

DEEP  
LEARNING  
WORKSHOP

Dublin City University  
27-28 April 2017

Day 2 Lecture 12

Attention Models



Amaia Salvador

[amaia.salvador@upc.edu](mailto:amaia.salvador@upc.edu)

PhD Candidate

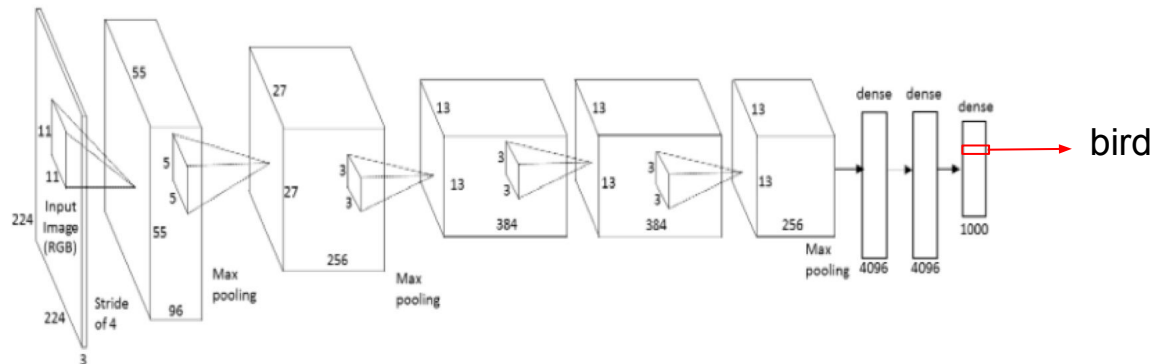
Universitat Politècnica de Catalunya



# Attention Models: Motivation



Image:  
 $H \times W \times 3$



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

# Attention Models: Motivation

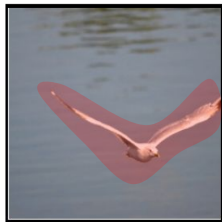
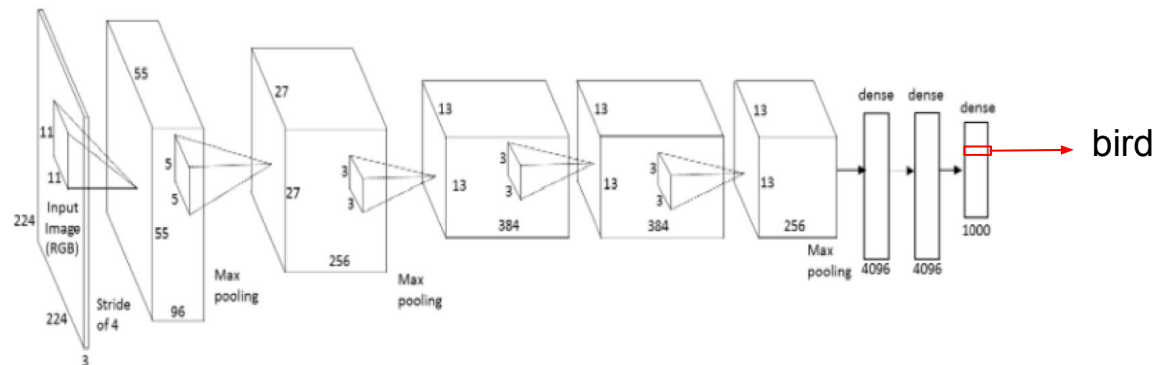


Image:  
 $H \times W \times 3$



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

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

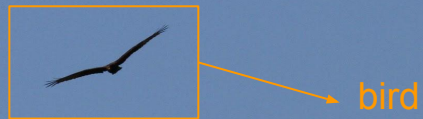
# Attention Models: Motivation



Attention models can  
relieve computational burden

Helpful when processing big  
images !

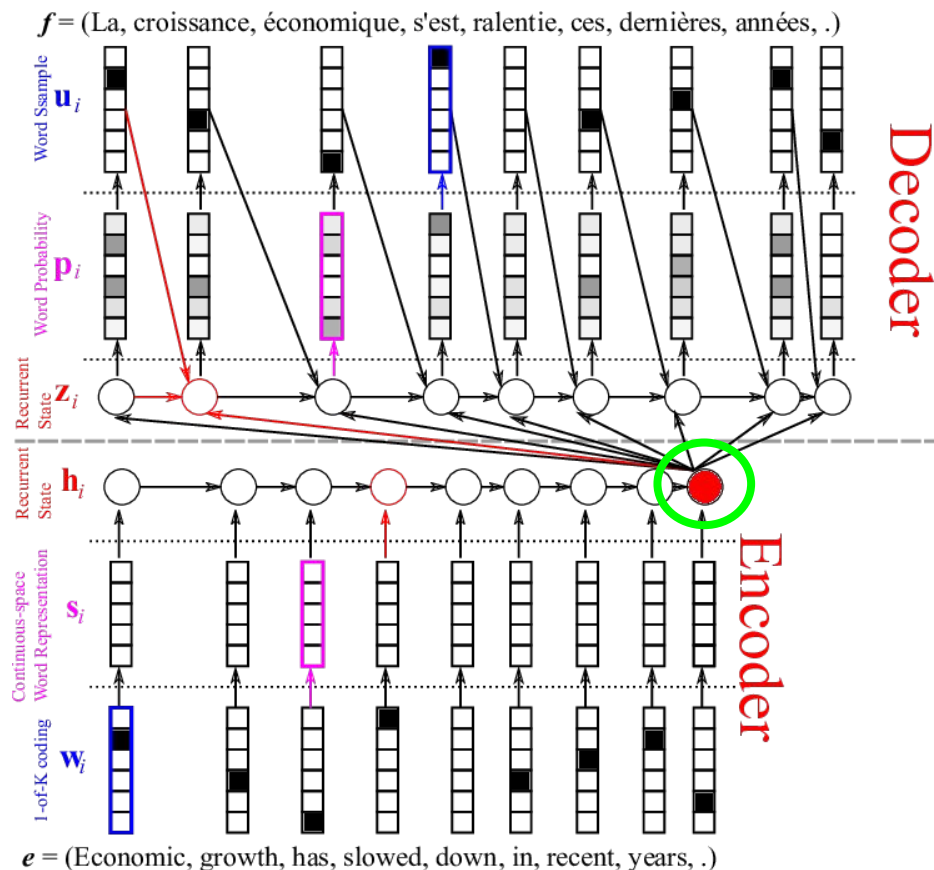
# Attention Models: Motivation



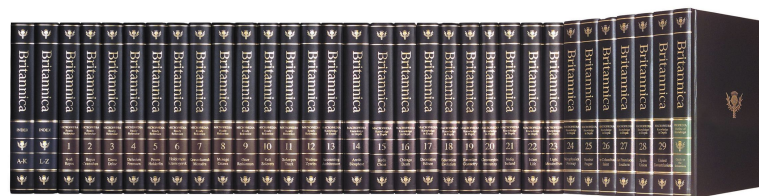
Attention models can  
relieve computational burden

Helpful when processing big  
images !

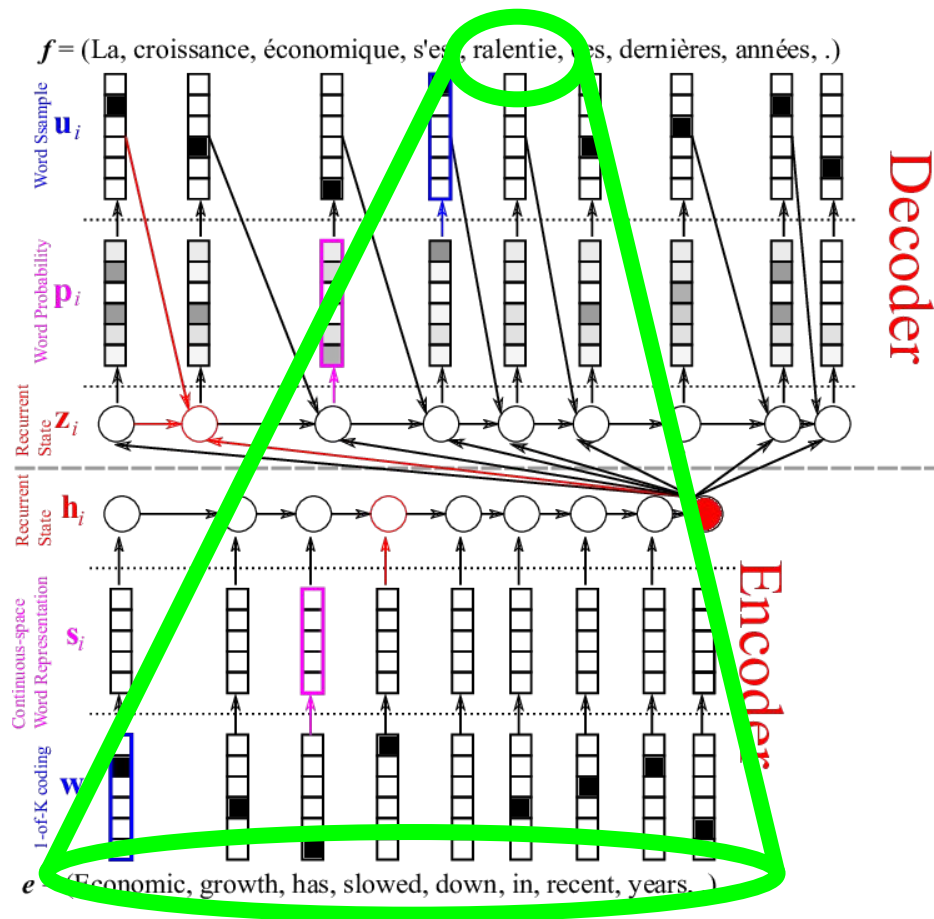
# Encoder & Decoder for Machine Translation



Limitation 1: The whole information is encoded in a **fixed-size vector**, no matter the length of the input sentence.



# Encoder & Decoder

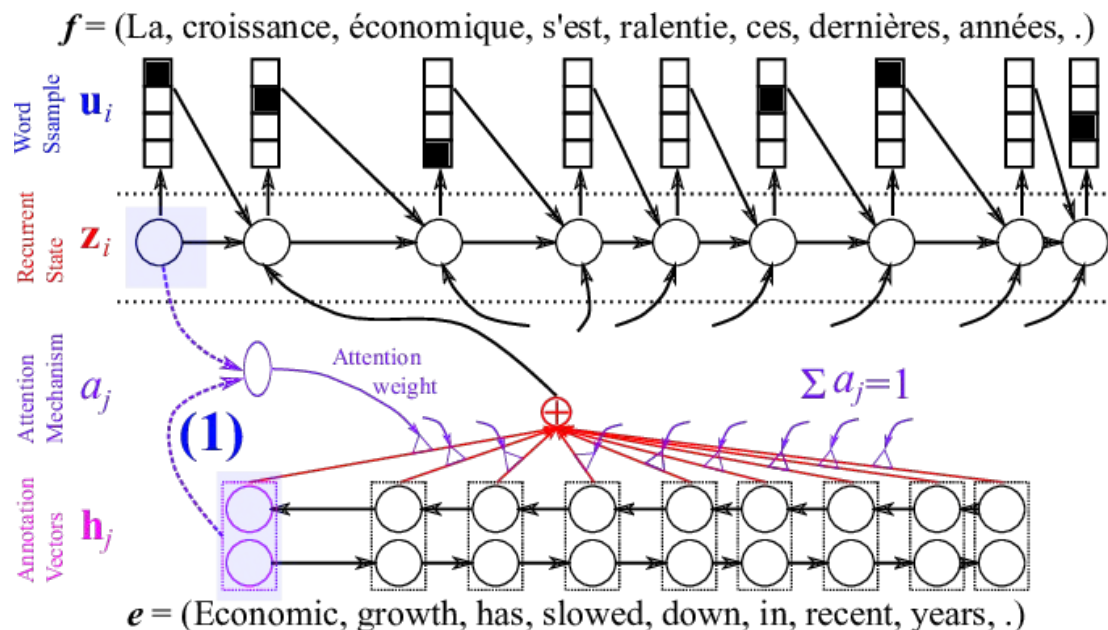


Limitation 2: All output predictions are based on the **final and static** recurrent state of the encoder ( $h_T$ )

No matter the output word being predicted at each time step, all input words are considered in an equal way.

# Attention Mechanism

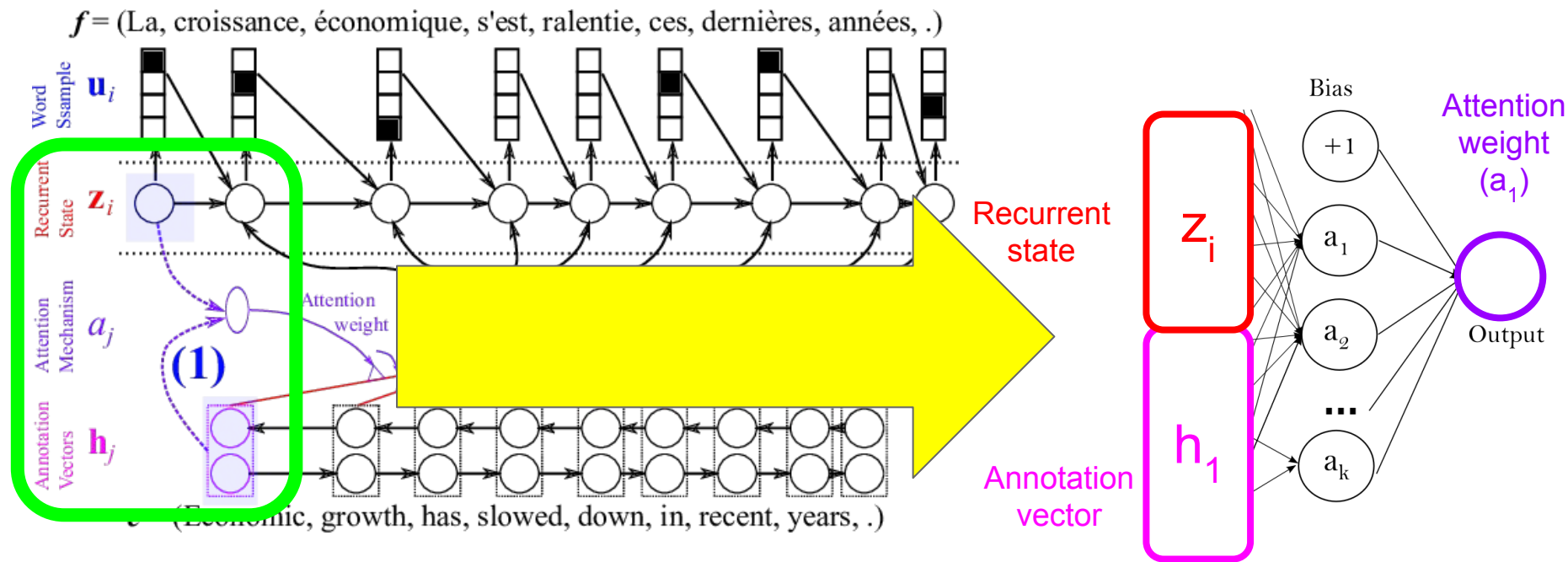
The vector to be fed to the RNN at each timestep is a weighted sum of all the annotation vectors.





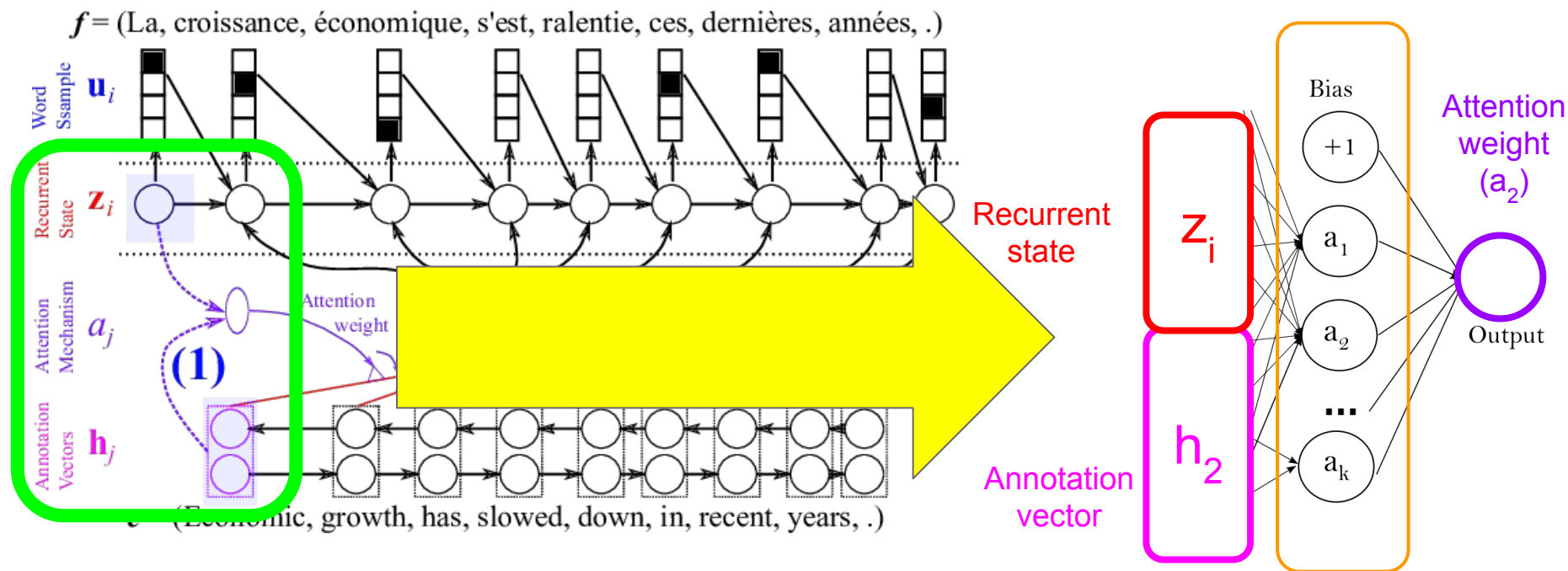
# Attention Mechanism

An attention weight (scalar) is predicted at each time-step for each annotation vector  $h_j$  with a simple fully connected neural network.



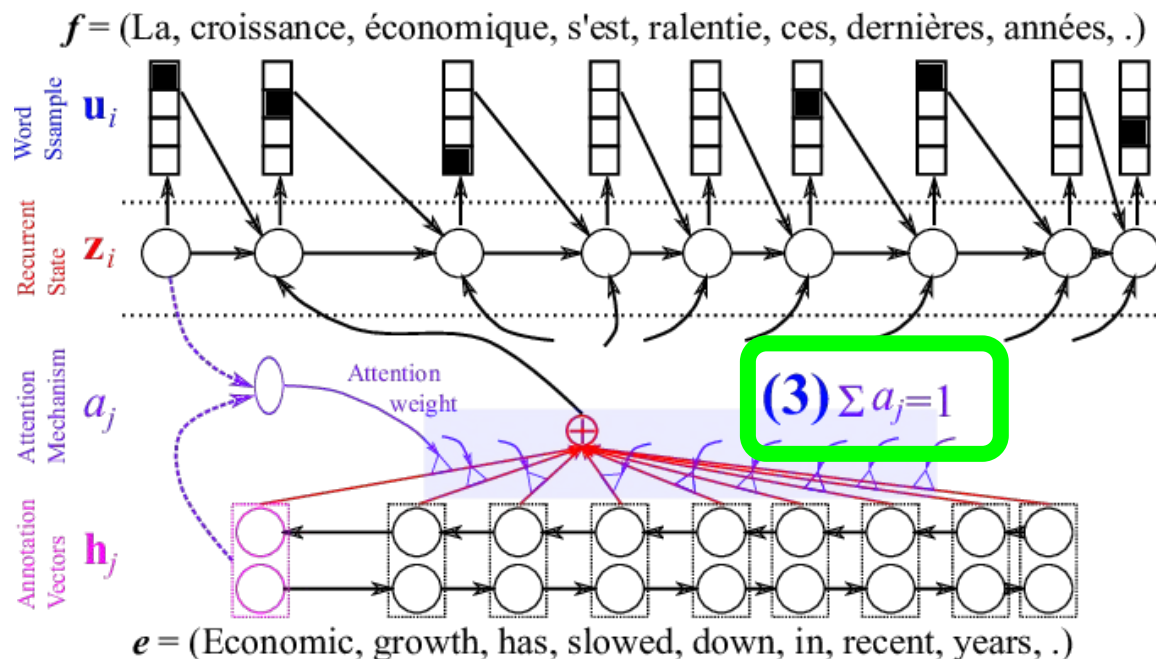
# Attention Mechanism

An attention weight (scalar) is predicted at each time-step for each annotation vector  $h_j$  with a simple fully connected neural network.



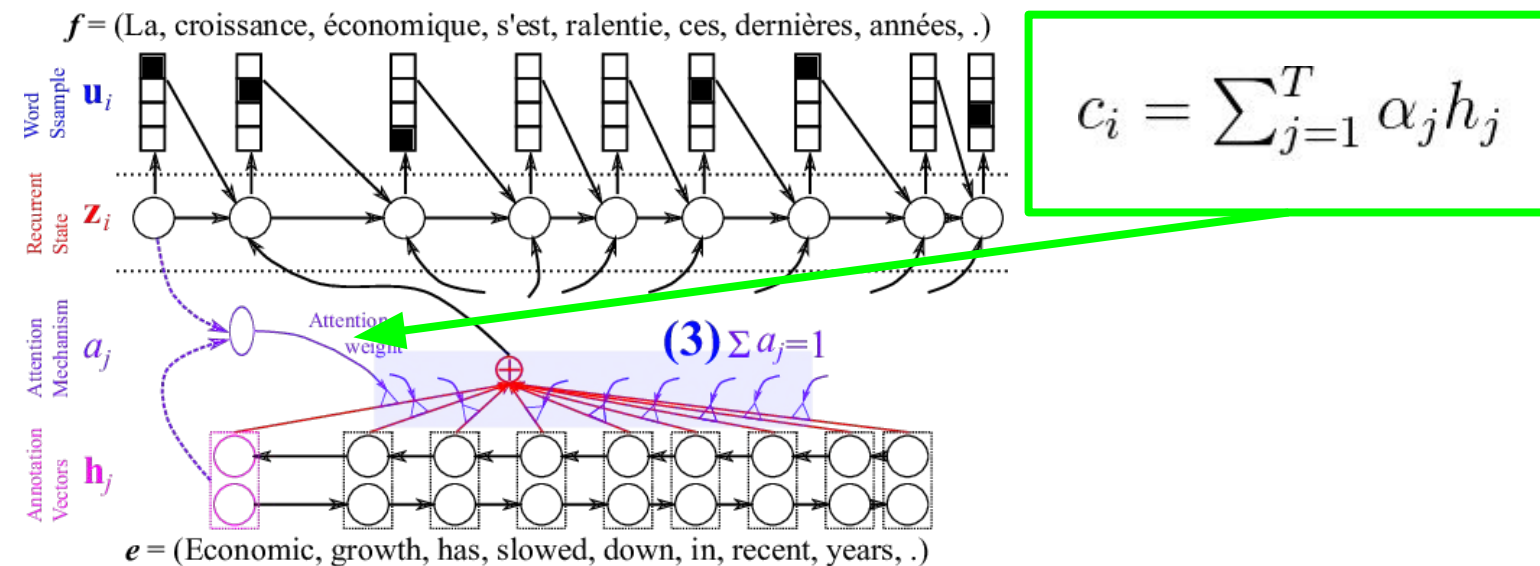
# Attention Mechanism

Once a relevance score (weight) is estimated for each word, they are normalized with a softmax function so they sum up to 1.



# Attention Mechanism

Finally, a context-aware representation  $c_i$  for the output word at timestep  $i$  can be defined as:



# Attention Mechanism

The model automatically finds the correspondence structure between two languages (alignment).



Economic growth has slowed down in recent years .



La croissance économique s' est ralentie ces dernières années .

(Edge thicknesses represent the attention weights found by the attention model)

# Attention Mechanism

The model automatically found the correspondence structure between two languages (alignment).



Economic growth has slowed down in recent years .



Das Wirtschaftswachstum hat sich in den letzten Jahren verlangsamt .

(Edge thicknesses represent the attention weights found by the attention model)

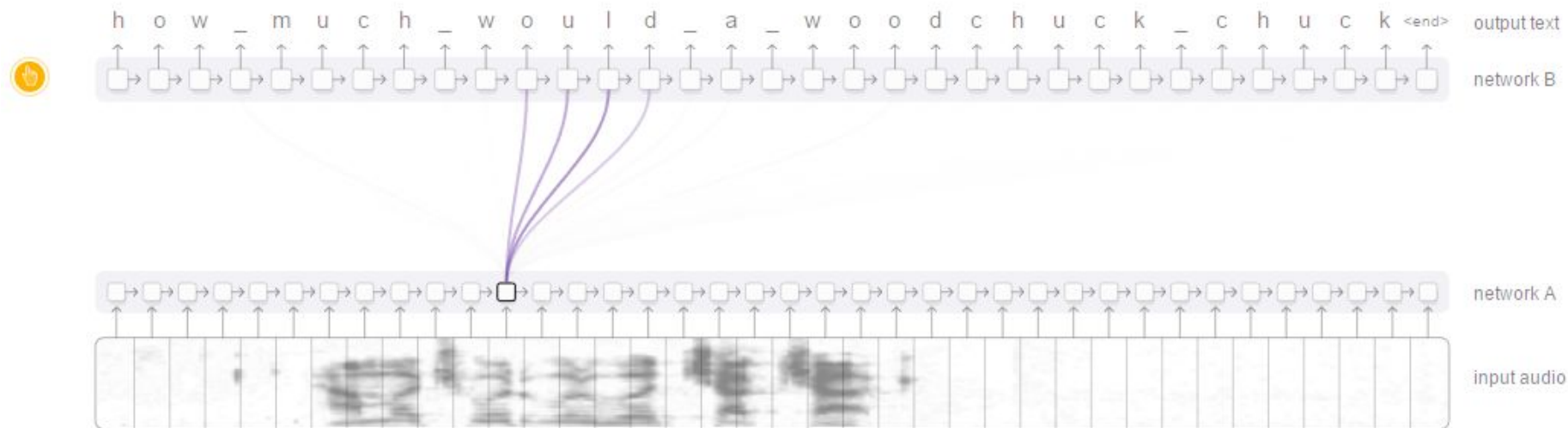
# Attention Models

Attend to different parts of the input to optimize a certain output

# Attention Models

Attend to different parts of the input to optimize a certain output

Input: Audio features; Output: Text



Chan et al. [Listen, Attend and Spell](#). ICASSP 2016

Source: [distill.pub](https://arxiv.org/abs/1603.01567)



# Attention Models

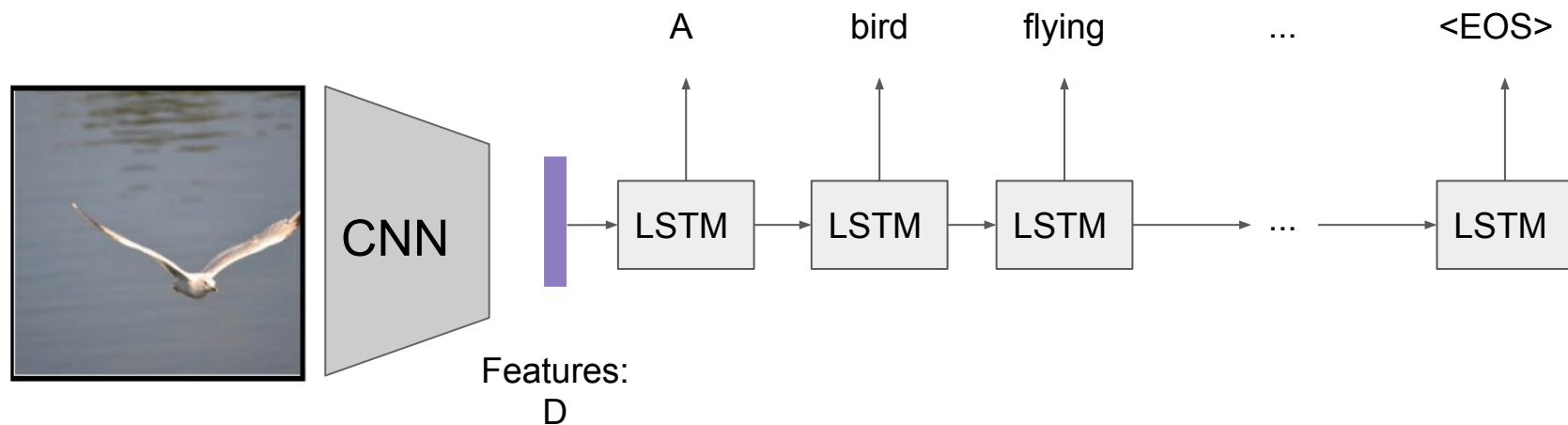
Attend to different parts of the input to optimize a certain output

Input: Image; Output: Text



A bird flying over a body of water

# LSTM Decoder for Image Captioning

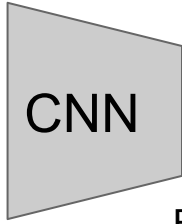


The LSTM decoder “sees” the input only at the beginning !

# Attention for Image Captioning

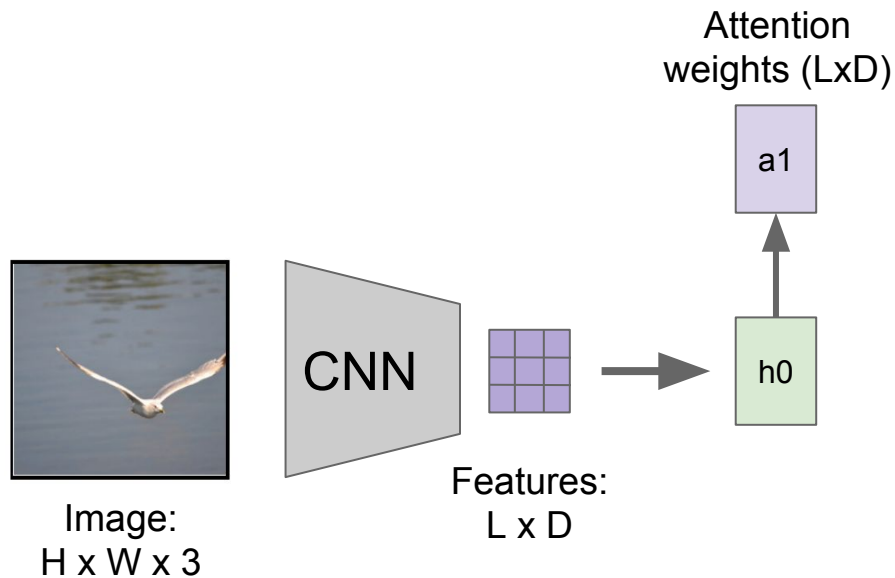


Image:  
 $H \times W \times 3$

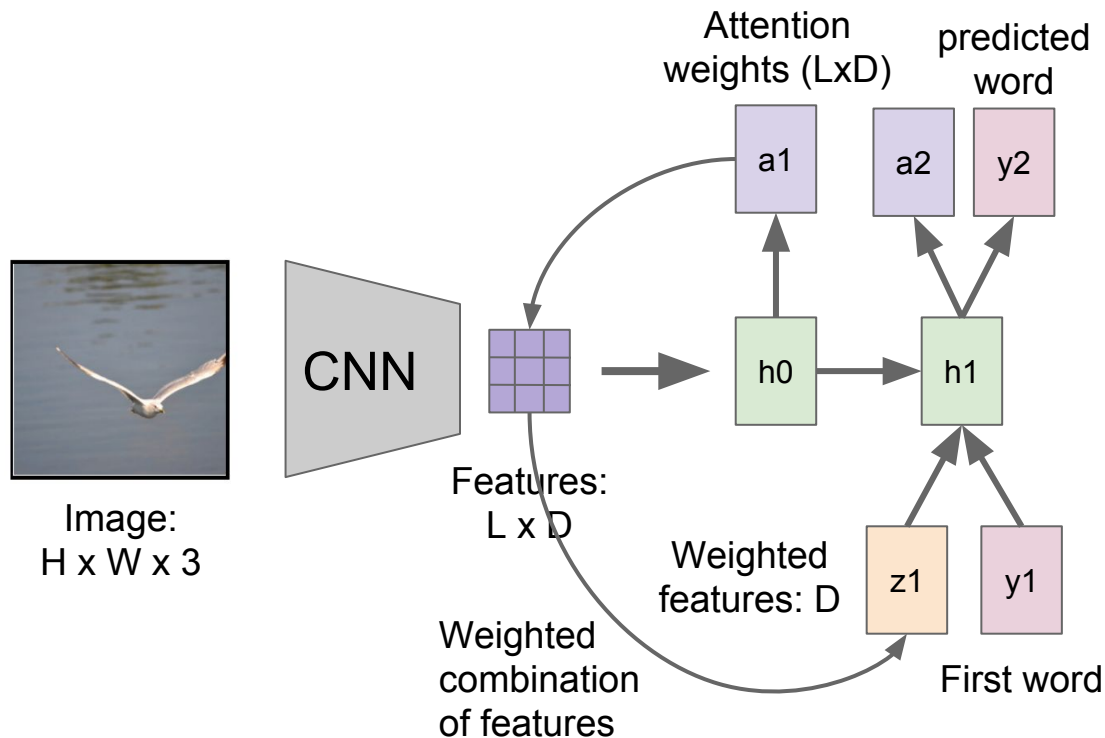


Features:  
 $L \times D$

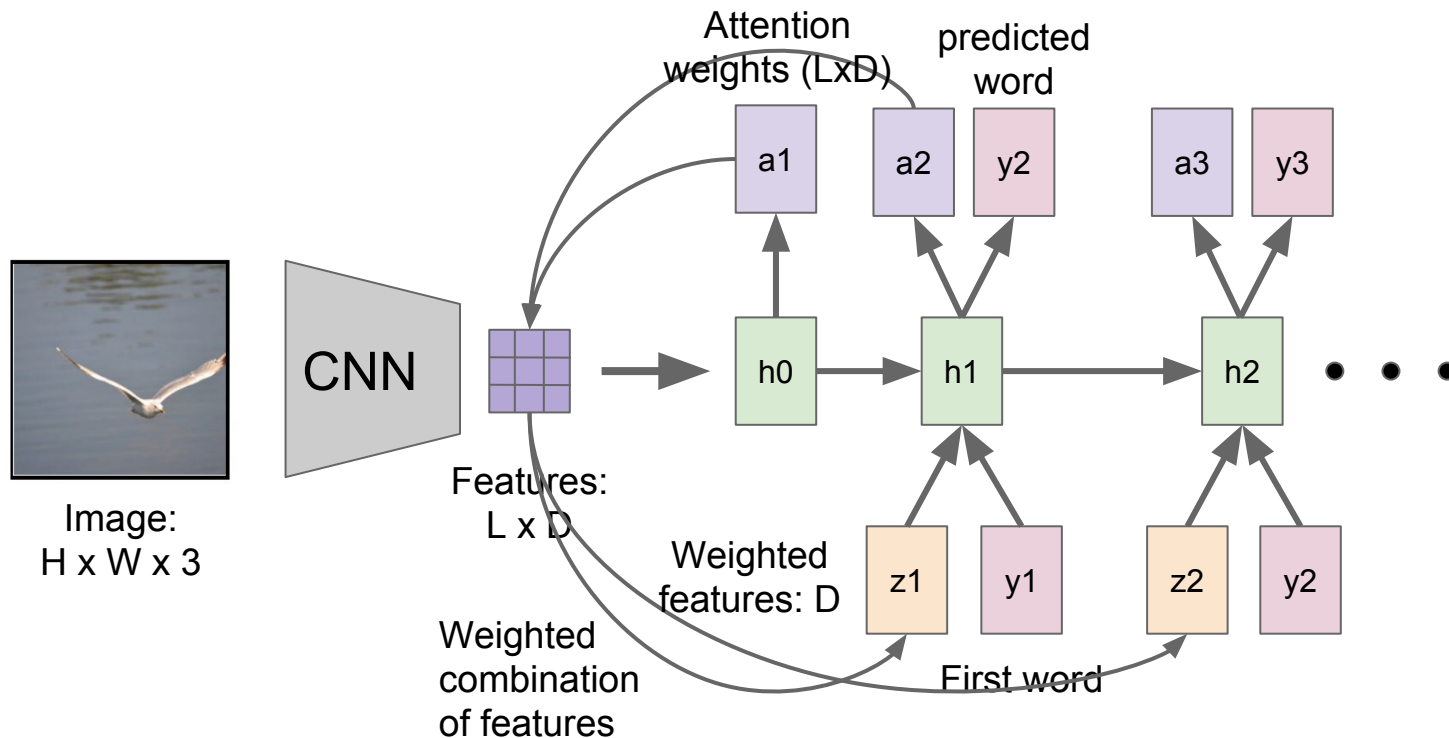
# Attention for Image Captioning



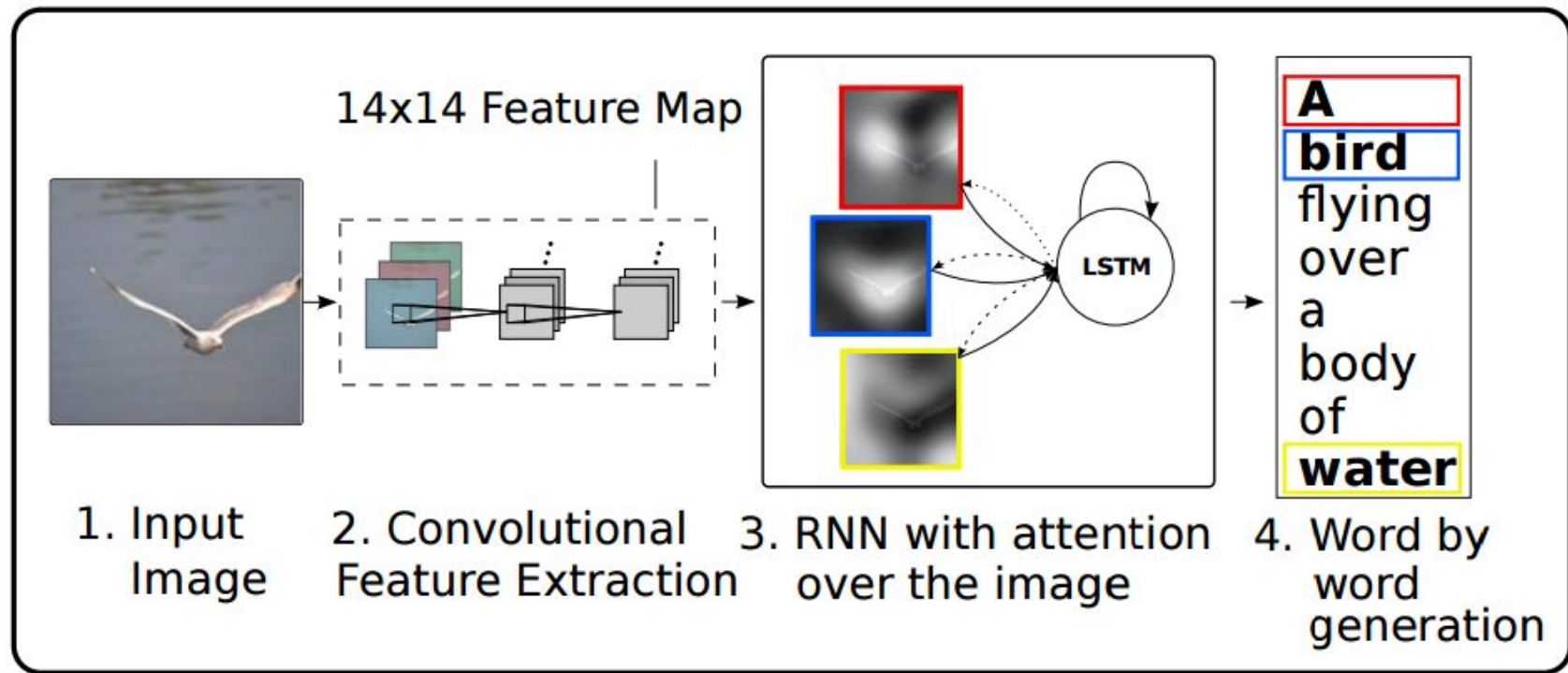
# Attention for Image Captioning



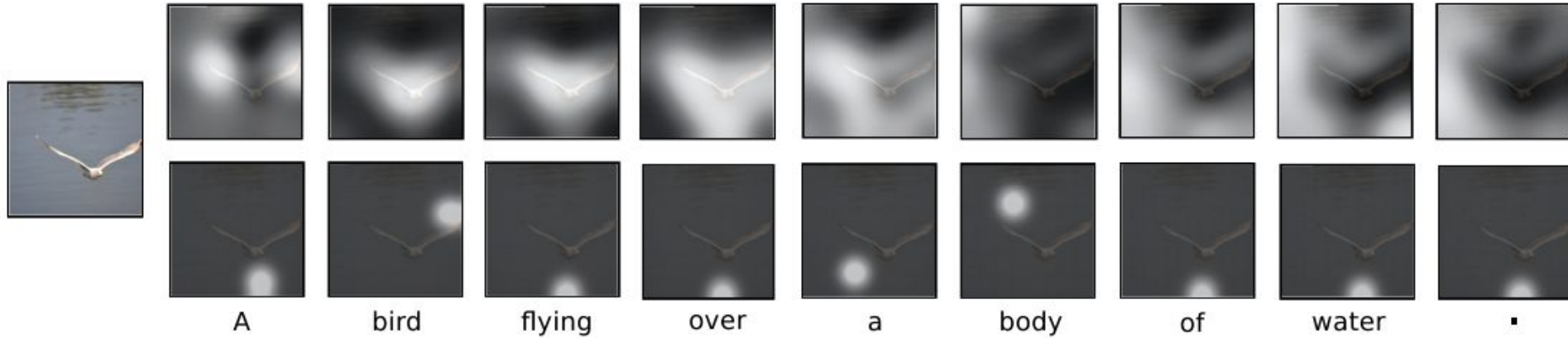
# Attention for Image Captioning



# Attention for Image Captioning



# Attention for Image Captioning





# Attention for Image Captioning



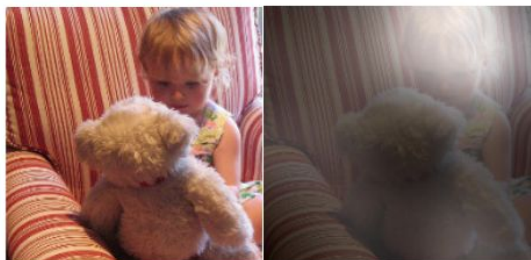
A woman is throwing a frisbee in a park.



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.



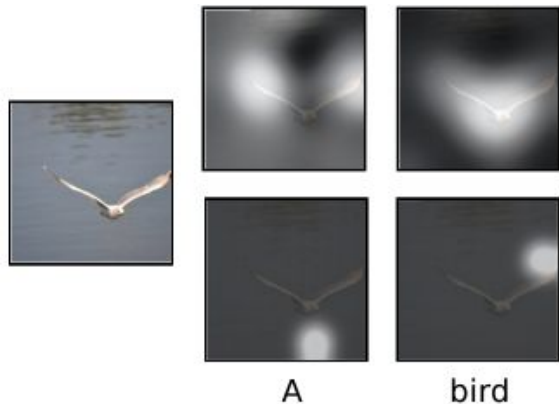
A group of people sitting on a boat in the water.



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

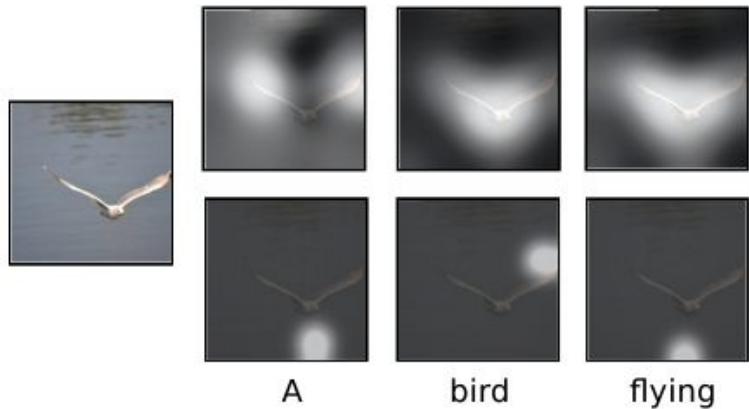
# Attention for Image Captioning

Some outputs can probably be predicted without looking at the image...



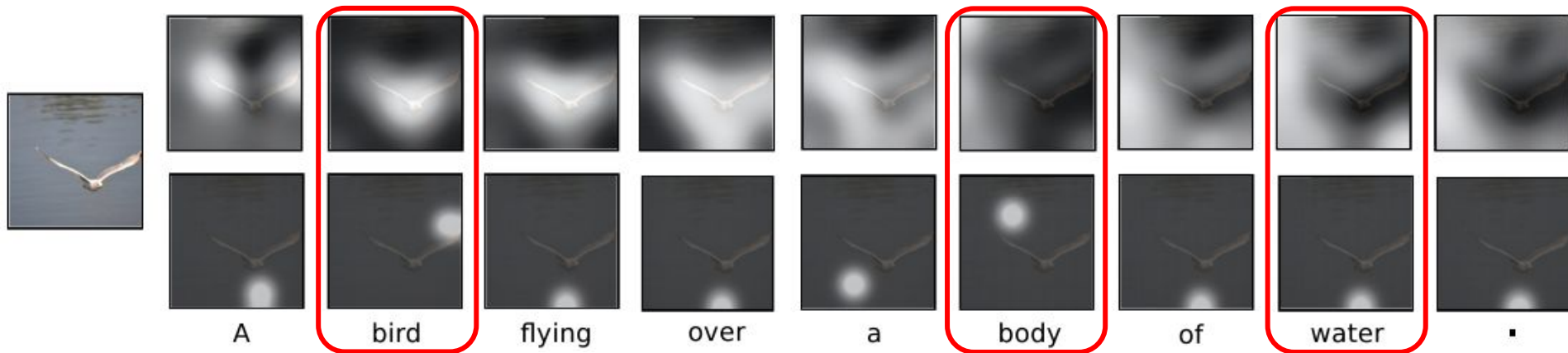
# Attention for Image Captioning

Some outputs can probably be predicted without looking at the image...

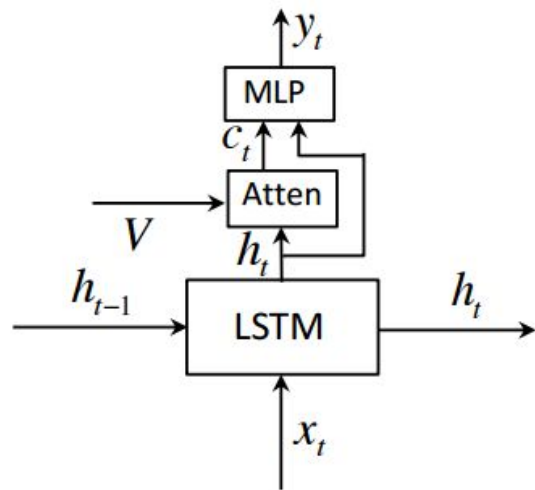


# Attention for Image Captioning

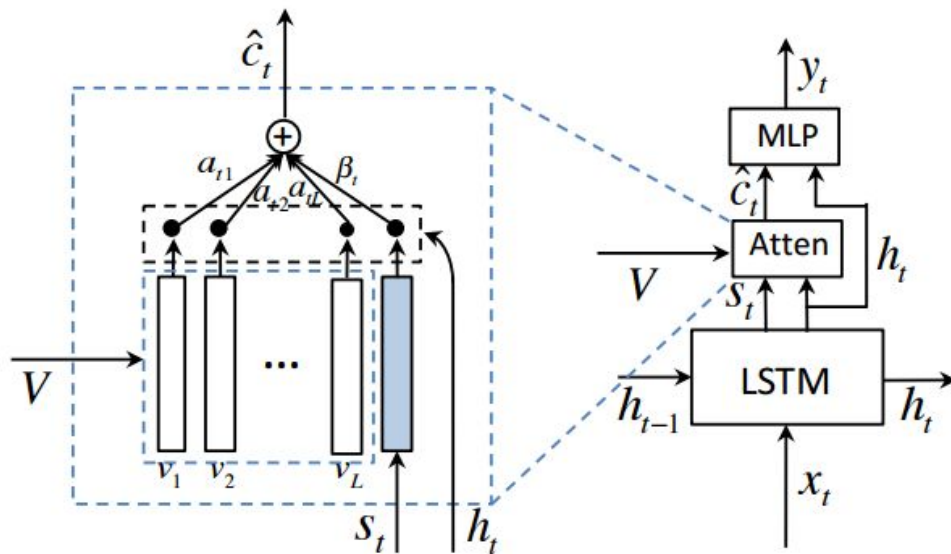
Can we focus on the image only when necessary?



# Attention for Image Captioning

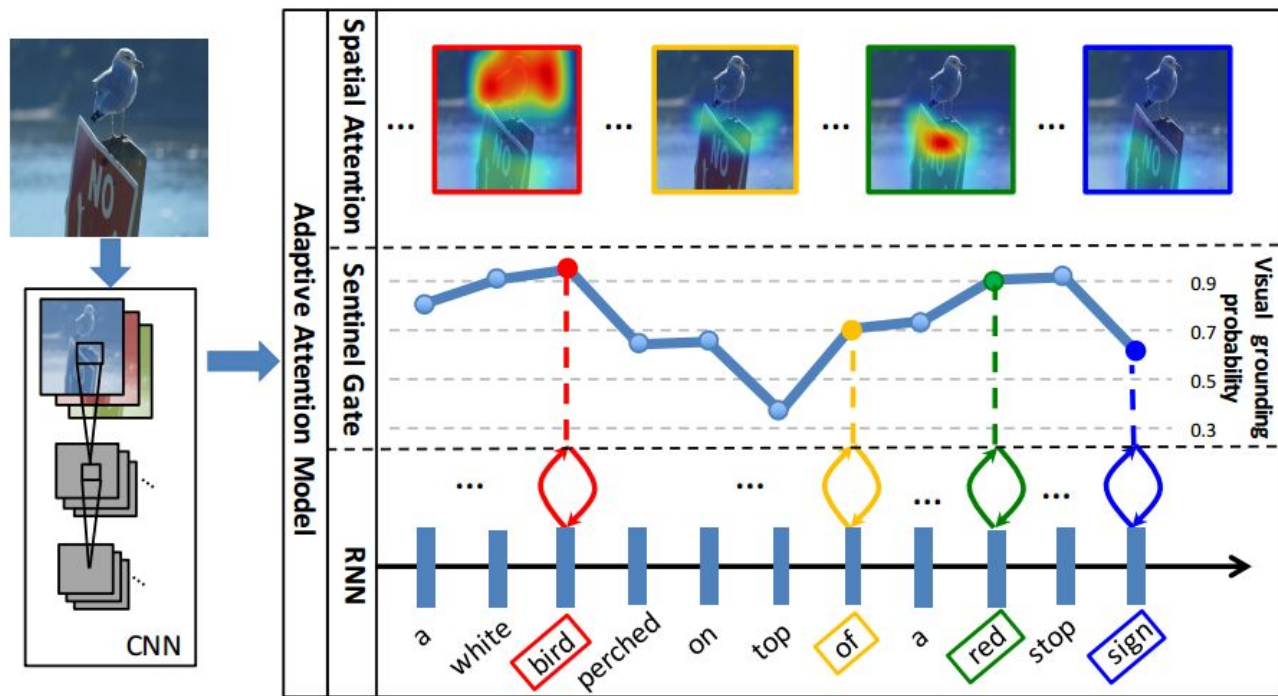


“Regular” Attention

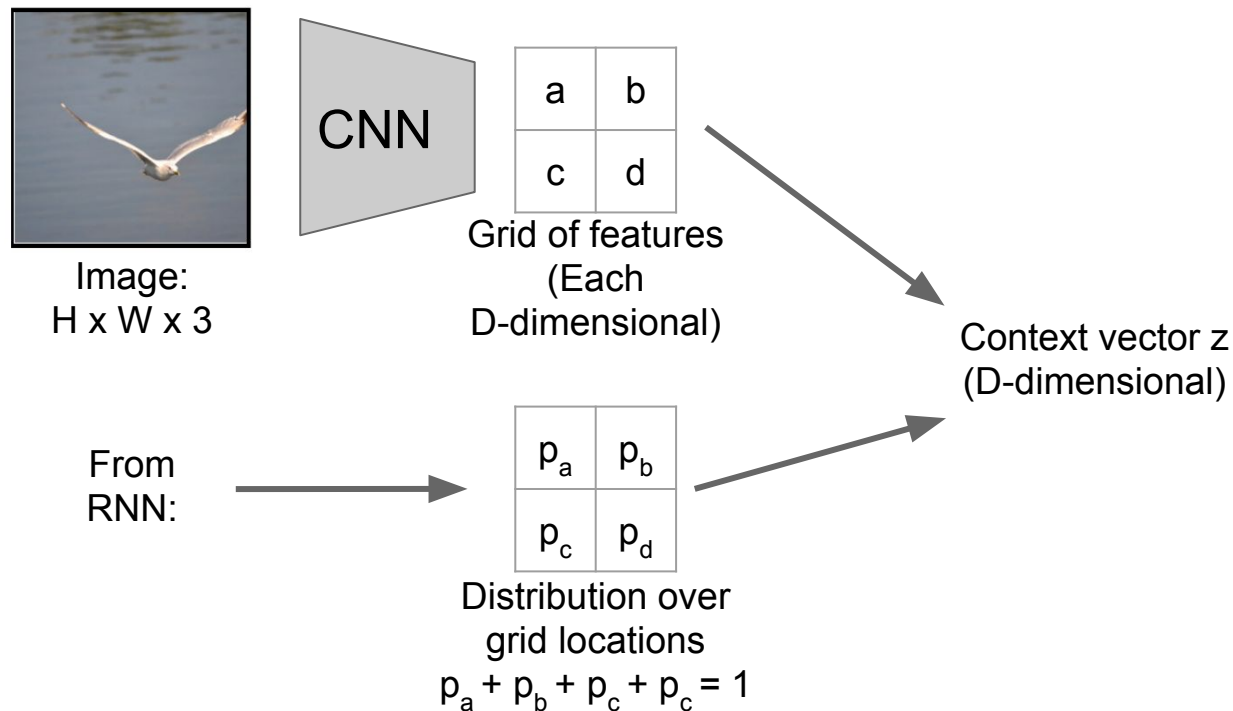


Attention with Sentinel

# Attention for Image Captioning



# Soft Attention



## Soft attention:

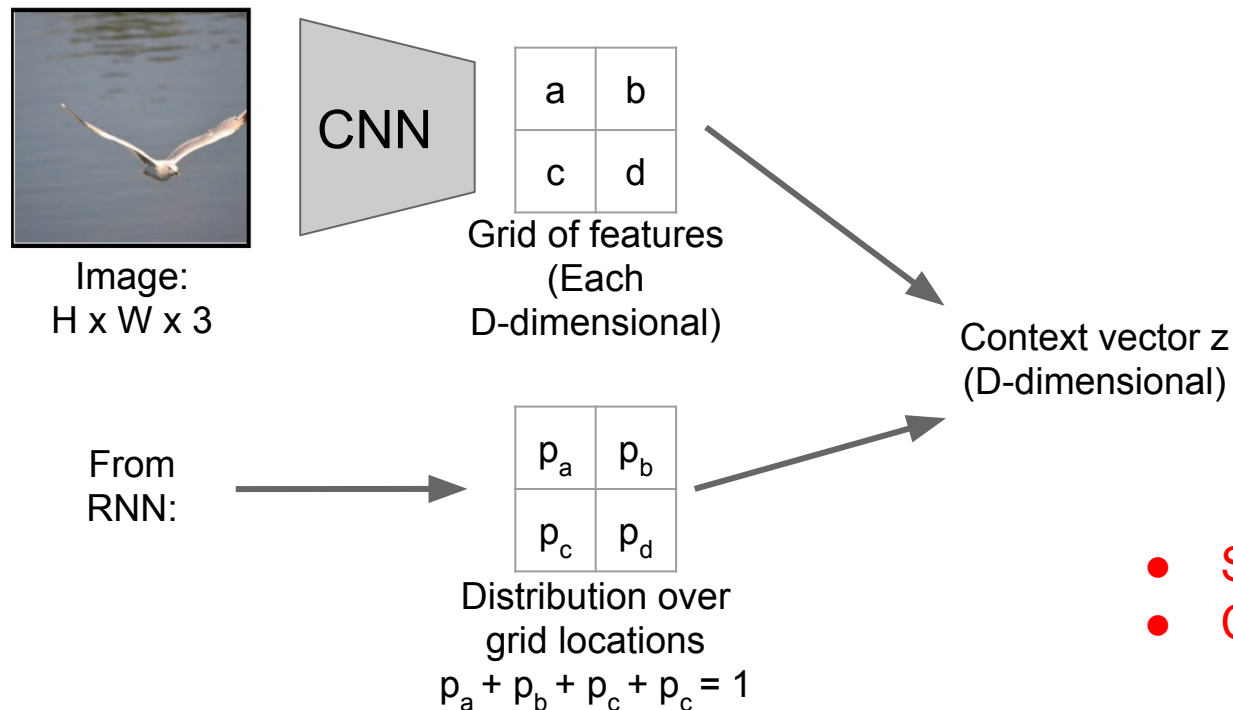
Summarize ALL locations

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

Derivative  $dz/dp$  is nice!

Train with gradient descent

# Soft Attention



## Soft attention:

Summarize ALL locations

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

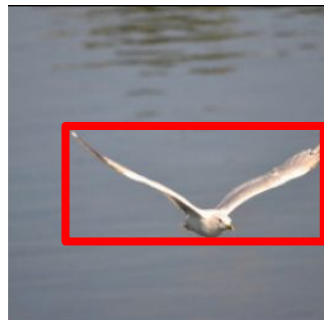
Differentiable function

Train with gradient descent

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



# Hard Attention

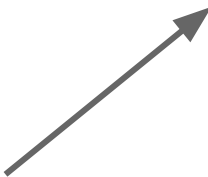


Input image:  
 $H \times W \times 3$

Box Coordinates:  
(xc, yc, w, h)



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



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



Not a differentiable function !

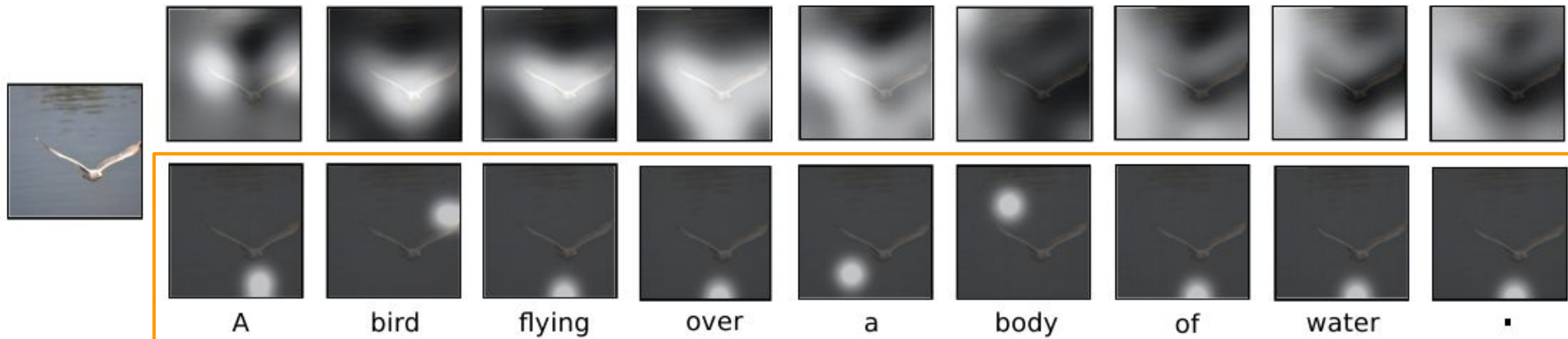


Can't train with backprop :(

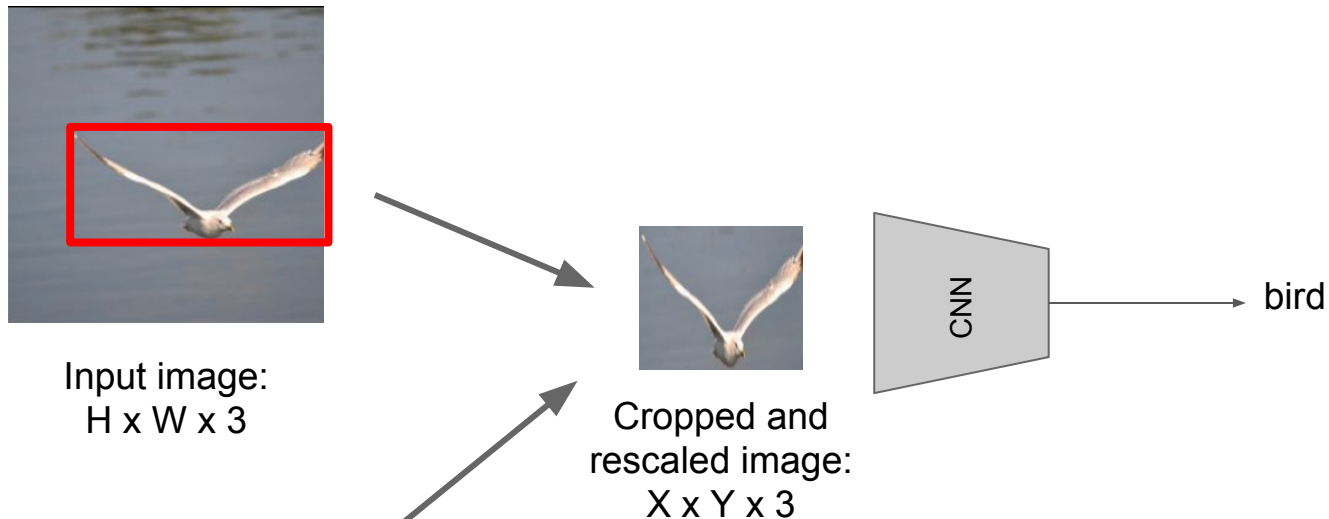


Need other optimization strategies  
e.g.: [reinforcement learning](#)

# Attention for Image Captioning



# Hard Attention



Not a differentiable function !



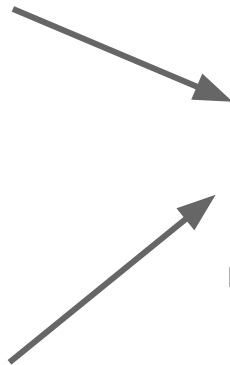
Can't train with backprop :(

# Spatial Transformer Networks

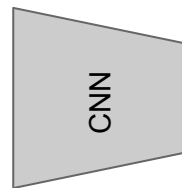


Input image:  
 $H \times W \times 3$

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



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



bird

Not a differentiable function !



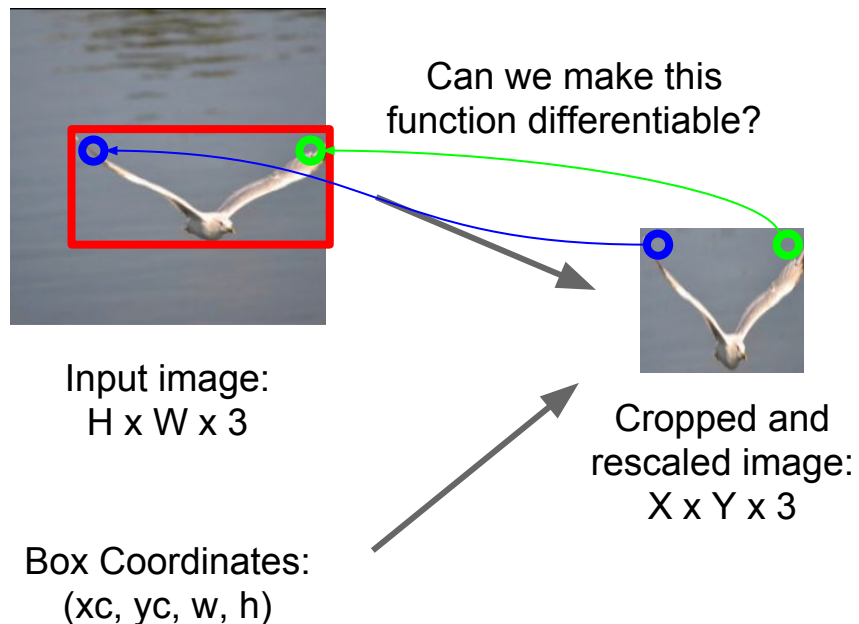
Can't train with backprop :(

Make it differentiable



Train with backprop :) 36

# Spatial Transformer Networks

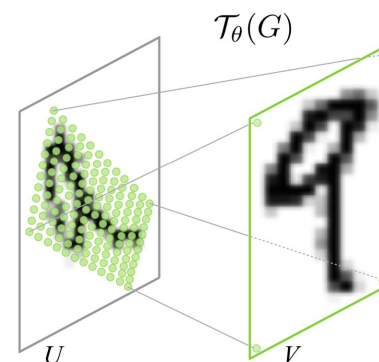


Mapping given by box coordinates  
(translation + scale)

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

Network  
attends to  
input by  
predicting  $\theta$

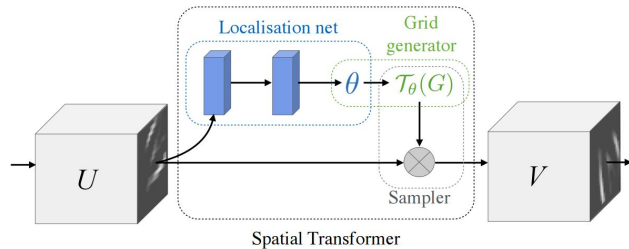
$$\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \begin{bmatrix} \theta_{11} & \theta_{12} & \theta_{13} \\ \theta_{21} & \theta_{22} & \theta_{23} \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix}$$



Repeat for all pixels  
in *output*

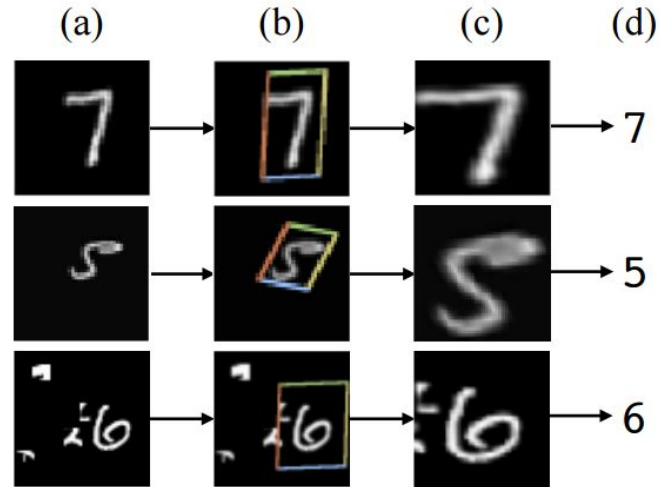
# Spatial Transformer Networks

Differentiable module

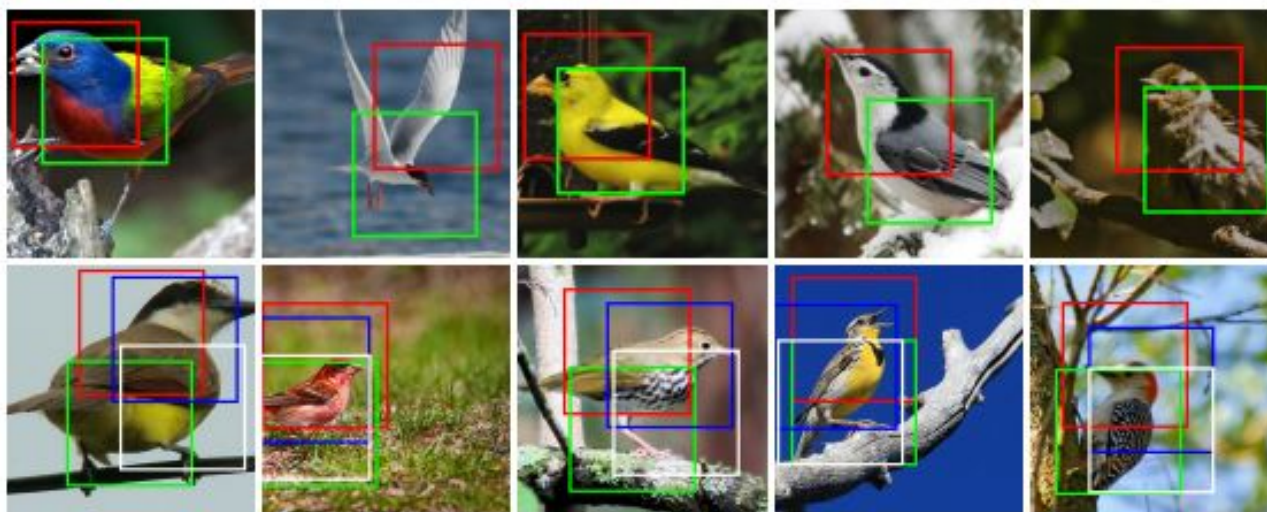


Easy to incorporate in any network, anywhere !

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



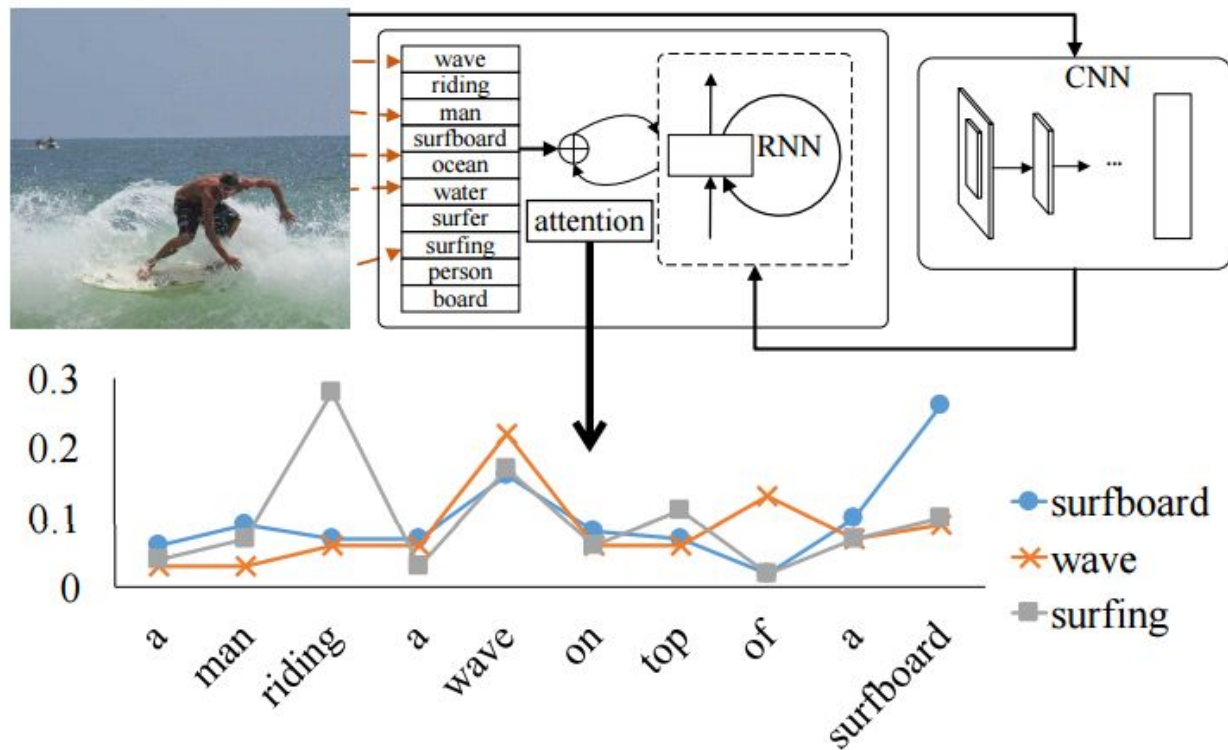
# Spatial Transformer Networks



Fine-grained classification

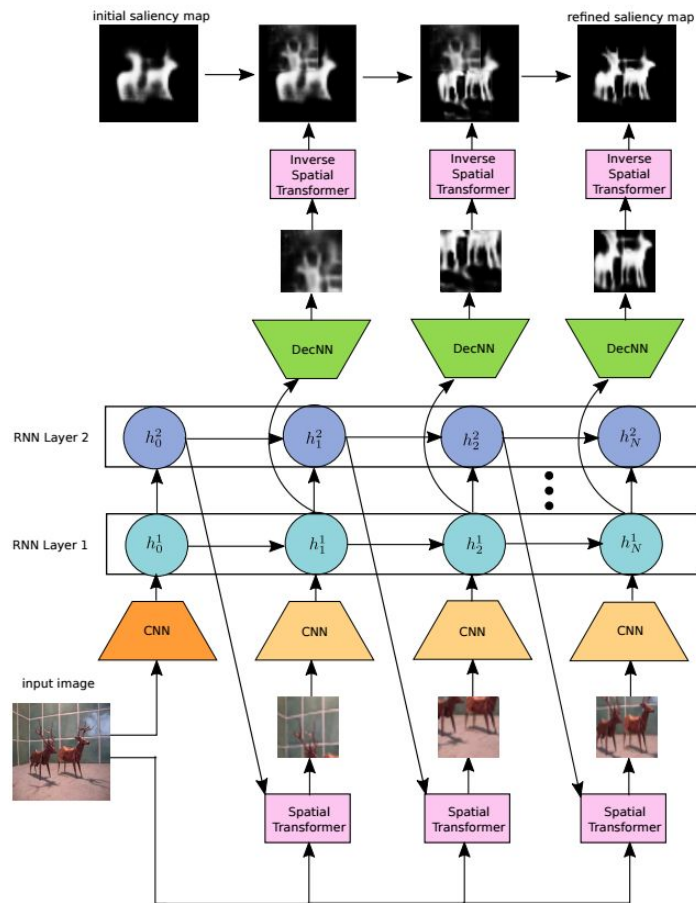
Also used as an alternative to RoI pooling in proposal-based detection & segmentation pipelines

# Semantic Attention: Image Captioning

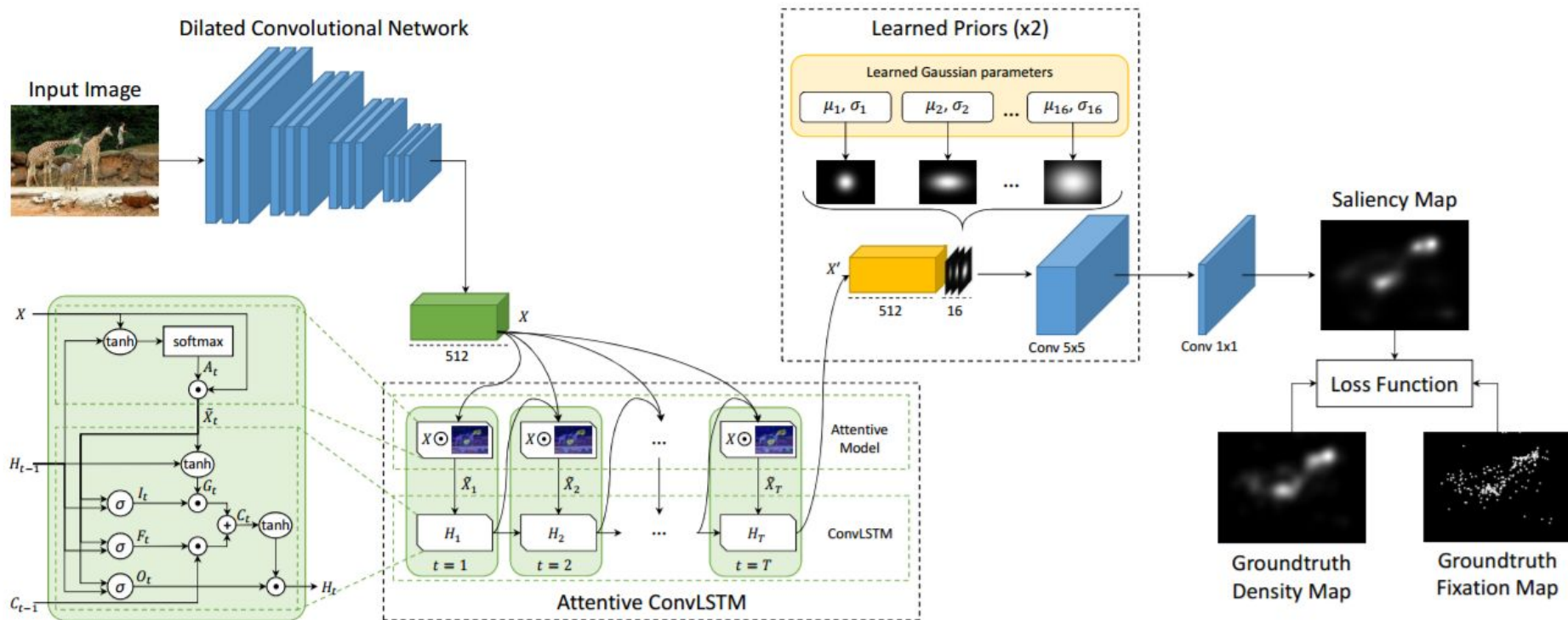




# Visual Attention: Saliency Detection

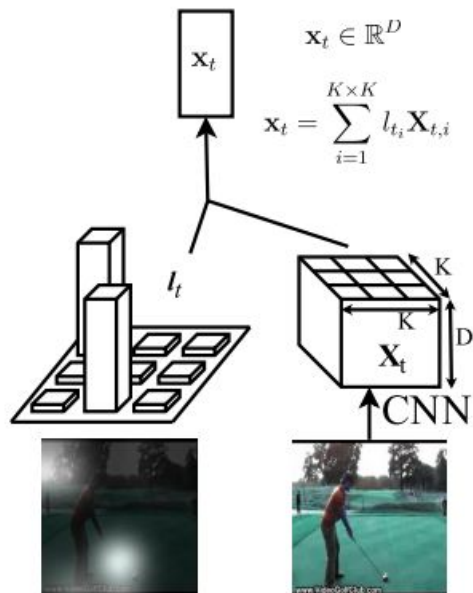


# Visual Attention: Fixation Prediction

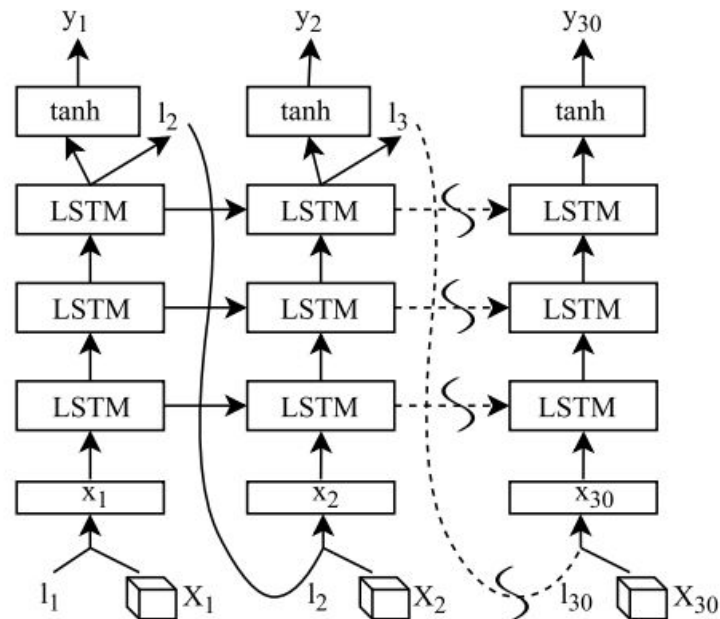


Cornia et al. [Predicting Human Eye Fixations via an LSTM-based Saliency Attentive Model.](#)

# Visual Attention: Action Recognition



(a) The soft attention mechanism



(b) Our recurrent model

# Deformable Convolutions

## Dynamic & learnable receptive field

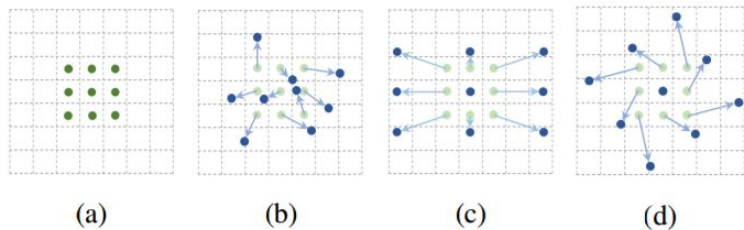
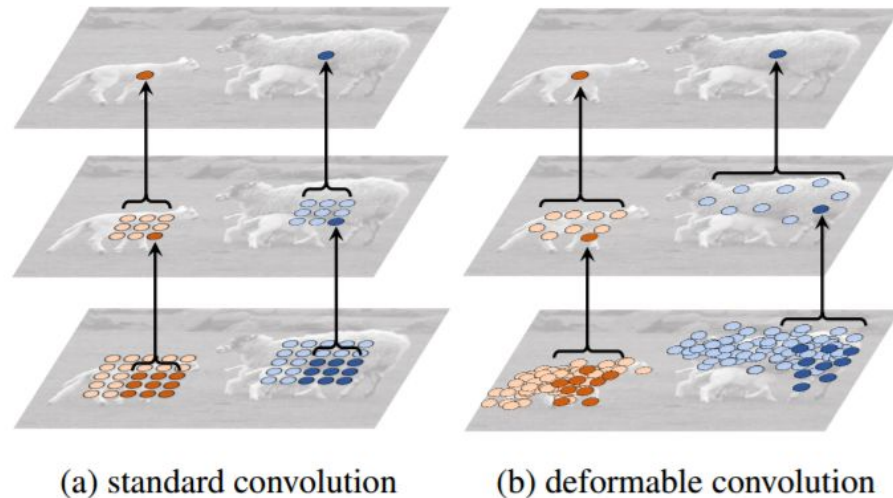


Figure 1: Illustration of sampling grids in  $3 \times 3$  regular and deformable convolutions. (a) regular sampling grid (green points) of standard convolution. (b) deformed sampling grid with augmented offsets (blue arrows) in deformable convolution. (c)(d) are special cases of (b), showing that the deformable convolution generalizes scale, aspect ratio and rotation transformations.



Dai, Qi, Xiong, Li, Zhang et al. [Deformable Convolutional Networks](#). arXiv Mar 2017

Questions?