

# DEEP LEARNING FOR COMPUTER VISION

Summer Seminar UPC TelecomBCN, 4 - 8 July 2016



## Instructors



Xavier  
Giró-i-Nieto



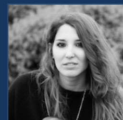
Elisa  
Sayrol



Amaia  
Salvador



Jordi  
Torres



Eva  
Mohedano



Kevin  
McGuinness

## Organizers



UNIVERSITAT POLITÈCNICA  
DE CATALUNYA  
BARCELONATECH



**Barcelona  
Supercomputing  
Center**  
Centro Nacional de Supercomputación



Dublin City University  
Oileici Chathair Bhaile Átha Cliath



Centre for Data Analytics



GPU  
CENTER OF  
EXCELLENCE

Co-funded by the  
Erasmus+ Programme  
of the European Union



+ info: [TelecomBCN.DeepLearning.Barcelona](https://TelecomBCN.DeepLearning.Barcelona)

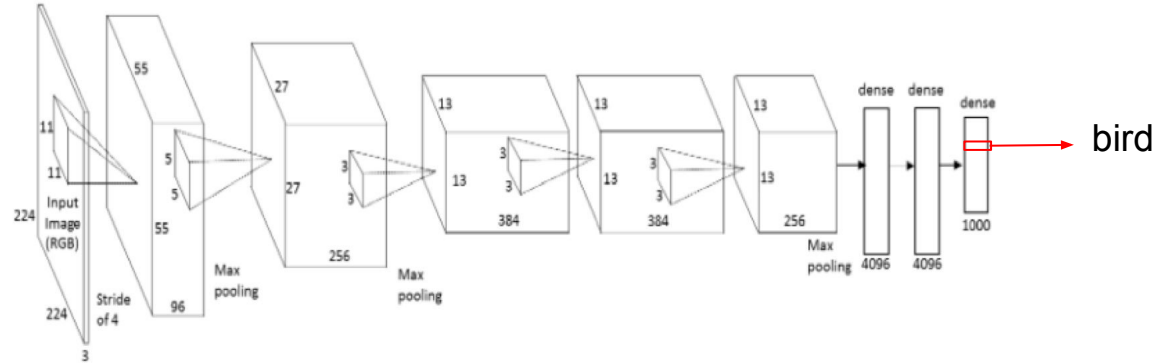
Day 4 Lecture 6

# Attention Models

# Attention Models: Motivation



Image:  
H x W x 3

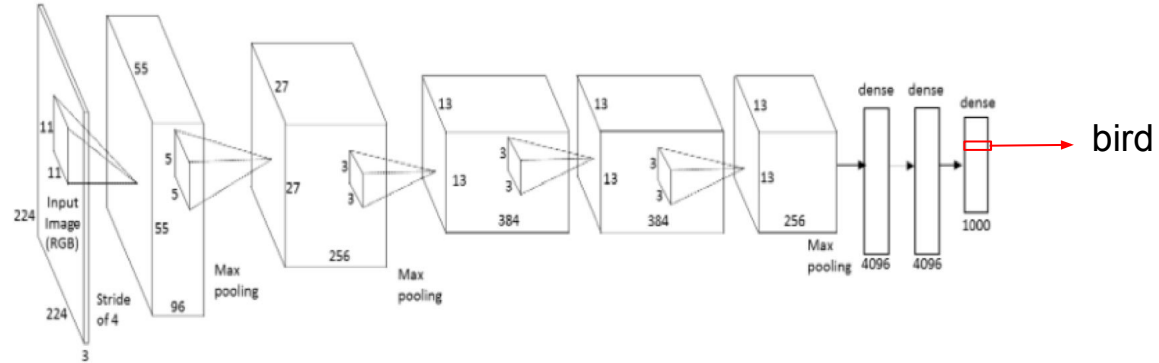


The whole input volume is used to predict the output...

# Attention Models: Motivation



Image:  
H x W x 3



The whole input volume is used to predict the output...

...despite the fact that not all pixels are equally important

# Attention Models: Motivation



Attention models can  
relieve computational burden

Helpful when processing big  
images !

# Attention Models: Motivation



Attention models can  
relieve computational burden

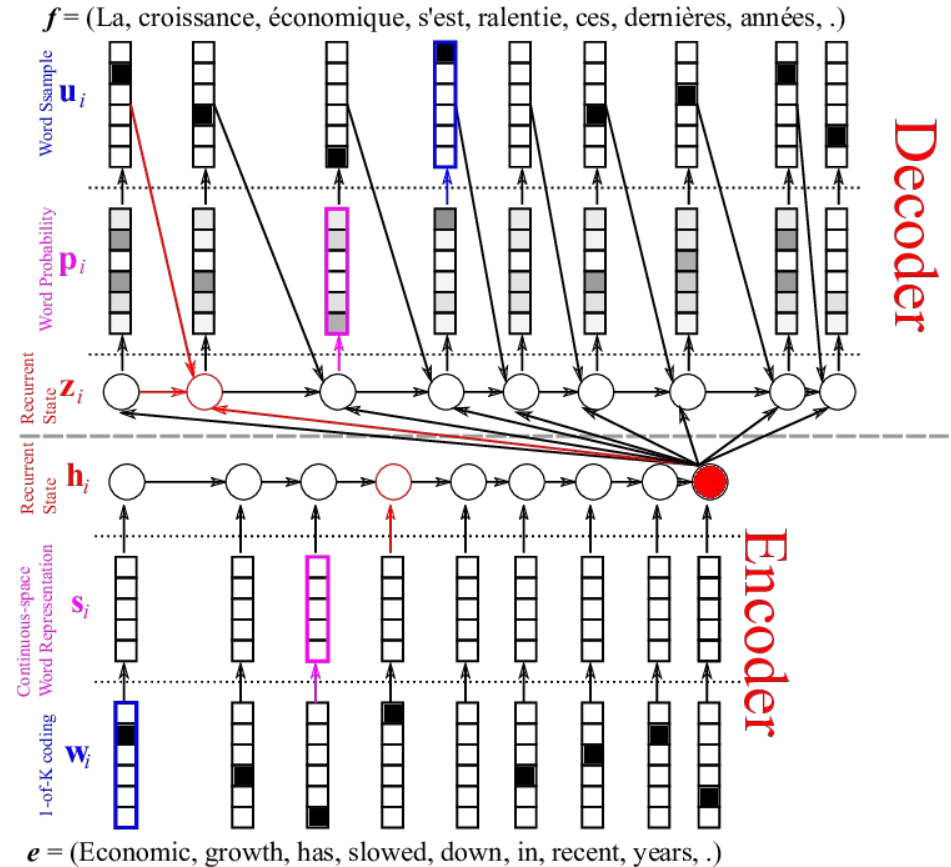
Helpful when processing big  
images !

# Encoder & Decoder

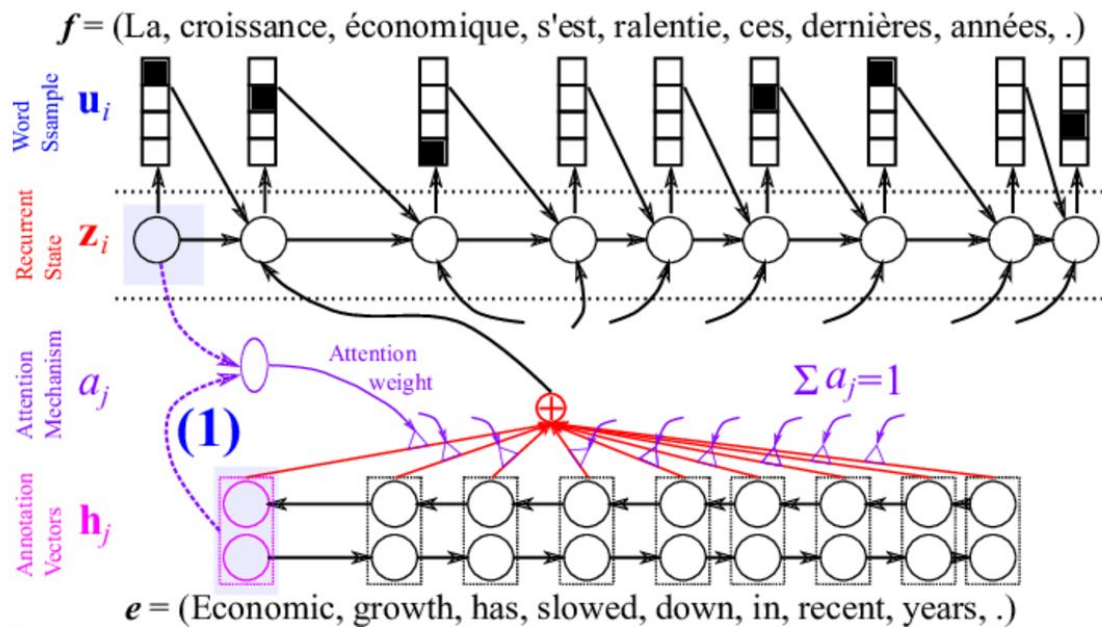
From previous lecture...

The whole input sentence  
is used to produce the translation

Kyunghyun Cho, [“Introduction to Neural Machine Translation with GPUs”](#) (2015)



# Attention Models



Bahdanau et al. [Neural Machine Translation by Jointly Learning to Align and Translate](#). ICLR 2015  
Kyunghyun Cho, [“Introduction to Neural Machine Translation with GPUs”](#) (2015)



# Attention Models

Idea: Focus in different parts of the input as you make/refine predictions in time

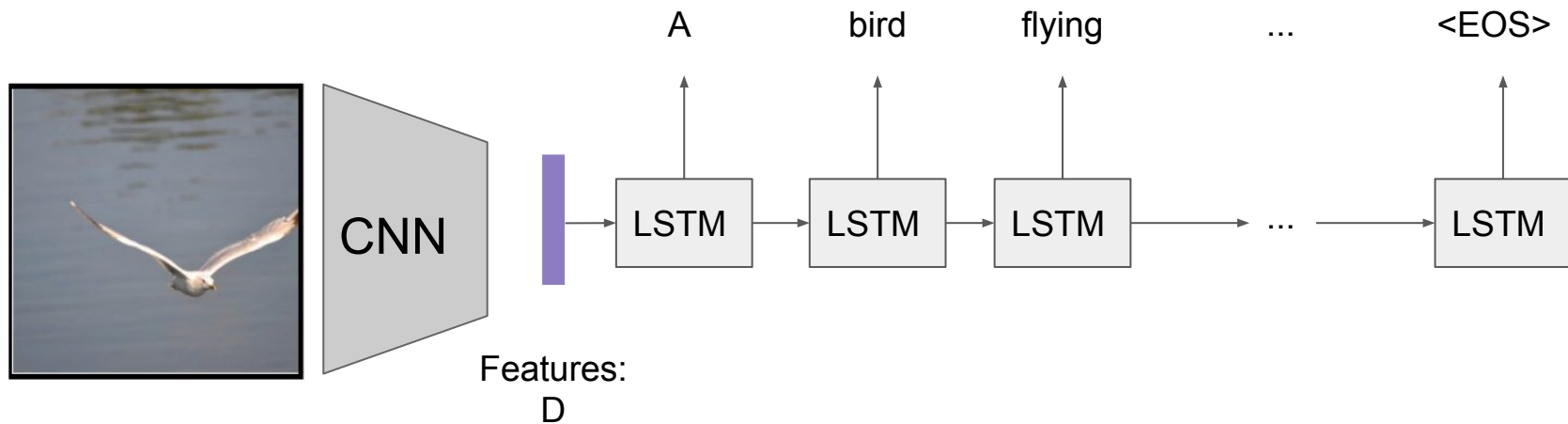
E.g.: Image Captioning



A bird flying over a body of water



# LSTM Decoder



The LSTM decoder “sees” the input only at the beginning !

# Attention for Image Captioning

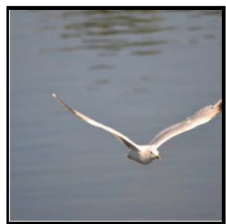
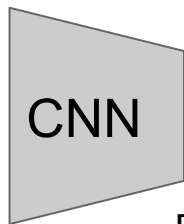
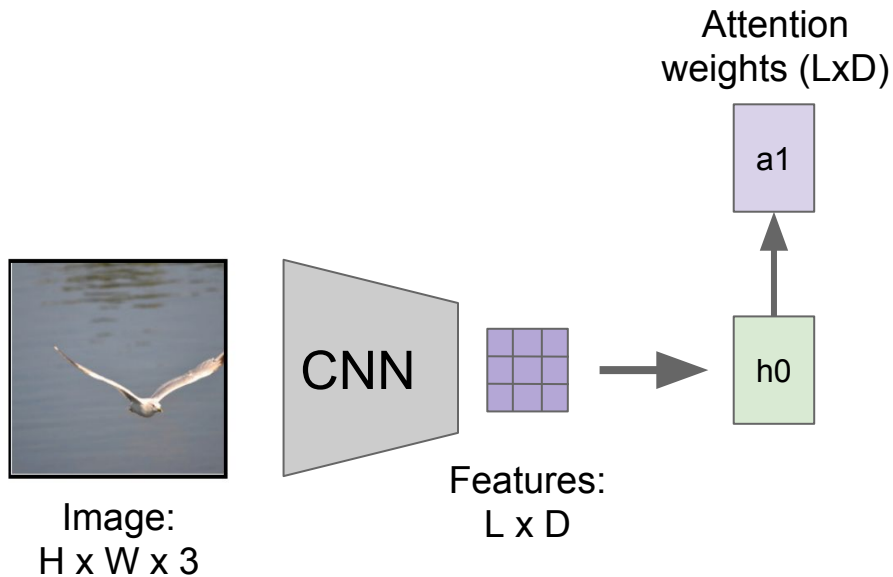


Image:  
 $H \times W \times 3$

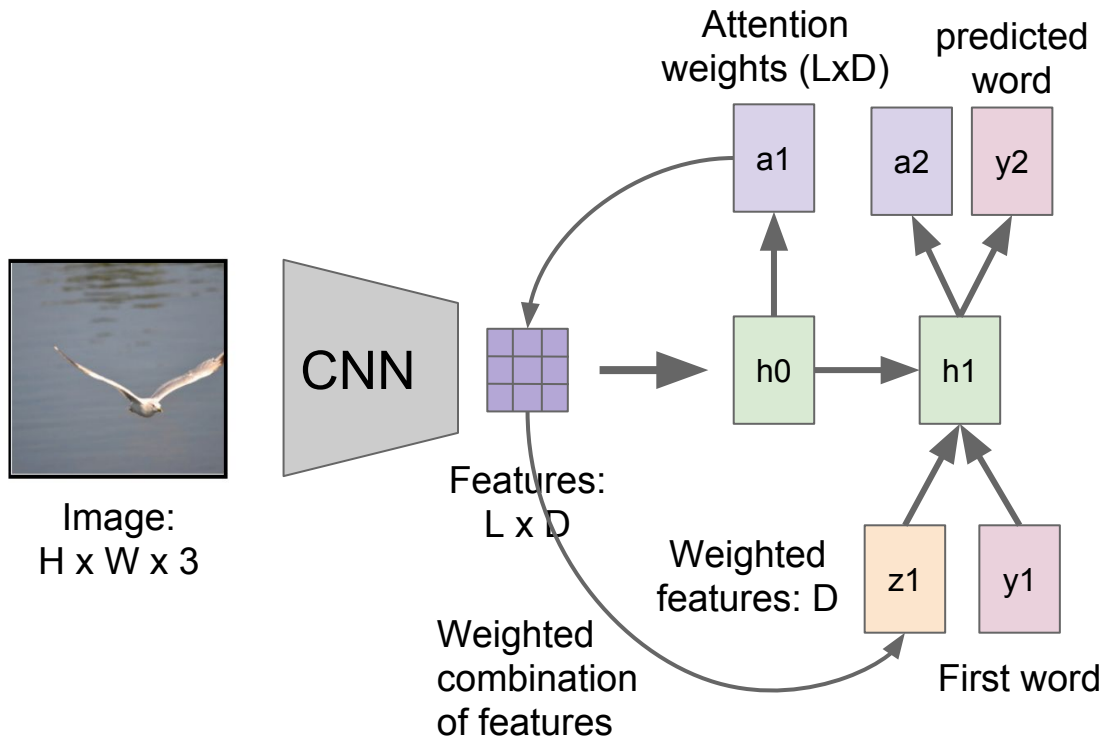


Features:  
 $L \times D$

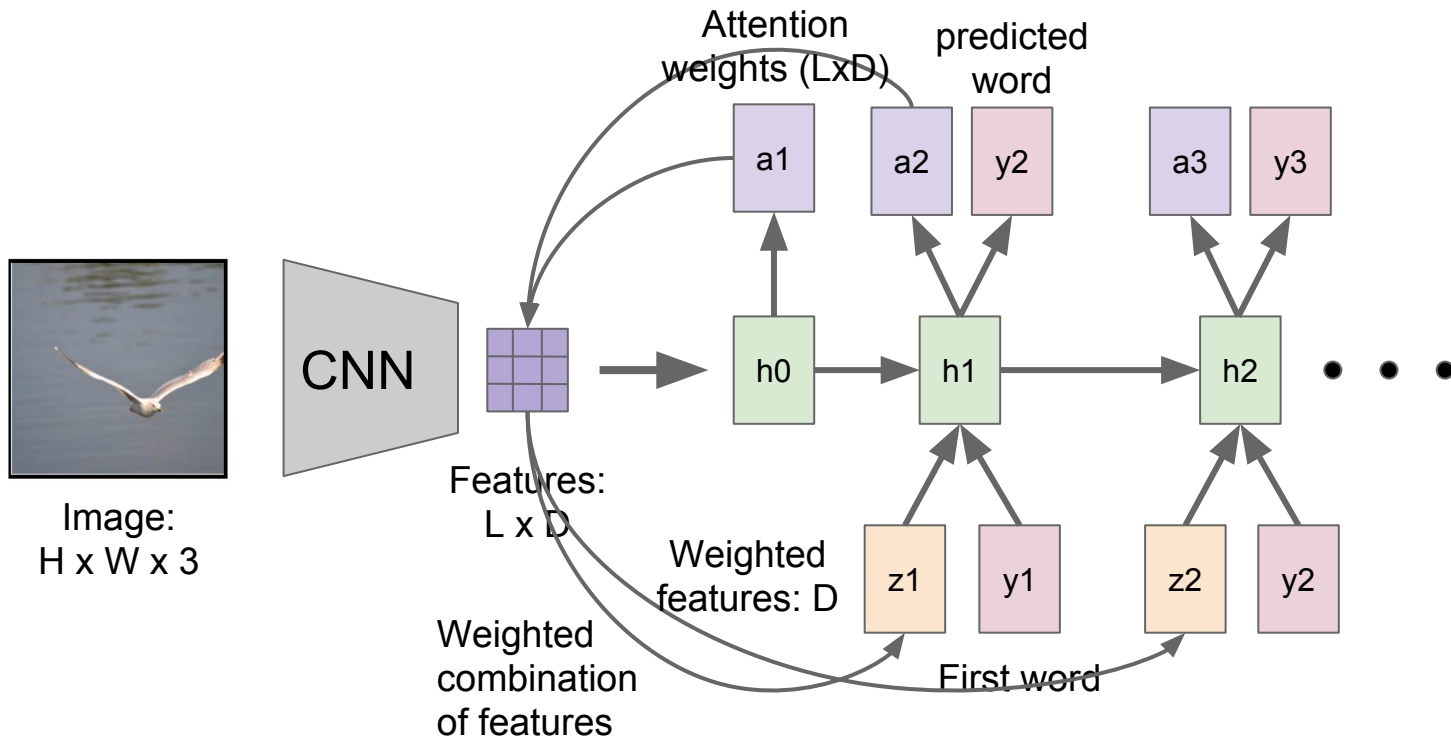
# Attention for Image Captioning



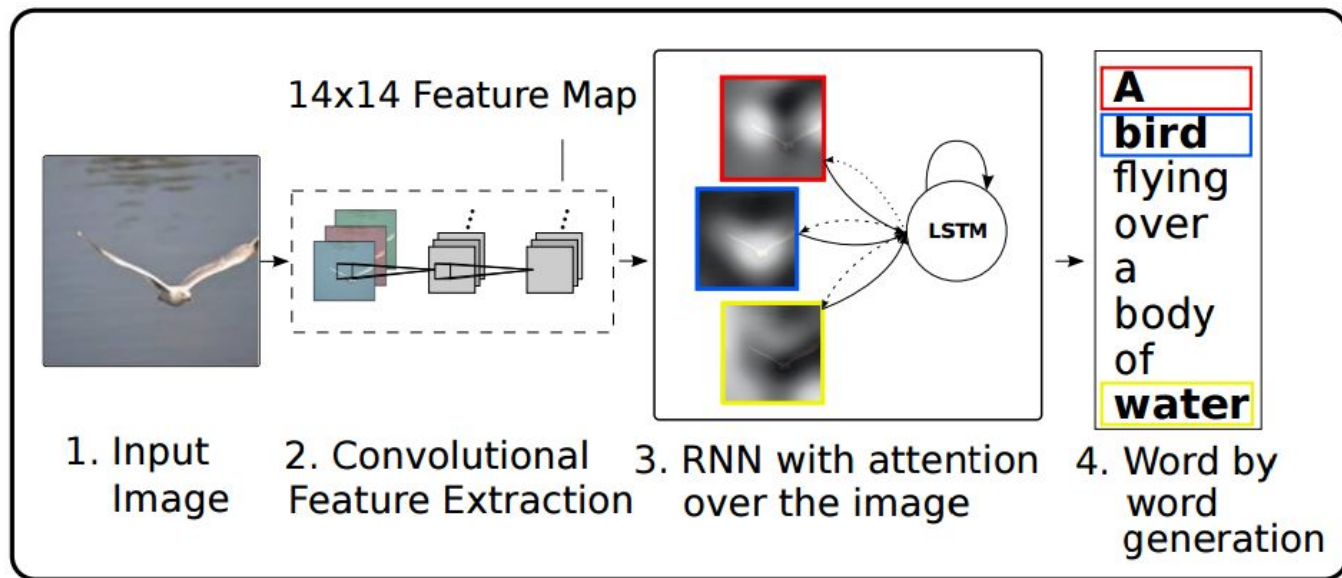
# Attention for Image Captioning



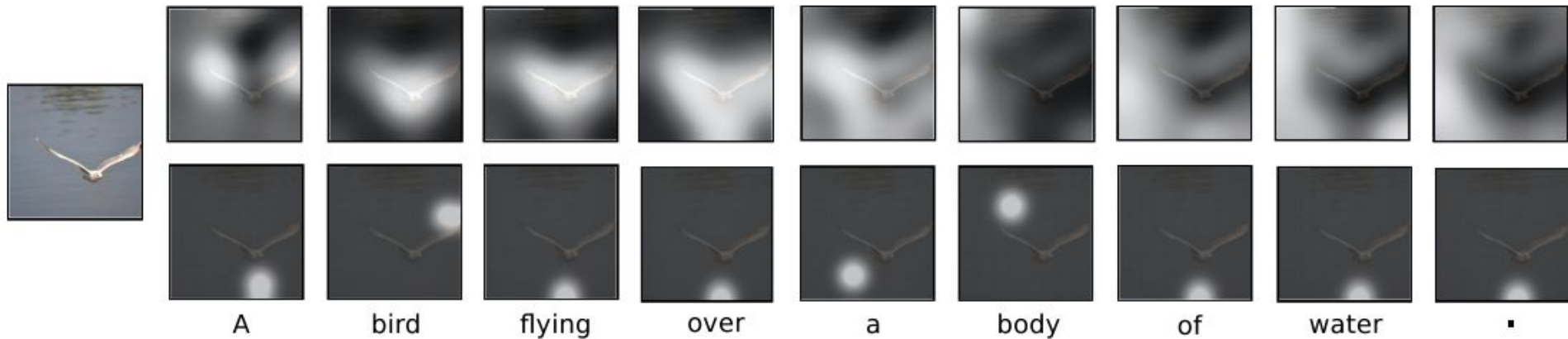
# Attention for Image Captioning



# Attention for Image Captioning



# Attention for Image Captioning





# Attention for Image Captioning



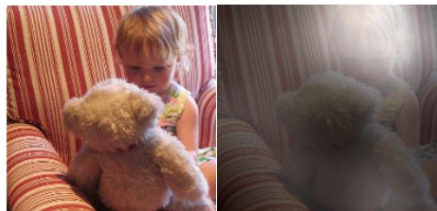
A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



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



A little girl sitting on a bed with a teddy bear.



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



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

# Soft Attention

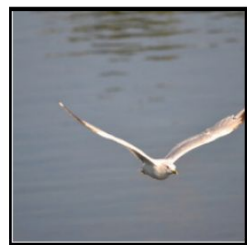
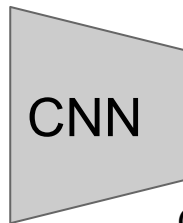


Image:  
 $H \times W \times 3$



a	b
c	d

Grid of features  
(Each D-dimensional)

From  
RNN:

$p_a$	$p_b$
$p_c$	$p_d$

Distribution over  
grid locations

$$p_a + p_b + p_c + p_c = 1$$

Context vector  $z$   
(D-dimensional)

## Soft attention:

Summarize ALL locations

$$z = p_a a + p_b b + p_c c + p_d d$$

Derivative  $dz/dp$  is nice!

Train with gradient descent

# Soft Attention

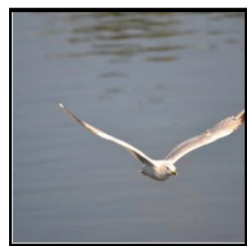
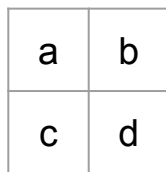
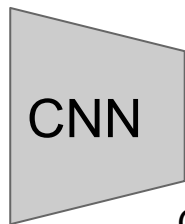
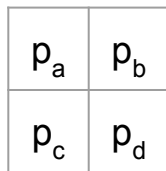


Image:  
 $H \times W \times 3$



Grid of features  
(Each D-dimensional)

From  
RNN:



Distribution over  
grid locations

$$p_a + p_b + p_c + p_c = 1$$

Context vector  $z$   
(D-dimensional)

## Soft attention:

Summarize ALL locations

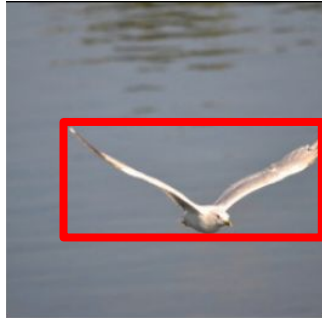
$$z = p_a a + p_b b + p_c c + p_d d$$

Differentiable function

Train with gradient descent

- Still uses the whole input !
- Constrained to fix grid

# Hard attention

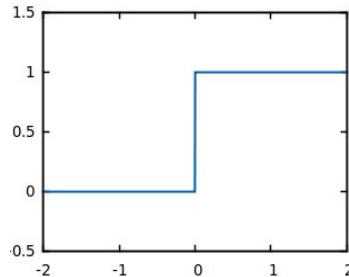
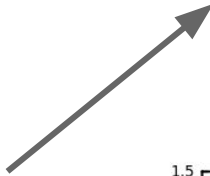


Input image:  
 $H \times W \times 3$

Box Coordinates:  
( $x_c, y_c, w, h$ )



Cropped and  
rescaled image:  
 $X \times Y \times 3$



Gradient is 0 almost everywhere  
Gradient is undefined at  $x = 0$

**Hard attention:**  
Sample a subset  
of the input



Not a differentiable function !



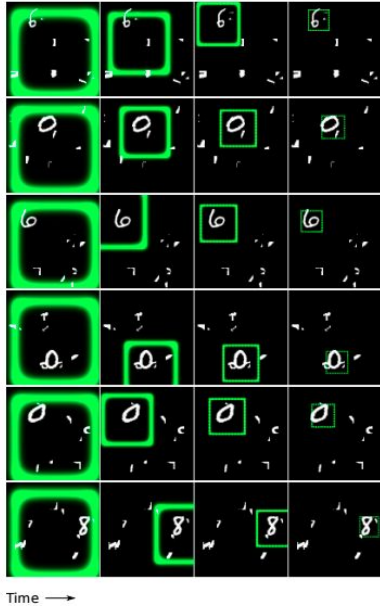
Can't train with backprop :(



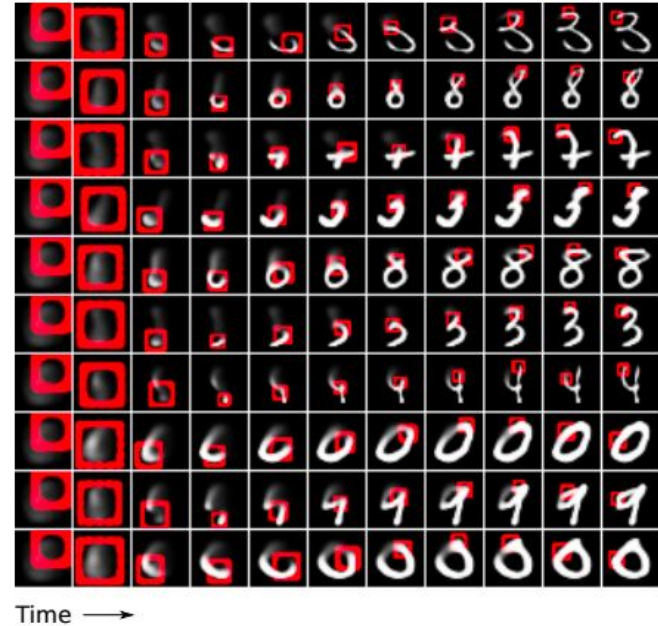
need [reinforcement learning](#)

# Hard attention

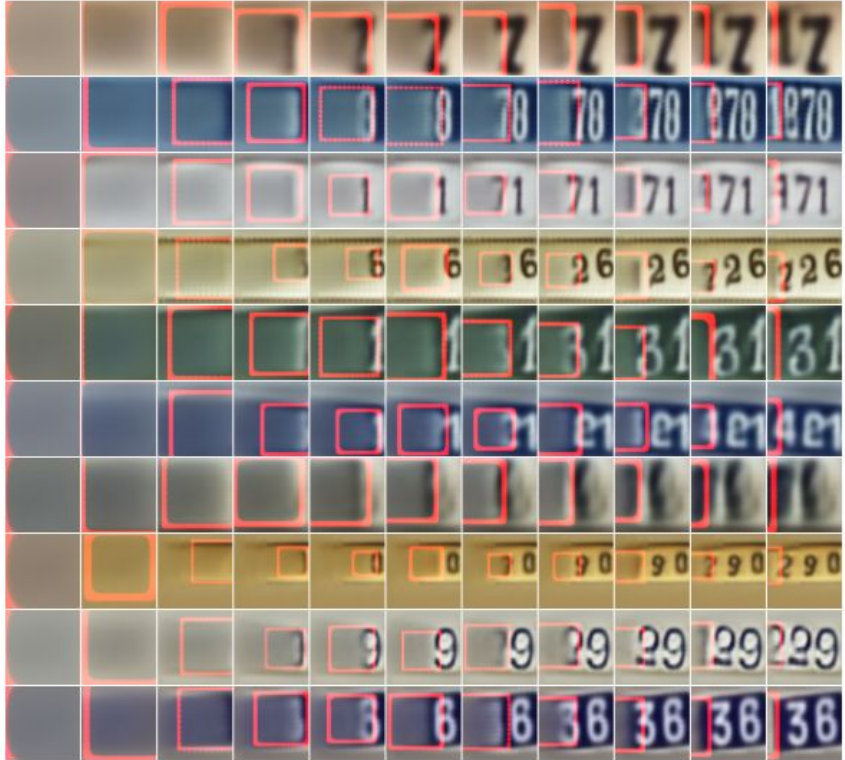
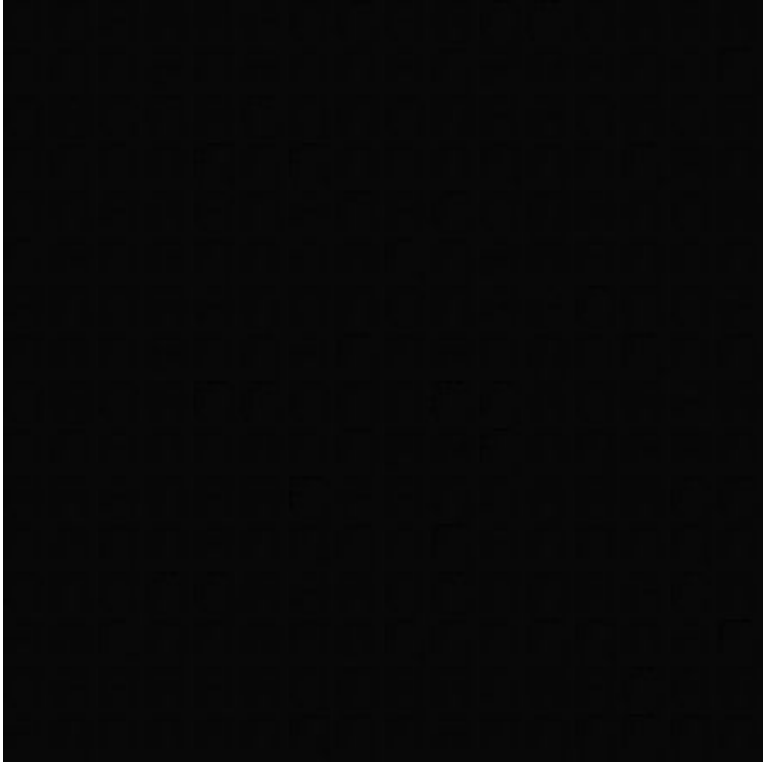
**Classify** images by attending to arbitrary regions of the *input*



**Generate** images by attending to arbitrary regions of the *output*



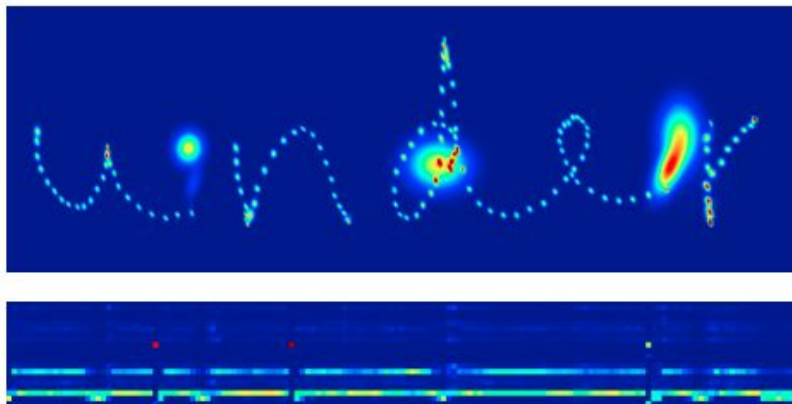
# Hard attention



Time →

# Hard attention

Read text, generate handwriting using an RNN that attends at different arbitrary regions over time



more of national temperament

more of national temperament  
more of national temperament  
more of national temperament  
more of national temperament  
more of national temperament

REAL

**GENERATED**

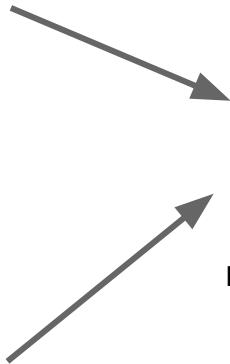


# Hard attention

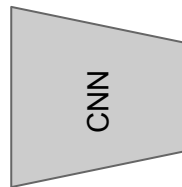


Input image:  
 $H \times W \times 3$

Box Coordinates:  
( $x_c, y_c, w, h$ )



Cropped and  
rescaled image:  
 $X \times Y \times 3$



bird

Not a differentiable function !



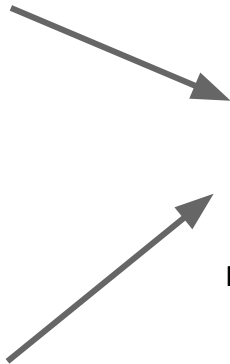
Can't train with backprop :(

# Spatial Transformer Networks

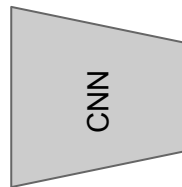


Input image:  
 $H \times W \times 3$

Box Coordinates:  
( $x_c, y_c, w, h$ )



Cropped and  
rescaled image:  
 $X \times Y \times 3$



bird

Not a differentiable function !

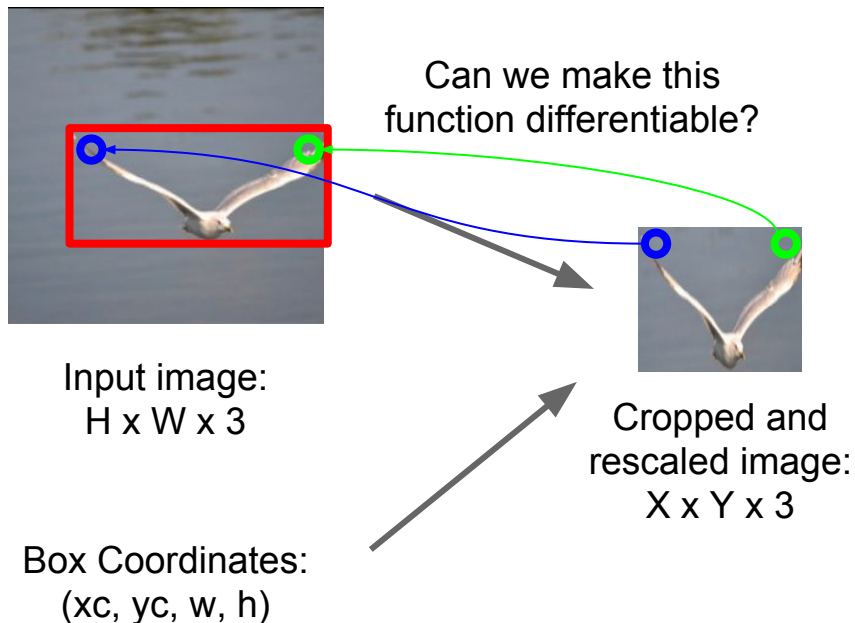
Make it differentiable



Can't train with backprop :(

Train with backprop :) 24

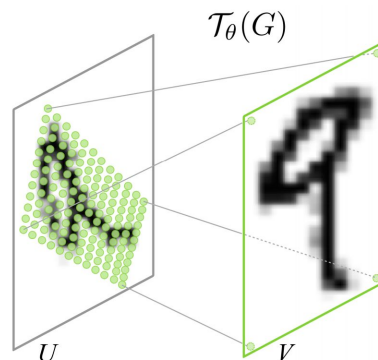
# Spatial Transformer Networks



**Idea:** Function mapping *pixel coordinates*  $(x_t, y_t)$  of output to *pixel coordinates*  $(x_s, y_s)$  of input

$$\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \begin{bmatrix} \theta_{11} & \theta_{12} & \theta_{13} \\ \theta_{21} & \theta_{22} & \theta_{23} \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix}$$

Network attends to input by predicting  $\theta$

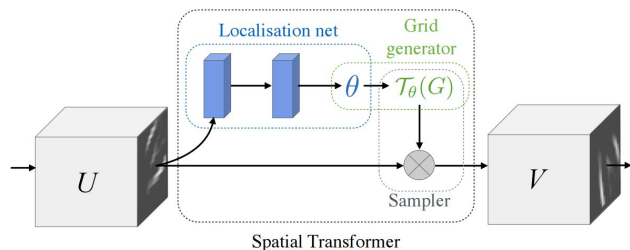


Repeat for all pixels in *output* to get a **sampling grid**

Then use **bilinear interpolation** to compute output

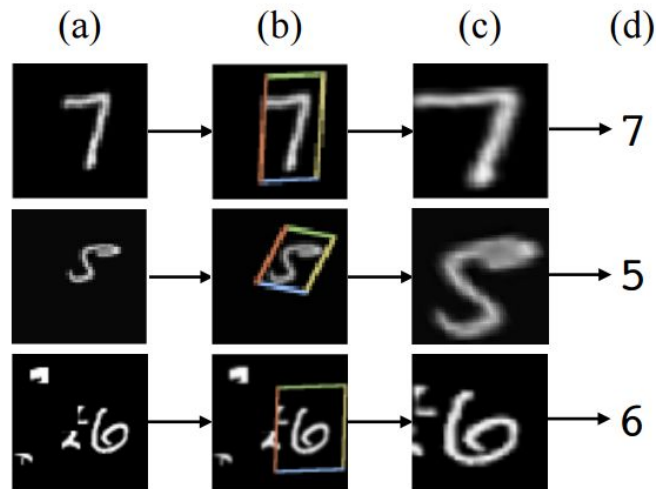
# Spatial Transformer Networks

Differentiable module

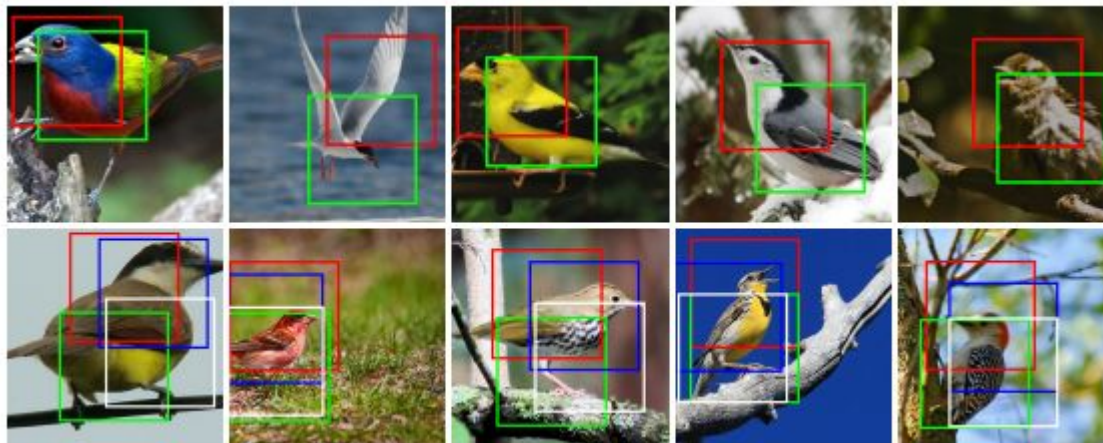


Easy to incorporate in any network, anywhere !

Insert spatial transformers into a classification network and it learns to attend and transform the input



# Spatial Transformer Networks

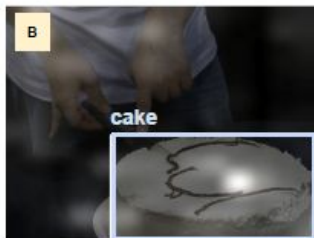


Fine-grained classification

# Visual Attention



What kind of animal is in the photo?  
A cat.



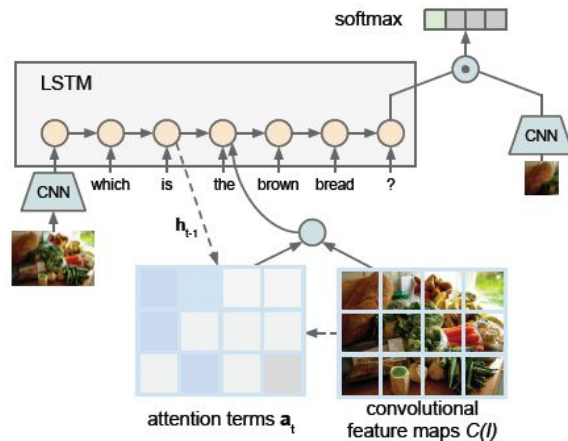
Why is the person holding a knife?  
To cut the cake with.



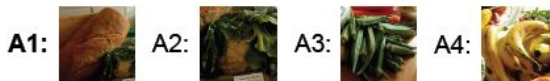
Where are the carrots?  
At the top.



How many people are there?  
Three.

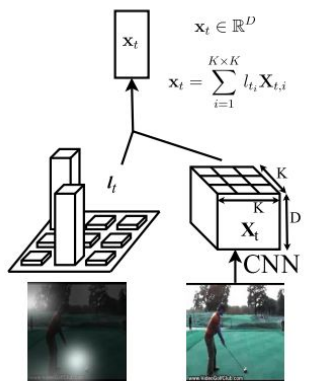


Q: Which is the brown bread?

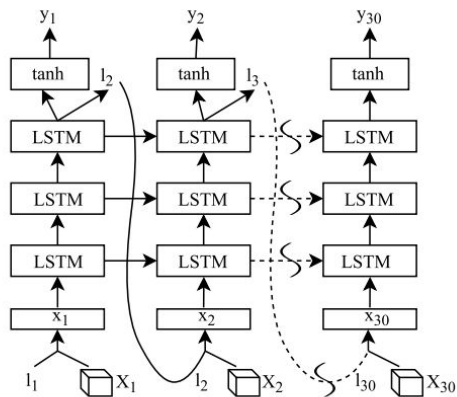


Visual Question Answering

# Visual Attention



(a) The soft attention mechanism

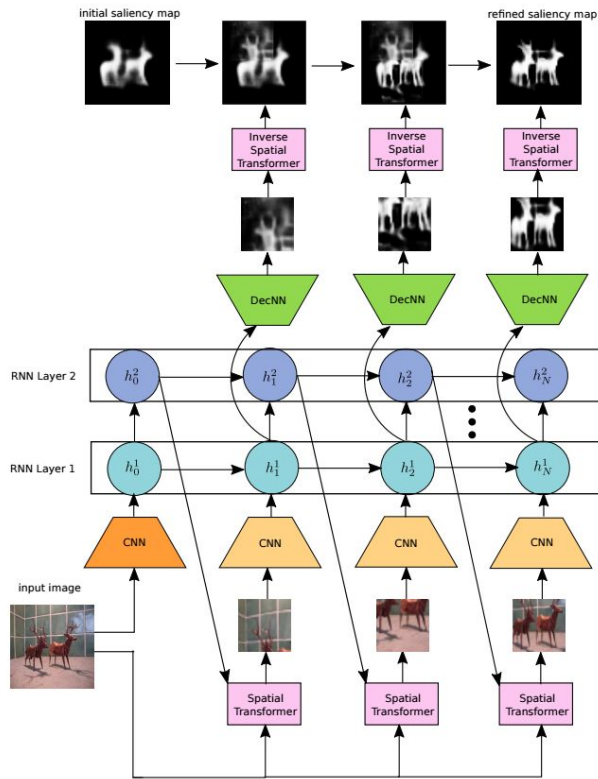


(b) Our recurrent model

## Action Recognition in Videos

Sharma et al. [Action Recognition Using Visual Attention](#). arXiv 2016

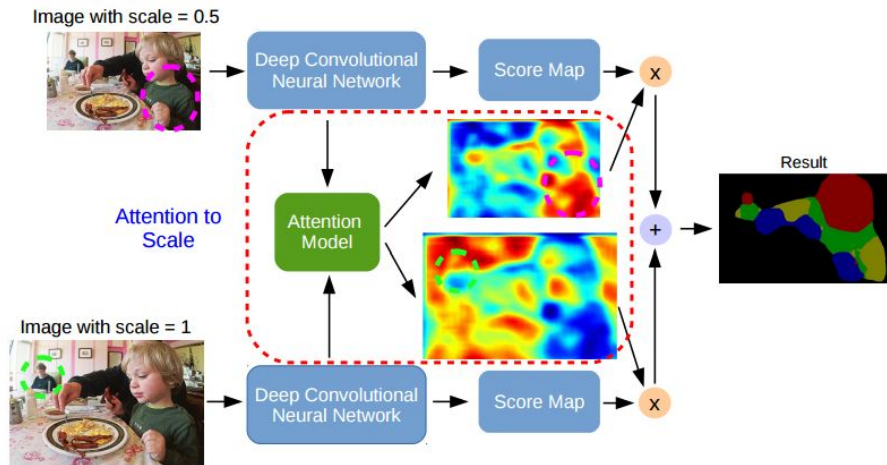
Kuen et al. [Recurrent Attentional Networks for Saliency Detection](#). CVPR 2016



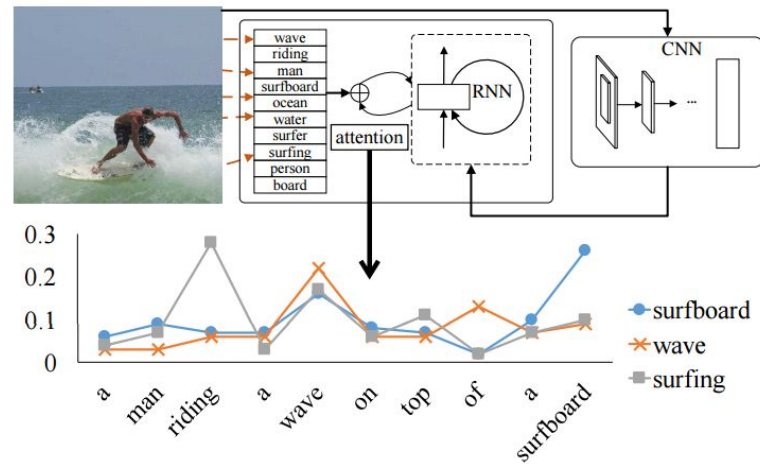
## Salient Object Detection



# Other examples



Attention to scale for semantic segmentation



Semantic attention For image captioning

Chen et al. [Attention to Scale: Scale-aware Semantic Image Segmentation](#). CVPR 2016

You et al. [Image Captioning with Semantic Attention](#). CVPR 2016

# Resources

- CS231n Lecture @ Stanford [[slides](#)][[video](#)]
- [More on Reinforcement Learning](#)
- [Soft vs Hard attention](#)
- [Handwriting generation demo](#)
- Spatial Transformer Networks - [Slides](#) & [Video](#) by Victor Campos
- Attention implementations:
  - [Seq2seq in Keras](#)
  - [DRAW & Spatial Transformers in Keras](#)
  - [DRAW in Lasagne](#)
  - [DRAW in Tensorflow](#)