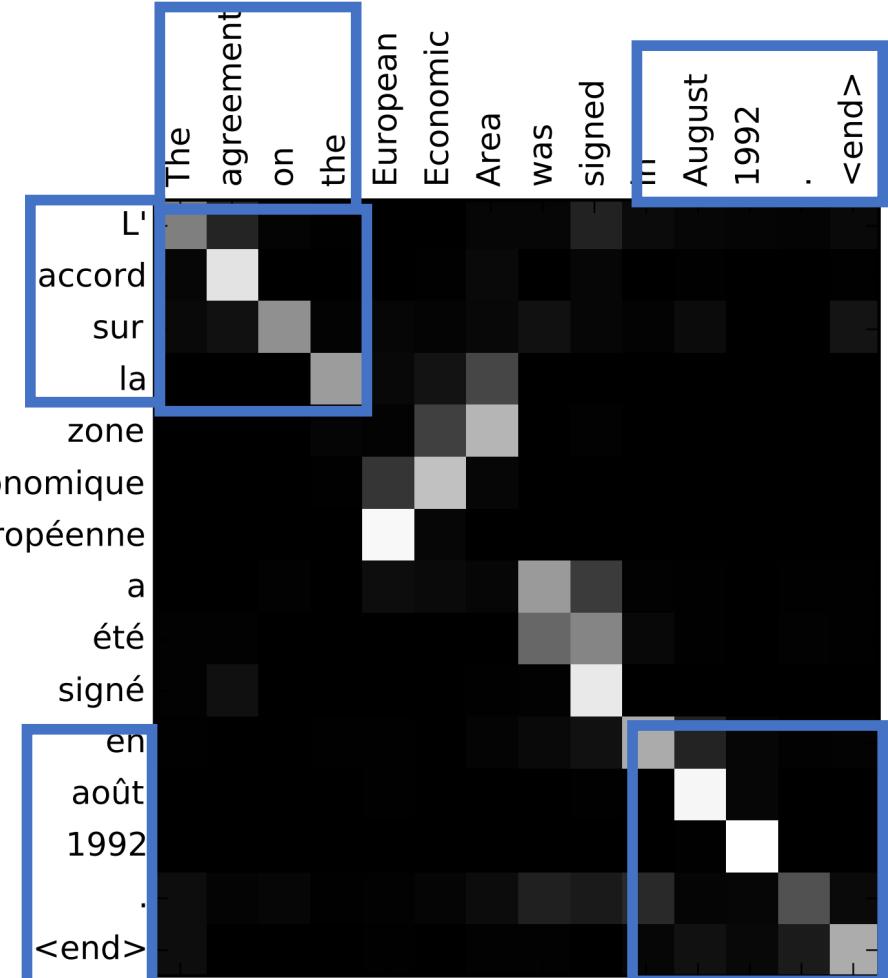
**Input:** “**The agreement on the European Economic Area was signed in August 1992.**”

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

Diagonal attention means words correspond in order

Diagonal attention means words correspond in order

Visualize attention weights  $a_{t,i}$



# 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.”

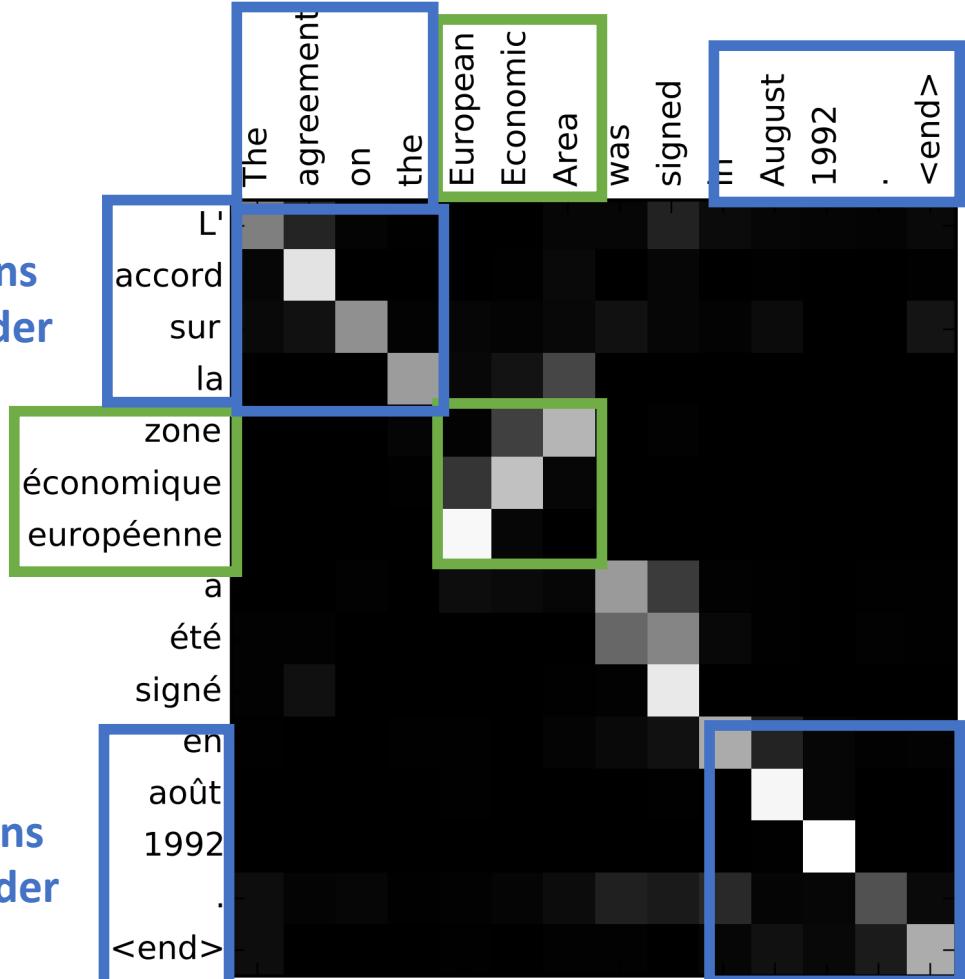
**Output:** “L'accord sur la zone économique européenne a été signé en août 1992.”

Diagonal attention means words correspond in order

Attention figures out different word orders

Diagonal attention means words correspond in order

Visualize attention weights  $a_{t,i}$



# 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.”

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

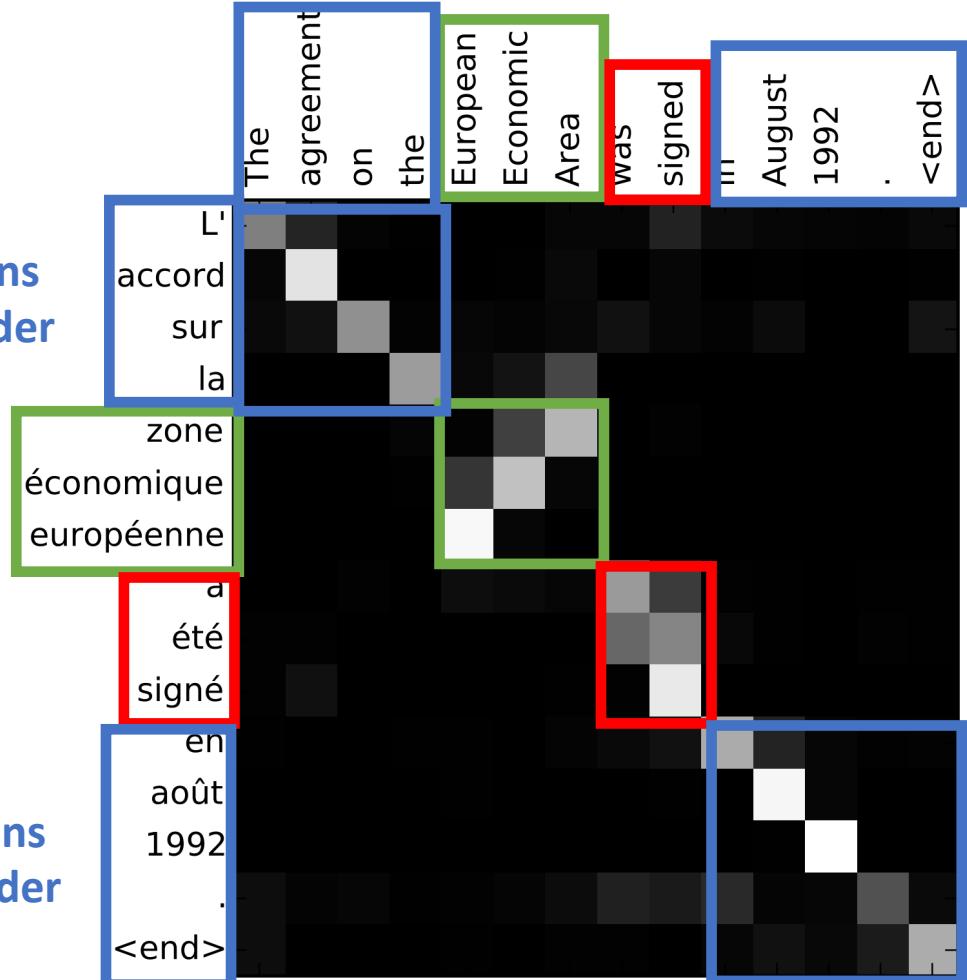
Diagonal attention means words correspond in order

Attention figures out different word orders

Verb conjugation

Diagonal attention means words correspond in order

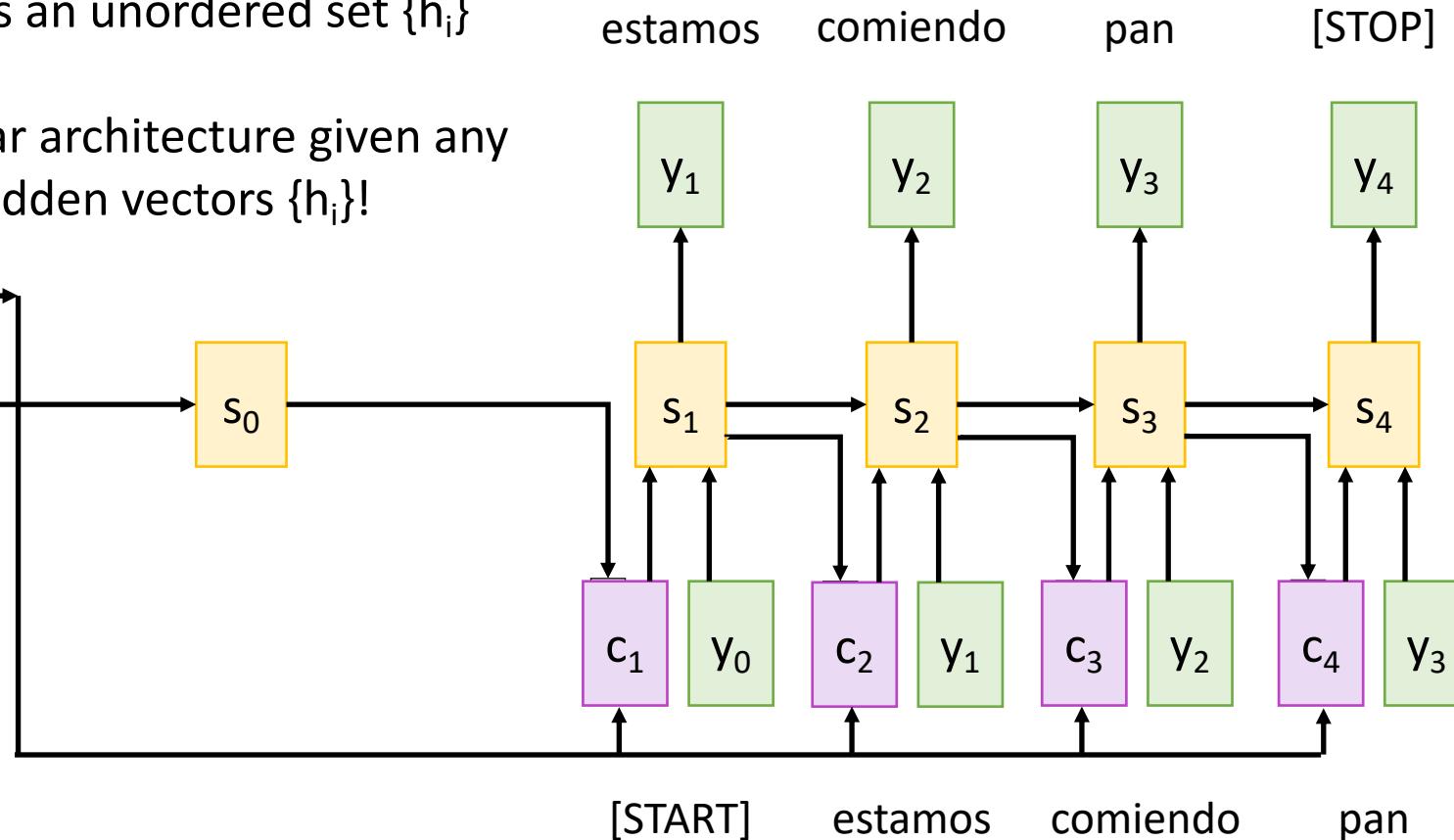
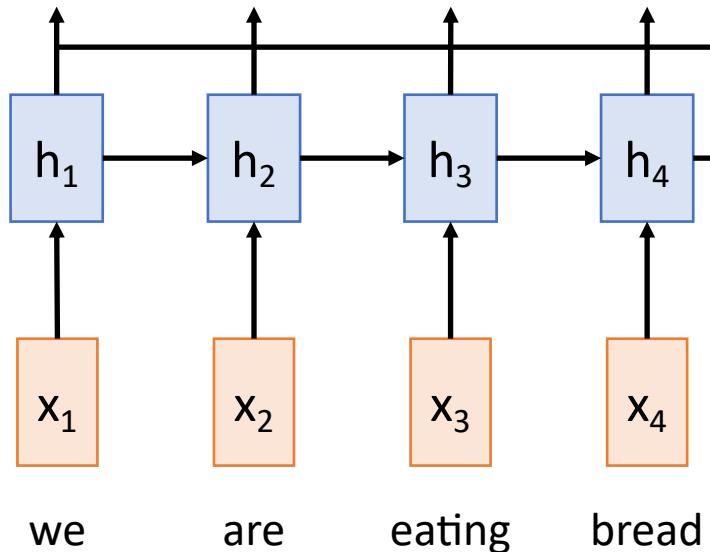
Visualize attention weights  $a_{t,i}$



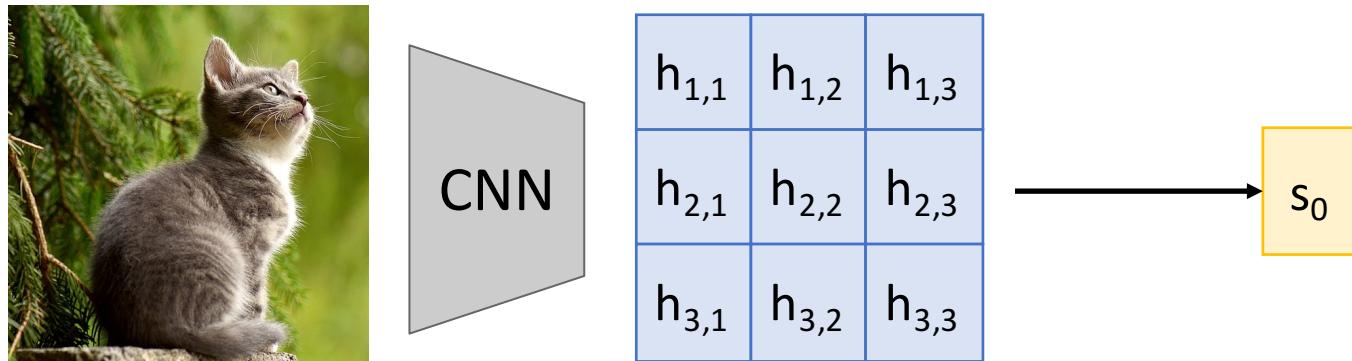
# Sequence-to-Sequence with RNNs and Attention

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  $\{h_i\}$ !



# 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

# Image Captioning with RNNs and Attention

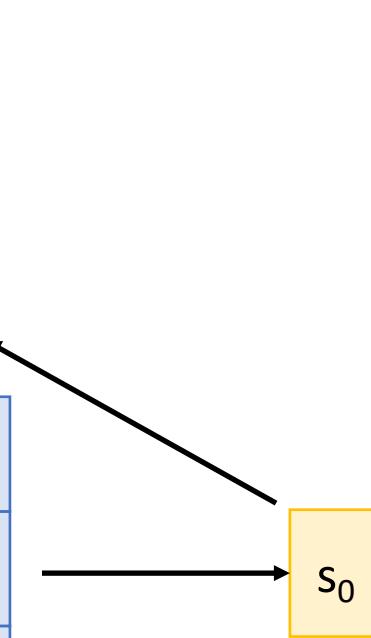
$$e_{t,i,j} = f_{\text{att}}(s_{t-1}, h_{i,j})$$

Alignment scores

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



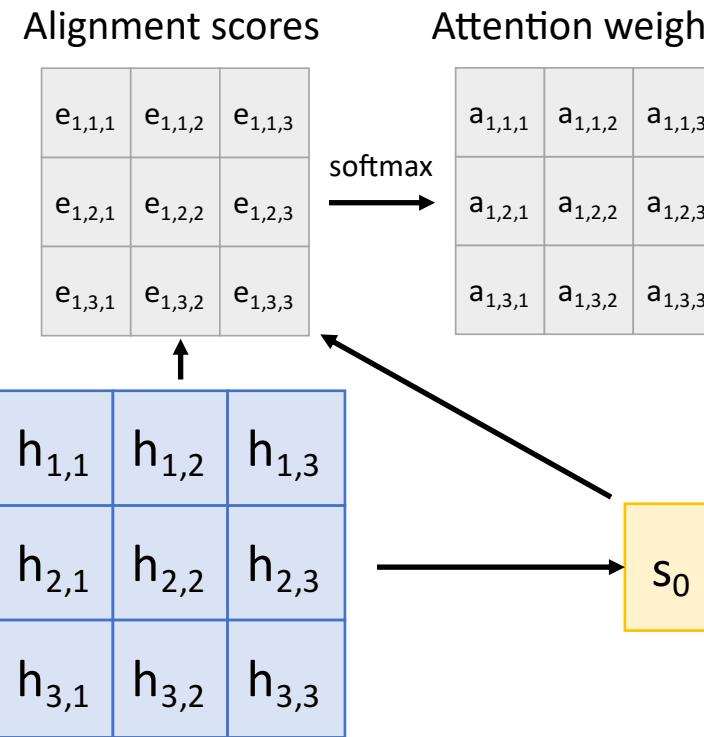
$h_{1,1}$	$h_{1,2}$	$h_{1,3}$
$h_{2,1}$	$h_{2,2}$	$h_{2,3}$
$h_{3,1}$	$h_{3,2}$	$h_{3,3}$



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

# Image Captioning with RNNs and Attention

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



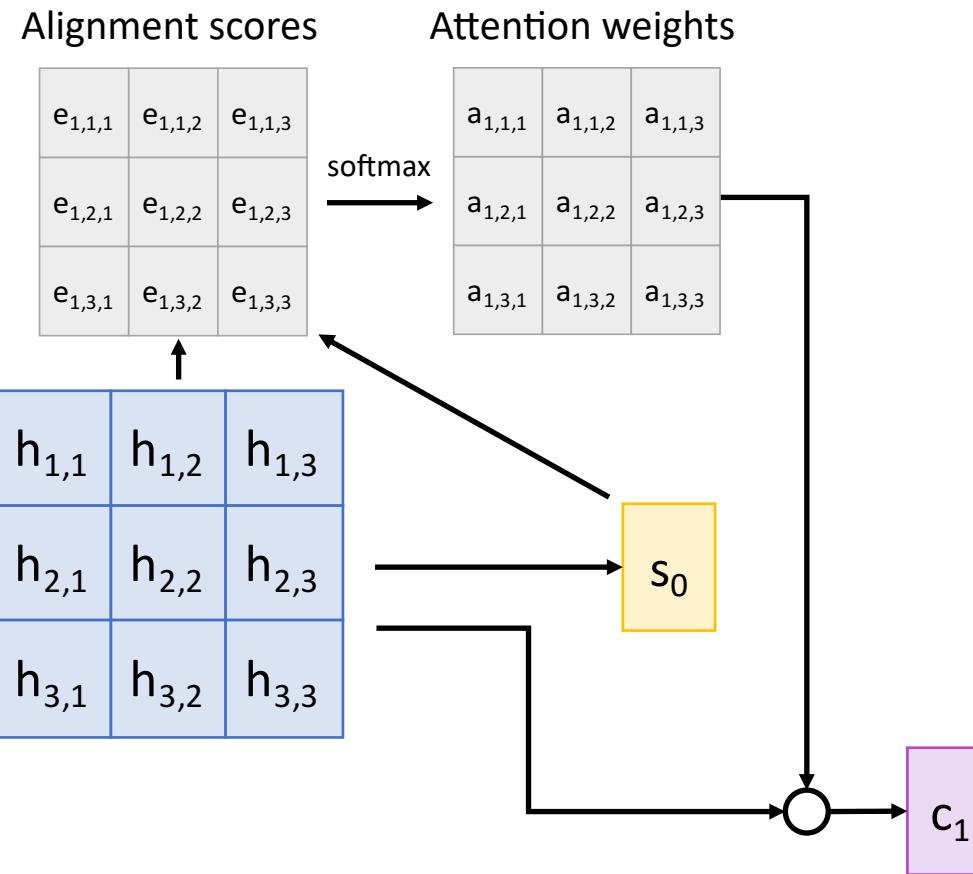
Use a CNN to compute a grid of features for an image

# Image Captioning with RNNs and Attention

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

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

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



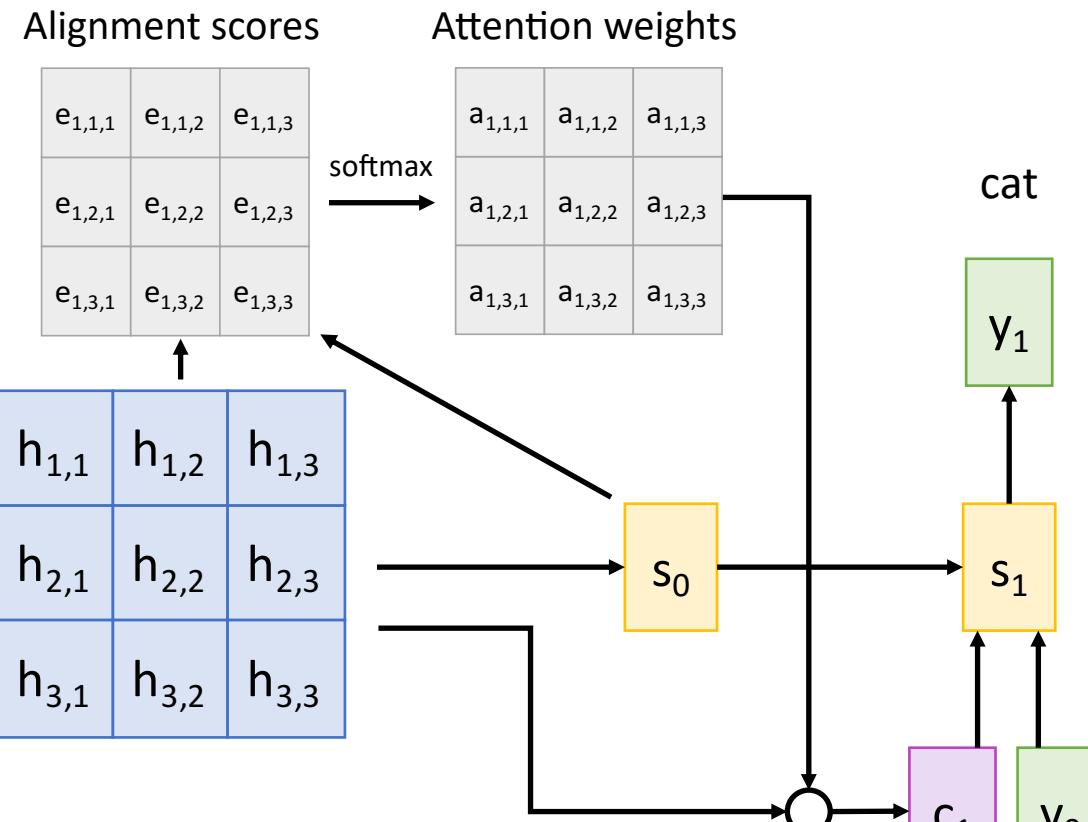
Use a CNN to compute a grid of features for an image

# Image Captioning with RNNs and Attention

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

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

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



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

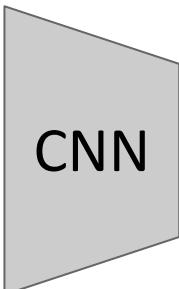
[START]

# Image Captioning with RNNs and Attention

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

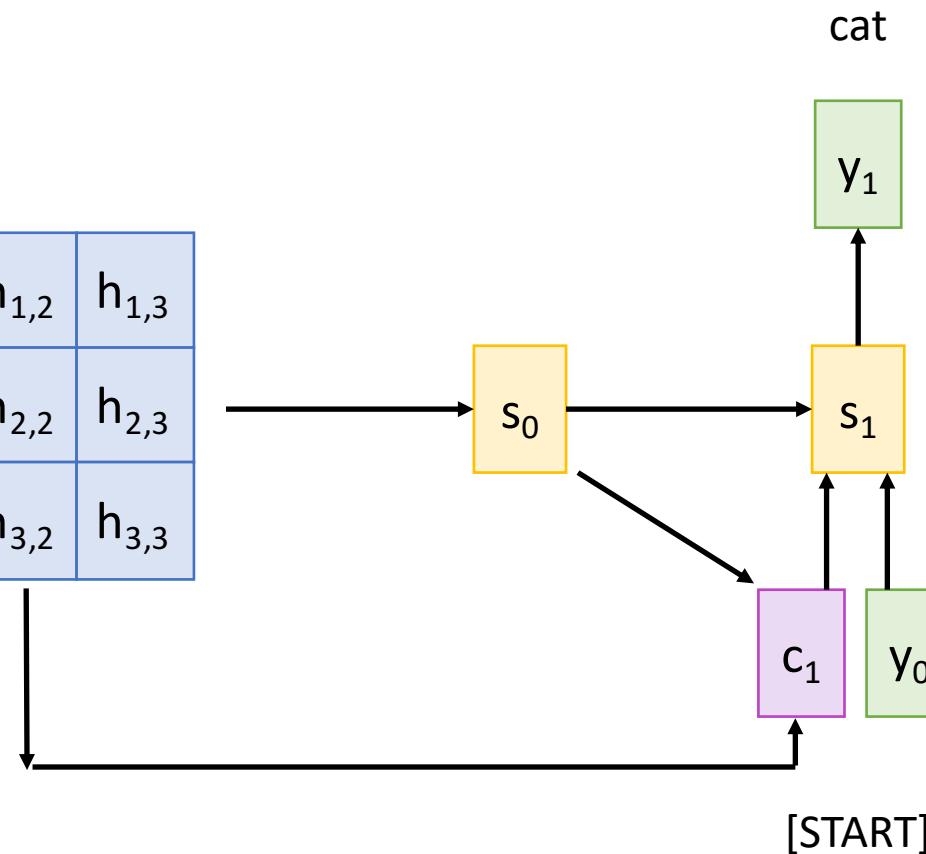
$$a_{t,:,:} = \text{softmax}(e_{t,:,:})$$

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



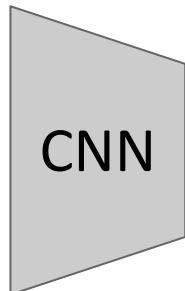
$h_{1,1}$	$h_{1,2}$	$h_{1,3}$
$h_{2,1}$	$h_{2,2}$	$h_{2,3}$
$h_{3,1}$	$h_{3,2}$	$h_{3,3}$

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

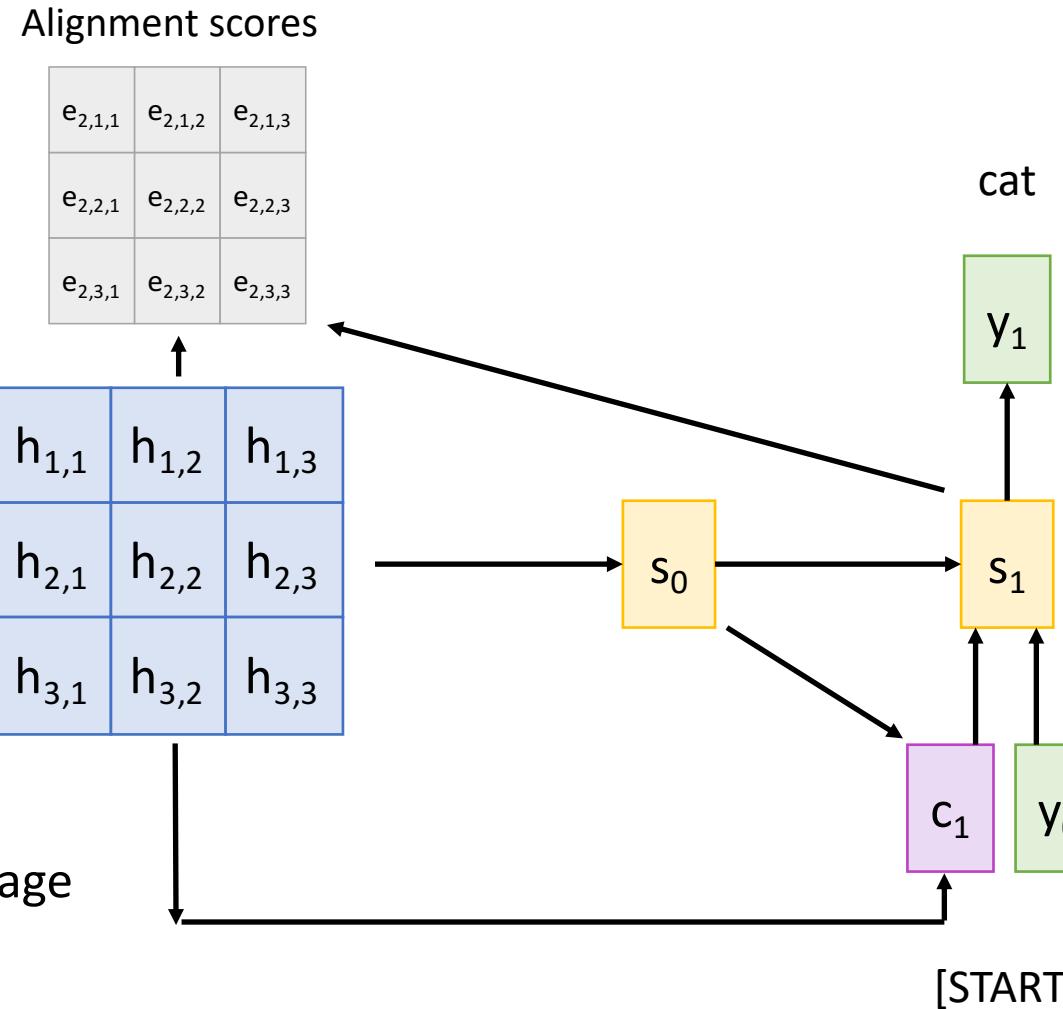


# Image Captioning with RNNs and Attention

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

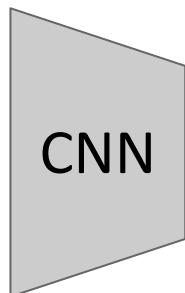


Use a CNN to compute a grid of features for an image

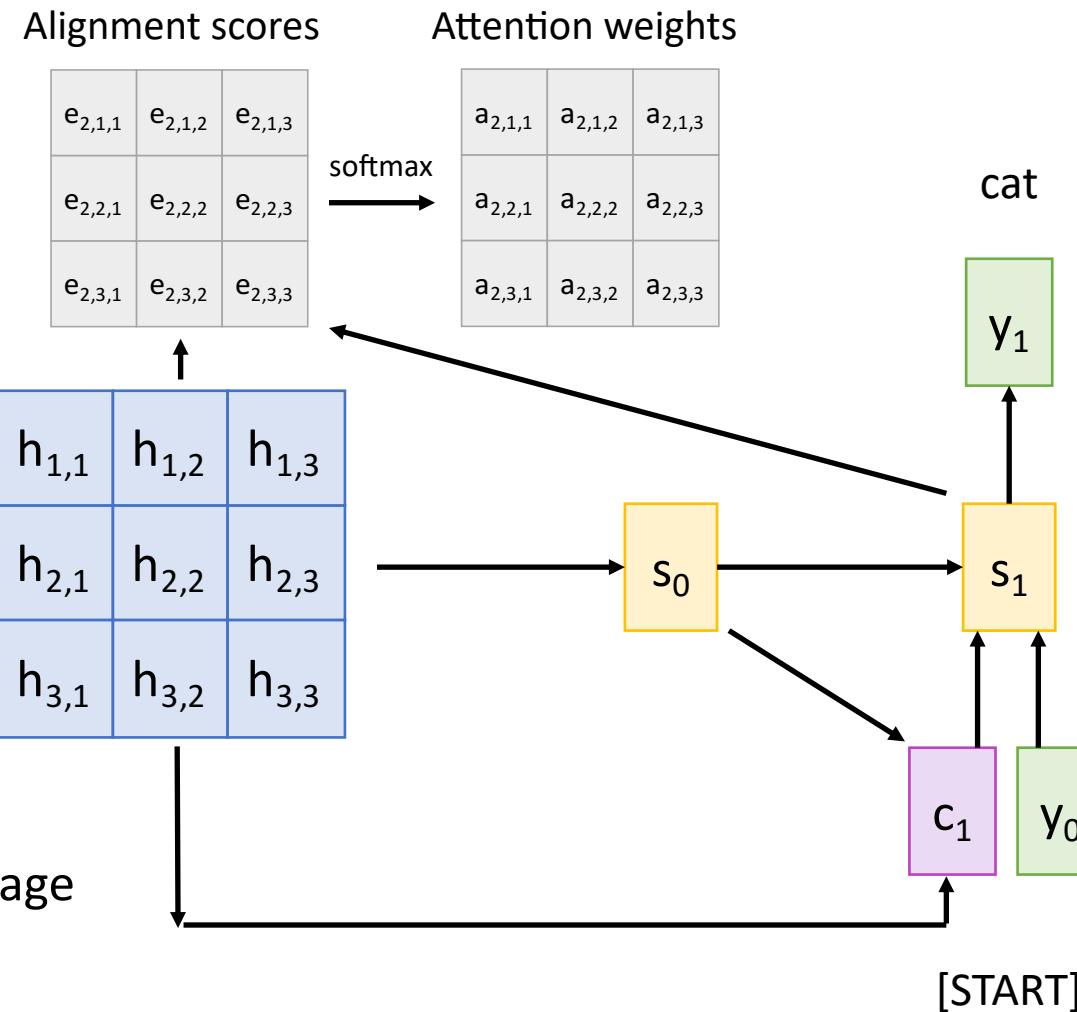


# Image Captioning with RNNs and Attention

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



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

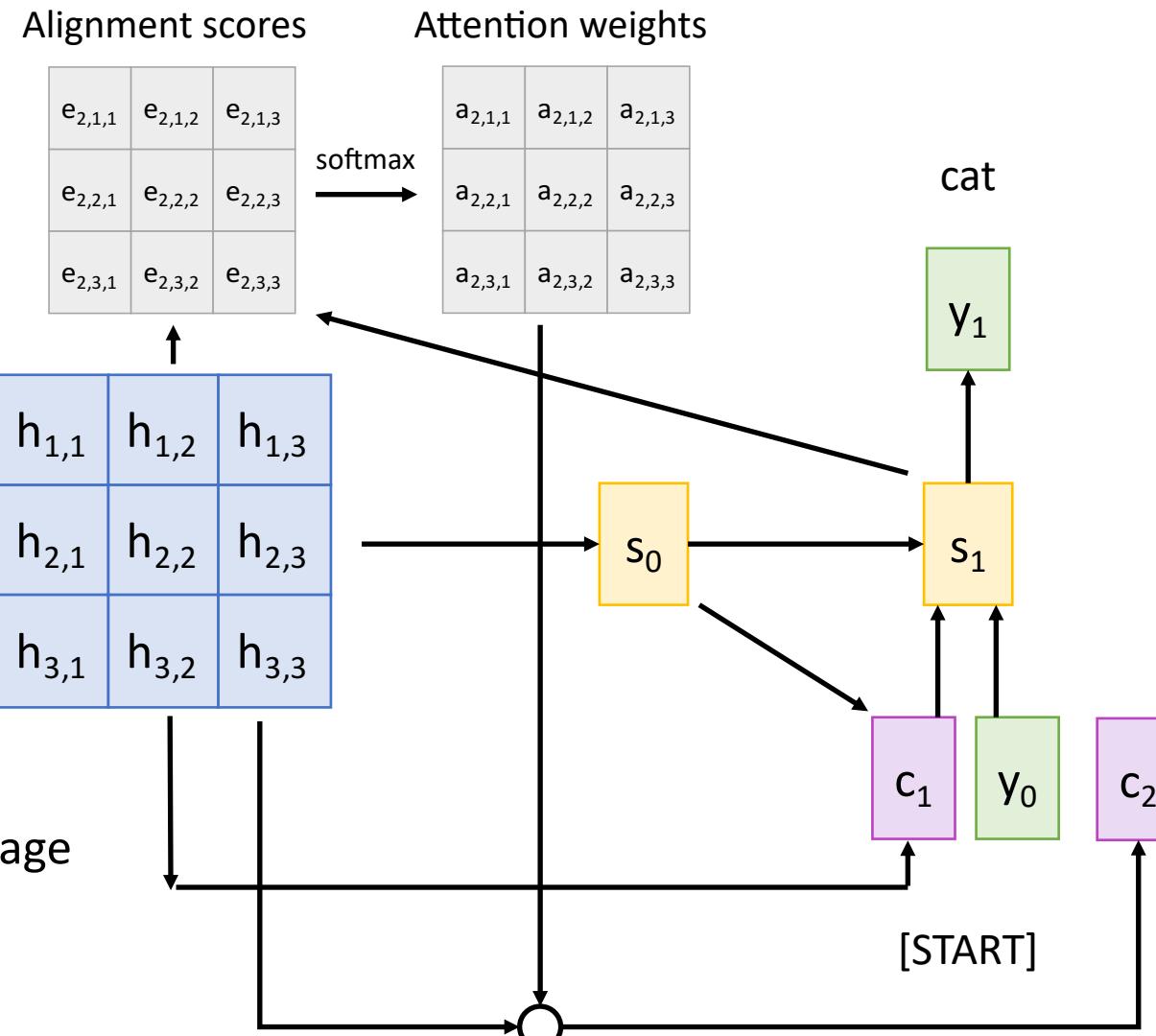


# Image Captioning with RNNs and Attention

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



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

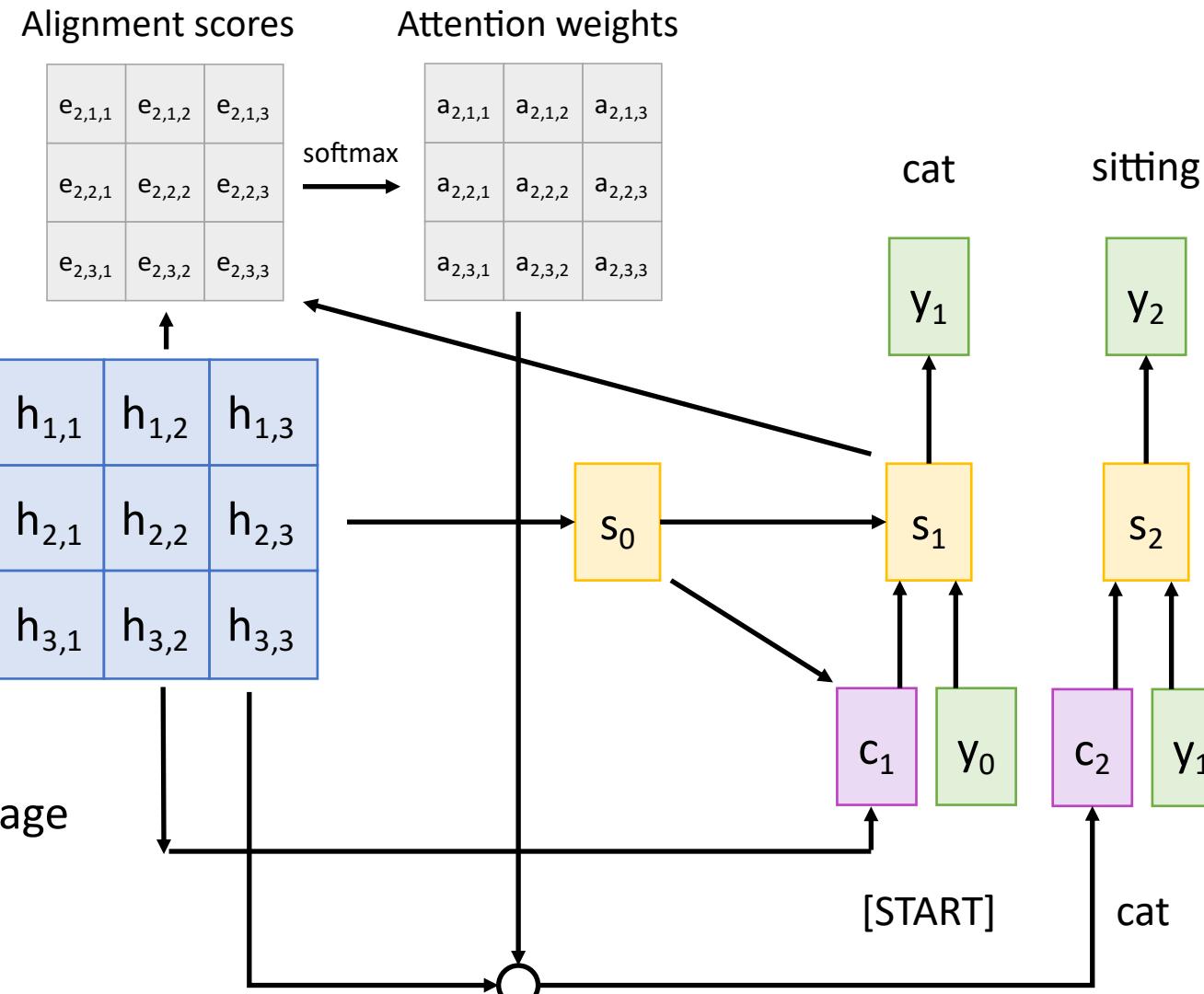


# Image Captioning with RNNs and Attention

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



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



# Image Captioning with RNNs and Attention

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

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

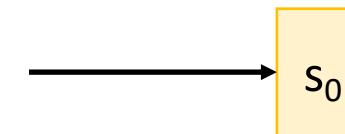
$$c_t = \sum_{i,j} a_{t,i,j} h_{i,j}$$



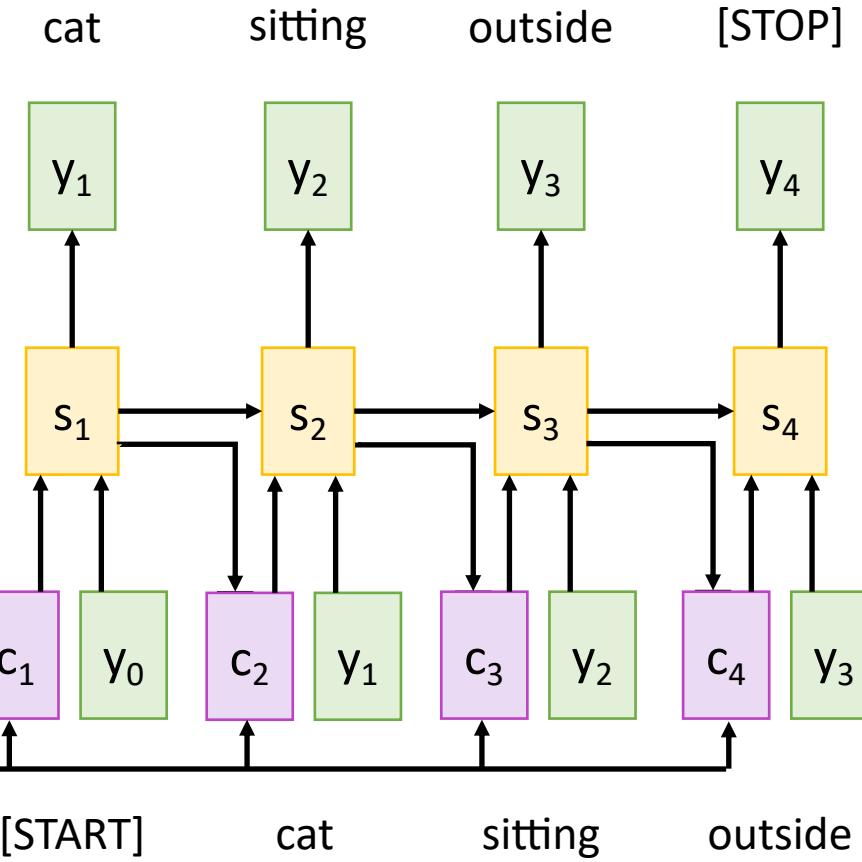
Use a CNN to compute a grid of features for an image

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

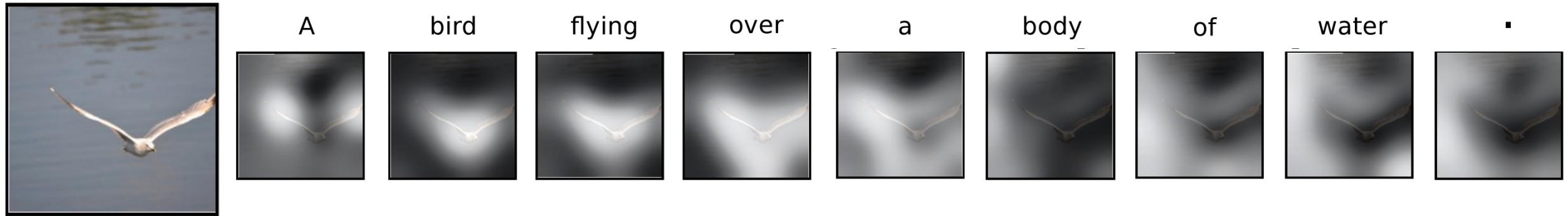
$h_{1,1}$	$h_{1,2}$	$h_{1,3}$
$h_{2,1}$	$h_{2,2}$	$h_{2,3}$
$h_{3,1}$	$h_{3,2}$	$h_{3,3}$



[START] cat sitting outside [STOP]



# Image Captioning with RNNs and Attention



# Image Captioning with RNNs and Attention



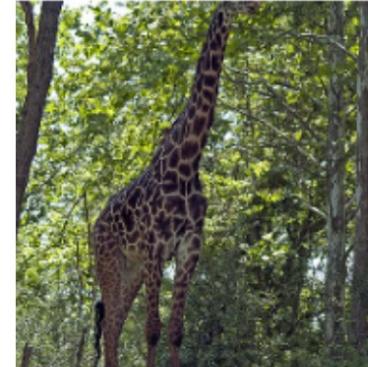
A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.

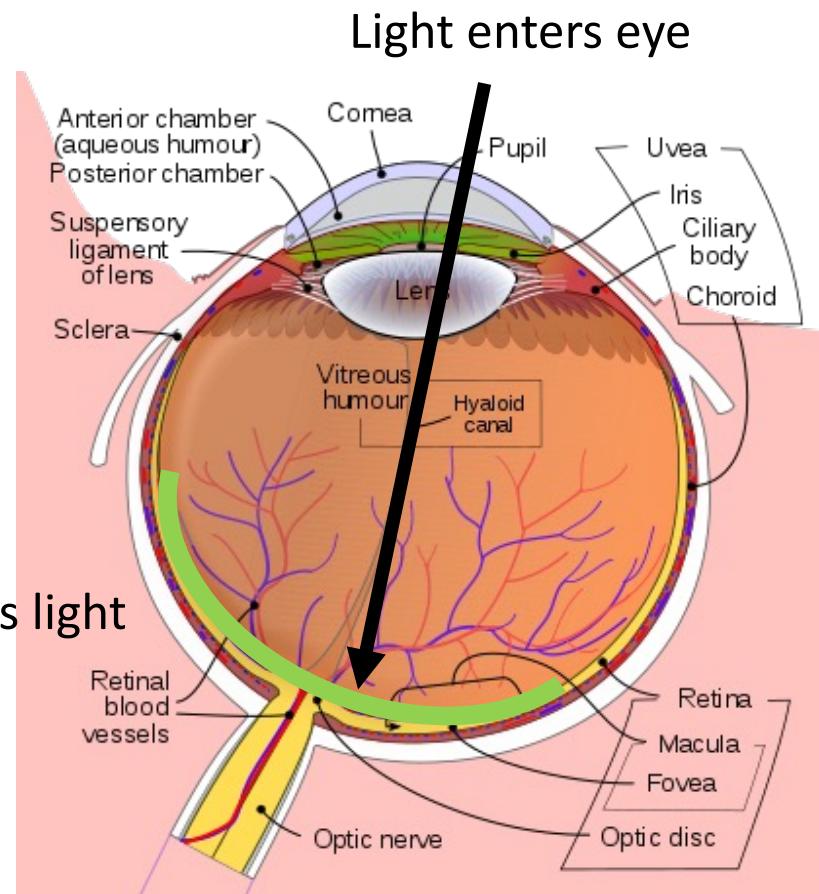


A group of people sitting on a boat in the water.



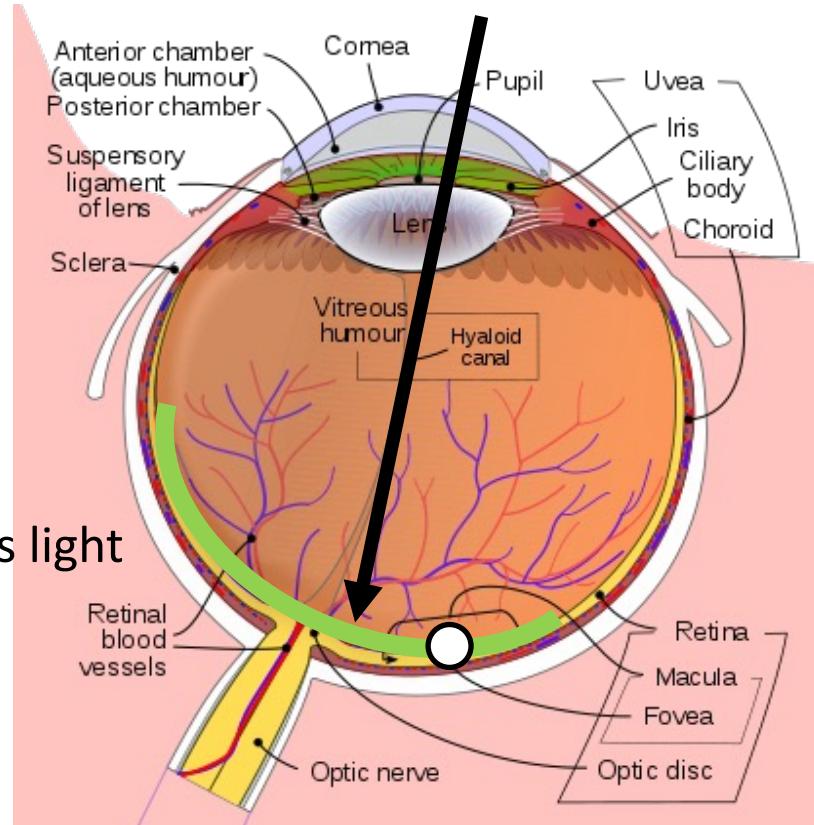
A giraffe standing in a forest with trees in the background.

# Human Vision: Fovea

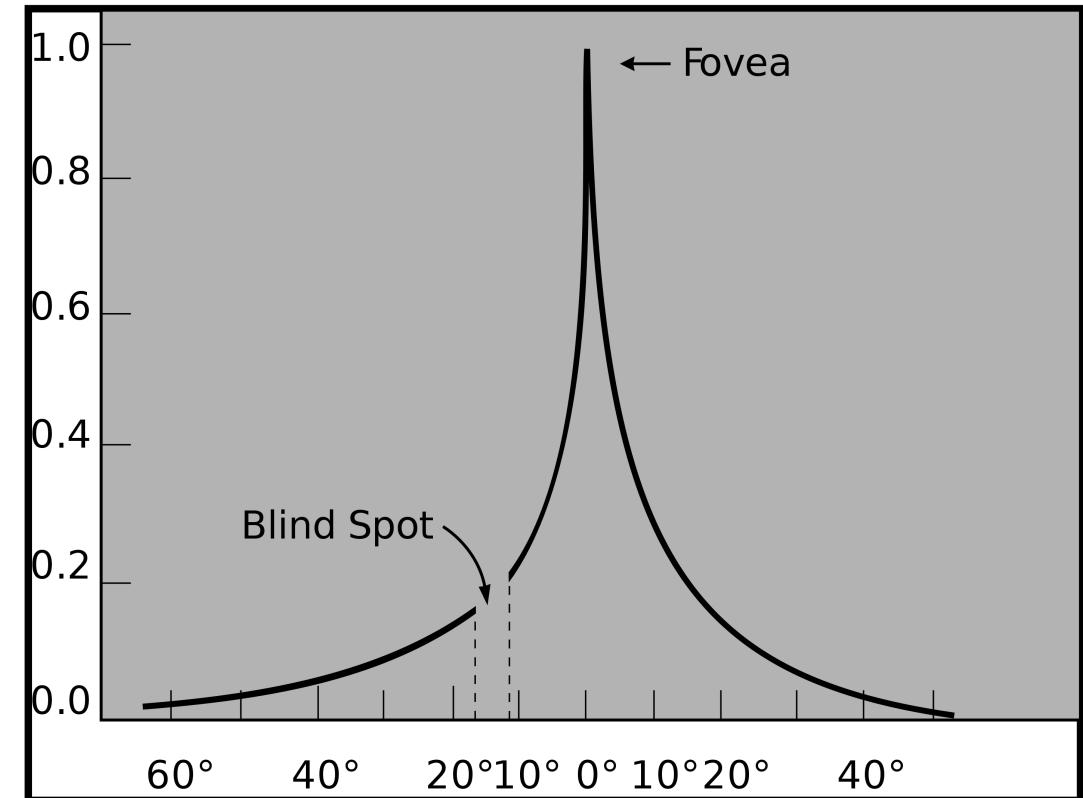


# Human Vision: Fovea

Light enters eye



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

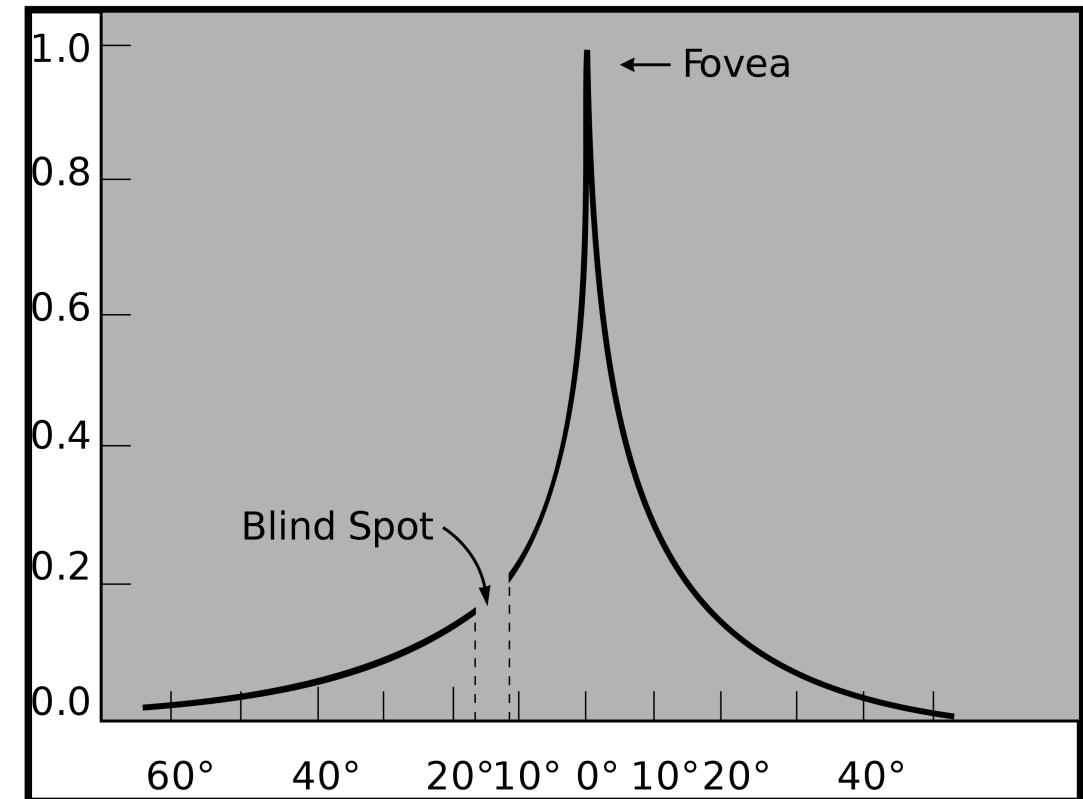


# Human Vision: Saccades

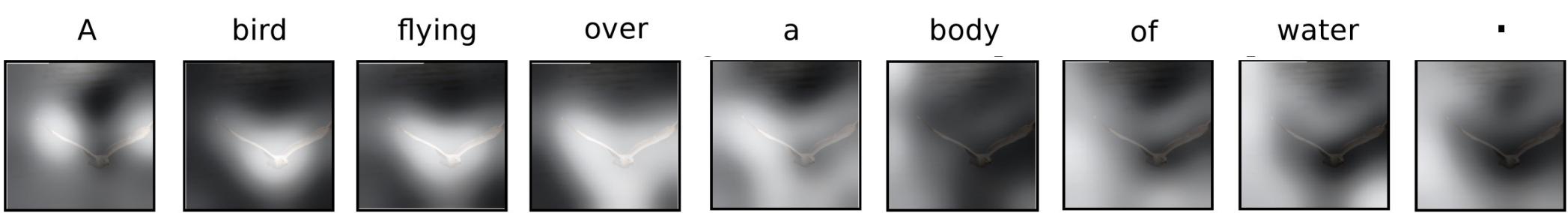
Human eyes are constantly moving so we don't notice



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



# Image Captioning with RNNs and Attention



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



# 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

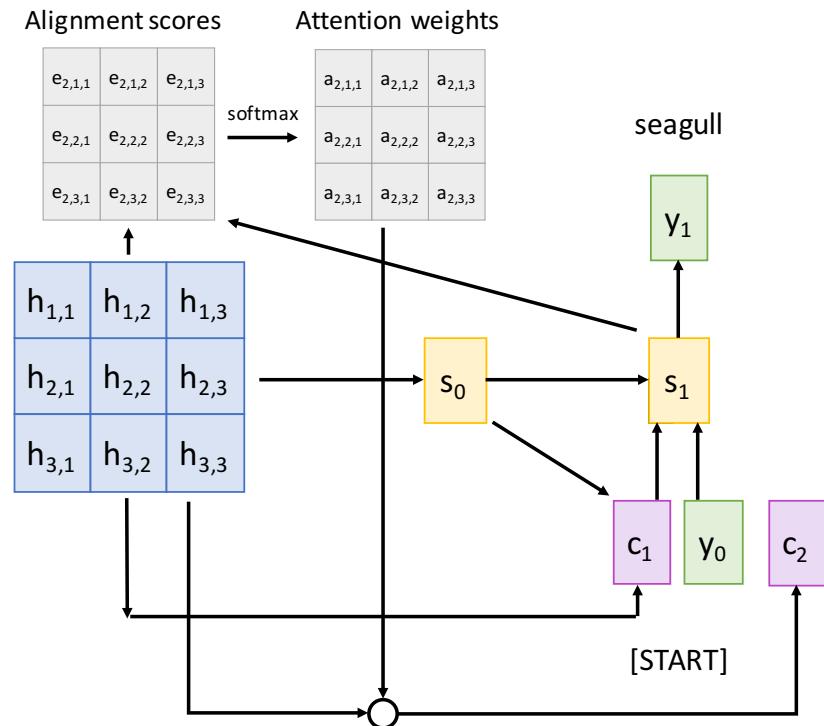
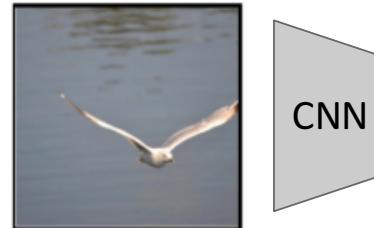
## Inputs:

Query vector:  $q$  (Shape:  $D_Q$ )

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

Similarity function:  $f_{att}$

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



## Computation:

Similarities:  $e$  (Shape:  $N_x$ )  $e_i = f_{att}(q, X_i)$

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

Output vector:  $y = \sum_i a_i X_i$  (Shape:  $D_x$ )

# Attention Layer

## Inputs:

Query vector:  $q$  (Shape:  $D_Q$ )

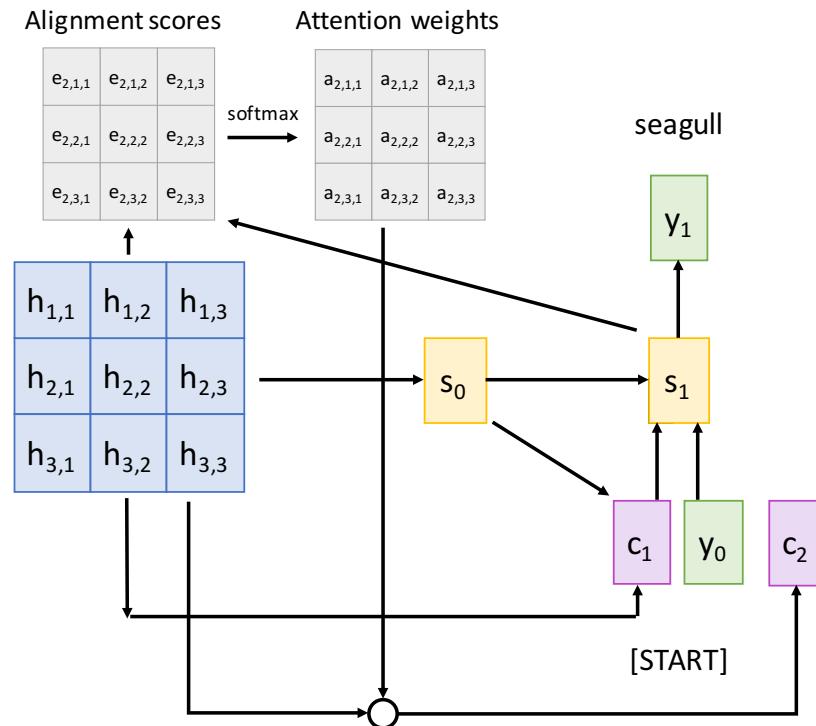
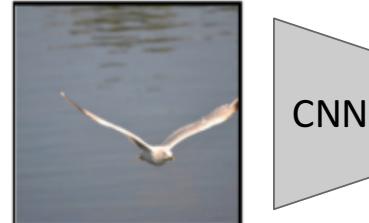
Input vectors:  $X$  (Shape:  $N_x \times D_Q$ )

Similarity function: dot product

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

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

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



## Computation:

Similarities:  $e$  (Shape:  $N_x$ )  $e_i = q \cdot X_i$

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

Output vector:  $y = \sum_i a_i X_i$  (Shape:  $D_X$ )

## Changes:

- Use dot product for similarity

# Attention Layer

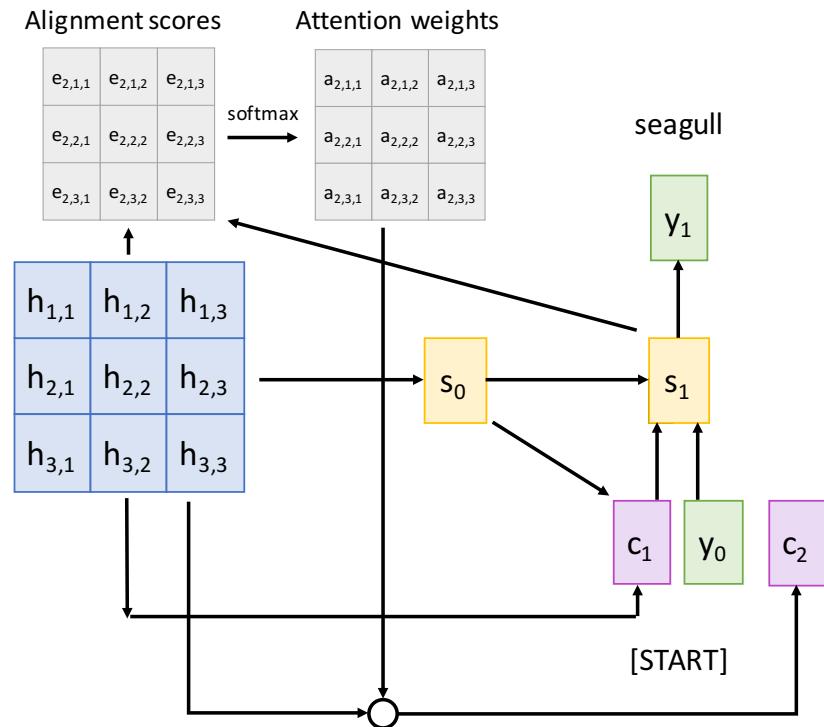
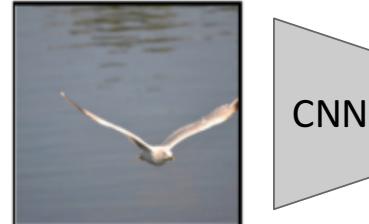
## Inputs:

Query vector:  $q$  (Shape:  $D_Q$ )

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

Similarity function: scaled dot product

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



## Computation:

Similarities:  $e$  (Shape:  $N_x$ )  $e_i = q \cdot X_i / \sqrt{D_Q}$

Attention weights:  $a = \text{softmax}(e)$  (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 vector:  $q$  (Shape:  $D_Q$ )

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

Similarity function: scaled dot product

Large similarities will cause softmax to saturate and give vanishing gradients

Recall  $a \cdot b = |a| |b| \cos(\text{angle})$

Suppose that  $a$  and  $b$  are constant vectors of dimension  $D$

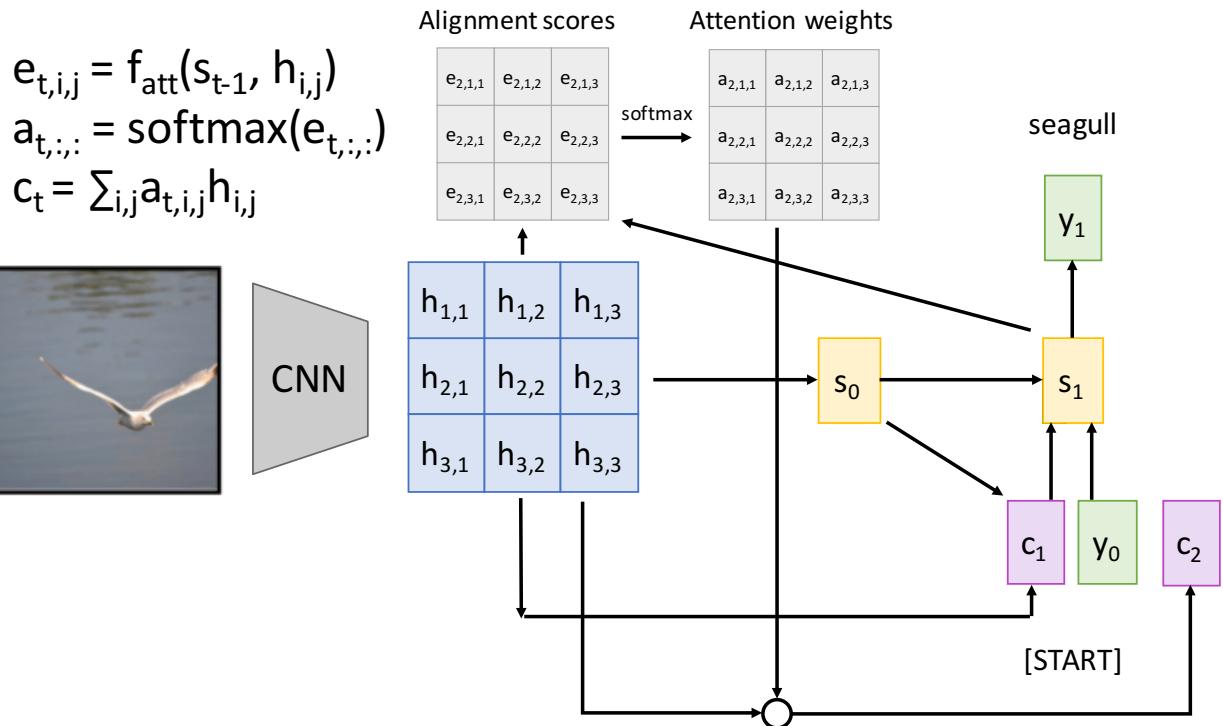
Then  $|a| = (\sum_i a_i^2)^{1/2} = a \sqrt{D}$

## Computation:

Similarities:  $e$  (Shape:  $N_x$ )  $e_i = q \cdot X_i / \sqrt{D_Q}$

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

Output vector:  $y = \sum_i a_i X_i$  (Shape:  $D_X$ )



## Changes:

- Use **scaled** dot product for similarity

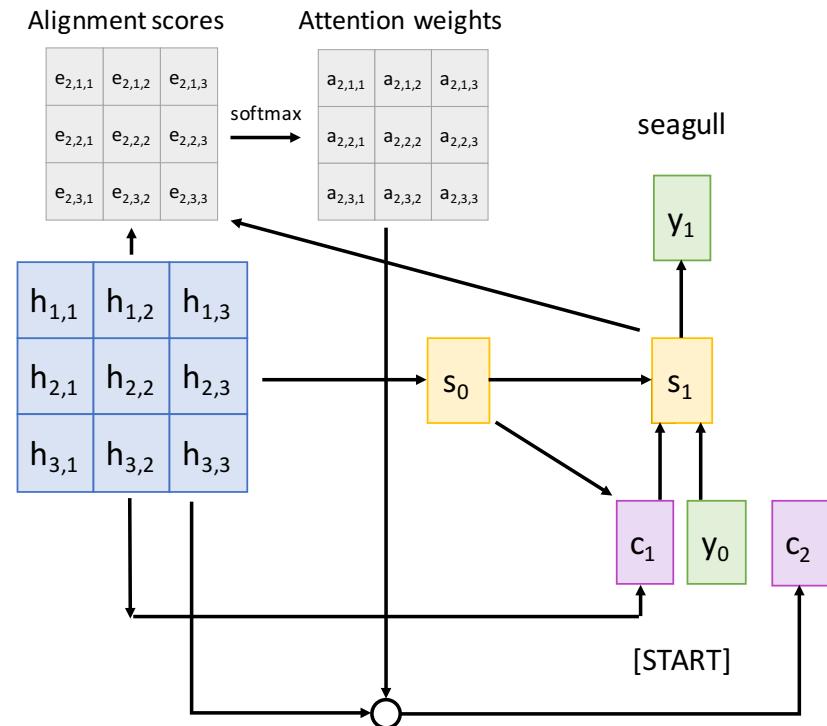
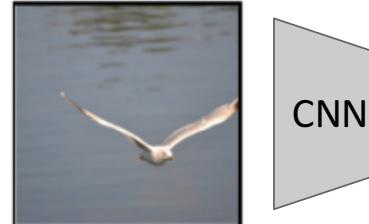
# Attention Layer

## Inputs:

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

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

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



## Computation:

**Similarities:**  $E = \mathbf{Q}\mathbf{X}^T / \sqrt{D_Q}$  (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

## Inputs:

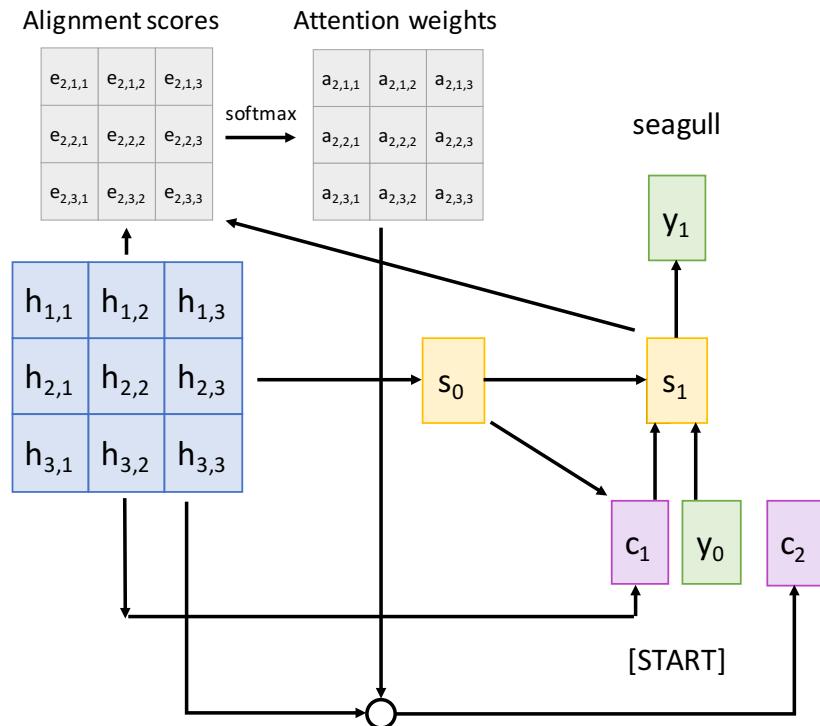
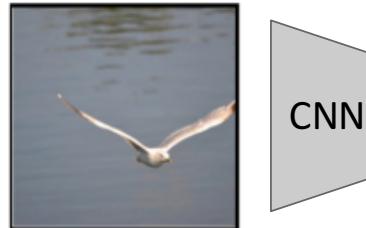
**Query vectors:**  $\mathbf{Q}$  (Shape:  $N_Q \times D_Q$ )

**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:  $D_X \times D_V$ )

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



## 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:**  $E = \mathbf{Q}\mathbf{K}^T / \sqrt{D_Q}$  (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:**  $\mathbf{Y} = A\mathbf{V}$  (Shape:  $N_Q \times D_V$ )  $Y_i = \sum_j A_{i,j} V_j$

## 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:**  $E = \mathbf{Q}\mathbf{K}^T / \sqrt{D_Q}$  (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:**  $\mathbf{Y} = A\mathbf{V}$  (Shape:  $N_Q \times D_V$ )  $\mathbf{Y}_i = \sum_j A_{i,j} \mathbf{V}_j$

$X_1$

$X_2$

$X_3$

$Q_1$

$Q_2$

$Q_3$

$Q_4$

# Attention Layer

## Inputs:

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

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

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

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

## Computation:

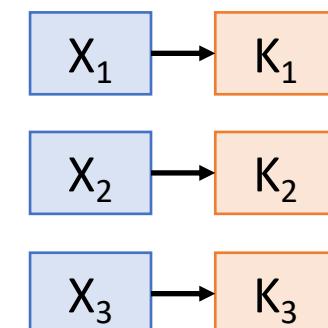
**Key vectors:**  $\mathbf{K} = \mathbf{X}\mathbf{W}_K$  (Shape:  $N_X \times D_Q$ )

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

**Similarities:**  $E = \mathbf{Q}\mathbf{K}^T / \sqrt{D_Q}$  (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:**  $\mathbf{Y} = A\mathbf{V}$  (Shape:  $N_Q \times D_V$ )  $\mathbf{Y}_i = \sum_j A_{i,j} \mathbf{V}_j$



# 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:

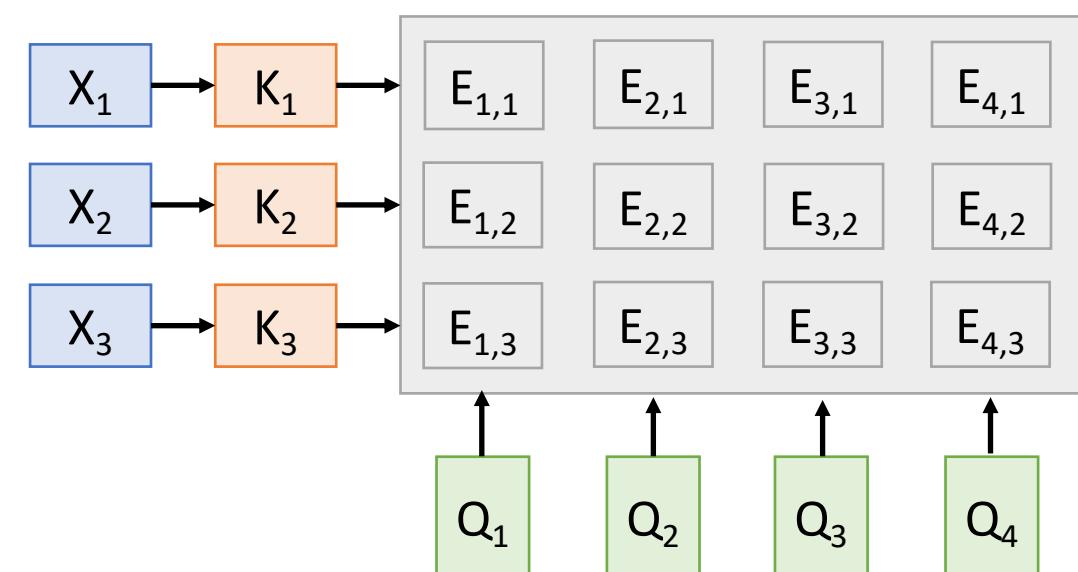
**Key vectors:**  $\mathbf{K} = \mathbf{X}\mathbf{W}_K$  (Shape:  $N_X \times D_Q$ )

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

**Similarities:**  $E = \mathbf{Q}\mathbf{K}^T / \sqrt{D_Q}$  (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:**  $\mathbf{Y} = A\mathbf{V}$  (Shape:  $N_Q \times D_V$ )  $Y_i = \sum_j A_{i,j} \mathbf{V}_j$



# Attention Layer

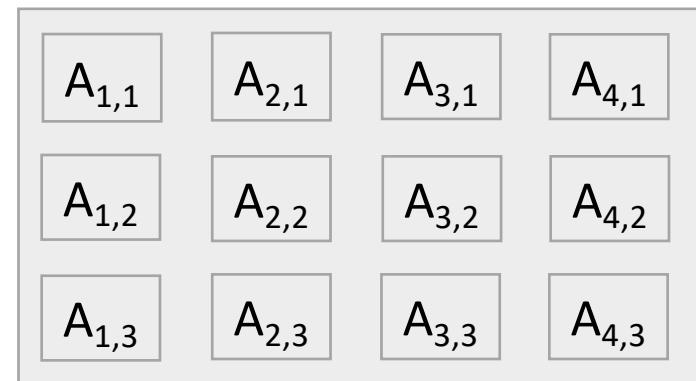
## Inputs:

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

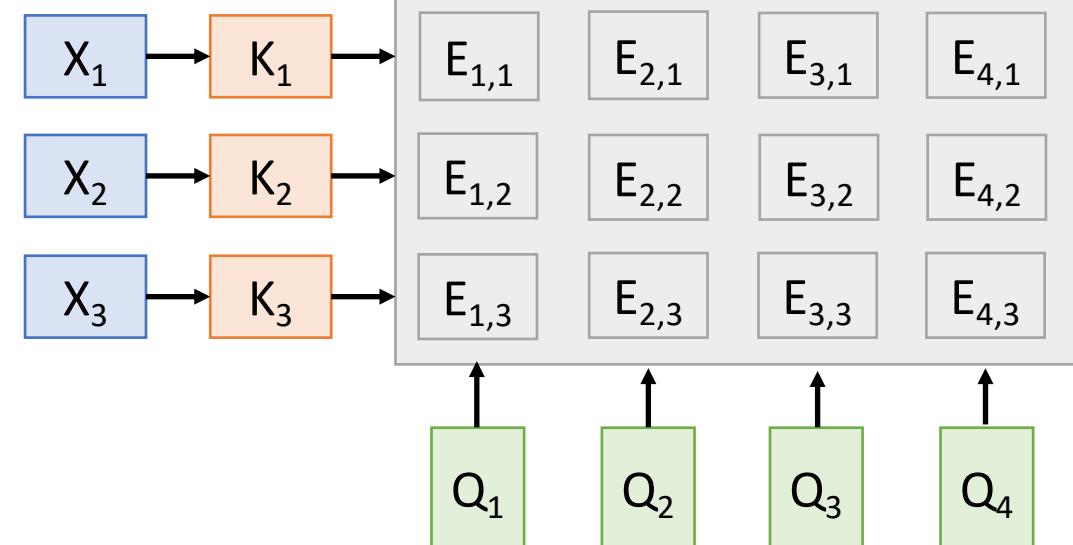
**Input vectors:**  $\mathbf{X}$  (Shape:  $N_X \times D_X$ )

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

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



Softmax(↑)



## 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:**  $E = \mathbf{Q}\mathbf{K}^T / \sqrt{D_Q}$  (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:**  $\mathbf{Y} = A\mathbf{V}$  (Shape:  $N_Q \times D_V$ )  $\mathbf{Y}_i = \sum_j A_{i,j} \mathbf{V}_j$

# 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:

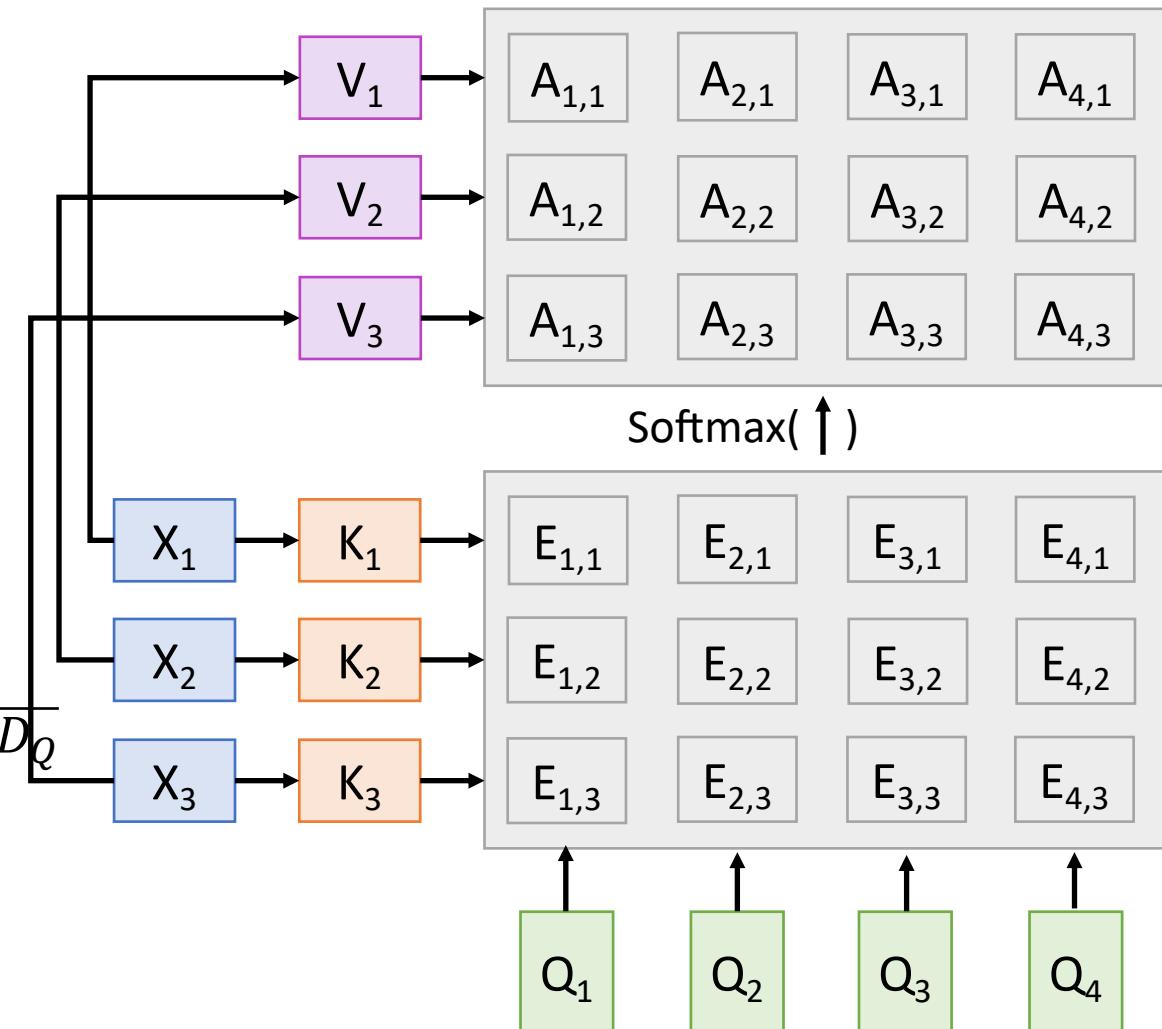
**Key vectors:**  $\mathbf{K} = \mathbf{X}\mathbf{W}_K$  (Shape:  $N_X \times D_Q$ )

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

**Similarities:**  $E = \mathbf{Q}\mathbf{K}^T / \sqrt{D_Q}$  (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:**  $\mathbf{Y} = A\mathbf{V}$  (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:

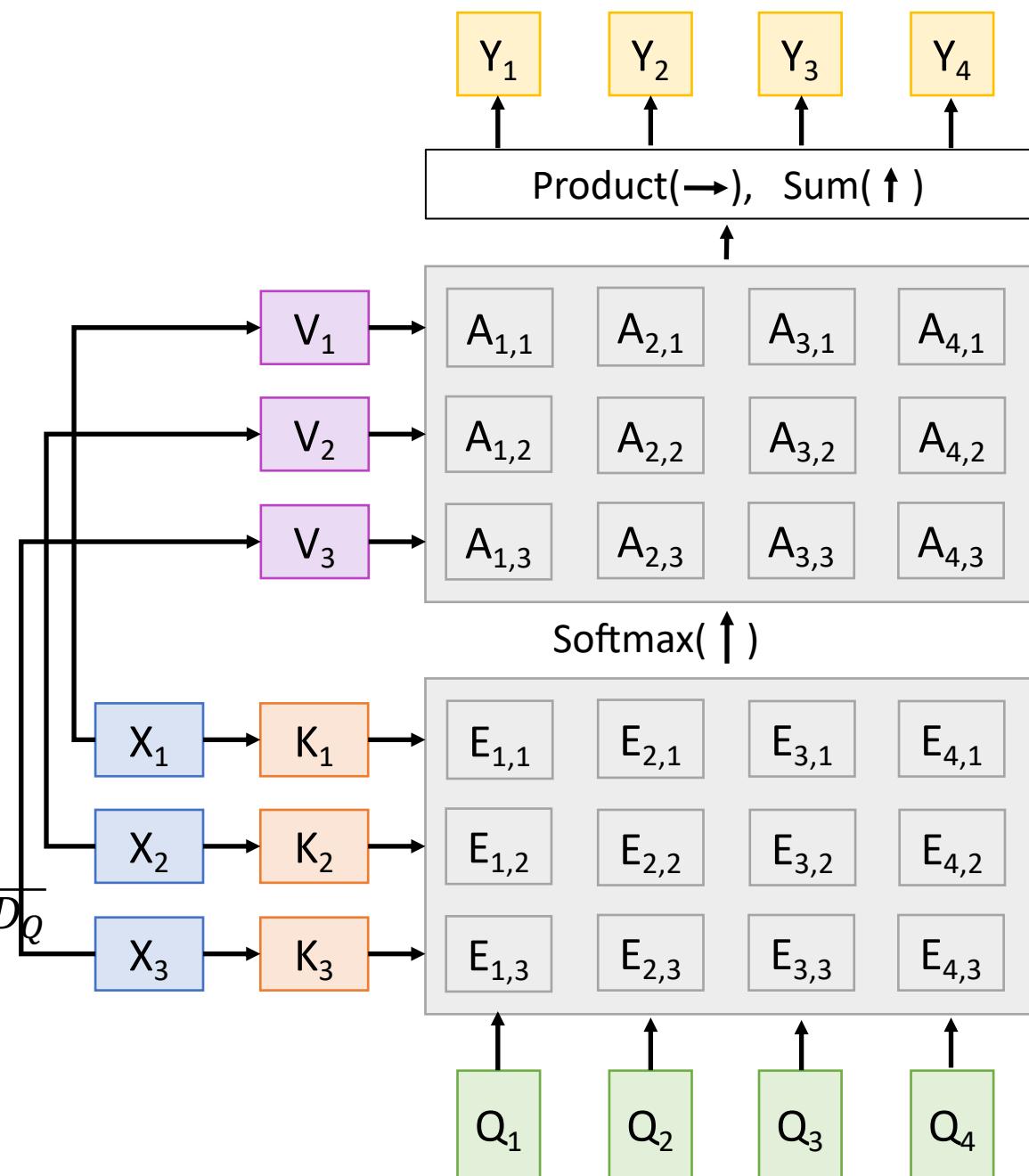
**Key vectors:**  $\mathbf{K} = \mathbf{X}\mathbf{W}_K$  (Shape:  $N_X \times D_Q$ )

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

**Similarities:**  $E = \mathbf{Q}\mathbf{K}^T / \sqrt{D_Q}$  (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:**  $\mathbf{Y} = A\mathbf{V}$  (Shape:  $N_Q \times D_V$ )  $Y_i = \sum_j A_{i,j} \mathbf{V}_j$



# Self-Attention Layer

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

## 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:**  $E = \mathbf{Q}\mathbf{K}^T / \sqrt{D_Q}$  (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:**  $\mathbf{Y} = \mathbf{A}\mathbf{V}$  (Shape:  $N_Q \times D_V$ )  $\mathbf{Y}_i = \sum_j A_{i,j} \mathbf{V}_j$

$x_1$

$x_2$

$x_3$

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

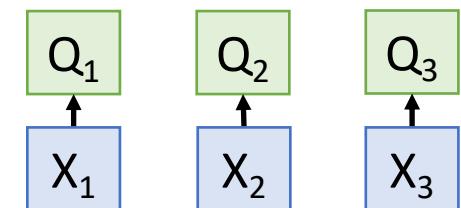
**Key vectors:**  $\mathbf{K} = \mathbf{X}\mathbf{W}_K$  (Shape:  $N_x \times D_Q$ )

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

**Similarities:**  $E = \mathbf{Q}\mathbf{K}^T / \sqrt{D_Q}$  (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:**  $\mathbf{Y} = \mathbf{A}\mathbf{V}$  (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:**  $X$  (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:**  $Q = XW_Q$

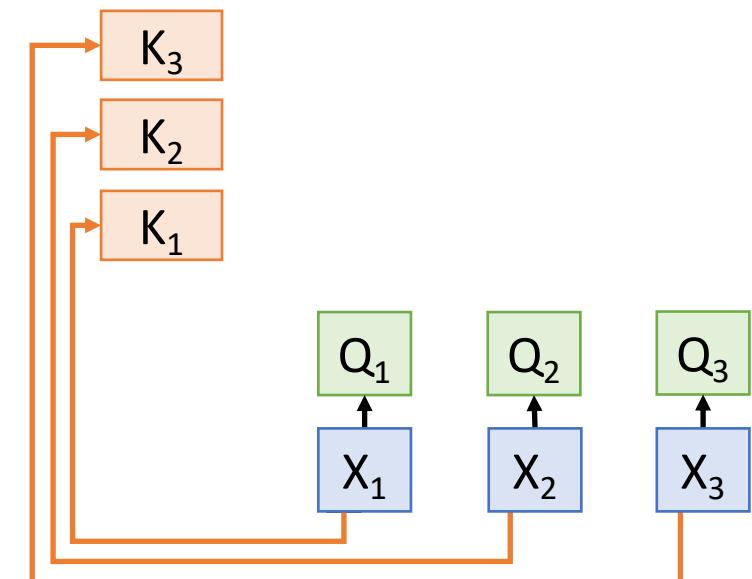
**Key vectors:**  $K = XW_K$  (Shape:  $N_x \times D_Q$ )

**Value Vectors:**  $V = XW_V$  (Shape:  $N_x \times D_V$ )

**Similarities:**  $E = QK^T / \sqrt{D_Q}$  (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

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

## 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:**  $Q = XW_Q$

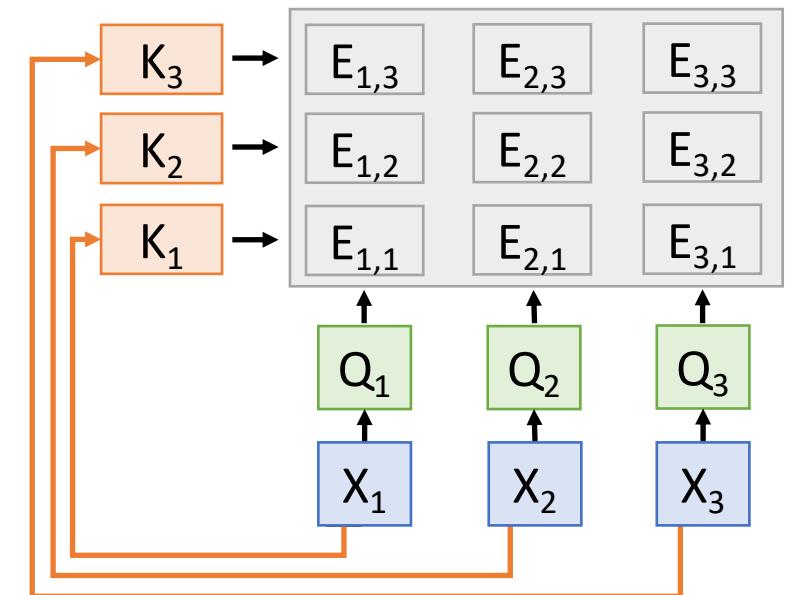
**Key vectors:**  $K = XW_K$  (Shape:  $N_x \times D_Q$ )

**Value Vectors:**  $V = XW_V$  (Shape:  $N_x \times D_V$ )

**Similarities:**  $E = QK^T / \sqrt{D_Q}$  (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

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

## 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:**  $Q = XW_Q$

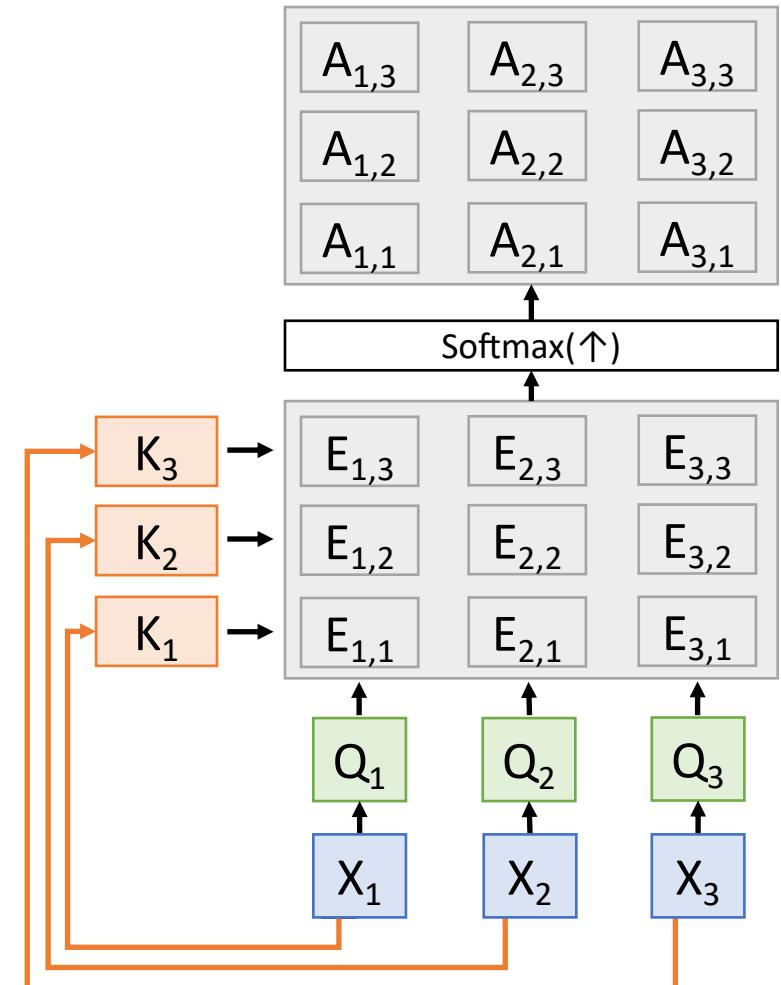
**Key vectors:**  $K = XW_K$  (Shape:  $N_x \times D_Q$ )

**Value Vectors:**  $V = XW_V$  (Shape:  $N_x \times D_V$ )

**Similarities:**  $E = QK^T / \sqrt{D_Q}$  (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

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

## 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:**  $Q = XW_Q$

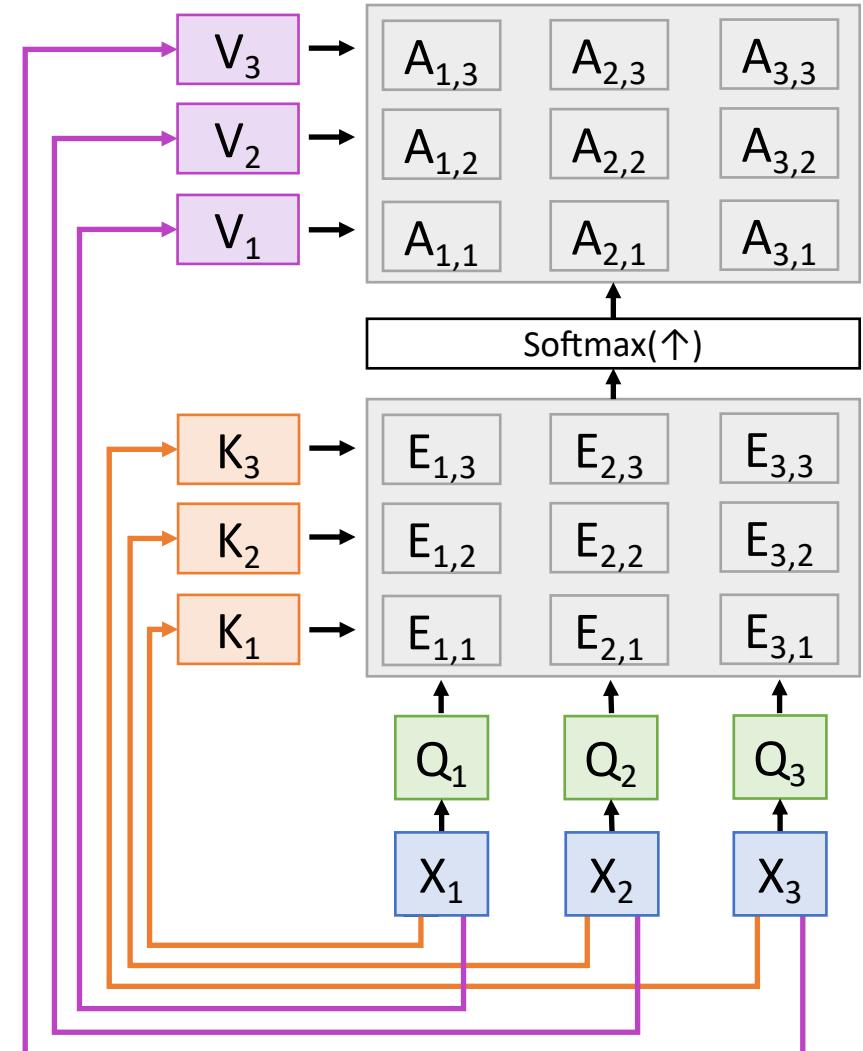
**Key vectors:**  $K = XW_K$  (Shape:  $N_x \times D_Q$ )

**Value Vectors:**  $V = XW_V$  (Shape:  $N_x \times D_V$ )

**Similarities:**  $E = QK^T / \sqrt{D_Q}$  (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:**  $Y = A V$  (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:**  $X$  (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:**  $Q = XW_Q$

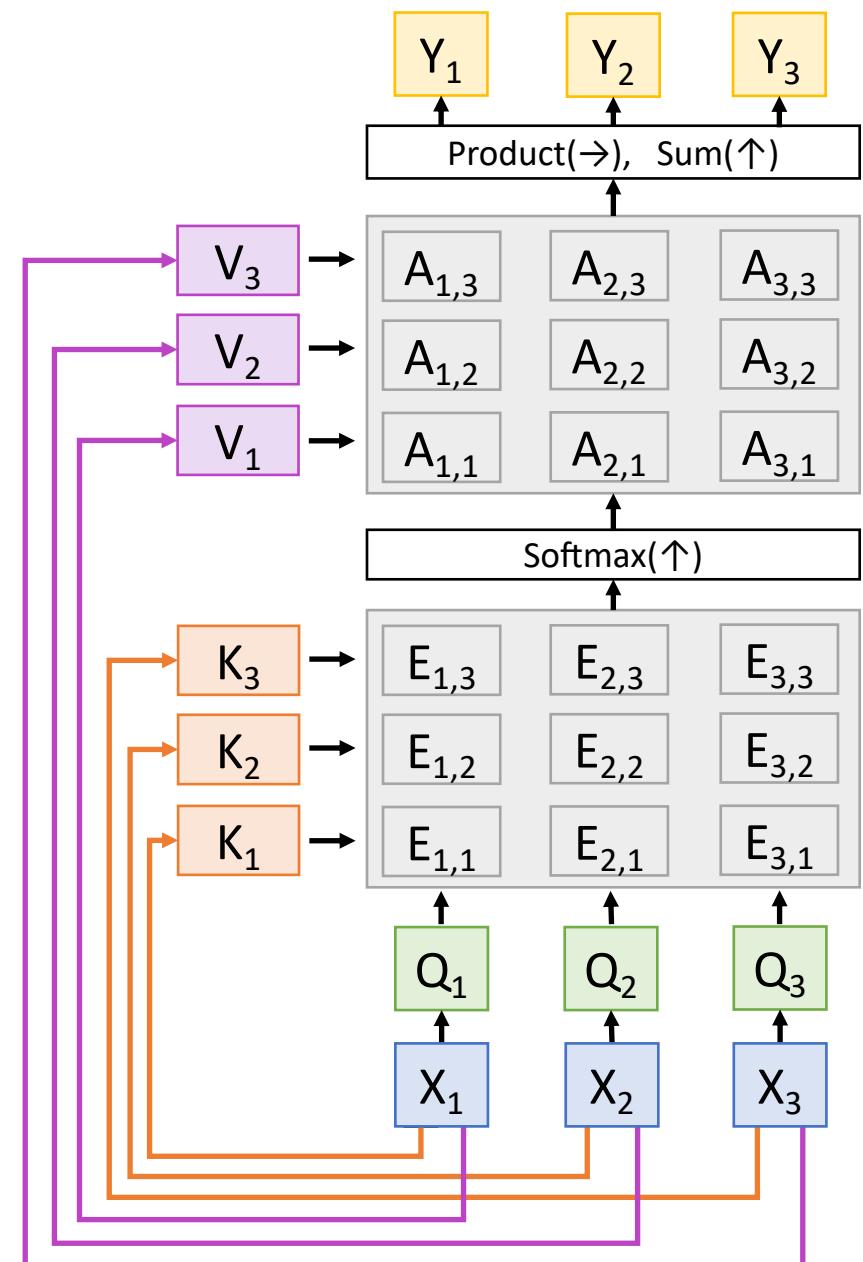
**Key vectors:**  $K = XW_K$  (Shape:  $N_x \times D_Q$ )

**Value Vectors:**  $V = XW_V$  (Shape:  $N_x \times D_V$ )

**Similarities:**  $E = QK^T / \sqrt{D_Q}$  (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:**  $Y = A V$  (Shape:  $N_x \times D_V$ )  $Y_i = \sum_j A_{i,j} V_j$



# Self-Attention Layer

Consider **permuting** the input vectors:

## 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:**  $Q = XW_Q$

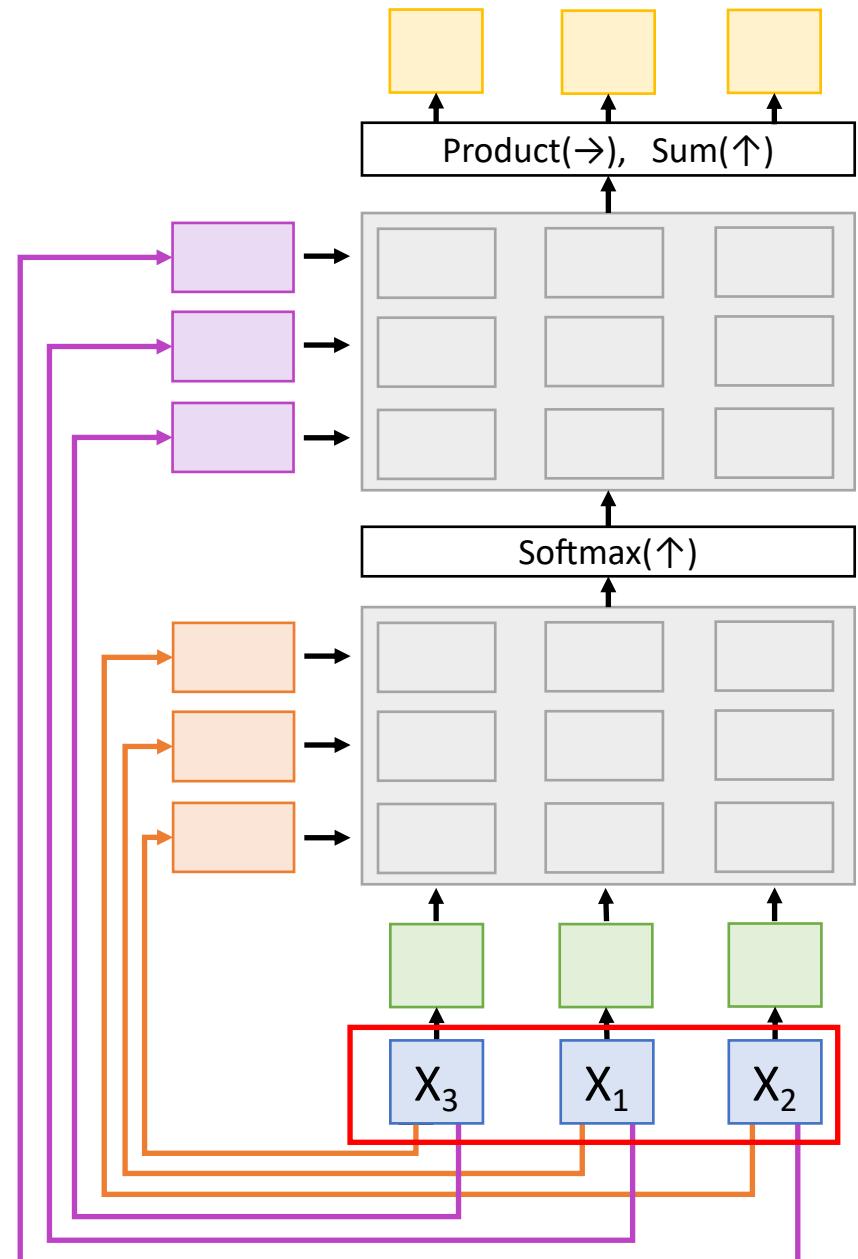
**Key vectors:**  $K = XW_K$  (Shape:  $N_x \times D_Q$ )

**Value Vectors:**  $V = XW_V$  (Shape:  $N_x \times D_V$ )

**Similarities:**  $E = QK^T / \sqrt{D_Q}$  (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

## 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:  $N_x \times D_V$ )

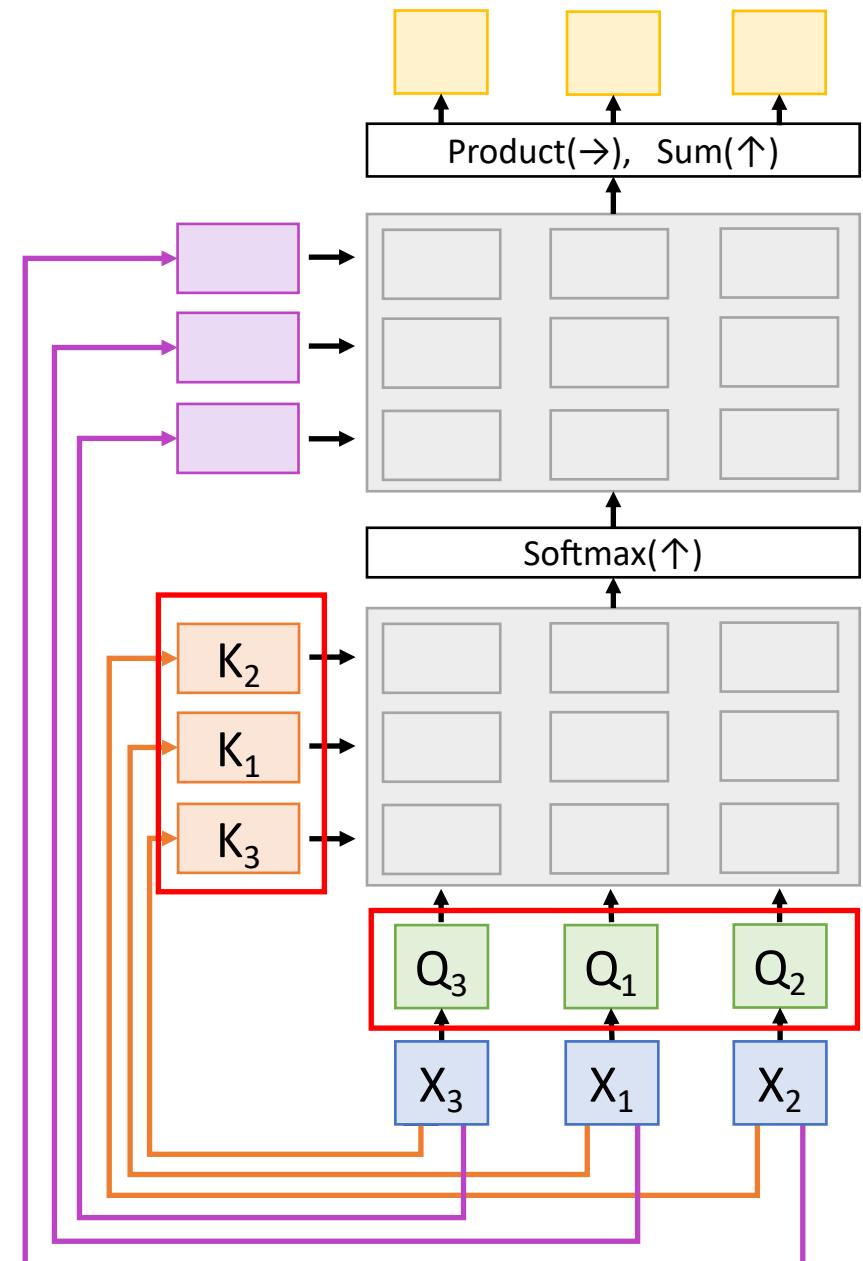
Similarities:  $E = QK^T / \sqrt{D_Q}$  (Shape:  $N_x \times N_x$ )  $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$

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

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

Consider **permuting** the input vectors:

Queries and Keys will be the same, but permuted



# 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:  $N_x \times D_V$ )

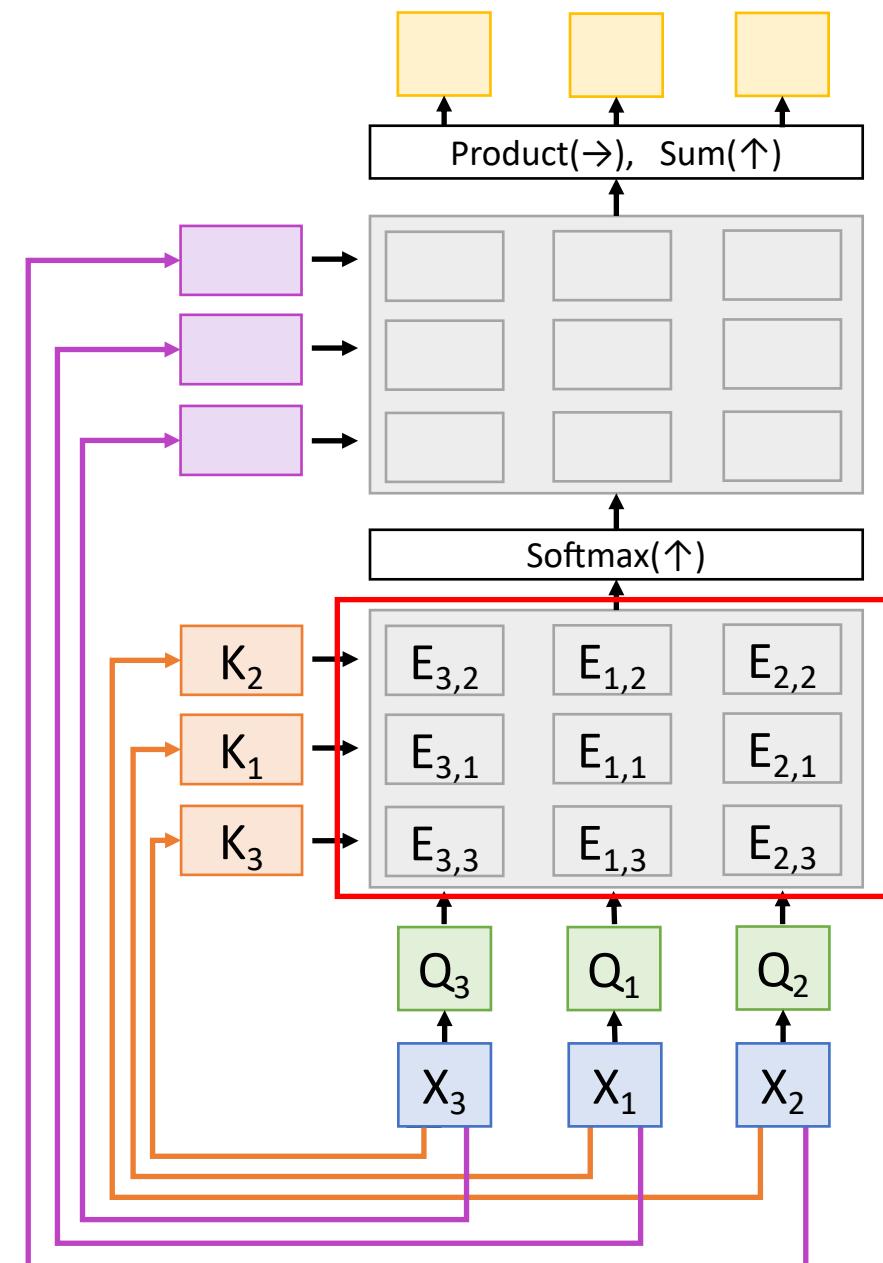
Similarities:  $E = \mathbf{Q}\mathbf{K}^T / \sqrt{D_Q}$  (Shape:  $N_x \times N_x$ )  $E_{i,j} = (\mathbf{Q}_i \cdot \mathbf{K}_j) / \sqrt{D_Q}$

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:

Similarities 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:  $D_x \times D_Q$ )

## Computation:

Query vectors:  $Q = XW_Q$

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

Value Vectors:  $V = XW_V$  (Shape:  $N_x \times D_V$ )

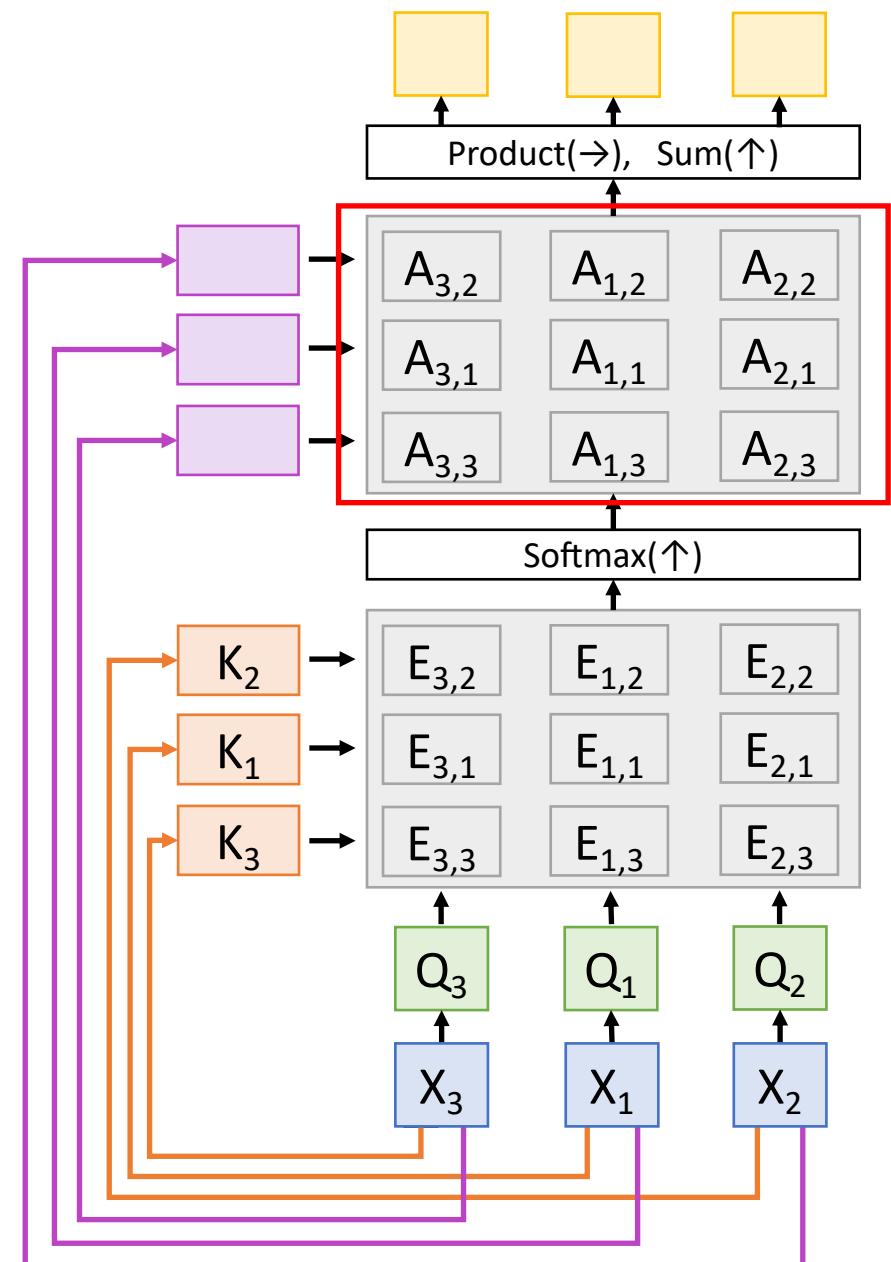
Similarities:  $E = QK^T / \sqrt{D_Q}$  (Shape:  $N_x \times N_x$ )  $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$

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

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

Consider **permuting** the input vectors:

Attention weights will be the same, but permuted



# Self-Attention Layer

## Inputs:

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

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:

Values will be the same, but permuted

## Computation:

Query vectors:  $Q = XW_Q$

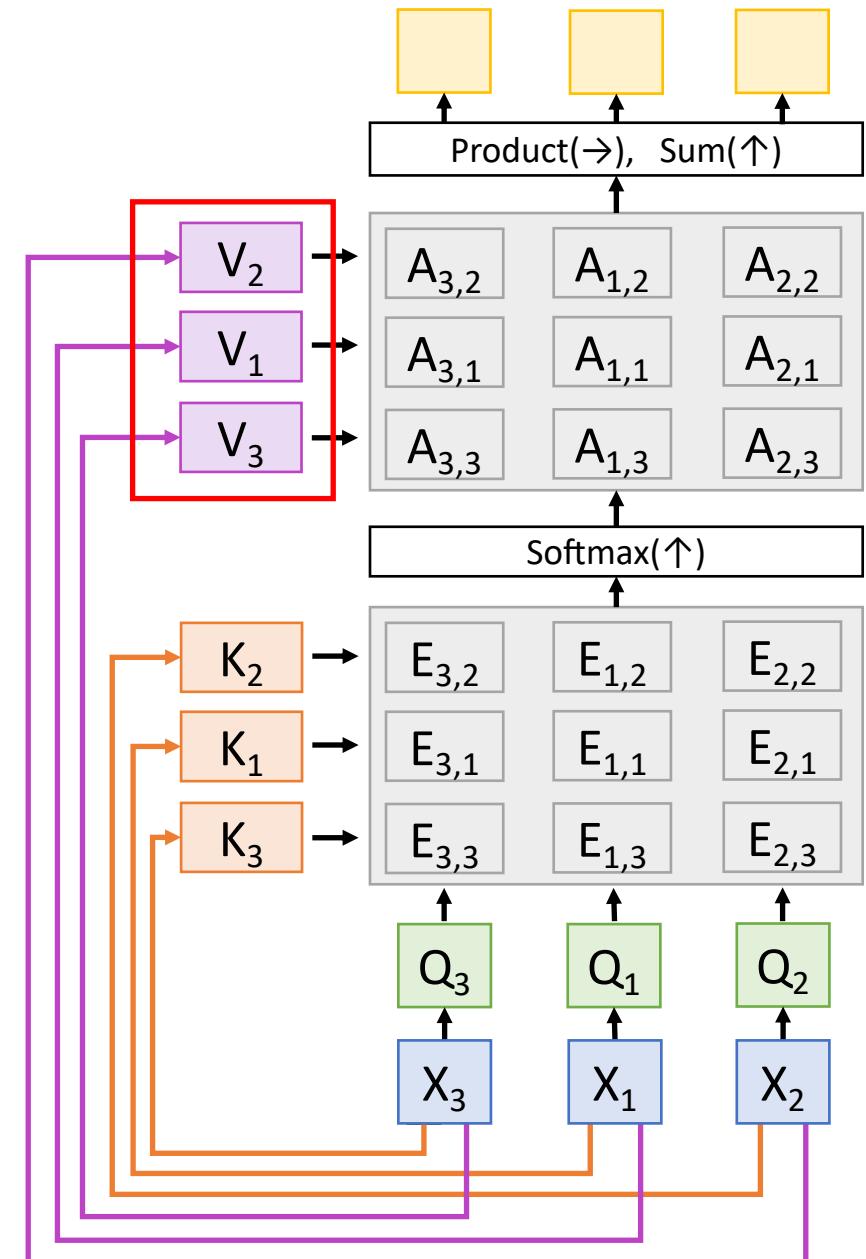
Key vectors:  $K = XW_K$  (Shape:  $N_x \times D_Q$ )

Value Vectors:  $V = XW_V$  (Shape:  $N_x \times D_V$ )

Similarities:  $E = QK^T / \sqrt{D_Q}$  (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:  $Y = A V$  (Shape:  $N_x \times D_V$ )  $Y_i = \sum_j A_{i,j} V_j$



# 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:  $N_x \times D_V$ )

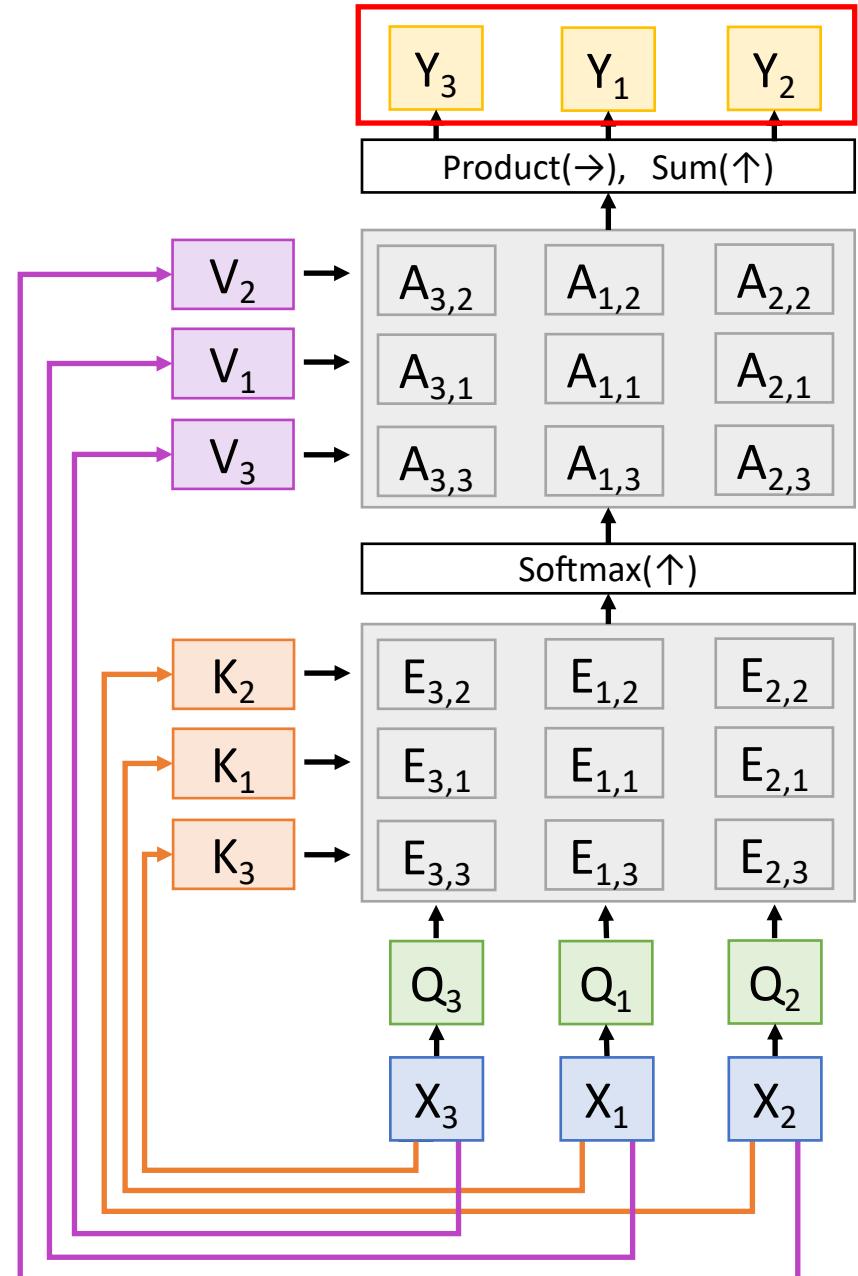
Similarities:  $E = QK^T / \sqrt{D_Q}$  (Shape:  $N_x \times N_x$ )  $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$

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

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

Consider **permuting** the input vectors:

Outputs will be the same, but permuted



# Self-Attention Layer

## 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:  $E = QK^T / \sqrt{D_Q}$  (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:  $N_x \times N_x$ )

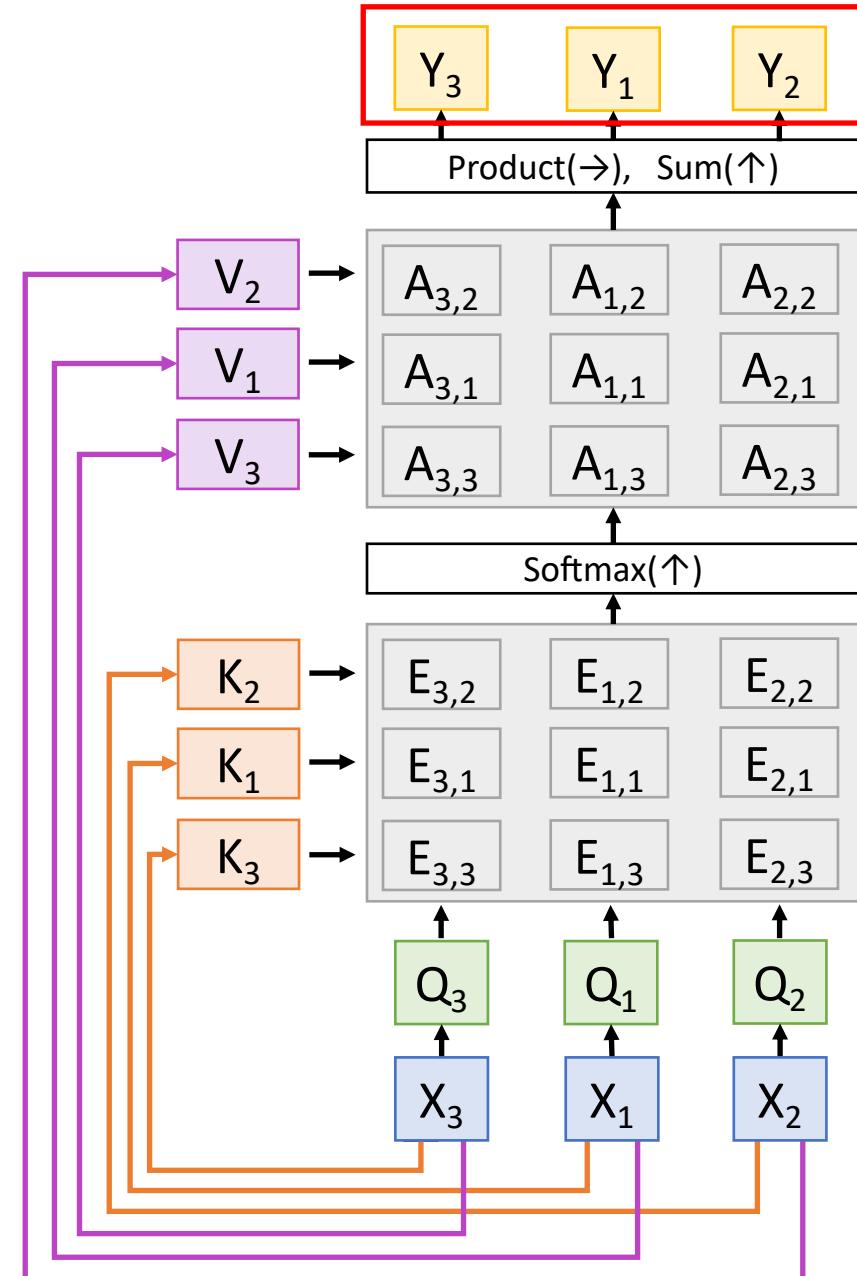
Output vectors:  $Y = A V$  (Shape:  $N_x \times D_V$ )  $Y_i = \sum_j A_{i,j} V_j$

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 works on **sets** of vectors



# Self-Attention Layer

Self attention doesn't  
“know” the order of the  
vectors it is processing!

## 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:**  $Q = XW_Q$

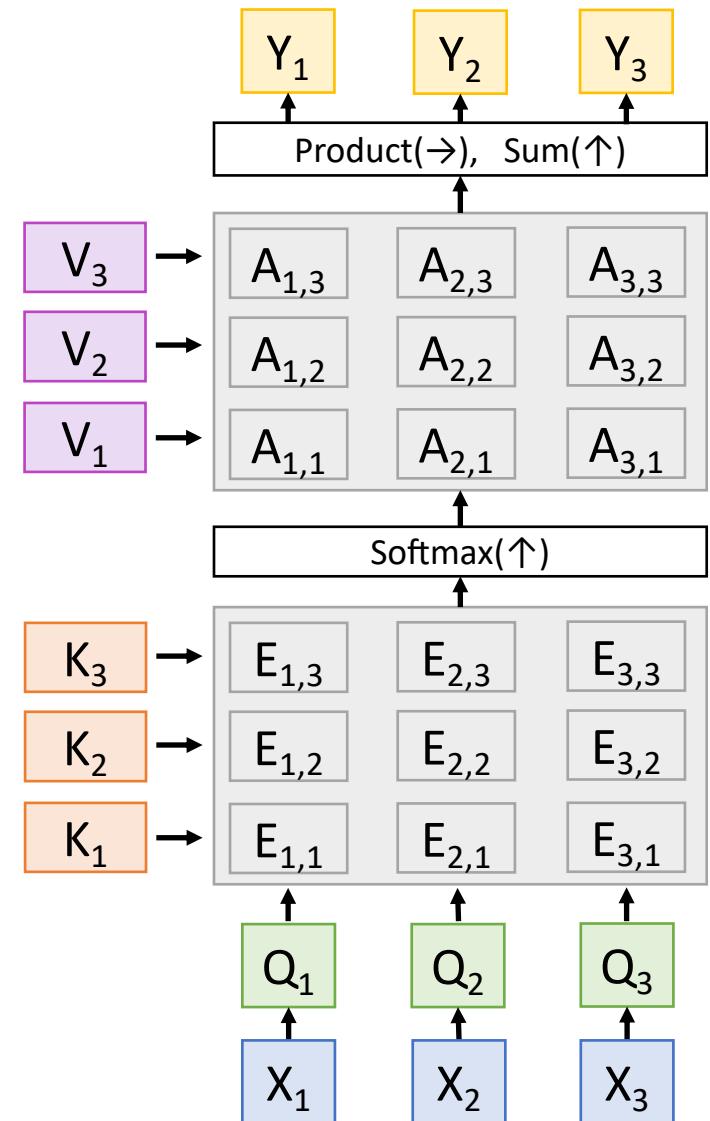
**Key vectors:**  $K = XW_K$  (Shape:  $N_x \times D_Q$ )

**Value Vectors:**  $V = XW_V$  (Shape:  $N_x \times D_V$ )

**Similarities:**  $E = QK^T / \sqrt{D_Q}$  (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

## 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:  $D_x \times D_Q$ )

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**

## Computation:

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

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

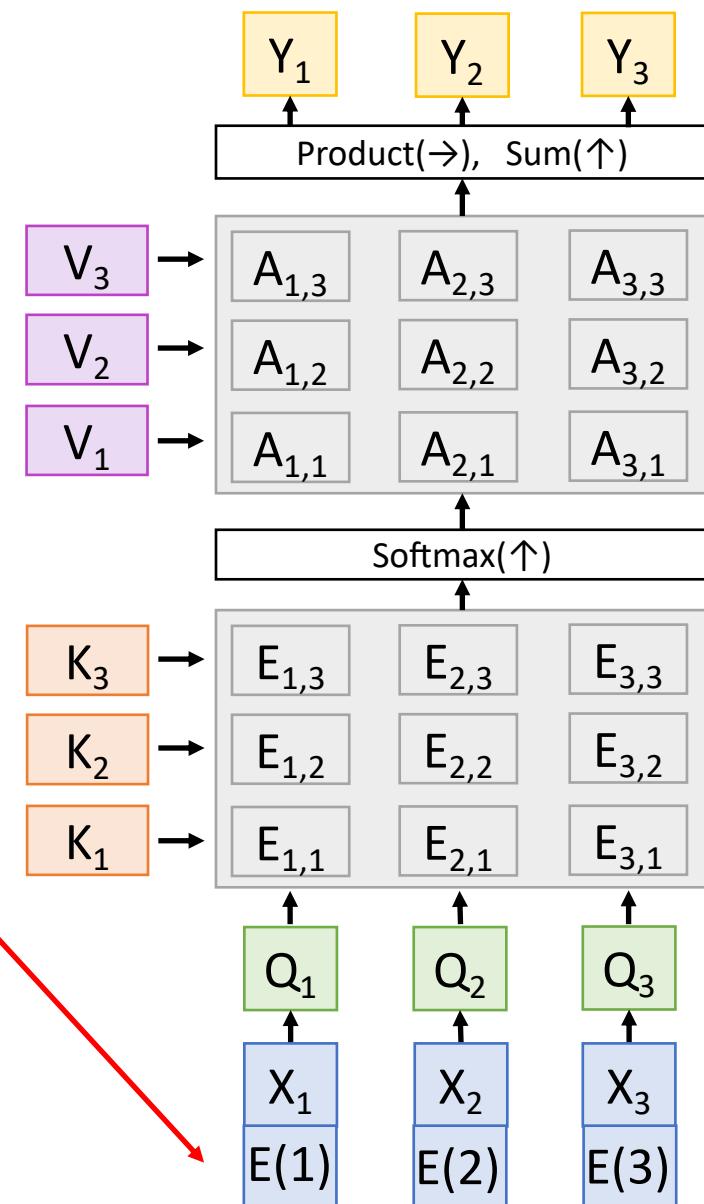
**Value Vectors:**  $V = XW_V$  (Shape:  $N_x \times D_V$ )

**Similarities:**  $E = QK^T / \sqrt{D_Q}$  (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:  $N_x \times D_V$ )  $Y_i = \sum_j A_{i,j} V_j$

$E$  can be learned lookup  
table, or fixed function



# Masked Self-Attention Layer

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

## 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:  $Q = XW_Q$

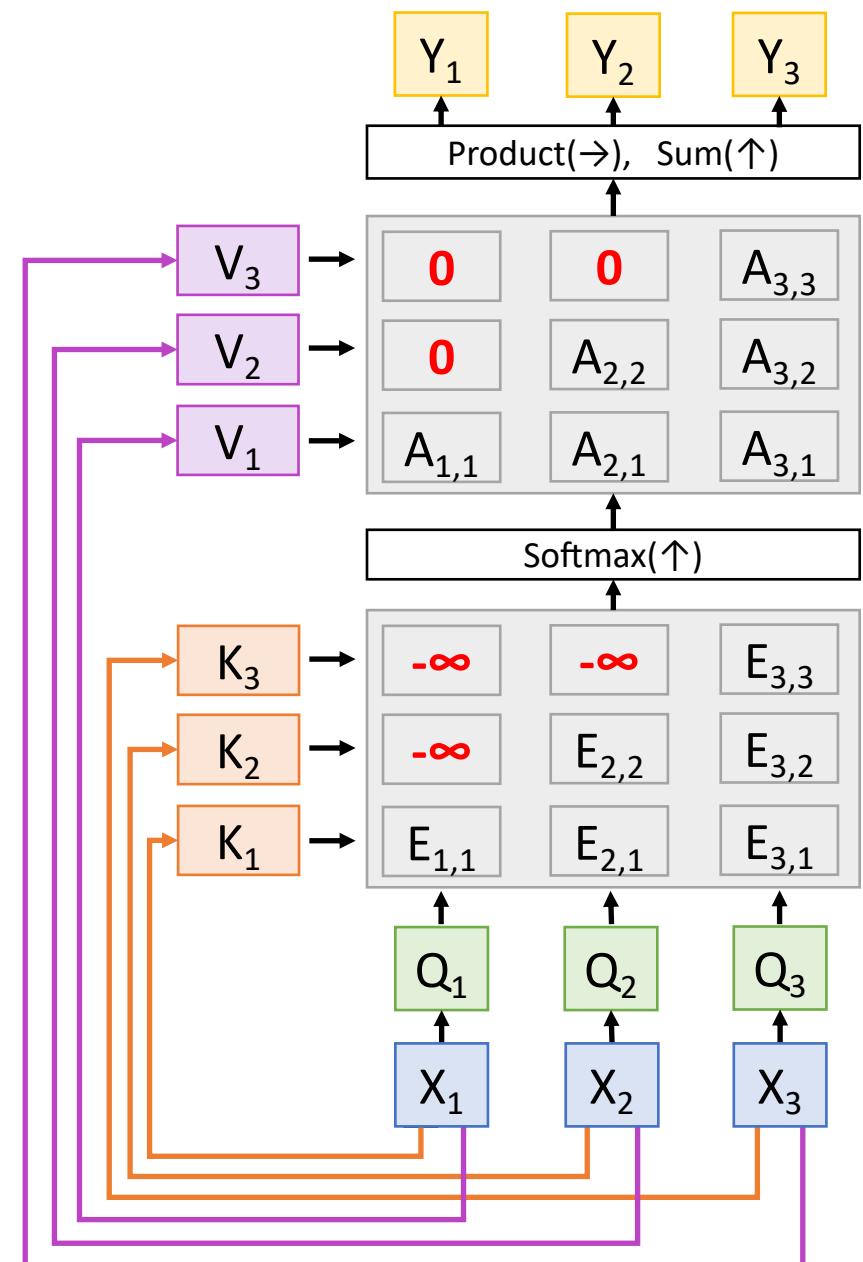
Key vectors:  $K = XW_K$  (Shape:  $N_x \times D_Q$ )

Value Vectors:  $V = XW_V$  (Shape:  $N_x \times D_V$ )

Similarities:  $E = QK^T / \sqrt{D_Q}$  (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:  $N_x \times D_V$ )  $Y_i = \sum_j A_{i,j} V_j$



# Masked Self-Attention Layer

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

Used for language modeling (predict next word)

## 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:**  $Q = XW_Q$

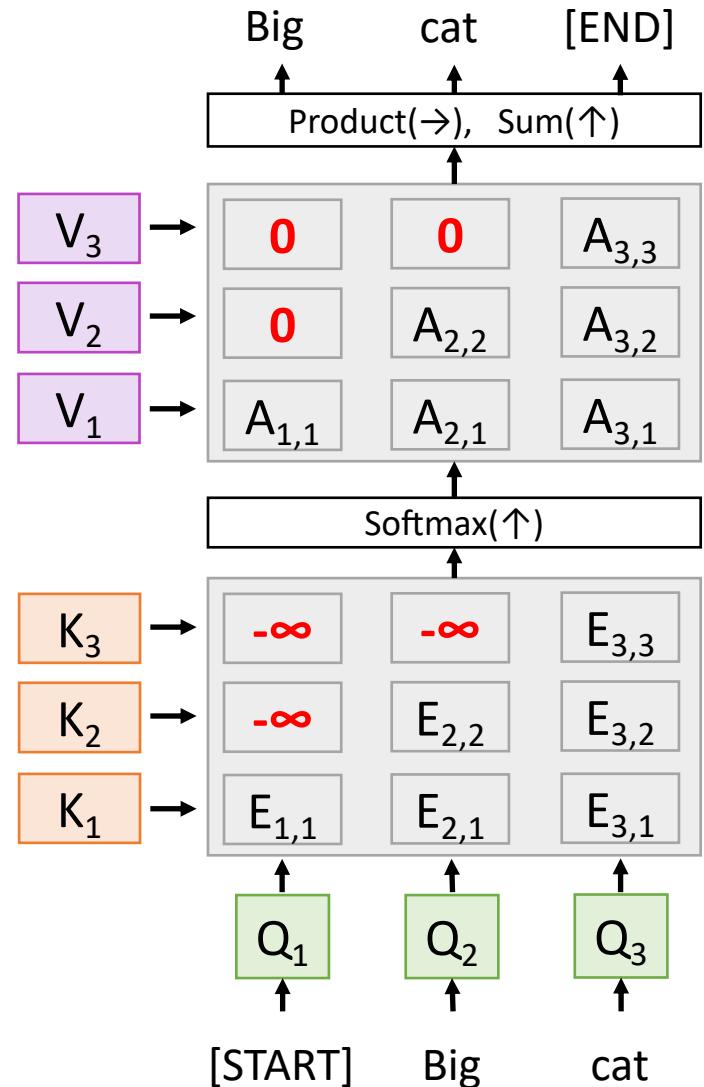
**Key vectors:**  $K = XW_K$  (Shape:  $N_x \times D_Q$ )

**Value Vectors:**  $V = XW_V$  (Shape:  $N_x \times D_V$ )

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

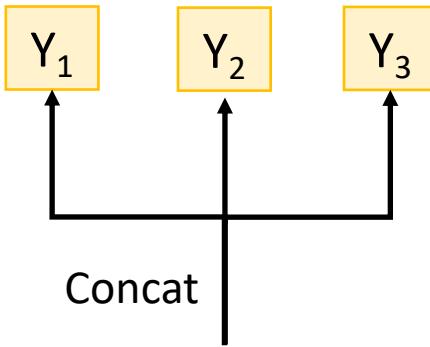
**Attention weights:**  $A = \text{softmax}(E, \text{dim}=1)$  (Shape:  $N_x \times N_x$ )

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



# Multihead Self-Attention Layer

Use H independent  
“Attention Heads” in parallel



## 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:  $N_x \times D_Q$ )

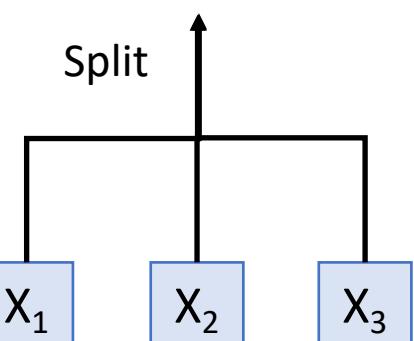
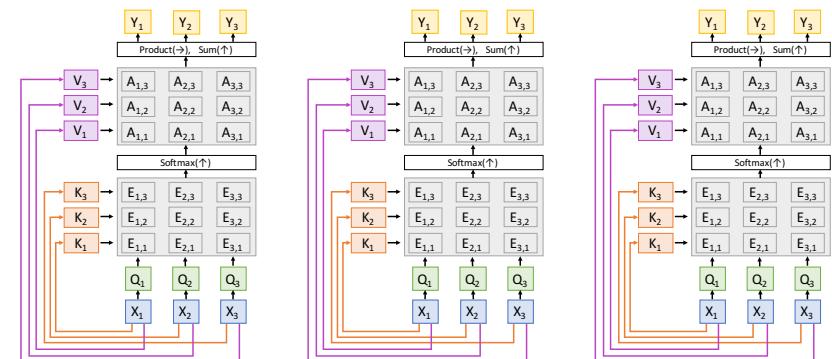
**Value Vectors:**  $V = XW_V$  (Shape:  $N_x \times D_V$ )

**Similarities:**  $E = QK^T / \sqrt{D_Q}$  (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:  $N_x \times D_V$ )  $Y_i = \sum_j A_{i,j} V_j$

**Hyperparameters:**  
Query dimension  $D_Q$   
Number of heads  $H$



Typically, if the dimension of the inputs  $X$  is  $D$  and there are  $H$  heads, the values, queries, and keys will all be of size  $D/H$ , as this allows for an efficient implementation.

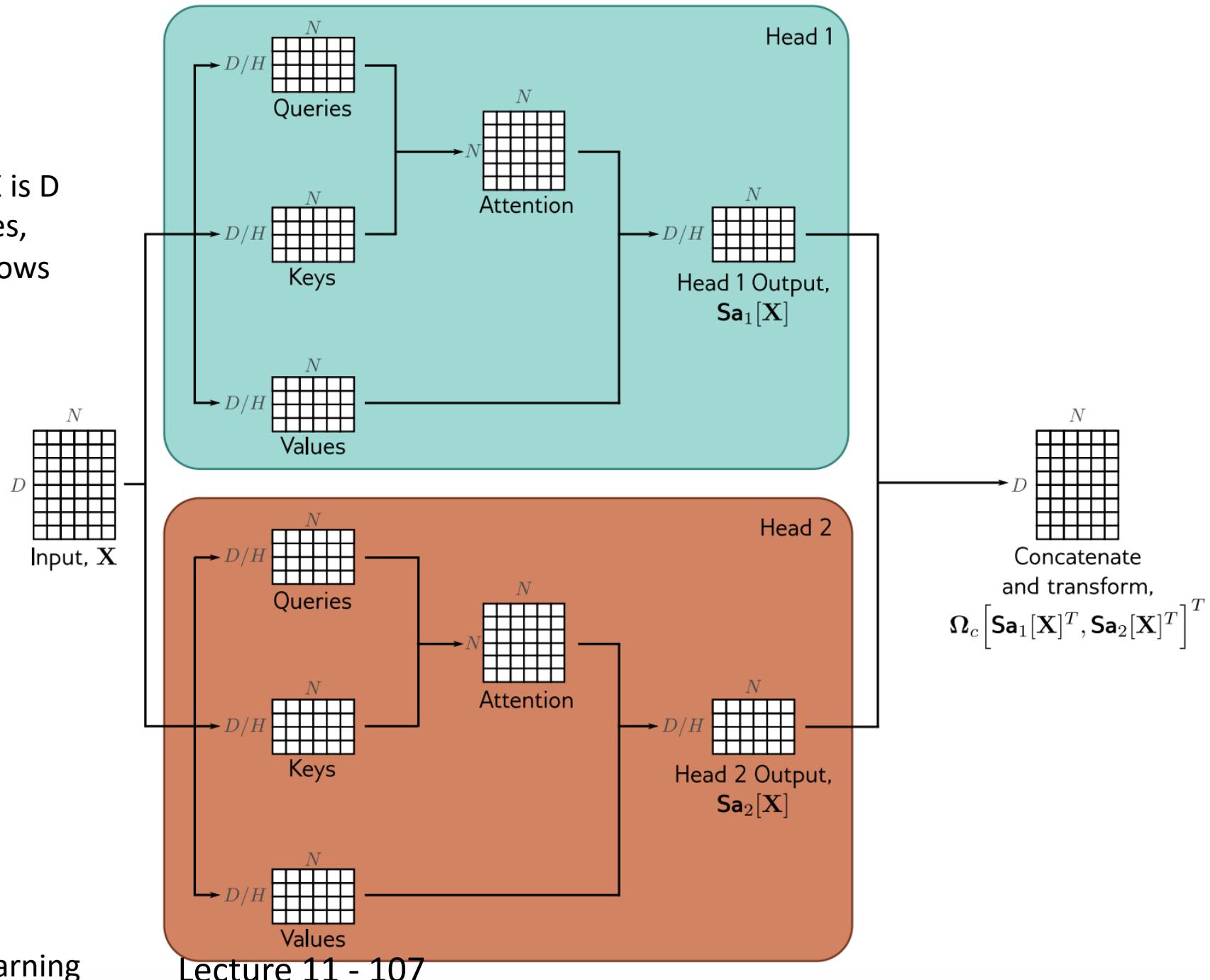
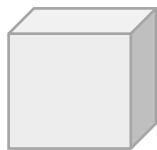
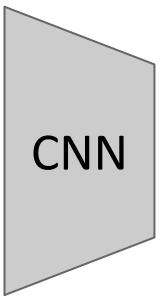


Figure 12.6 from Understanding Deep Learning

# Example: CNN with Self-Attention

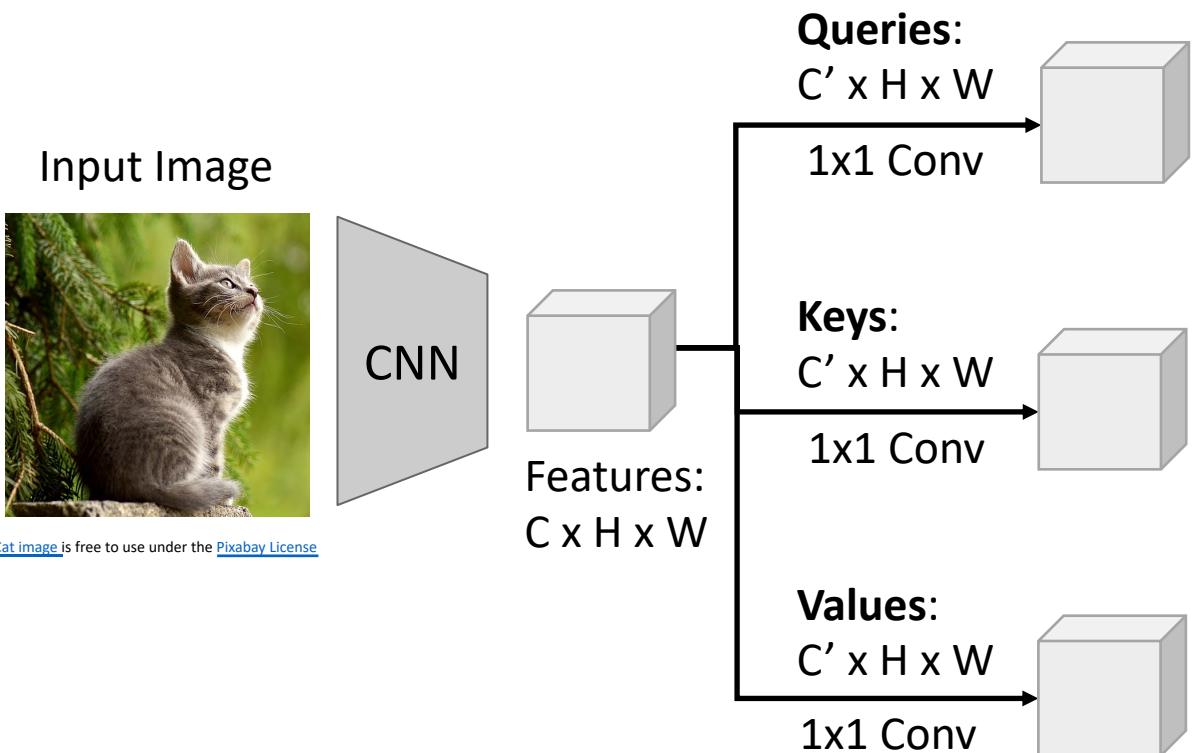
Input Image



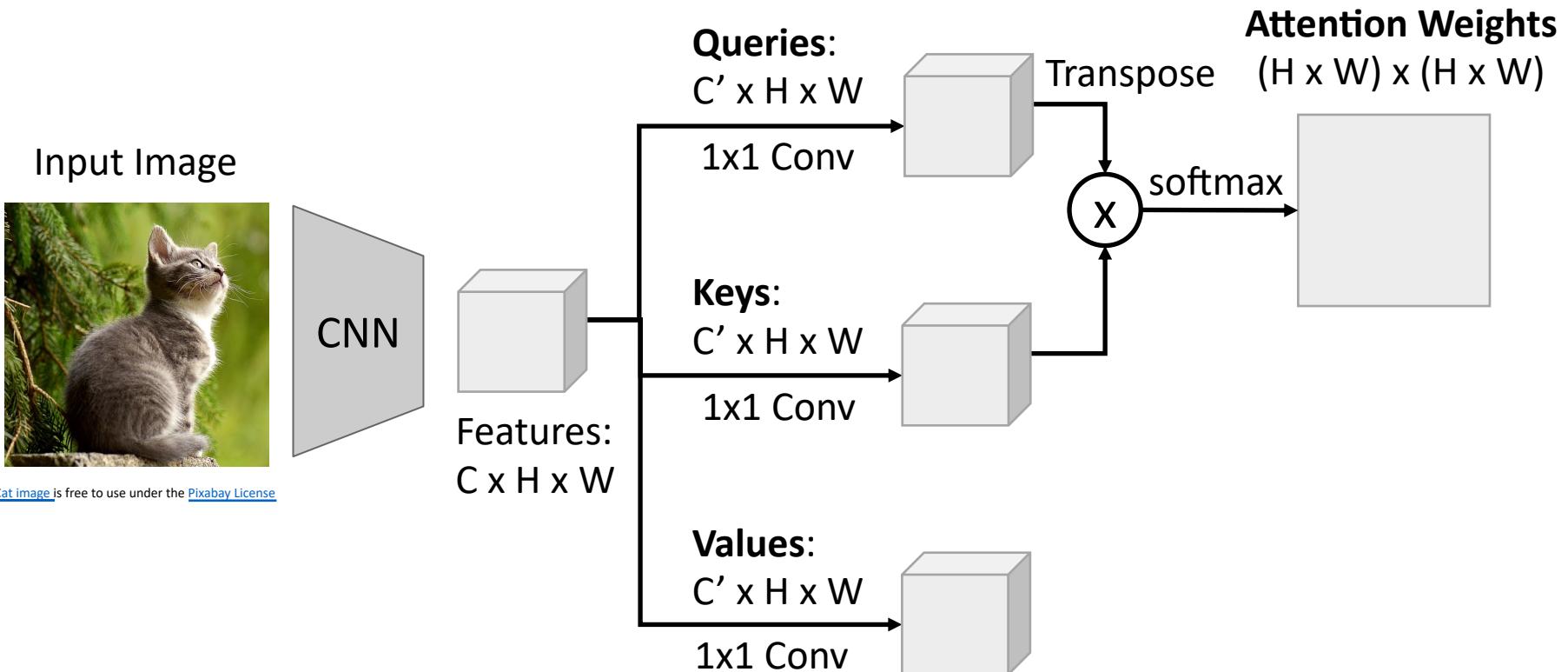
Features:  
 $C \times H \times W$

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

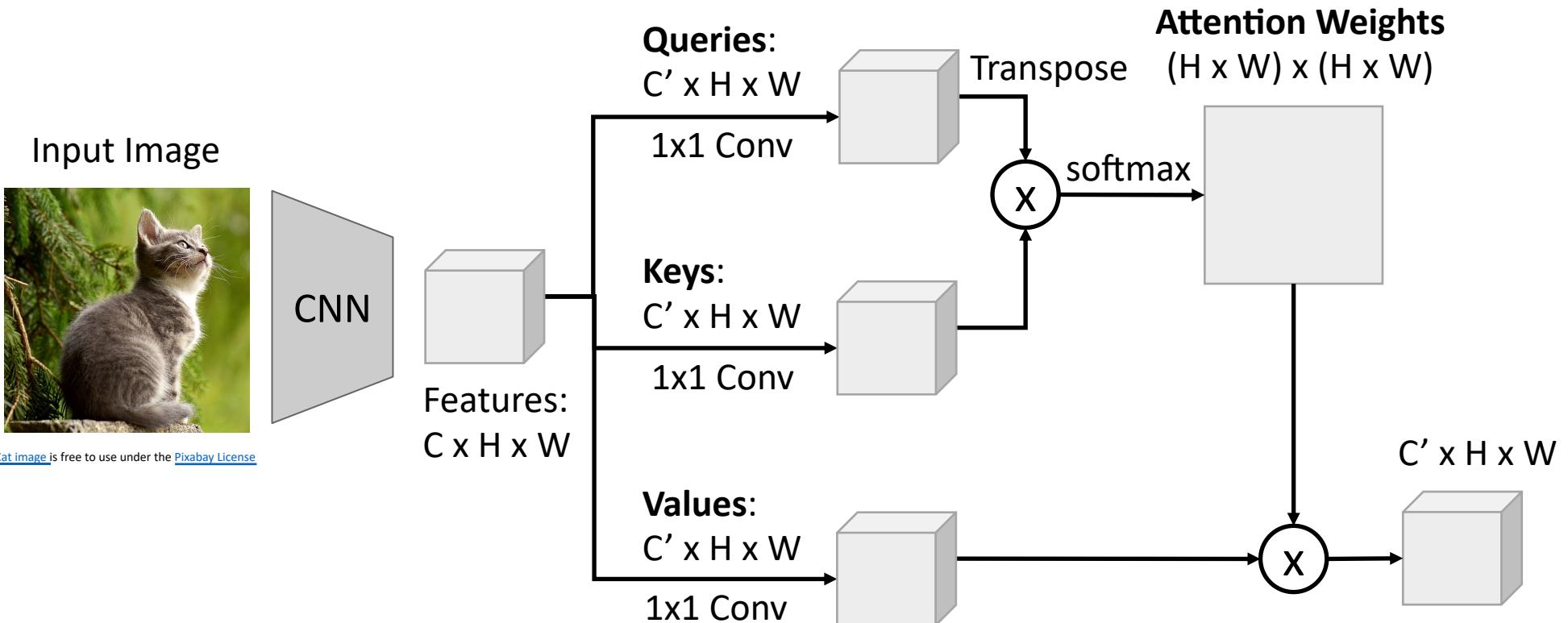
# Example: CNN with Self-Attention



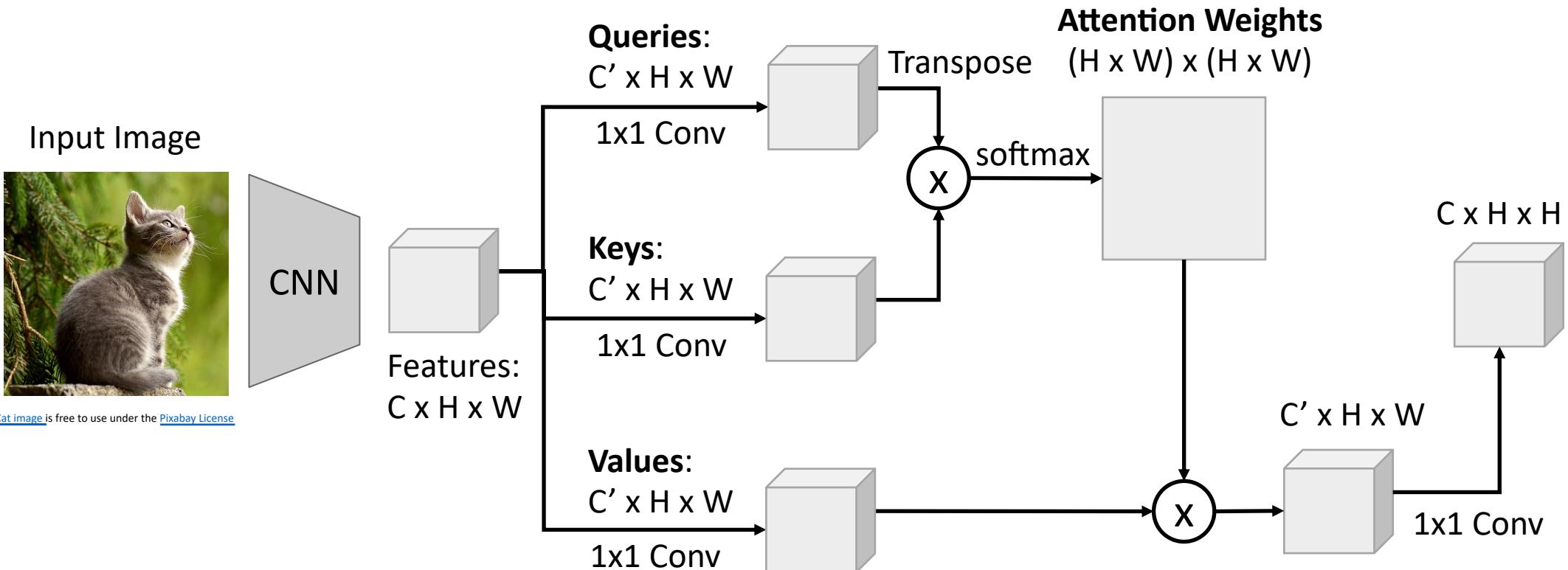
# Example: CNN with Self-Attention



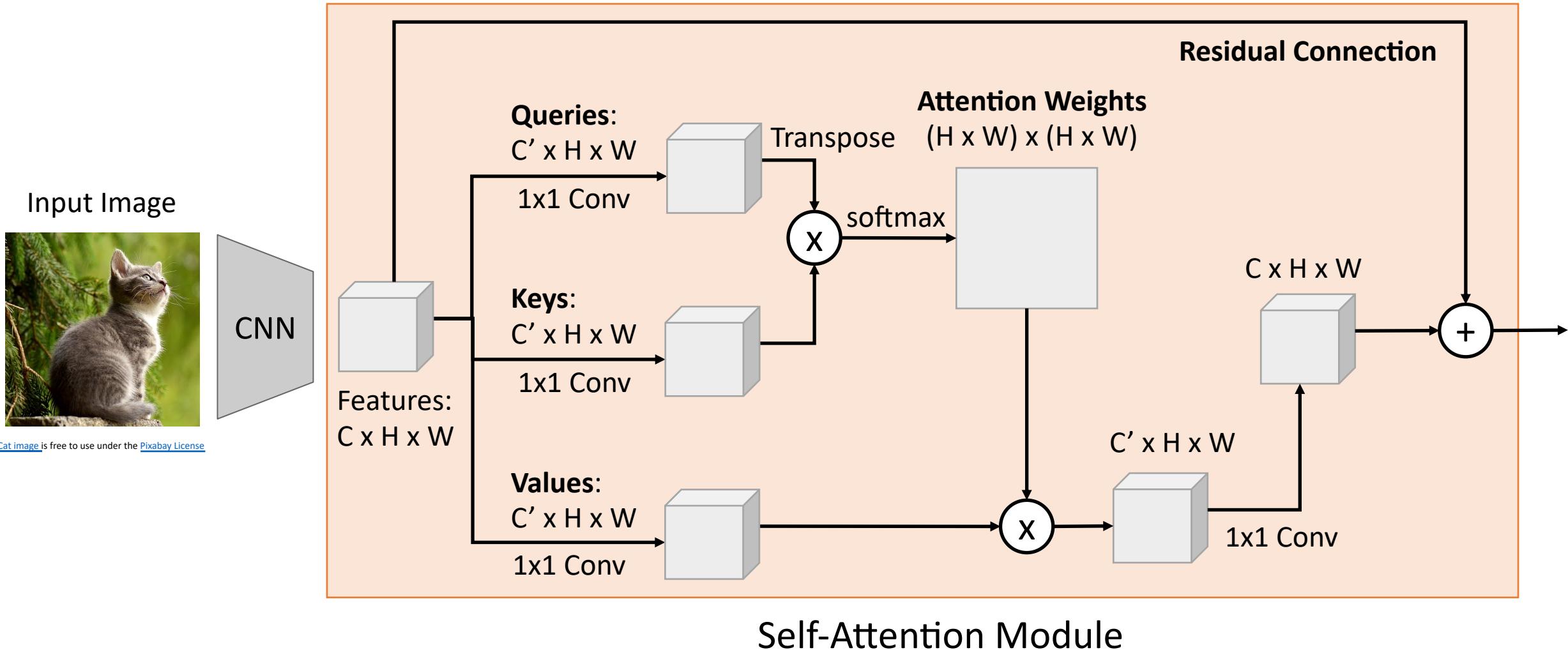
# Example: CNN with Self-Attention



# Example: CNN with Self-Attention

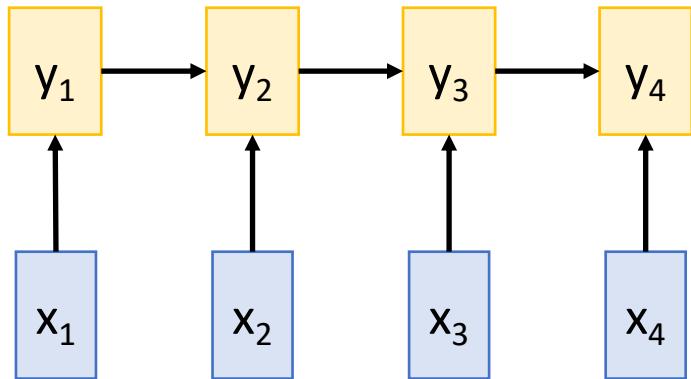


# Example: CNN with Self-Attention



# Three Ways of Processing Sequences

## Recurrent Neural Network



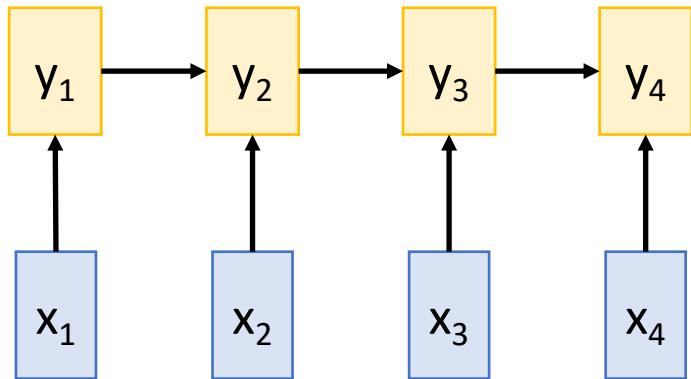
Works on **Ordered Sequences**

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

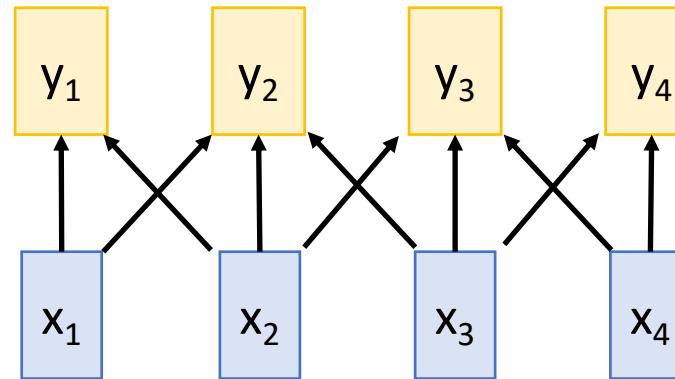
(-) **Not parallelizable:** need to compute hidden states sequentially

# 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

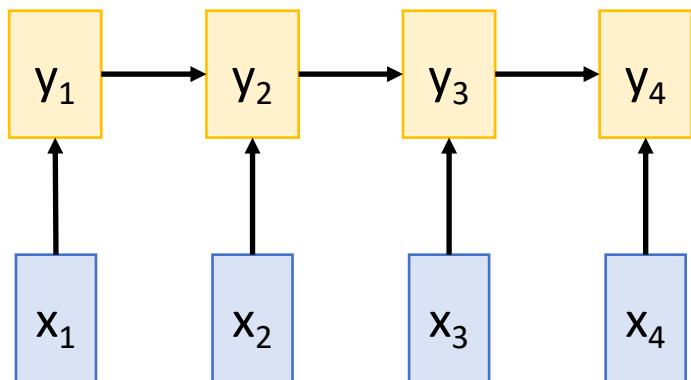
Works on **Multidimensional Grids**

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

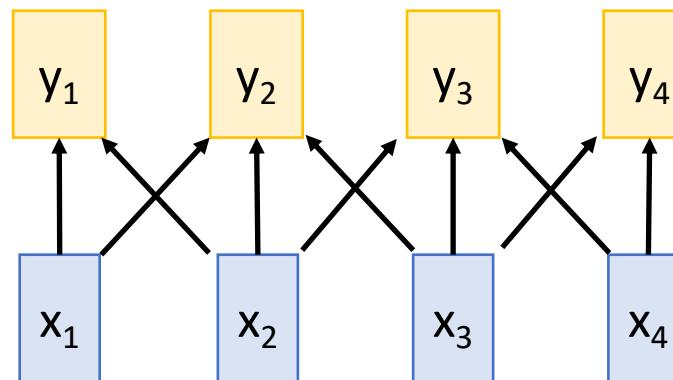
(+) **Highly parallel:** Each output can be computed in parallel

# Three Ways of Processing Sequences

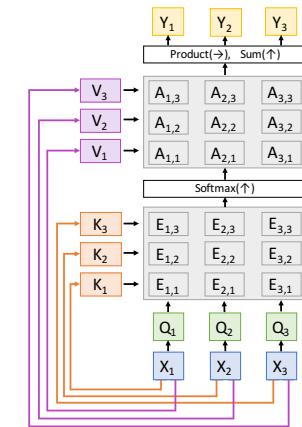
Recurrent Neural Network



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**

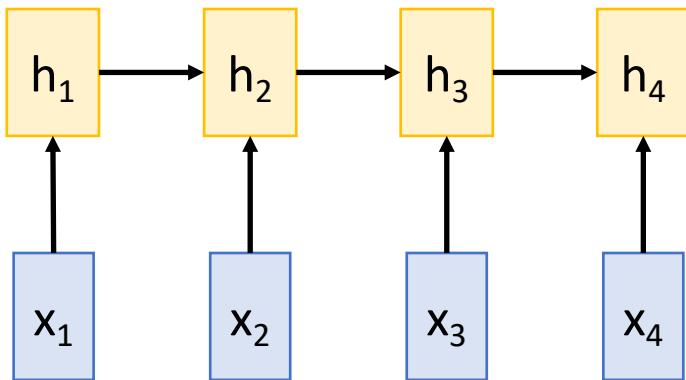
- (-) Bad at long sequences: Need to stack many conv layers for outputs to "see" the whole sequence
- (+) Highly parallel: Each output can be computed in parallel

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

# Summary

Recurrent Neural Network  
(RNN) for modeling sequence



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

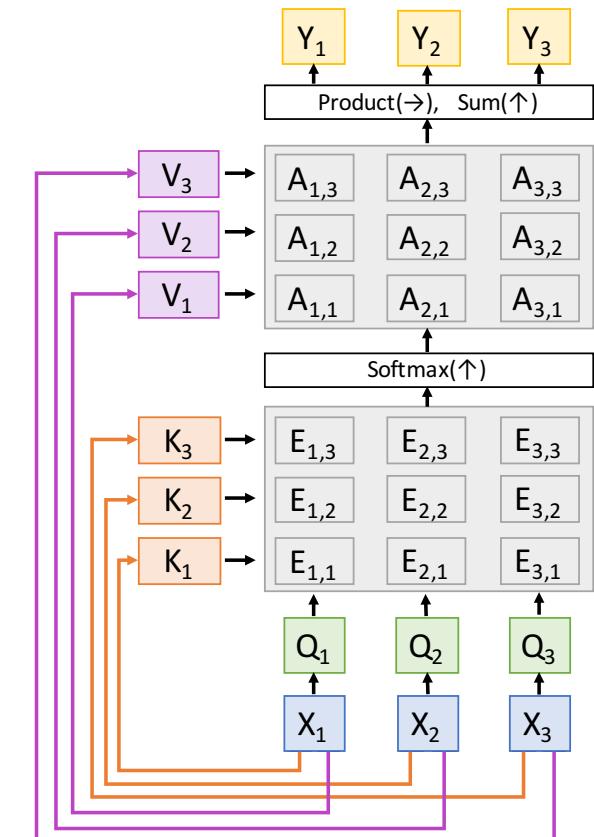
Adding **Attention** to RNN models lets them look at different parts of the input at each timestep



A dog is standing on a hardwood floor.

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

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



# Reading and Practice

Dive into Deep learning:

RNN: Chapter 9

Attention: Chapter 10.1, 10.2, 10.3, 10.5, 10.6

Understanding Deep Learning:

Chapter 12