# Lecture 16:
# Recurrent Neural Networks

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# Visual Question Answering



Q: What endangered animal is featured on the truck?
A: **A bald eagle.**
A: A sparrow.
A: A humming bird.
A: A raven.

Q: Where will the driver go if turning right?
A: **Onto 24 ¾ Rd.**
A: Onto 25 ¾ Rd.
A: Onto 23 ¾ Rd.
A: Onto Main Street.

Q: When was the picture taken?
A: **During a wedding.**
A: During a bar mitzvah.
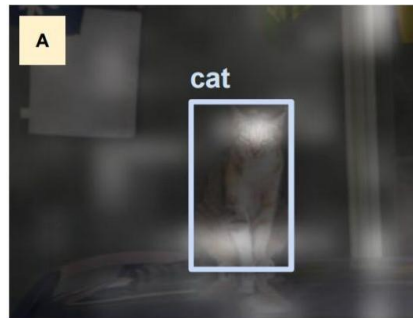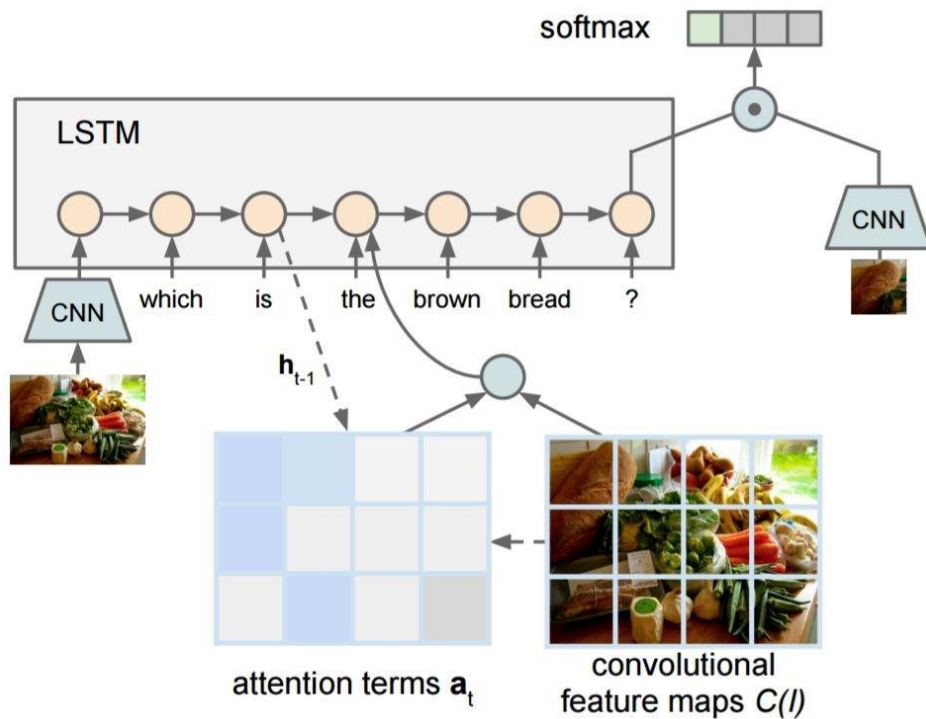A: During a funeral.
A: During a Sunday church service

Q: Who is under the umbrella?
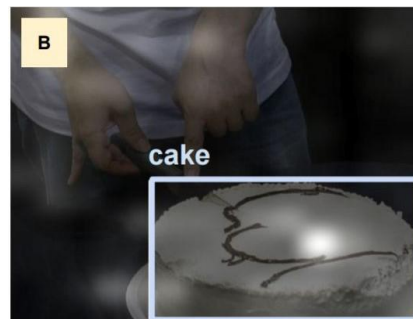A: **Two women.**
A: A child.
A: An old man.
A: A husband and a wife.

Agrawal et al, "VQA: Visual Question Answering", ICCV 2015
Zhu et al, "Visual 7W: Grounded Question Answering in Images", CVPR 2016
Figure from Zhu et al, copyright IEEE 2016. Reproduced for educational purposes.

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# Visual Question Answering: RNNs with Attention



Zhu et al, "Visual 7W: Grounded Question Answering in Images", CVPR 2016
Figures from Zhu et al, copyright IEEE 2016. Reproduced for educational purposes.
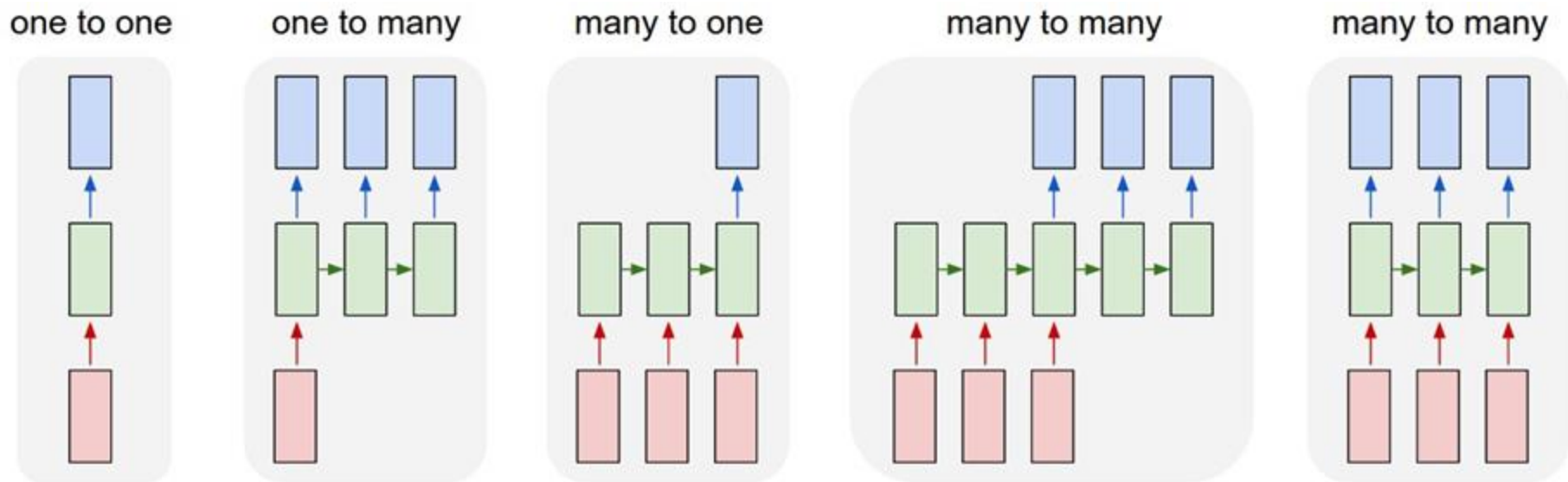
Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# Recurrent Networks offer a lot of flexibility:



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

**Vanilla Neural Networks**

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# Recurrent Networks offer a lot of flexibility:

one to one   one to many   many to one   many to many   many to many

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

# Recurrent Networks offer a lot of flexibility:

one to one     one to many     many to one     many to many     many to many
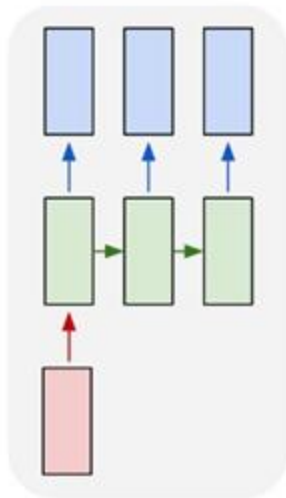
e.g. **Sentiment Classification**
sequence of words -> sentiment

# Recurrent Networks offer a lot of flexibility:

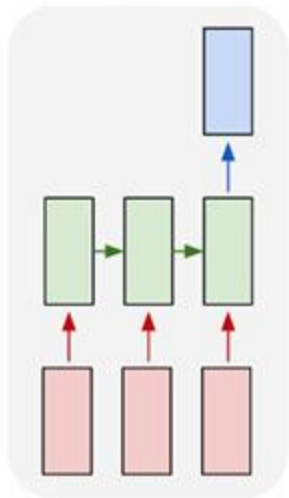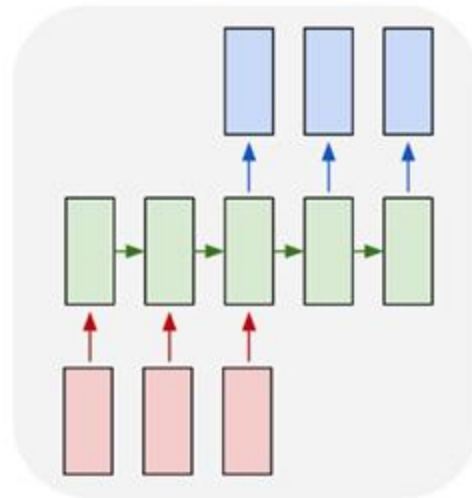one to one     one to many     many to one     many to many     many to many

e.g. **Machine Translation**
seq of words -> seq of words

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# Recurrent Networks offer a lot of flexibility:



e.g. **Video classification on frame level**

# Recurrent Neural Network

# Recurrent Neural Network



usually want to predict a vector at some time steps

10

# Recurrent Neural Network

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

some function with parameters W

old state

input vector at some time step

y

RNN

x

# Recurrent Neural Network

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

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

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

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# (Vanilla) Recurrent Neural Network

The state consists of a single *"hidden"* vector **h**:



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

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

$$y_t = W_{hy}h_t$$

**Character-level language model example**

Vocabulary:
[h,e,l,o]

Example training sequence:
**"hello"**

**Character-level language model example**

Vocabulary:
[h,e,l,o]

Example training sequence:
**"hello"**

# Character-level language model example

Vocabulary: [h,e,l,o]

Example training sequence: **"hello"**

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

# Character-level language model example

Vocabulary: [h,e,l,o]

Example training sequence: **"hello"**

# Image Captioning



Explain Images with Multimodal Recurrent Neural Networks, Mao et al.
Deep Visual-Semantic Alignments for Generating Image Descriptions, Karpathy and Fei-Fei
Show and Tell: A Neural Image Caption Generator, Vinyals et al.
Long-term Recurrent Convolutional Networks for Visual Recognition and Description, Donahue et al.
Learning a Recurrent Visual Representation for Image Caption Generation, Chen and Zitnick

# **Recurrent Neural Network**



# **Convolutional Neural Network**

test image

| 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 |

test image

image ← test 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

X

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

test image

x0
<START>

<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

v

**Wih**

test image

y0

h0

x0
<STA
RT>

<START>

**before:**

h = tanh(Wxh * x + Whh * h)

**now:**

h = tanh(Wxh * x + Whh * h **+ Wih * v**)

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

test image

y0

h0

x0
<START>

straw

sample!

<START>

test image

sample!

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

y0   y1

h0 → h1

x0
<START>   straw   hat

<START>

test image

sample <END> token => finish.

# Image Sentence Datasets

a man riding a bike on a dirt path through a forest.
bicyclist raises his fist as he rides on desert dirt trail.
this dirt bike rider is smiling and raising his fist in triumph.
a man riding a bicycle while pumping his fist in the air.
a mountain biker pumps his fist in celebration.



Microsoft COCO
*[Tsung-Yi Lin et al. 2014]*
mscoco.org


currently:
~120K images
~5 sentences each

"man in black shirt is playing guitar."

"construction worker in orange safety vest is working on road."

"two young girls are playing with lego toy."

"boy is doing backflip on wakeboard."

"man in black shirt is playing guitar."

"construction worker in orange safety vest is working on road."

"two young girls are playing with lego toy."

"boy is doing backflip on wakeboard."

"a young boy is holding a baseball bat."

"a cat is sitting on a couch with a remote control."

"a woman holding a teddy bear in front of a mirror."

"a horse is standing in the middle of a road."

# RNN:

$$h_t^l = \tanh W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$h \in \mathbb{R}^n$, $\qquad W^l \ [n \times 2n]$



depth

time

# Recurrent Neural Networks have loops



Figure credit: *Understanding LSTM Networks* on colah's blog

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# An unrolled recurrent neural network



Figure credit: *Understanding LSTM Networks* on colah's blog

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# Problem of Long-Term Dependencies

"the clouds are in the *sky*"



Figure credit: *Understanding LSTM Networks* on colah's blog

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# Problem of Long-Term Dependencies

"I grew up in France… I speak fluent *French*."



Figure credit: *Understanding LSTM Networks* on colah's blog

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# RNN:

$$h_t^l = \tanh W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$$h \in \mathbb{R}^n \qquad W^l \; [n \times 2n]$$

# LSTM:

$$W^l \; [4n \times 2n]$$

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

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

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



depth

time

# Summary

- RNNs allow a lot of flexibility in architecture design
- Vanilla RNNs are simple but don't work very well
- Common to use LSTM or GRU: their additive interactions improve gradient flow
- Backward flow of gradients in RNN can explode or vanish. Exploding is controlled with gradient clipping. Vanishing is controlled with additive interactions (LSTM)
- Better/simpler architectures are a hot topic of current research
- Better understanding (both theoretical and empirical) is needed.

# Sequence to Sequence with RNNs

**Input**: Sequence $x_1, \dots x_T$
**Output**: Sequence $y_1, \dots, y_{T'}$

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

Subhransu Maji, Chuang Gan and TAs

# Sequence to Sequence with RNNs

**Input**: Sequence $x_1, \ldots x_T$
**Output**: Sequence $y_1, \ldots, y_{T'}$

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

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



| $h_1$ | → | $h_2$ | → | $h_3$ | → | $h_4$ | → | $s_0$ |

c

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

we          are          eating          bread

Subhransu Maji, Chuang Gan and TAs
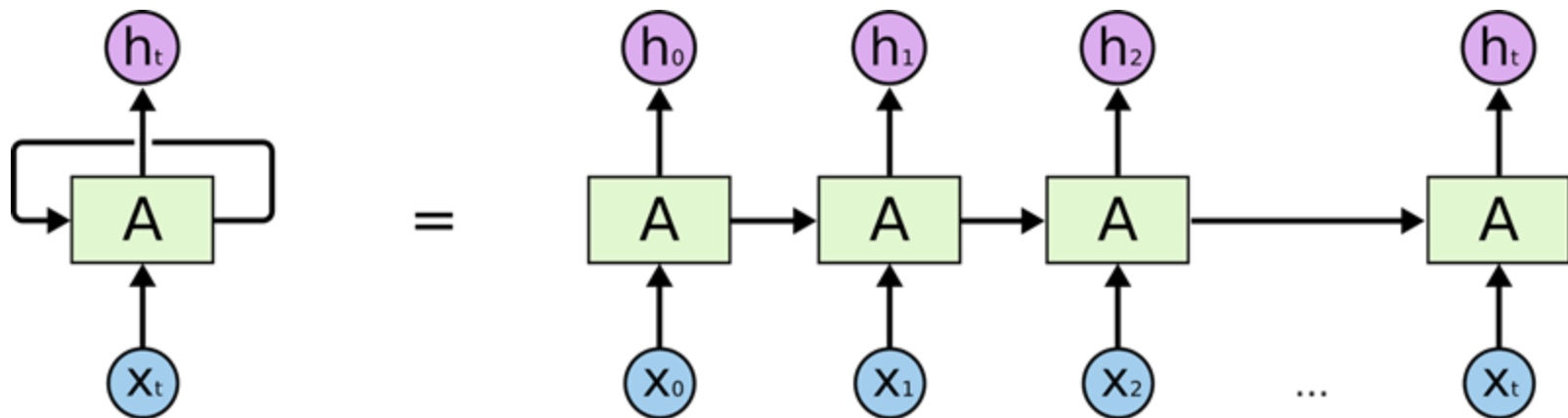Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# Sequence to Sequence with RNNs

**Input**: Sequence $x_1, \ldots x_T$
**Output**: Sequence $y_1, \ldots, y_{T'}$

**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 $c=h_T$)

estamos

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# Sequence to Sequence with RNNs

**Input**: Sequence $x_1, \ldots x_T$
**Output**: Sequence $y_1, \ldots, y_{T'}$

**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 $c = h_T$)



Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# Sequence to Sequence with RNNs

**Input**: Sequence $x_1, \ldots x_T$
**Output**: Sequence $y_1, \ldots, y_{T'}$

**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 $x_1, \ldots x_T$
**Output**: Sequence $y_1, \ldots, y_{T'}$

**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.**

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# Sequence to Sequence with RNNs

**Input**: Sequence $x_1, \ldots x_T$
**Output**: Sequence $y_1, \ldots, y_{T'}$

**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.**

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

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# Sequence to Sequence with RNNs and Attention
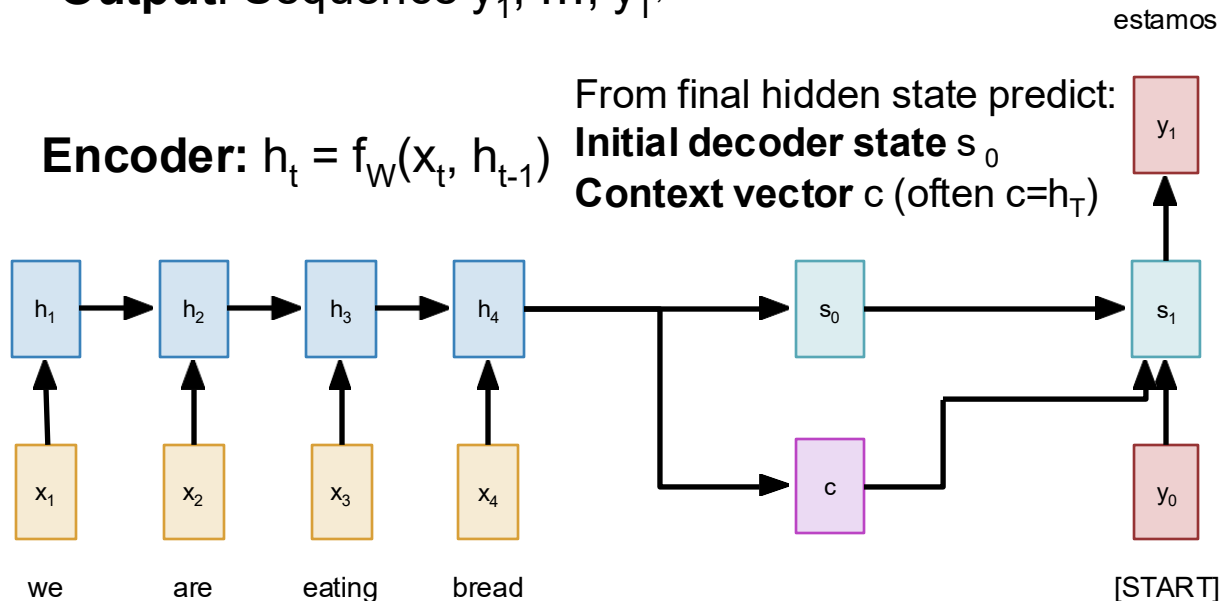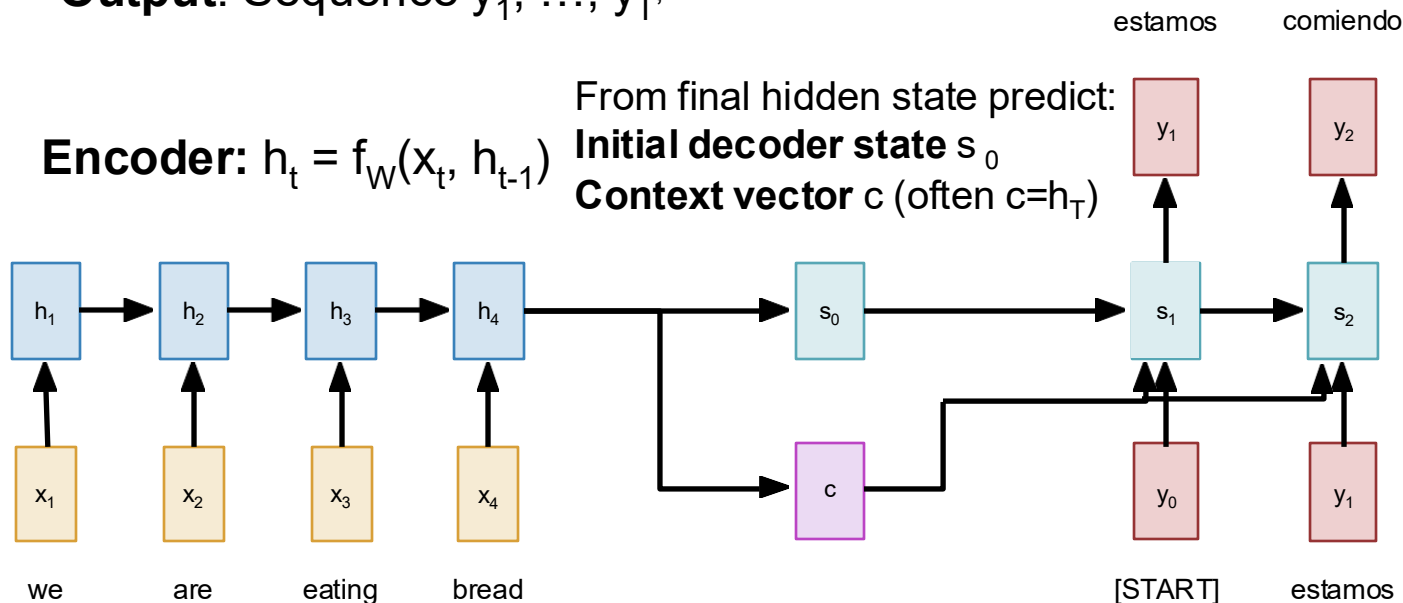
**Input**: Sequence $x_1, \ldots x_T$
**Output**: Sequence $y_1, \ldots, y_{T'}$

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

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

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# Sequence to Sequence with RNNs and Attention

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

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

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# Sequence to Sequence with RNNs and Attention

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

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

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

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# Sequence to Sequence with RNNs and Attention

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

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

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

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

softmax

estamos



we    are    eating    bread

[START]

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# Sequence to Sequence with RNNs and Attention



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

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

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

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

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

**Intuition**: Context vector attends to the relevant part of the input sequence
*"estamos"* = *"we are"*
so maybe $a_{11}=a_{12}=0.45$, $a_{13}=a_{14}=0.05$

Subhransu Maji, Chuang Gan and TAs
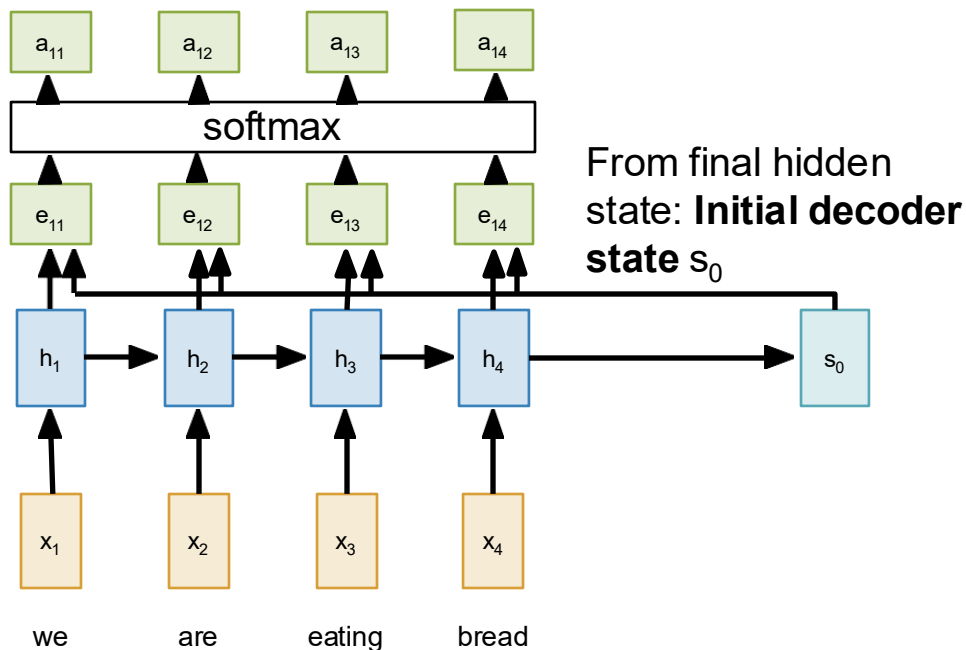Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# Sequence to Sequence with RNNs and Attention
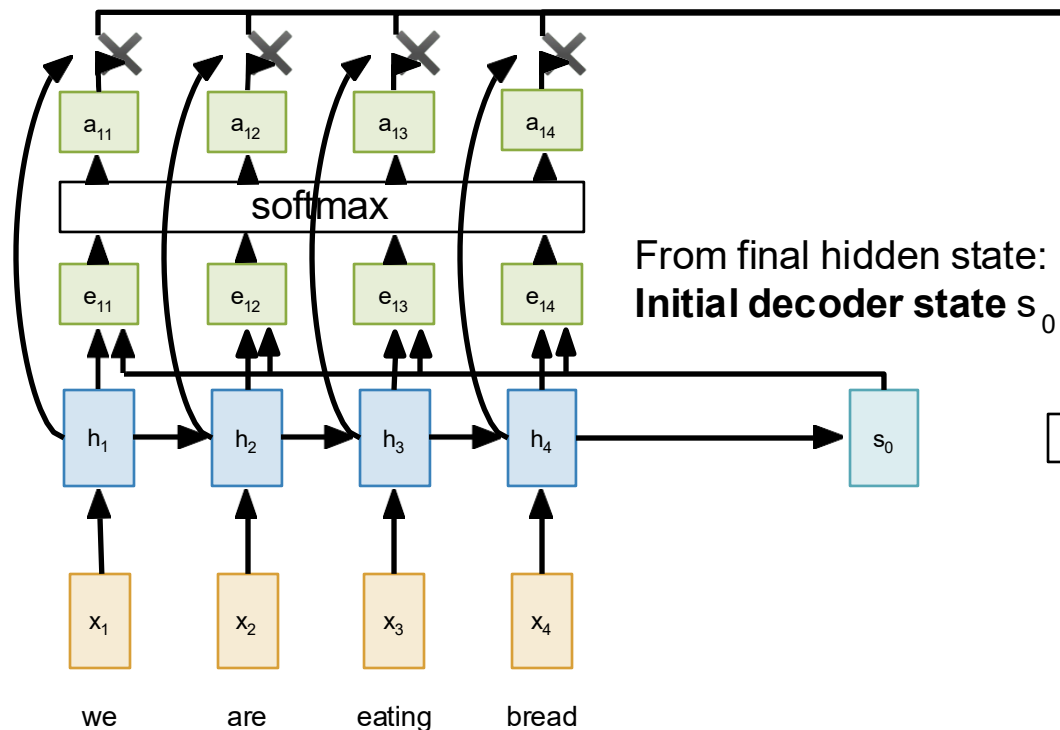


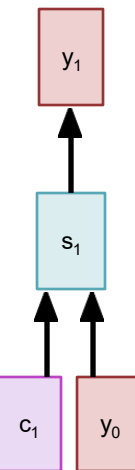Compute (scalar) **alignment scores**
$e_{t,i} = f_{att}(s_{t-1}, h_i)$    ($f_{att}$ is an MLP)

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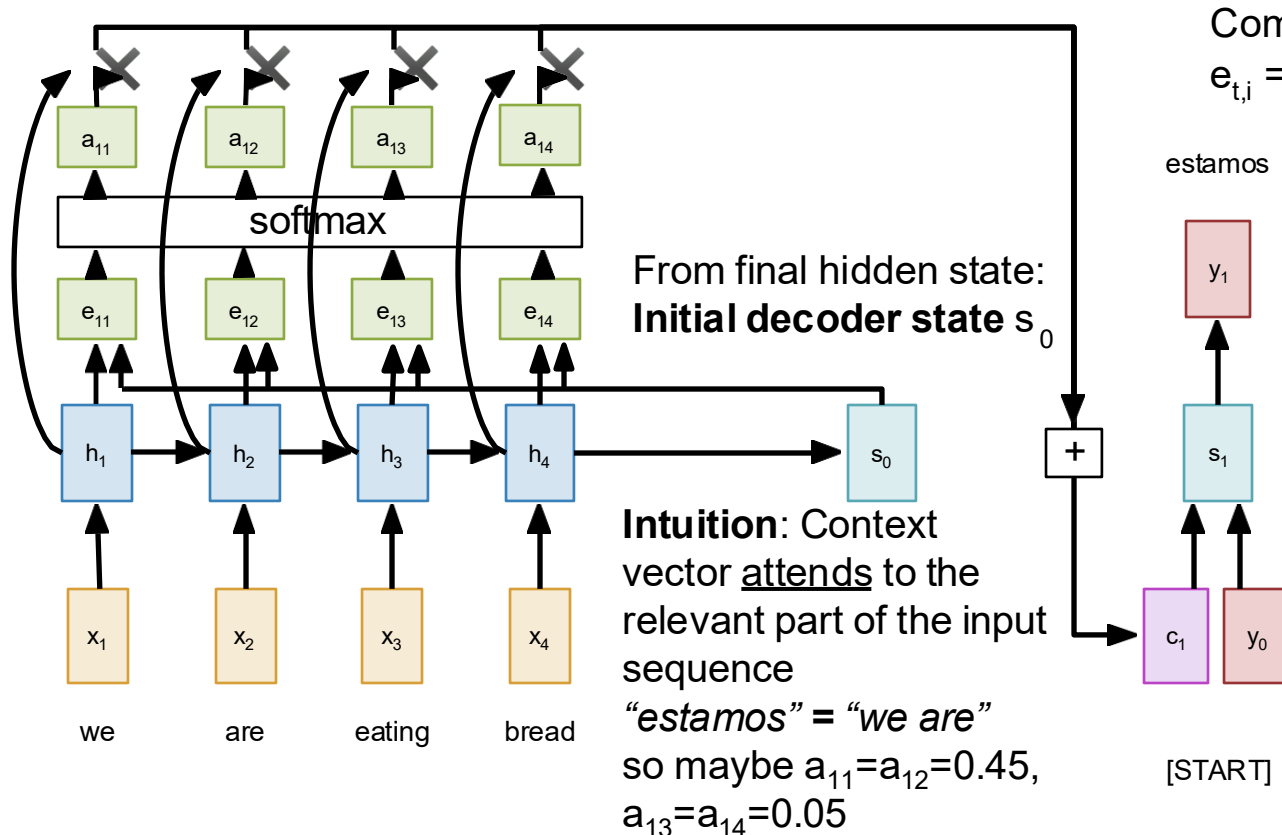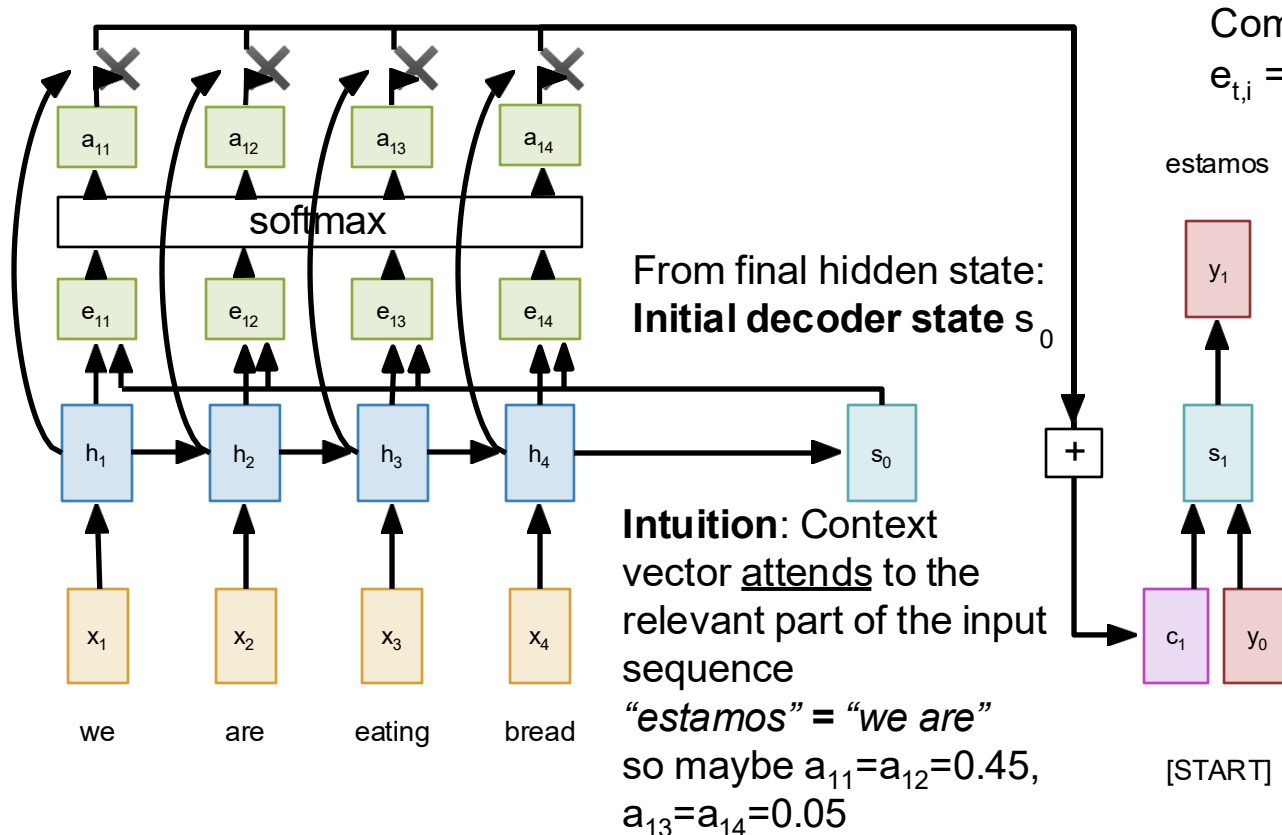$0 < a_{t,i} < 1$    $\sum_i a_{t,i} = 1$

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! No supervision on attention weights – backprop through everything**

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

**Intuition**: Context vector <u>attends</u> to the relevant part of the input sequence
*"estamos" = "we are"*
so maybe $a_{11}=a_{12}=0.45$, $a_{13}=a_{14}=0.05$

softmax

estamos

[START]

we    are    eating    bread

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# Sequence to Sequence with RNNs and Attention



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

estamos

$y_1$

softmax

$a_{21}$ $a_{22}$ $a_{23}$ $a_{24}$

$e_{21}$ $e_{22}$ $e_{23}$ $e_{24}$

$+$

$h_1$ $h_2$ $h_3$ $h_4$ $s_0$ $s_1$

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

$c_1$ $y_0$ $c_2$

we    are    eating    bread

[START]

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# Sequence to Sequence with RNNs and Attention

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

$a_{21}$  $a_{22}$  $a_{23}$  $a_{24}$

softmax

$e_{21}$  $e_{22}$  $e_{23}$  $e_{24}$

$h_1$ → $h_2$ → $h_3$ → $h_4$ → $s_0$

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

we  are  eating  bread

estamos  comiendo

$y_1$  $y_2$

$+$

Use $c_2$ to compute $s_2$, $y_2$

$s_1$  $s_2$

$c_1$  $y_0$  $c_2$  $y_1$

[START]  estamos

# Sequence to Sequence with RNNs and Attention

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



Use $c_2$ to compute $s_2$, $y_2$

**Intuition**: Context vector <u>attends</u> to the relevant part of the input sequence
*"comiendo" = "eating"*
so maybe $a_{21} = a_{24} = 0.05$, $a_{22} = 0.1$, $a_{23} = 0.8$

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# 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**
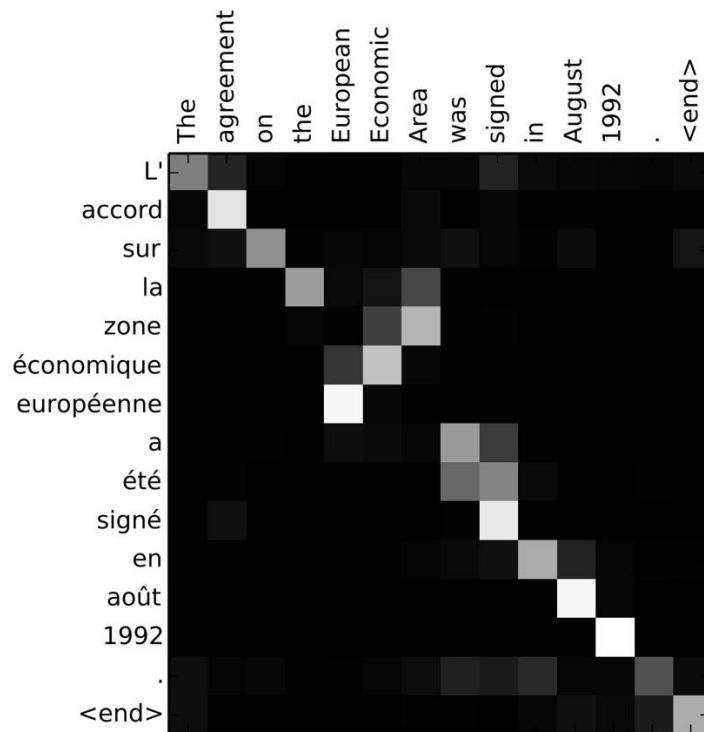
Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# 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

Visualize attention weights $a_{t,i}$
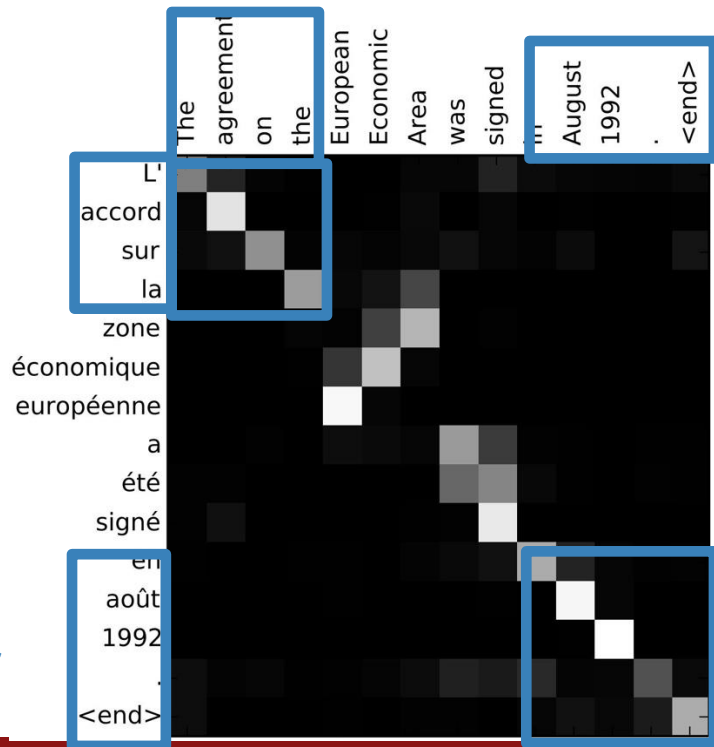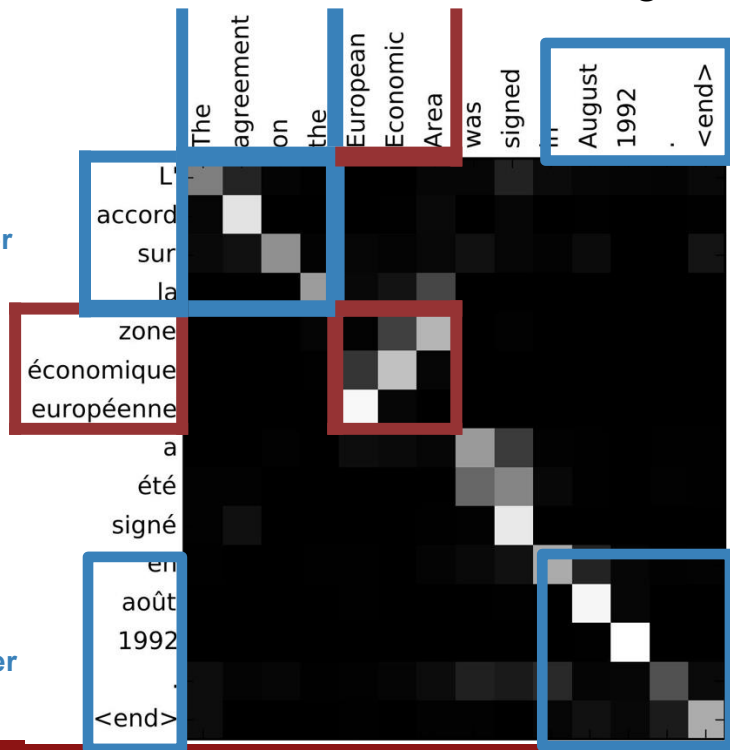
**Example**: English to French translation

**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**."



Diagonal attention means words correspond in order

Diagonal attention means words correspond in order

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# 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**."

**Diagonal attention means words correspond in order**

**Attention figures out different word orders**
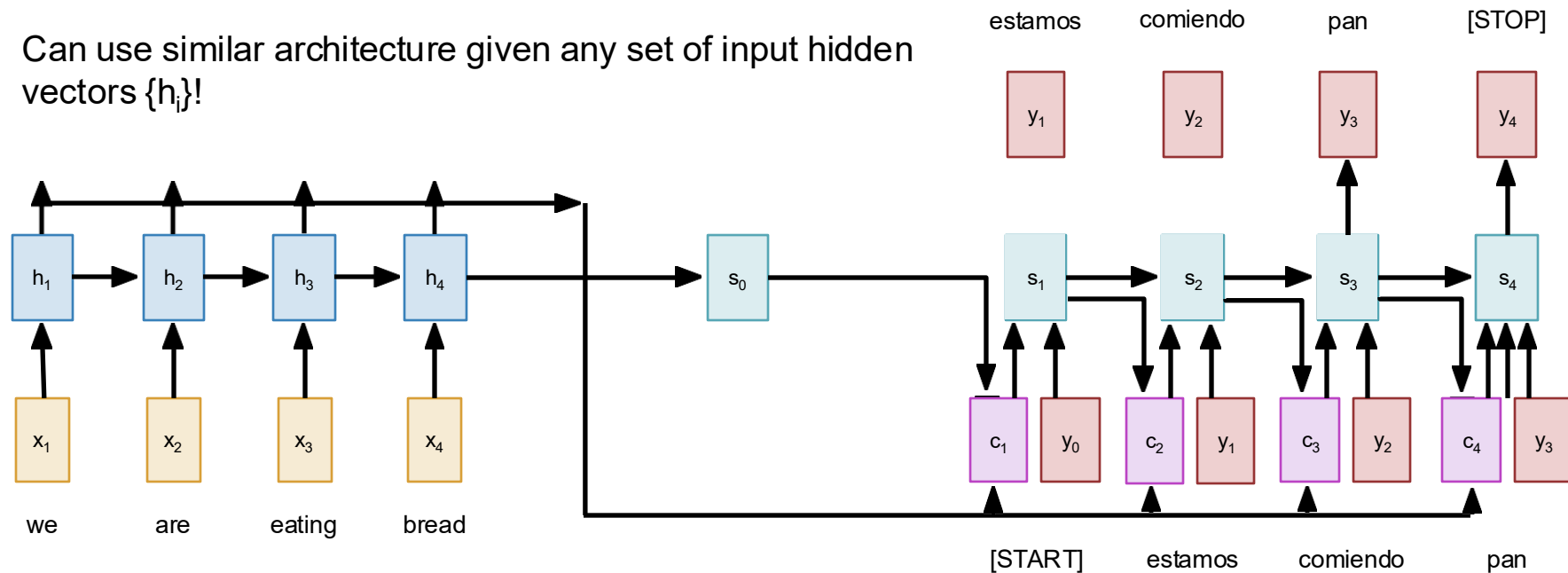
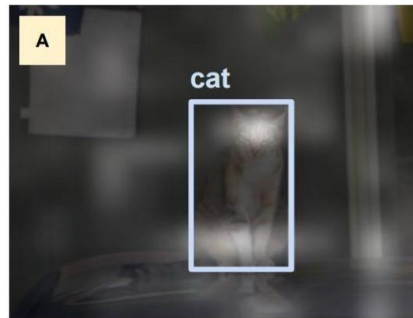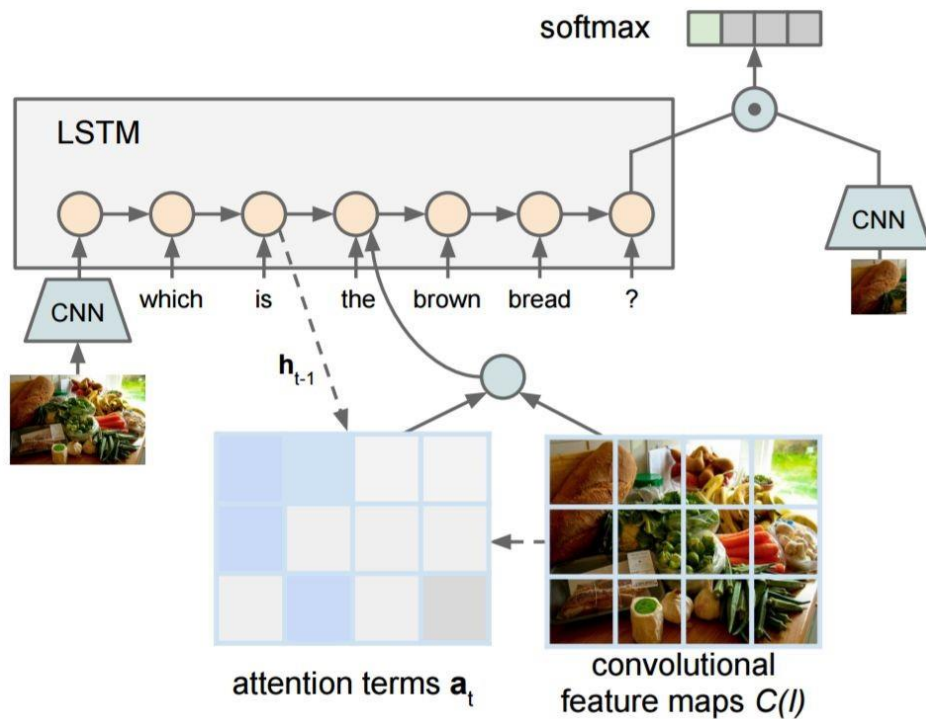**Diagonal attention means words correspond in order**

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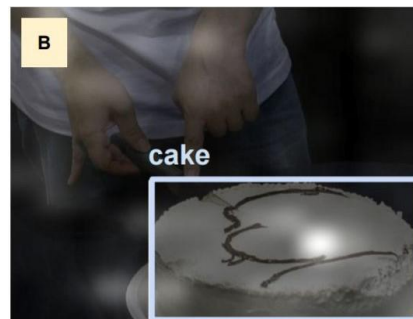
Can use similar architecture given any set of input hidden vectors $\{h_i\}$!

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller

# Visual Question Answering: RNNs with Attention



Zhu et al, "Visual 7W: Grounded Question Answering in Images", CVPR 2016
Figures from Zhu et al, copyright IEEE 2016. Reproduced for educational purposes.

Subhransu Maji, Chuang Gan and TAs
Some slides kindly provided by Fei-Fei Li, Jiajun Wu, Erik Learned-Miller