

Advanced Computer Vision

Tanveer Hussain

htanveer3797@gmail.com



جامعة الملك عبد الله
لعلوم والتكنولوجيا
King Abdullah University of
Science and Technology

KAUST Academy
King Abdullah University of Science and Technology

Course designed by **Naeem Ullah Khan** (naeemullah.khan@kaust.edu.sa), updated by Tanveer Hussain

Computer Vision Applications

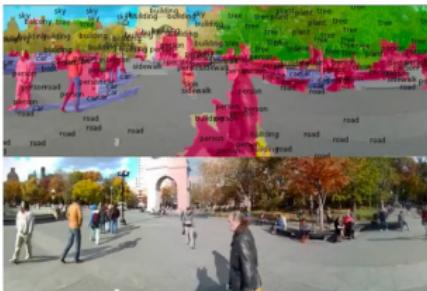
Image Classification



Image Retrieval



Image Segmentation



Fully Connected Neural Networks

Deep Neural Network

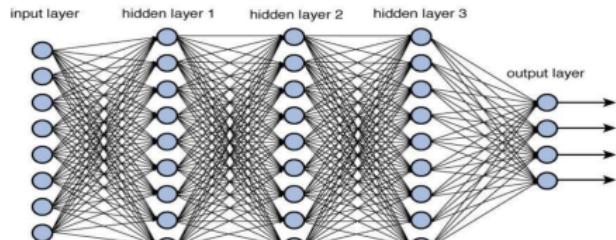
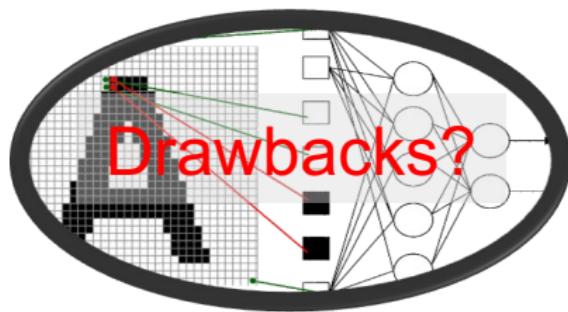
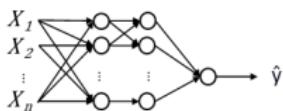
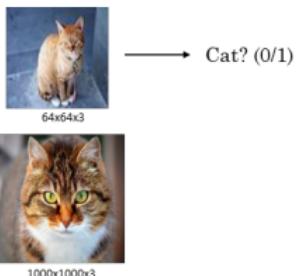
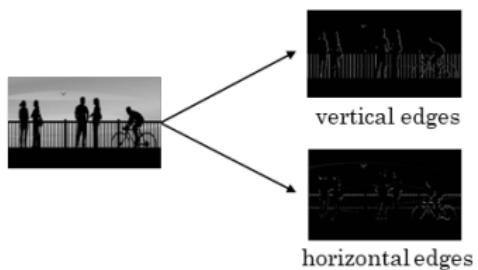


Figure 12.2 Deep network architecture with multiple layers.



Features Extraction Techniques



3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	3	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

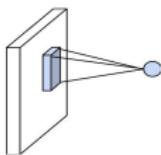
$$* \begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline & & \\ \hline \end{array}$$

$$\begin{array}{|c|c|c|c|c|c|c|} \hline 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ \hline 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ \hline 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ \hline 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ \hline 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ \hline 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ \hline \end{array} * \begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 0 & 30 & 30 & 0 \\ \hline 0 & 30 & 30 & 0 \\ \hline 0 & 30 & 30 & 0 \\ \hline 0 & 30 & 30 & 0 \\ \hline \end{array}$$

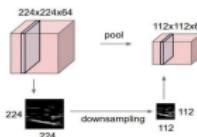
Previously

Components of CNNs

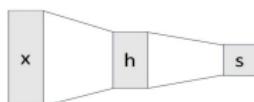
Convolution Layers



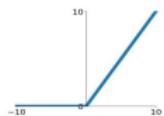
Pooling Layers



Fully-Connected Layers



Activation Function

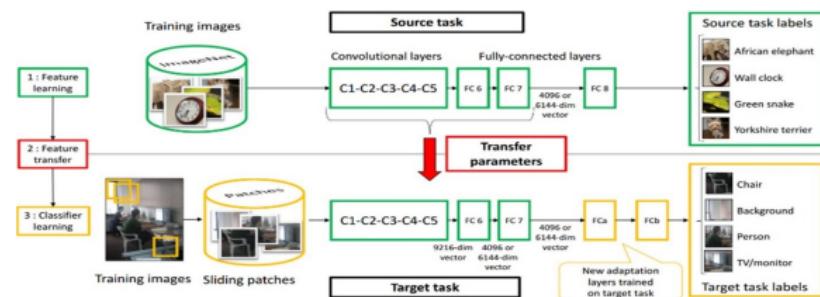
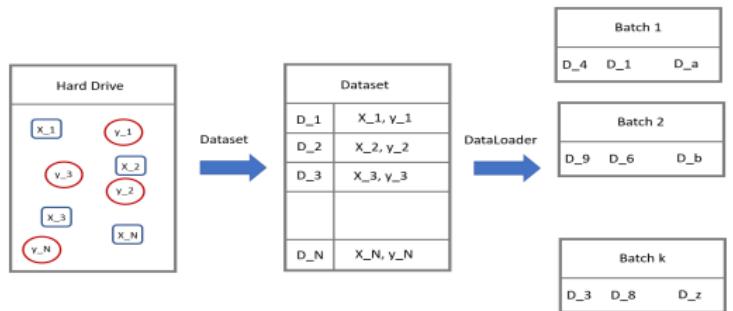


Normalization

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

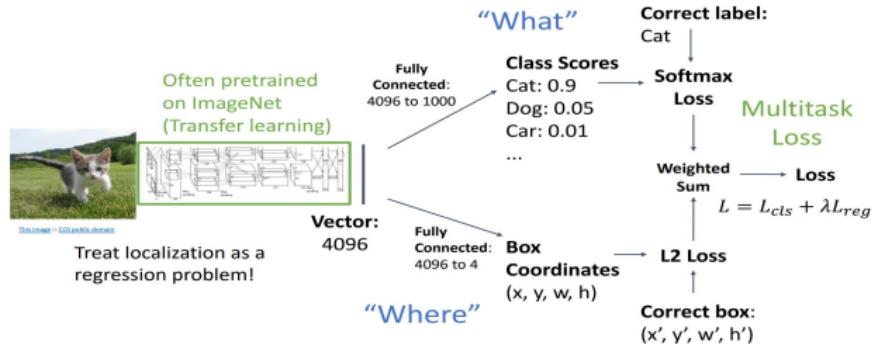
Previously

Data Loaders and Transfer Learning

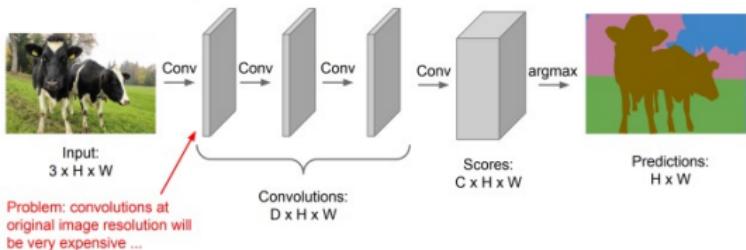


Previously

Objects Detection and Segmentation



Design a network with only convolutional layers without downsampling operators to make predictions for pixels all at once!



1. Recurrent Neural Networks and Self Attention
2. Attention and ViT
3. Video Classification
4. Generative AI (Part 1)
5. Generative AI (Part 2), CLIP model

Examples of Sequence Data

Music generation

\emptyset



Sentiment classification

"There is nothing to like
in this movie."



DNA sequence analysis

AGCCCCCTGTGAGGAACCTAG



AG**CCCCTGTGAGGAAC**TAG

Machine translation

Voulez-vous chanter avec
moi?



Do you want to sing with
me?

Video activity recognition



Running

Name entity recognition

Yesterday, Harry Potter
met Hermione Granger.



Yesterday, **Harry Potter**
met **Hermione Granger**.
Andrew Ng

Named Entity Recognition

x: Tom and Jerry went to the store to buy food.

y: 1 0 1 0

Vocabulary Indexes One-hot

A

Aa

...

...

...

And

...

Jerry

...

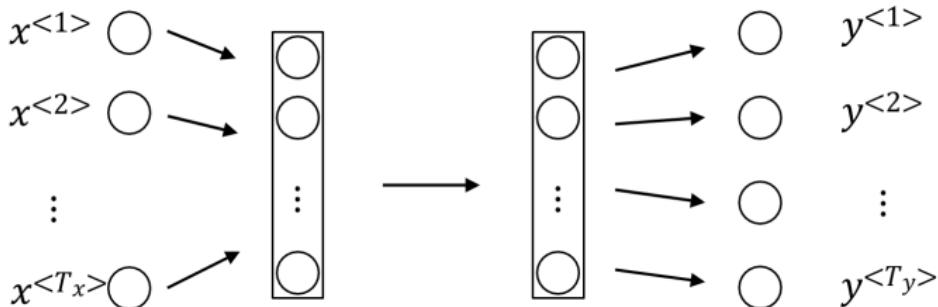
Tom

...

Unknown words?
End of sentence?

<unk>
<EOS>

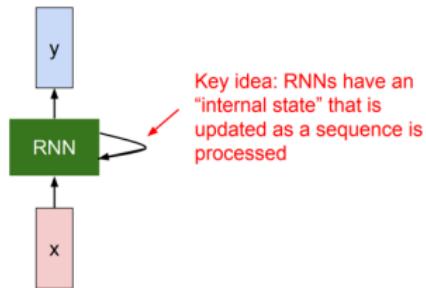
Why Sequence Models?



Problems?

- Inputs, outputs can be different lengths in different examples.
- Doesn't share features learned across different positions of text.

Rolled Diagram of RNNs



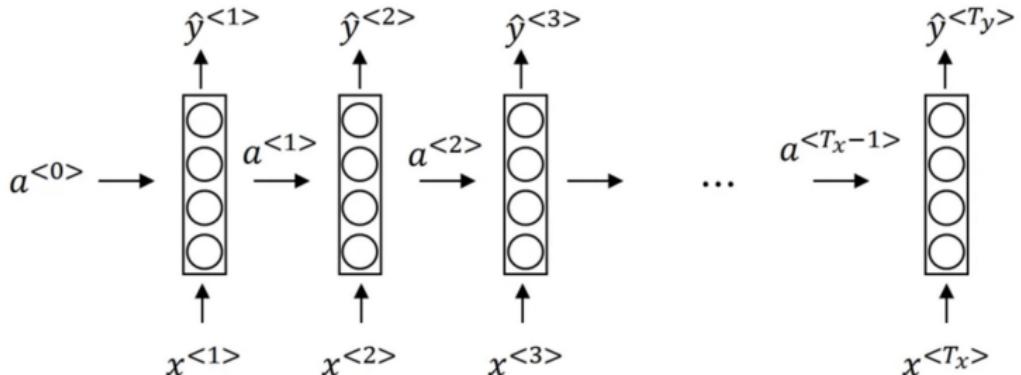
$$h_t = f_W(h_{t-1}, x_t)$$

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

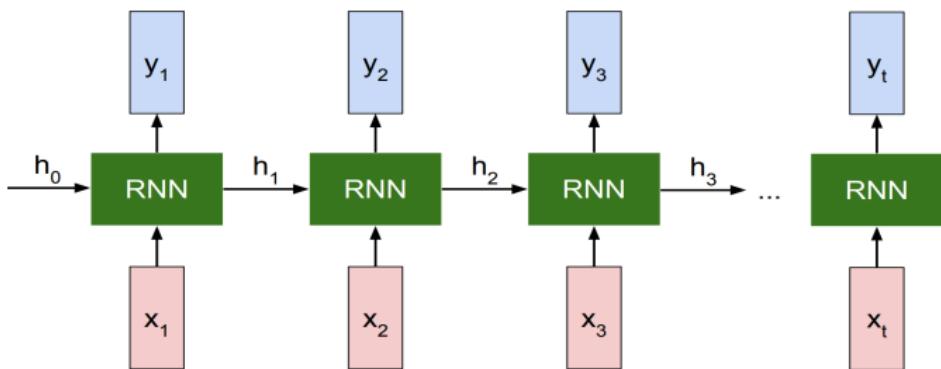
Recurrent Neural Networks (RNN)

Lets elaborate/unroll it.

Recurrent Neural Networks (RNN)

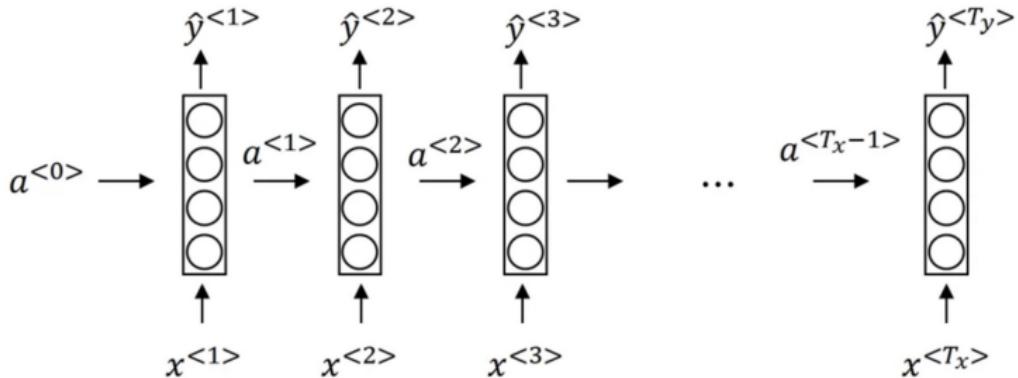


Another Representation

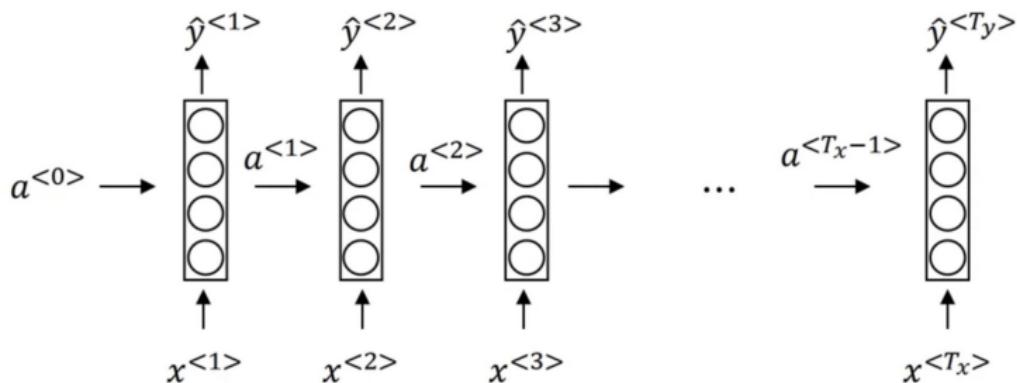


Recurrent Neural Networks (RNN)

Forward propagation



Backward propagation



Recurrent Neural Networks (RNN)

► RNN Advantages:

- Can process any length input
- Computation for step t can (in theory) use information from many steps back
- Model size doesn't increase for longer input
- Same weights applied on every timestep, so there is symmetry in how inputs are processed.

Recurrent Neural Networks (RNN)

► RNN Advantages:

- Can process any length input
- Computation for step t can (in theory) use information from many steps back
- Model size doesn't increase for longer input
- Same weights applied on every timestep, so there is symmetry in how inputs are processed.

► RNN Disadvantages:

- Recurrent computation is slow
- In practice, difficult to access information from many steps back

Types of RNNs

Many-to-many ($T_x = T_y$)

Many-to-one

Types of RNNs

One-to-many

Types of RNNs

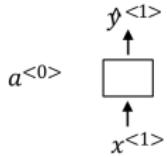
One-to-one

Types of RNNs

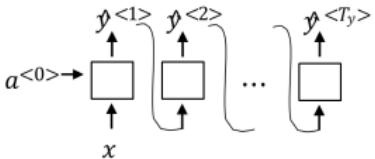
Many-to-many
See the whole input first.

Types of RNNs

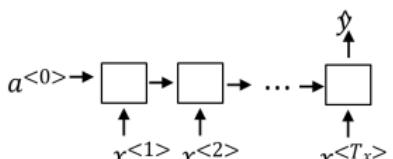
-1/**



One to one



One to many

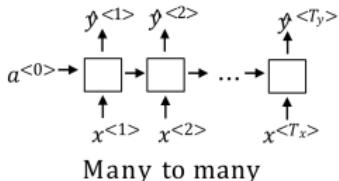


Many to one

Standard generic NN.

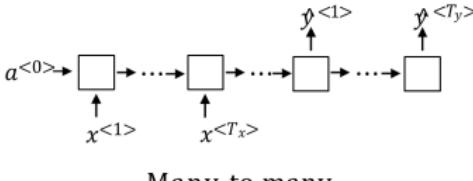
Music notes generation

The food is worst in this restaurant.
Sentiment classification



Many to many

Named entity recognition



Many to many

First: Read the input, then generate output.
Machine translation

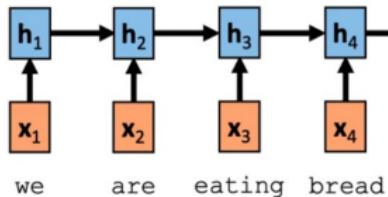
Sequence to Sequence with RNNs

Seq2seq: Many-to-one + one-to-many

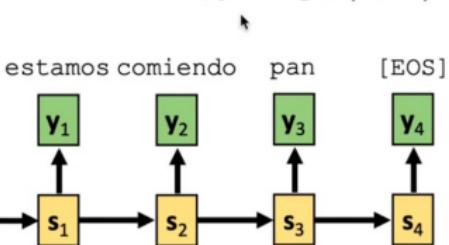
Input: Sequence $\mathbf{x}_1, \dots, \mathbf{x}_T$

Output: Sequence $\mathbf{y}_1, \dots, \mathbf{y}_{T'}$

Encoder: $\mathbf{h}_t = f_{\theta}(\mathbf{h}_{t-1}, \mathbf{x}_t)$



Decoder: $\mathbf{s}_t, \mathbf{y}_t = g_{\theta'}(\mathbf{s}_{t-1})$



Machine translation: English to French

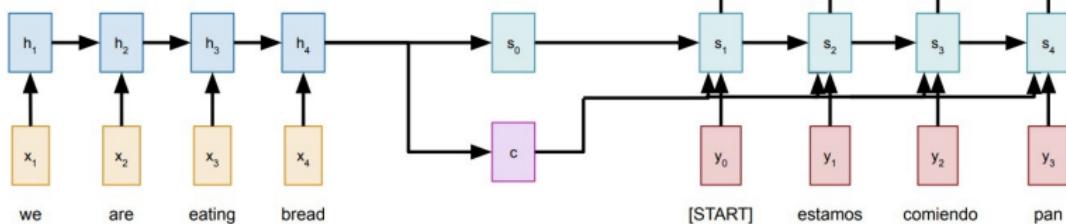
Sequence to Sequence with RNNs

Seq2seq: Many-to-one + one-to-many

Input: Sequence x_1, \dots, x_T
Output: Sequence y_1, \dots, y_T

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

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



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

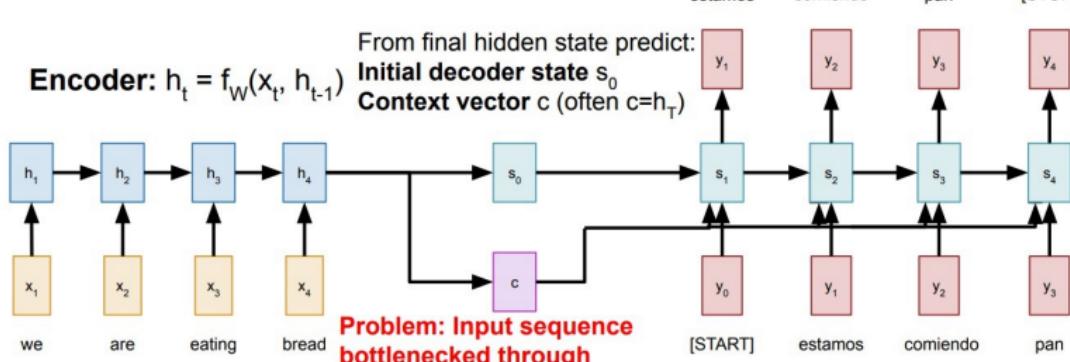
Sequence to Sequence with RNNs

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 $c=h_T$)



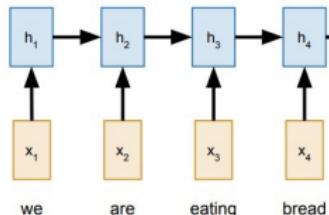
Problem: Input sequence bottlenecked through fixed-sized vector. What if $T=1000$?

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

Sequence to Sequence with RNNs

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

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

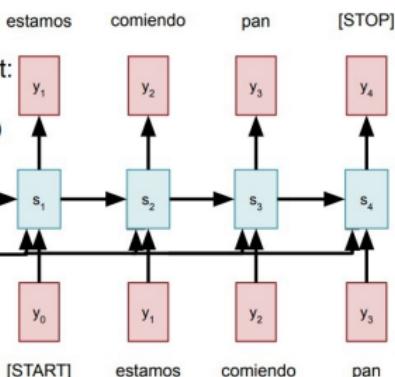


From final hidden state predict:
Initial decoder state s_0
Context vector c (often $c=h_T$)

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

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

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



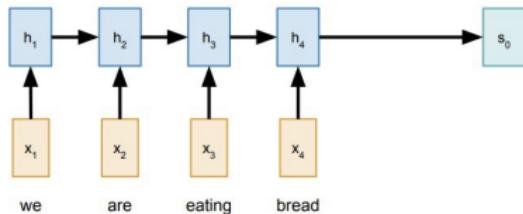
Idea: use new context vector at each step of decoder!

Sequence to Sequence with RNNs and Attention

Input: Sequence x_1, \dots, x_T

Output: Sequence y_1, \dots, y_T

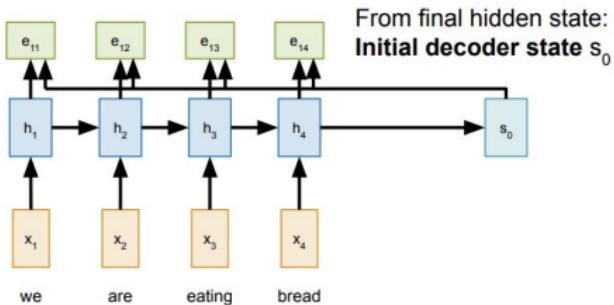
Encoder: $h_t = f_W(x_t, h_{t-1})$ From final hidden state:
Initial decoder state s_0



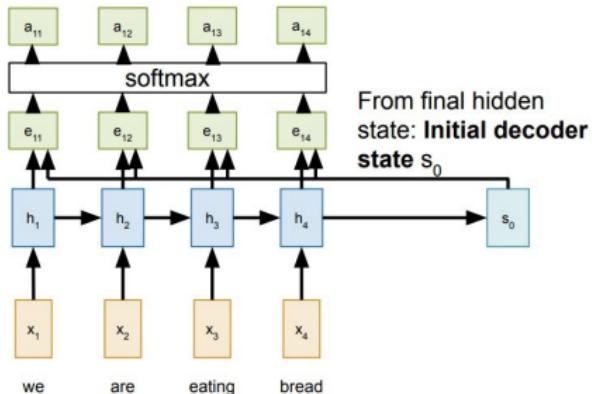
Sequence to Sequence with RNNs and Attention

How important is h_i [from encoder] for predicting output at S_t ?

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



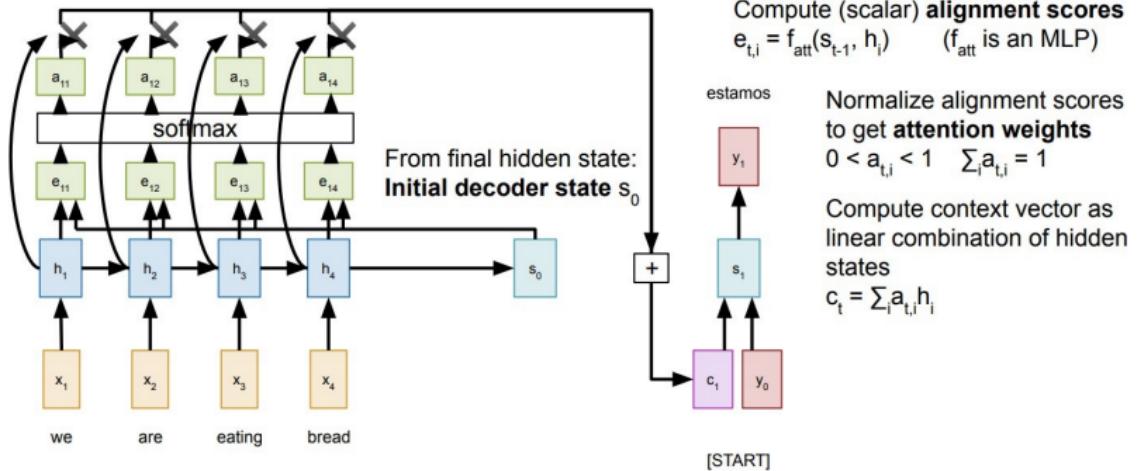
Sequence to Sequence with RNNs and Attention



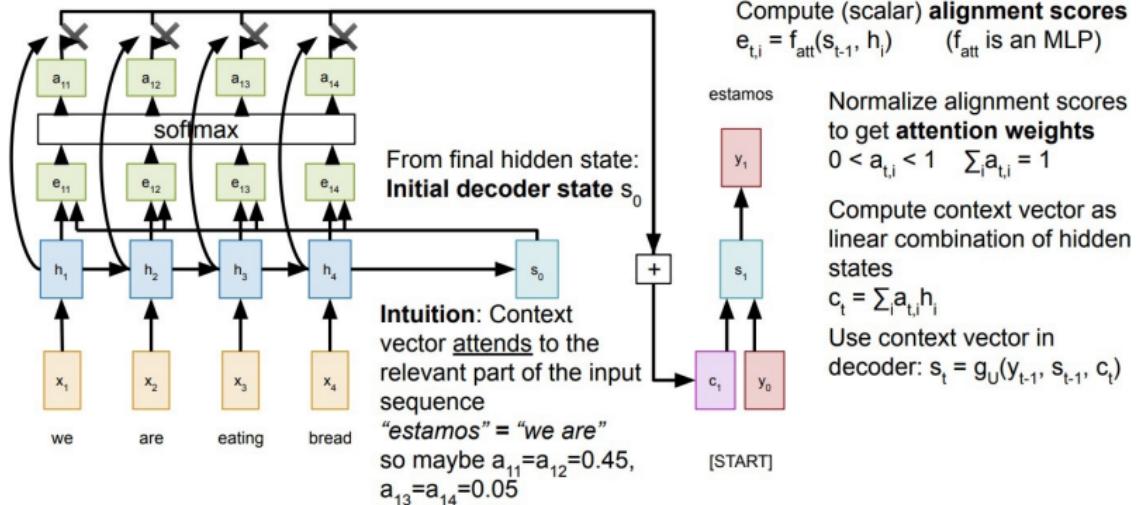
Compute (scalar) **alignment scores**
 $e_{t,i} = f_{att}(s_{t-1}, h_i)$ (f_{att} is an MLP)

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

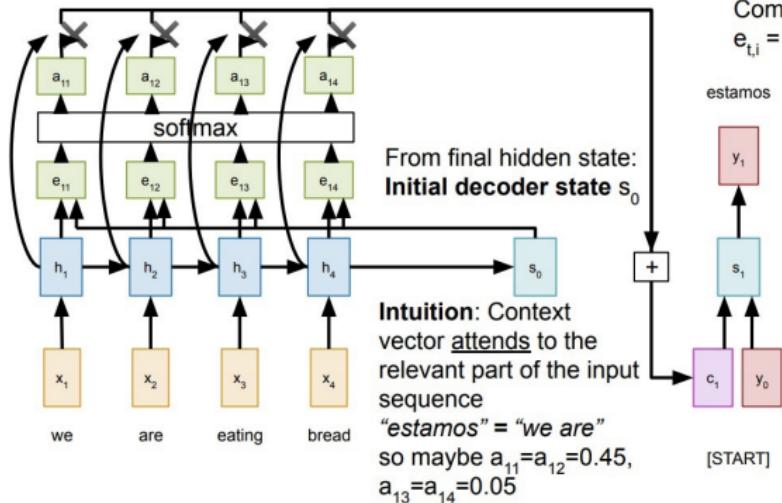
Sequence to Sequence with RNNs and Attention



Sequence to Sequence with RNNs and Attention



Sequence to Sequence with RNNs and Attention



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

estamos

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

Compute context vector as
linear combination of hidden
states

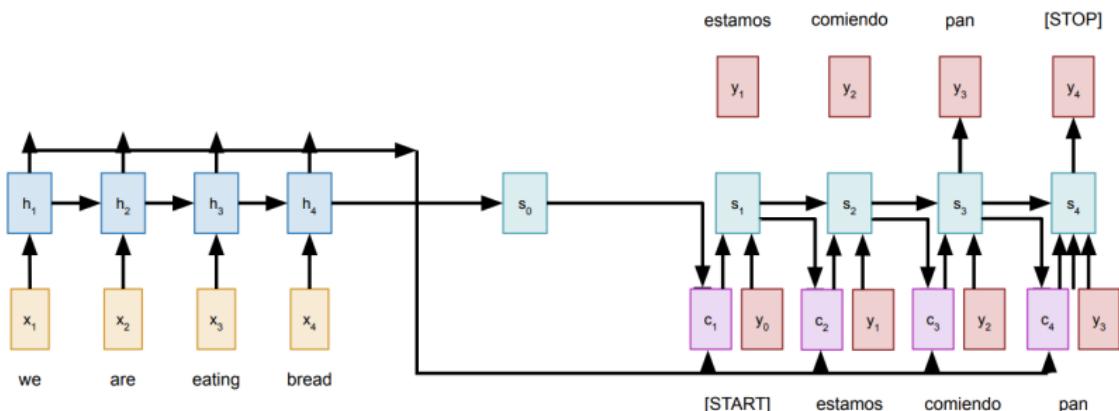
$$c_t = \sum_i a_{t,i} h_i$$

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

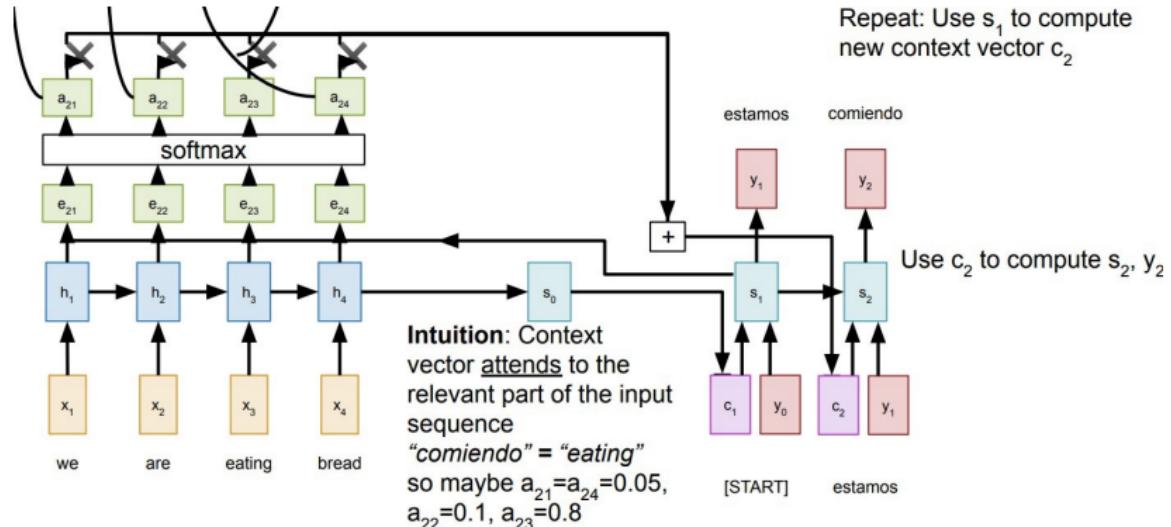
This is all differentiable! No
supervision on attention
weights – backprop through
everything

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

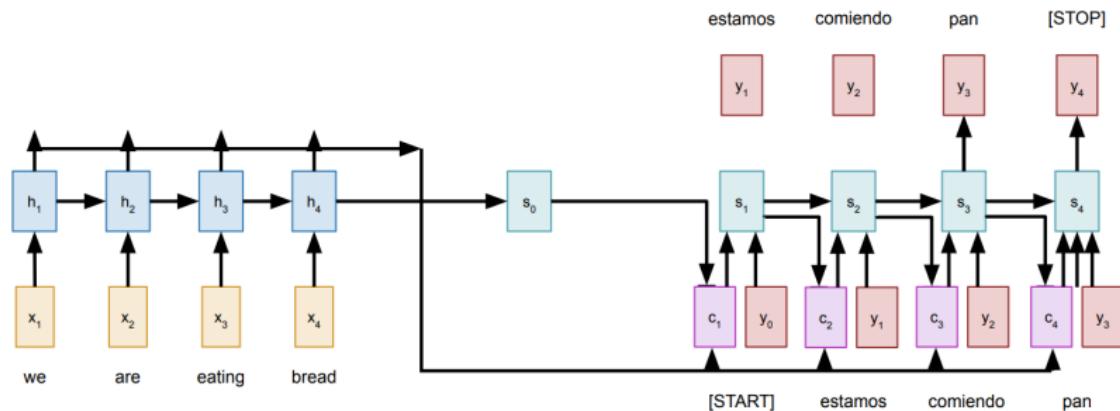


Sequence to Sequence with RNNs and Attention

► What is the cost of generating attention?

► Quadratic cost, at acceptable level.

Research is going to reduce the cost.



► Can we use this architecture for any type of hidden vectors?

Image Captioning

- ▶ A computer Vision task in which we generate captions for images



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

Sequence to Sequence with RNNs and Attention

Similar task in Computer Vision: Image Captioning Example

Input: Image I

Output: Sequence $y = y_1, y_2, \dots, y_T$

Decoder: $y_t = g_v(y_{t-1}, h_{t-1}, c)$

where context vector c is often $c = h_0$

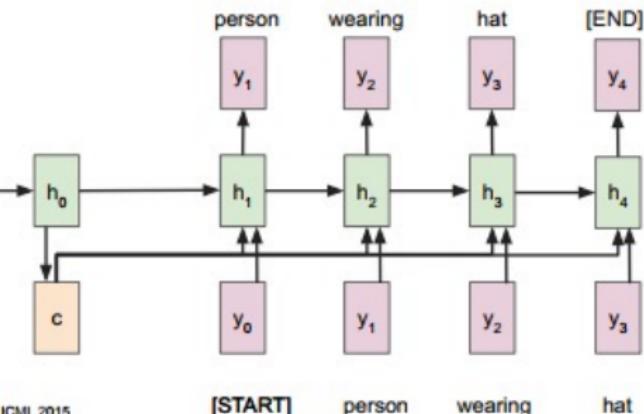
Encoder: $h_0 = f_w(z)$

where z is spatial CNN features

$f_w(\cdot)$ is an MLP



MLP



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

Problem: Input is "bottlenecked" through c

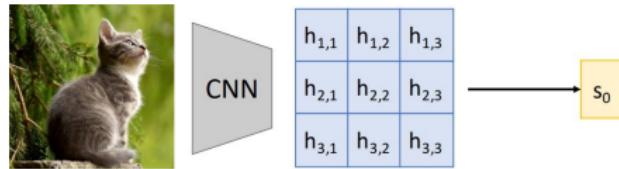
- Model needs to encode everything it wants to say within c

This is a problem if we want to generate really long descriptions? 100s of words long

Context vector at each time step

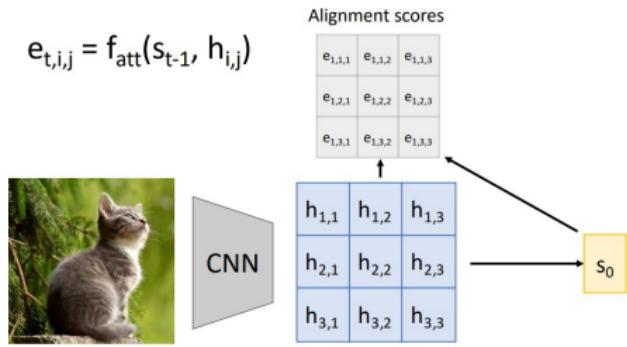
Each context vector attending different image regions

Image Captioning with RNNs and Attention



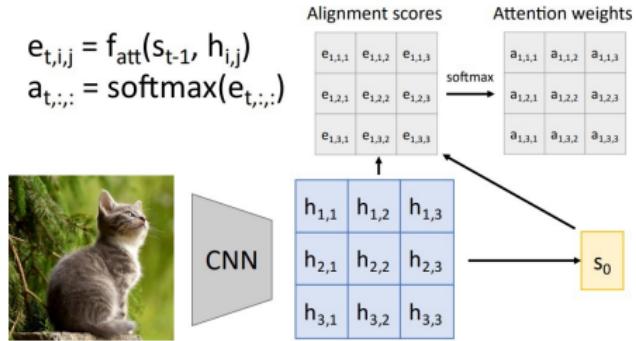
Use a CNN to compute a
grid of features for an image

Image Captioning with RNNs and Attention



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

Image Captioning with RNNs and Attention



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

Image Captioning with RNNs and Attention

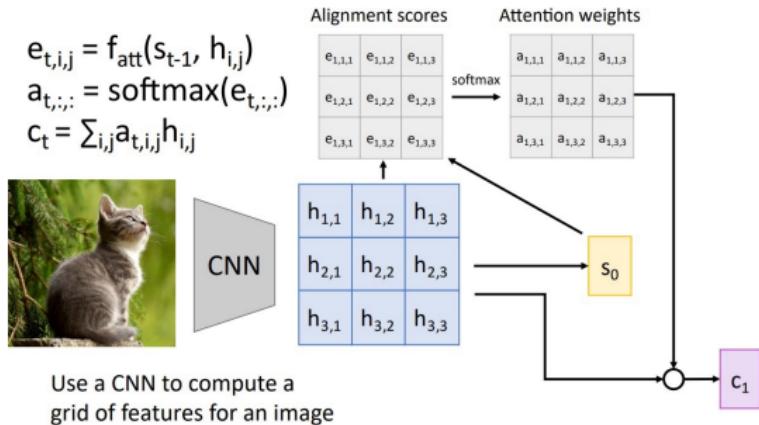


Image Captioning with RNNs and Attention

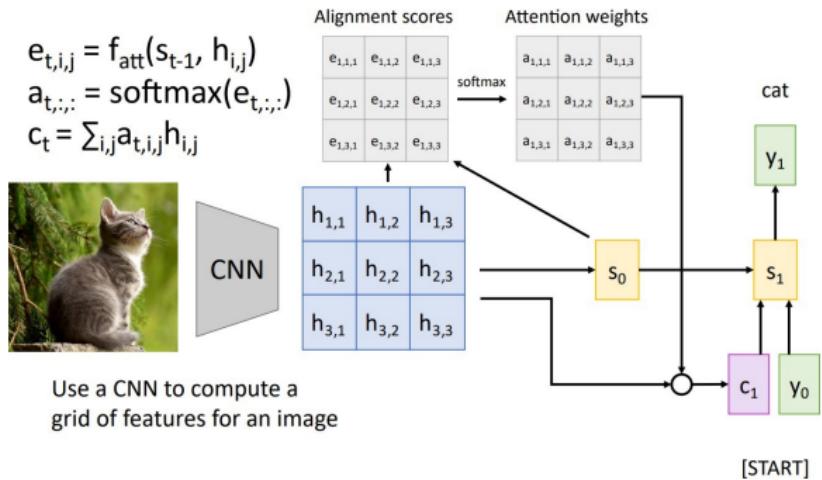


Image Captioning with RNNs and Attention

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



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

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

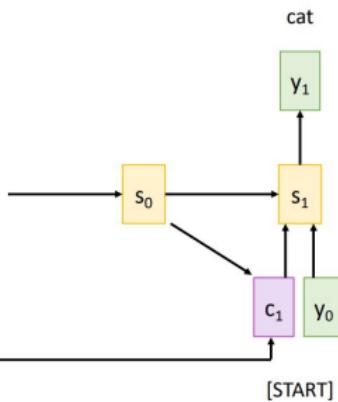


Image Captioning with RNNs and Attention

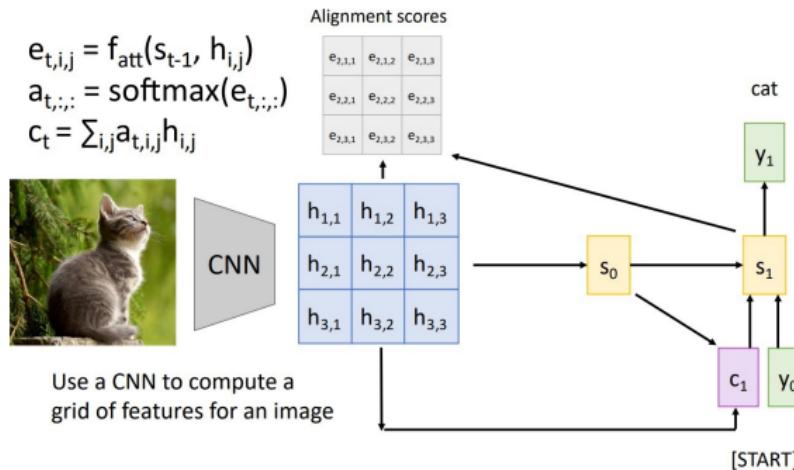


Image Captioning with RNNs and Attention

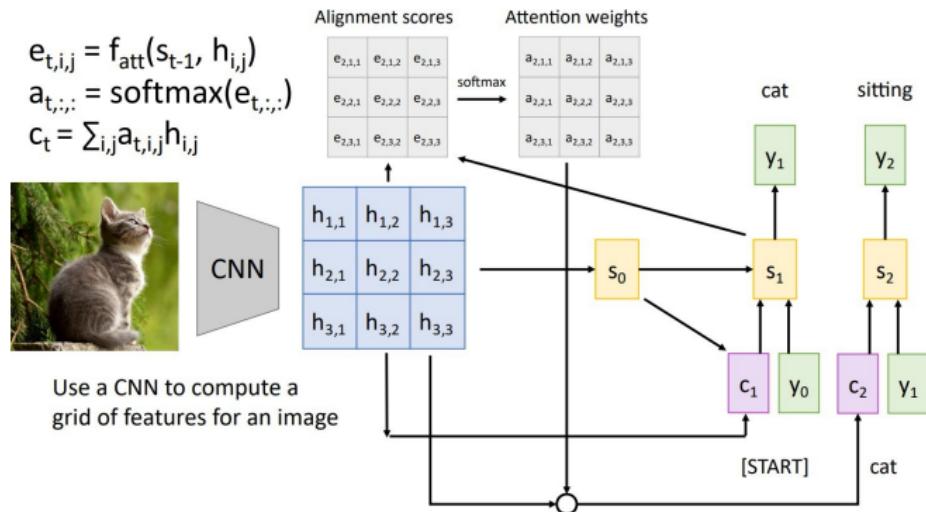


Image Captioning with RNNs and Attention

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

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



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

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

$$s_0$$

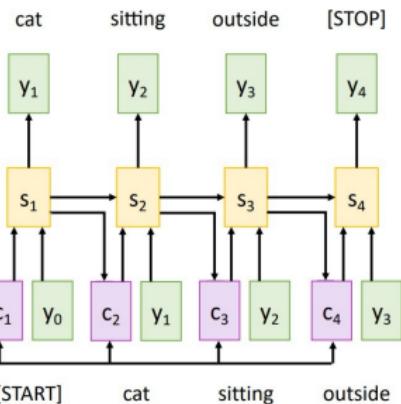


Image Captioning with RNNs and Attention



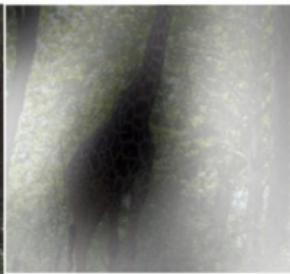
A dog is standing on a hardwood floor.



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



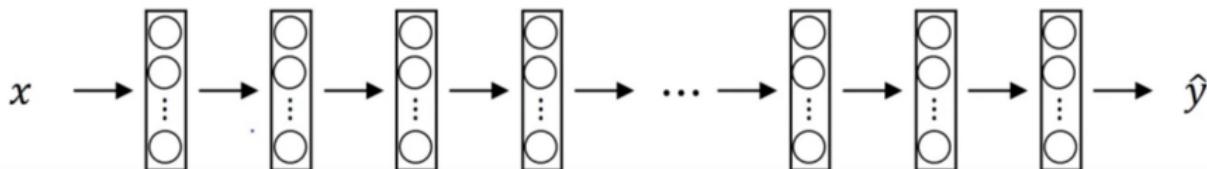
A group of people sitting on a boat in the water.



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

⁰Xu et al, "Show, Ask, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Long Sequences, Hard to Remember



Forward Prop: Left to right

Backward Prop: Right to left

Gradients from output > Hard times propagating back to effect the weights in the earlier layers.

Gated Recurrent Neural Networks

1. Recurrent Neural Networks and Self Attention
2. Attention and ViT
3. Video Classification
4. Generative AI (Part 1)
5. Generative AI (Part 2), CLIP model

These slides have been adapted from

- ▶ Fei-Fei Li, Yunzhu Li & Ruohan Gao, Stanford CS231n: [Deep Learning for Computer Vision](#)
- ▶ Assaf Shocher, Shai Bagon, Meirav Galun & Tali Dekel, WAIC DL4CV [Deep Learning for Computer Vision: Fundamentals and Applications](#)
- ▶ Justin Johnson, UMich EECS 498.008/598.008: [Deep Learning for Computer Vision](#)
- ▶ Andrew Ng, [Sequence Models](#)