

LECTURE 3: ATTENTION IN MACHINE LEARNING

CSC5991: Introduction to LLMs

OUTLINE



What is attention?



RNN + Attention in machine
translation (NLP)



CNN architecture

Attention in
CNNs



RNN + Attention in image
captioning (CV)

WHAT IS ATTENTION?

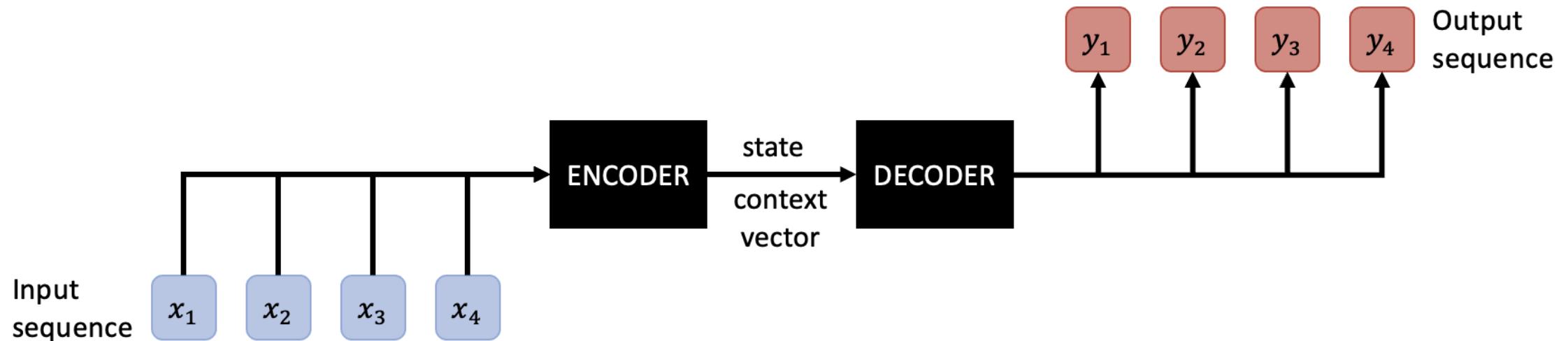
- The notion of **exploiting context** is not new
 - **RNN** – context from **temporal locality** (useful for sequences/time-series data)
 - **CNN** – context from **spatial locality** (useful for images)
 - Embedding priors into models forces them to pay “attention” to relevant features for a given problem
- What we now call “attention” in DL
 - The idea of paying “attention” to the most relevant parts of I/P
 - Very useful in **sequence-to-sequence** modelling
 - Ideally, we’d like to **learn** this!

WHAT IS A LEARNED ATTENTION MECHANISM?

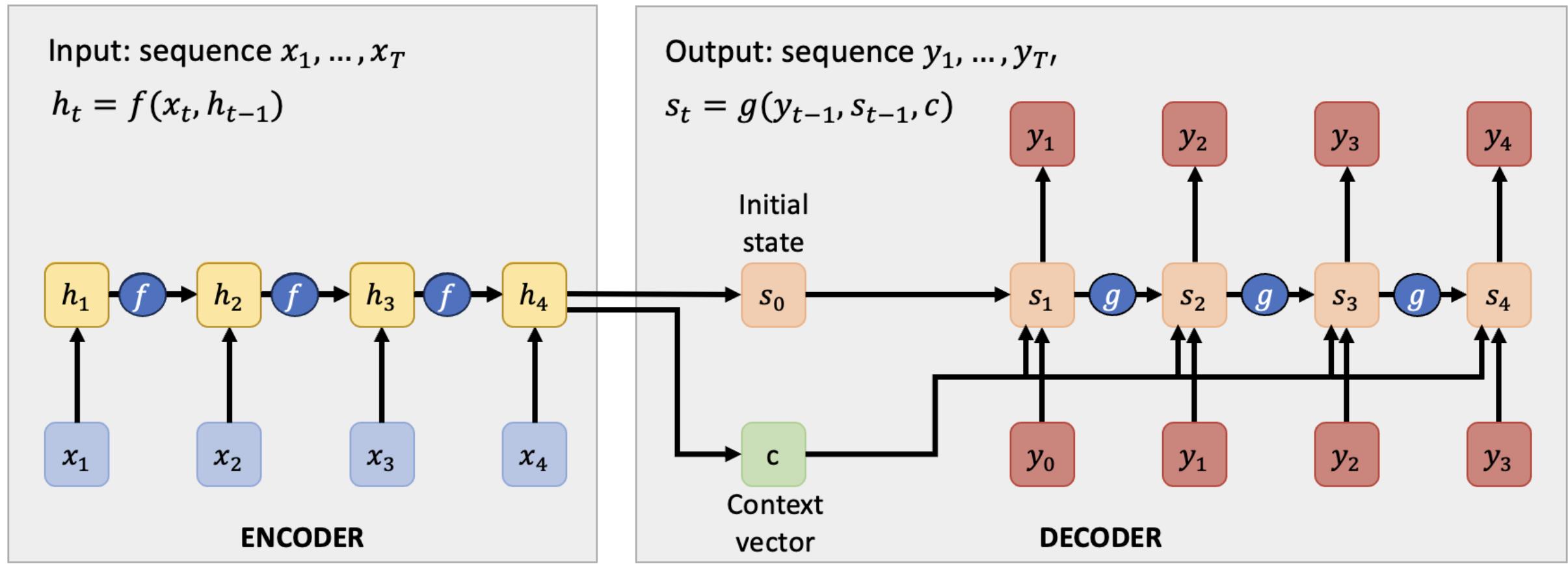
- An attention mechanism typically refers to **function** that allows a model to attend to different content
- There are many forms of attention mechanisms
 - **Additive**
 - **Dot-product**
- We have names to distinguish attention based on what is attended to
 - **Self-attention** (intra attention)
 - **Cross-attention** (encoder-decoder attention/inter attention)

SEQ2SEQ: RNNs + ATTENTION

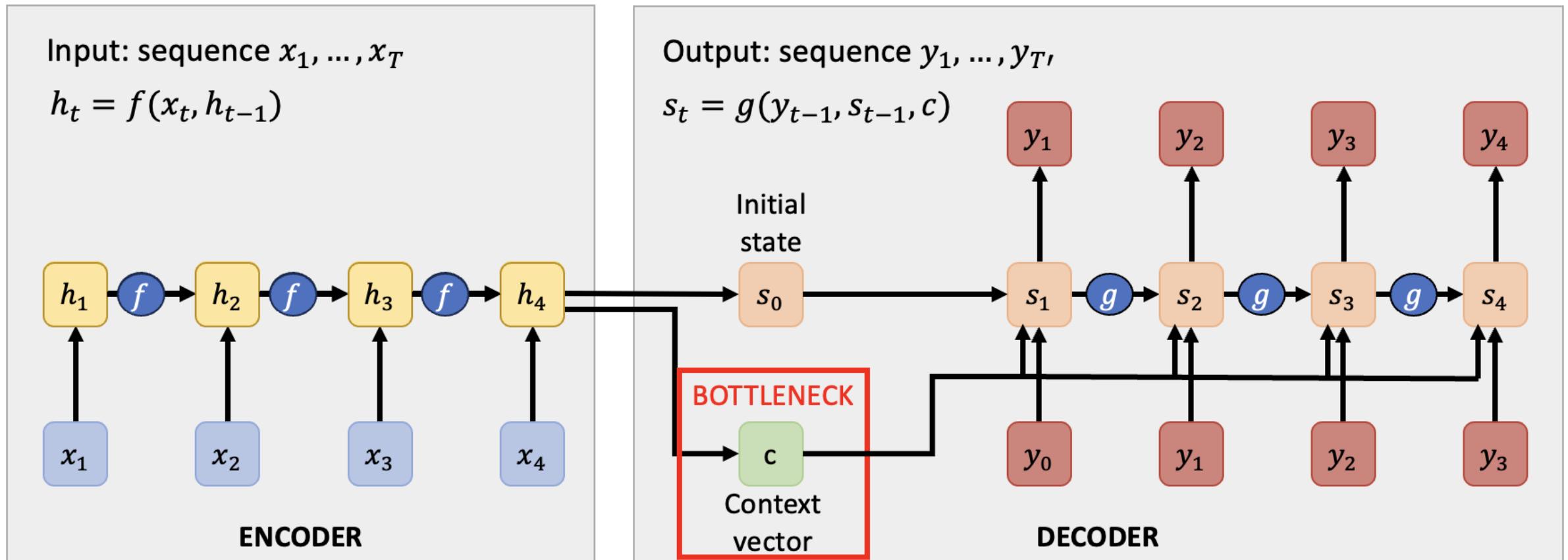
- **Example Scenarios**
 - Text → Text (e.g. Q/A, translation, text summarization)
 - Image → Text (e.g. image captioning)
- **How?** Usually Encoder-Decoder models
 - e.g. RNNs, transformers



FIXED CONTEXT



BOTTLENECK



“you can’t cram the meaning
of a whole %&@#&ing
sentence into a single
\$*&@ing vector!”

– Ray Mooney (NLP professor at UT Austin)

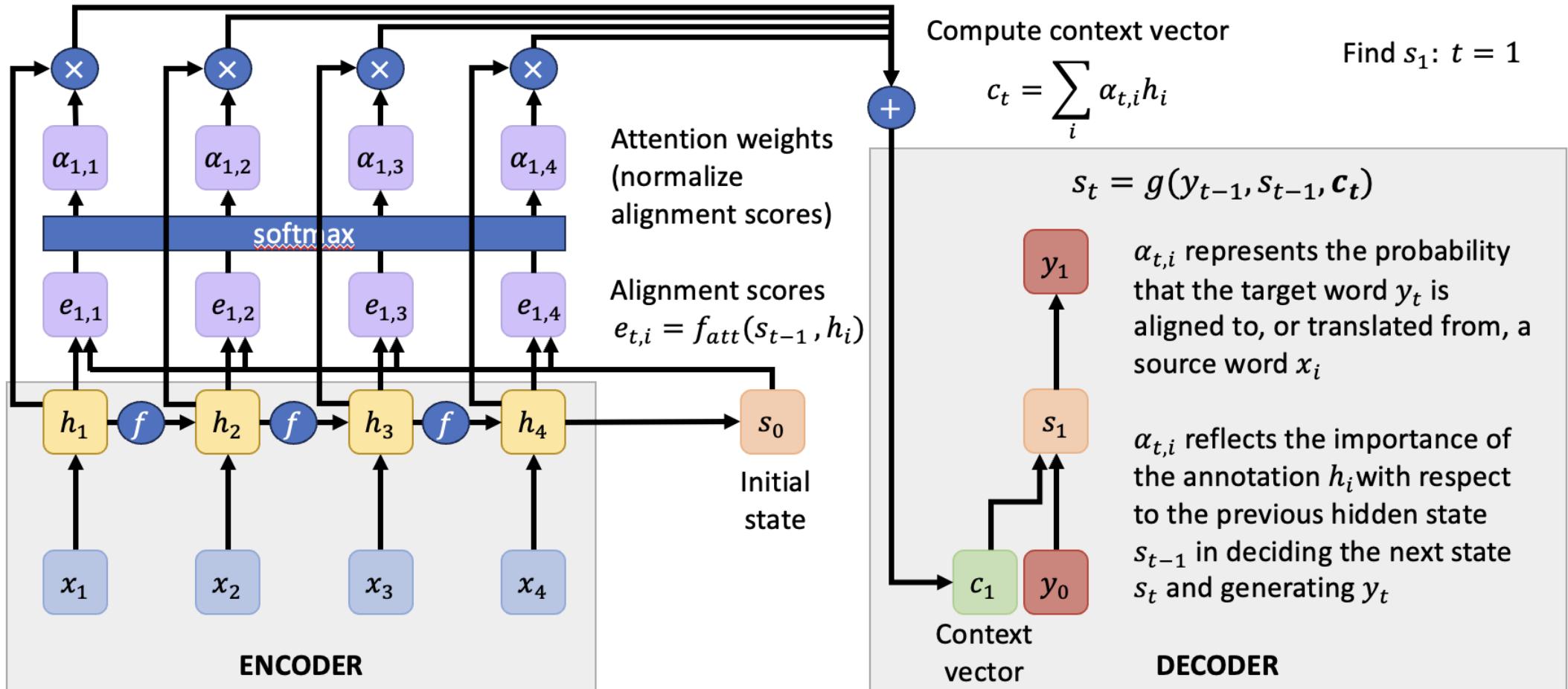
DESIGN CONTEXT USING ATTENTION

- **Idea!** Use a different context vector for each timestep in the decoder

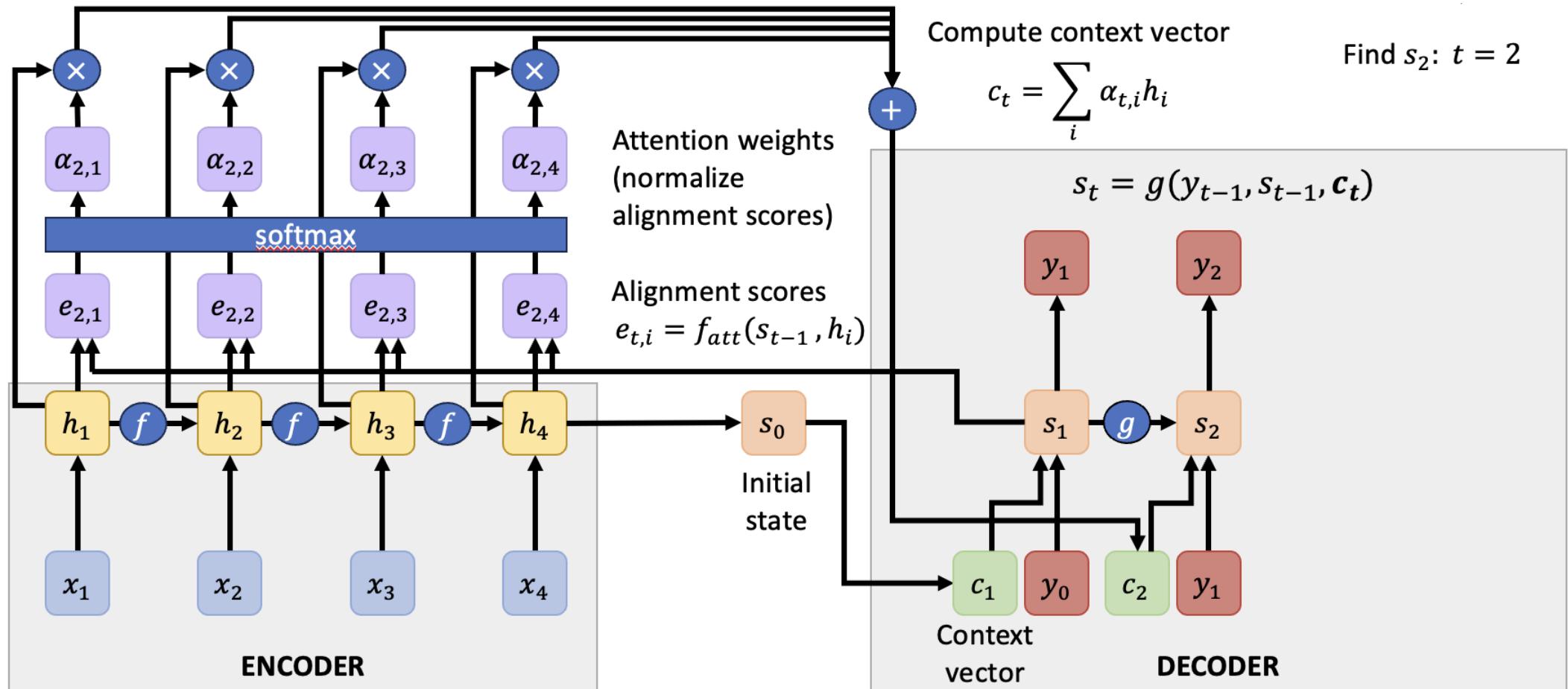
$$s_t = g(y_{t-1}, s_{t-1}, \mathbf{c}_t)$$

- No more bottleneck through a single vector
- Craft the context vector so that it “looks at” different parts of the input sequence for each decoder timestep

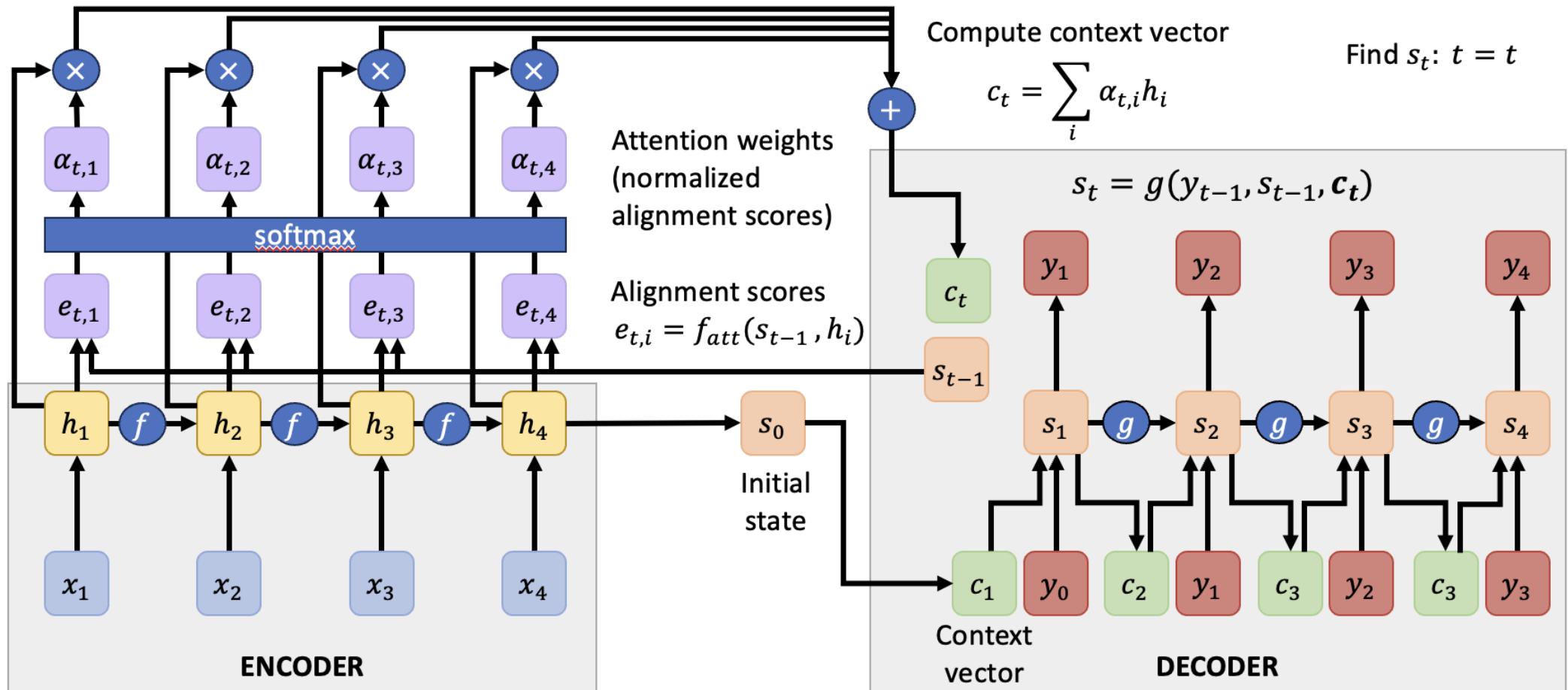
SEQ2SEQ: RNNs + ATTENTION



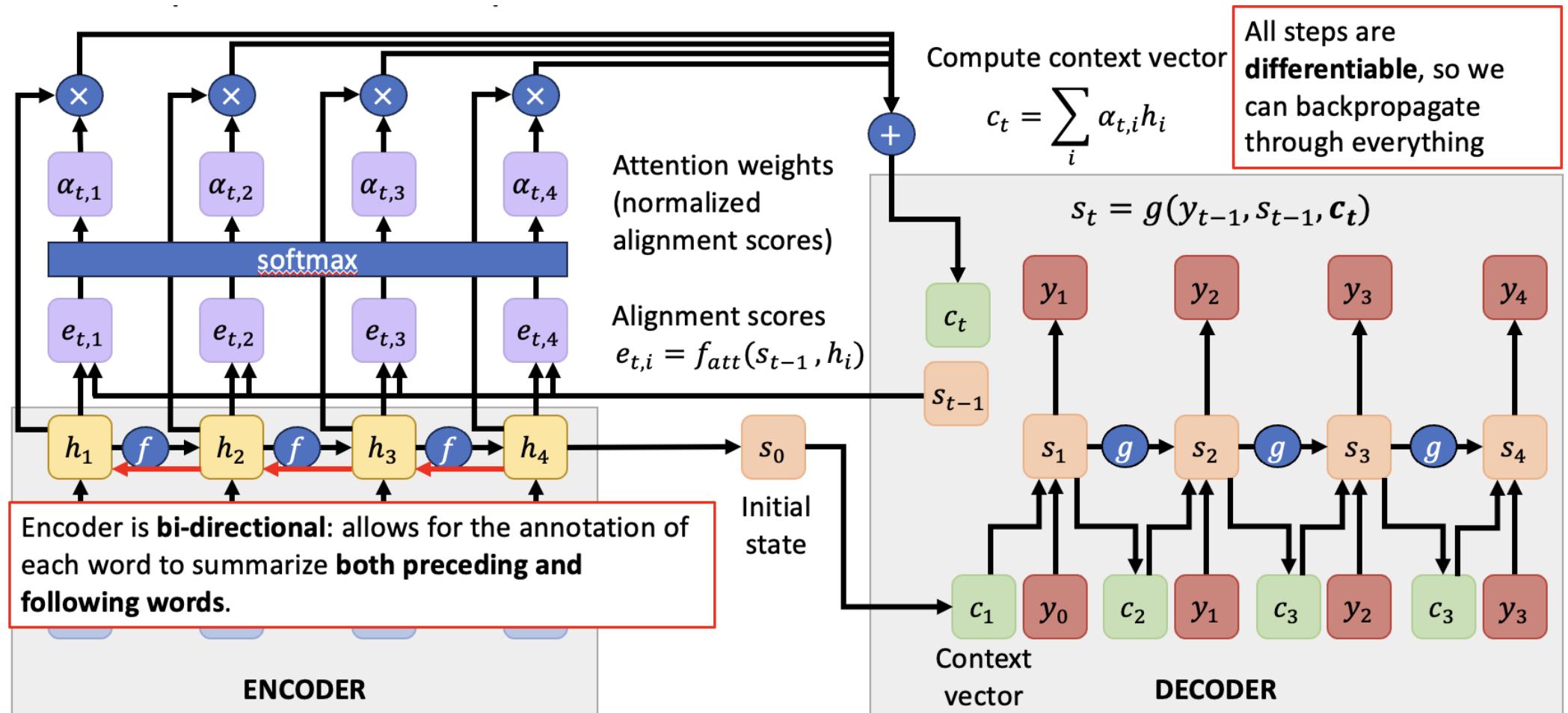
SEQ2SEQ: RNNs + ATTENTION



SEQ2SEQ: RNNs + ATTENTION



SEQ2SEQ: RNNs + ATTENTION



ATTENTION MAP

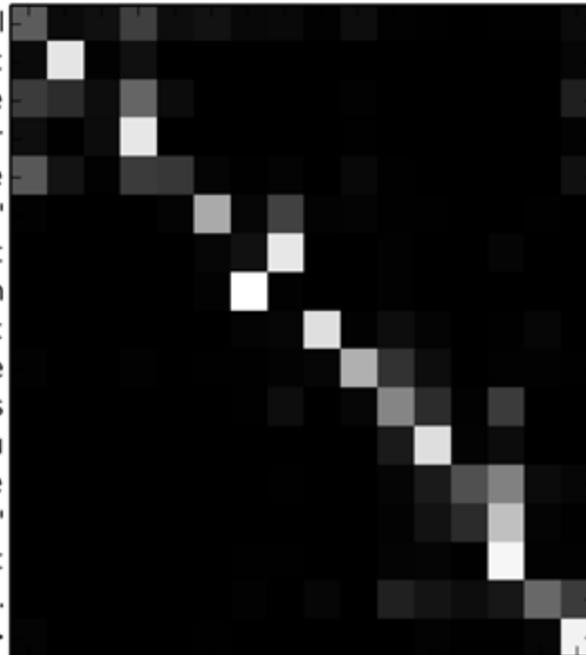
The
agreement
on
the
European
Economic
Area
was
signed
in
August
1992
. <end>



This attention map visualizes the importance of words in the English sentence. The words are represented by white and light gray squares, while other tokens like punctuation and end-of-sequence markers are black. The words 'The', 'agreement', 'on', 'the', 'European', 'Economic', 'Area', 'was', 'signed', 'in', 'August', '1992', and the final period and end marker are highlighted.

L'
accord
sur
la
zone
économique
européenne
a
été
signé
en
août
1992
. <end>

Il
convient
de
noter
que
l'
environnement
marin
est
le
moins
connu
de
l'
environnement
. <end>



This attention map visualizes the importance of words in the French sentence. The words are represented by white and light gray squares, while other tokens like punctuation and end-of-sequence markers are black. The words 'L', 'accord', 'sur', 'la', 'zone', 'économique', 'européenne', 'a', 'été', 'signé', 'en', 'août', '1992', the final period, and the end marker are highlighted. The second part of the sentence, starting with 'Il convient de noter que ...', also shows high attention weights.

CONVOLUTIONAL NEURAL NETWORKS (CNNs)

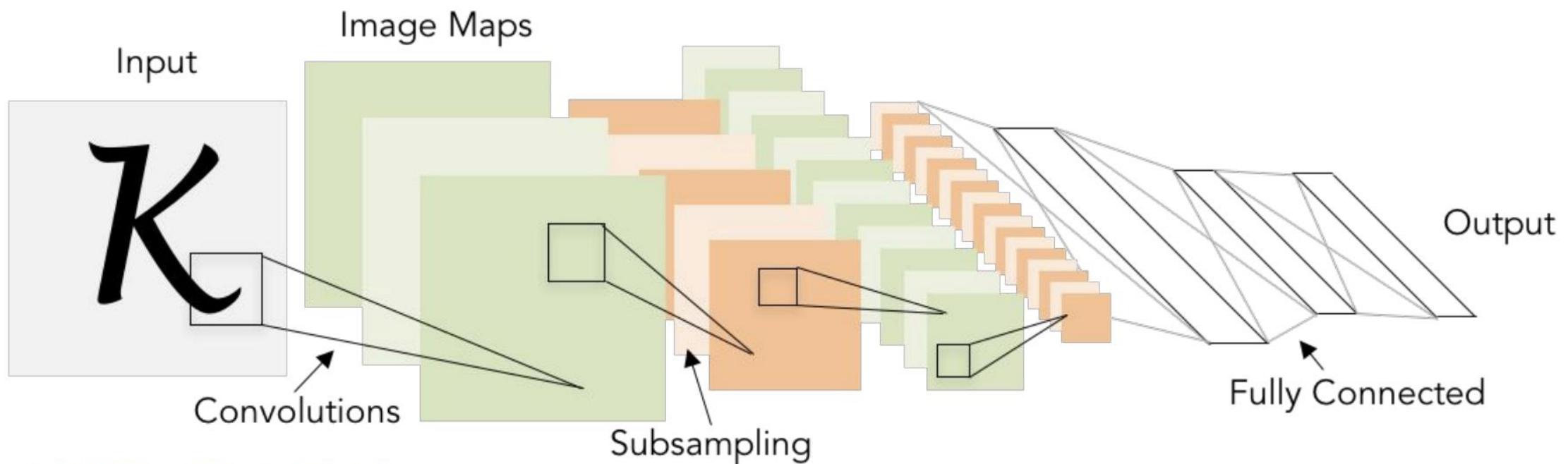
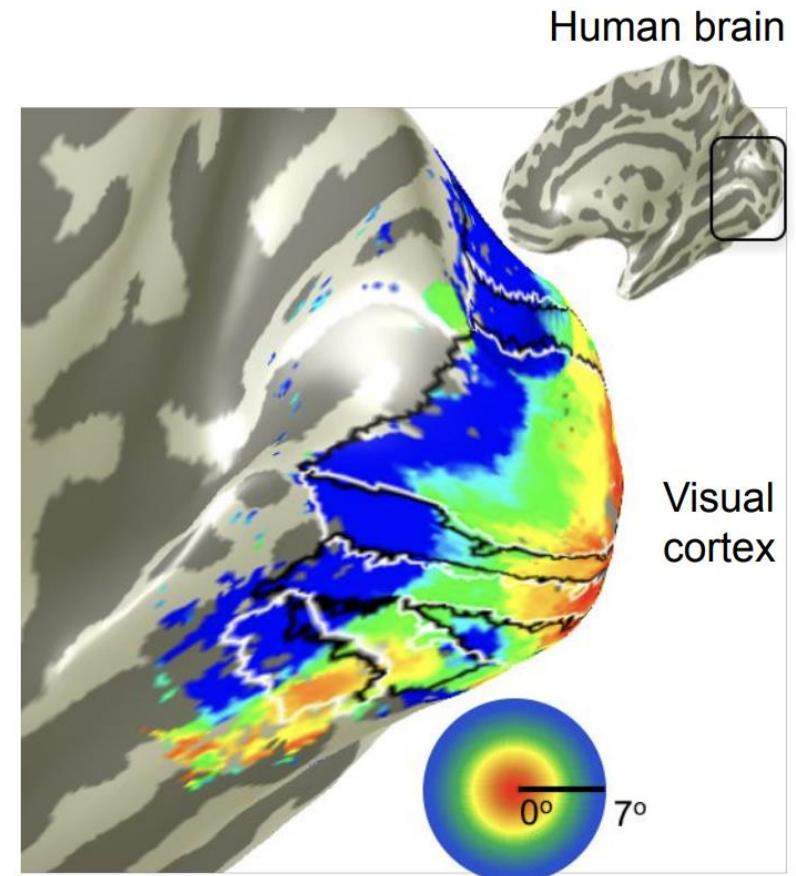
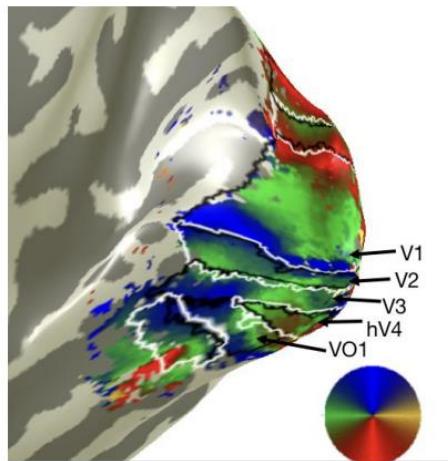


Illustration of LeCun et al. 1998 from CS231n 2017 Lecture 1

A bit of history

Topographical mapping in the cortex:
nearby cells in cortex represent
nearby regions in the visual field



Retinotopy images courtesy of Jesse Gomez in the Stanford Vision & Perception Neuroscience Lab.

Hierarchical organization

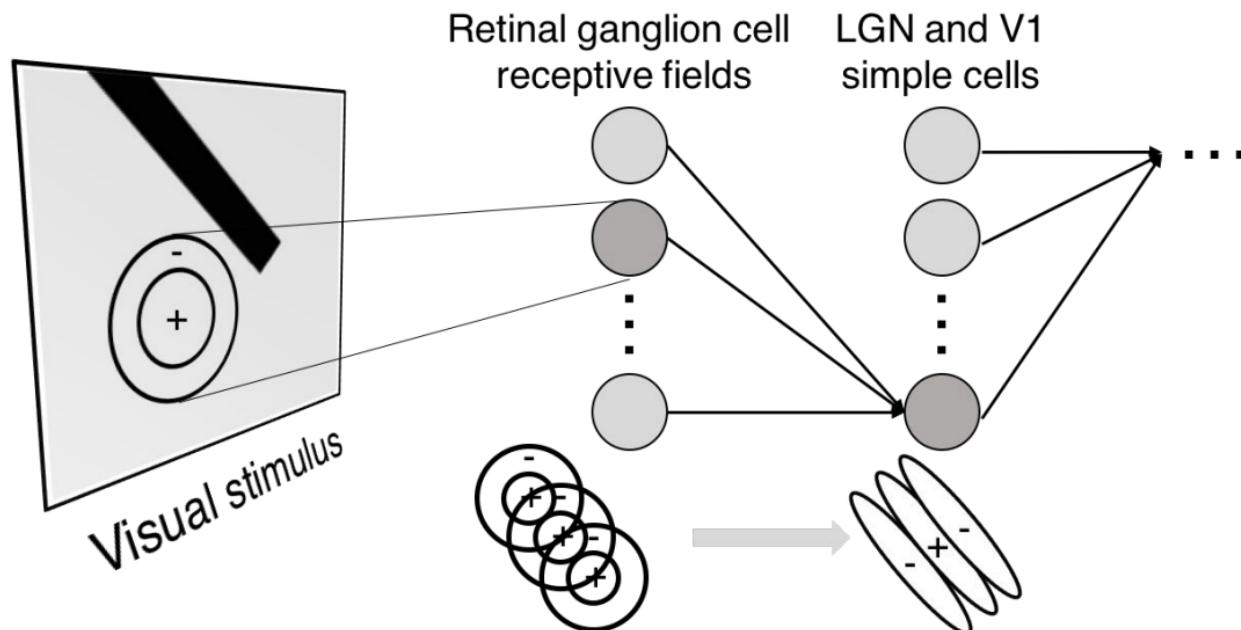
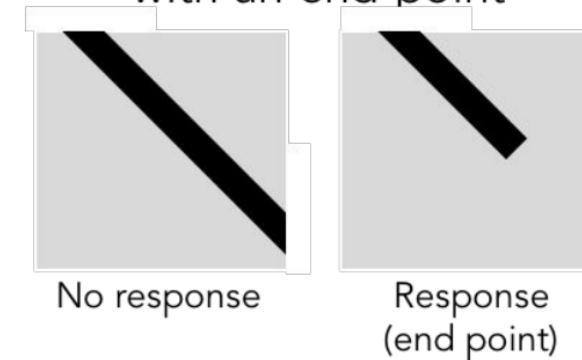


Illustration of hierarchical organization in early visual pathways by Lane McIntosh, copyright CS231n 2017

Simple cells:
Response to light orientation

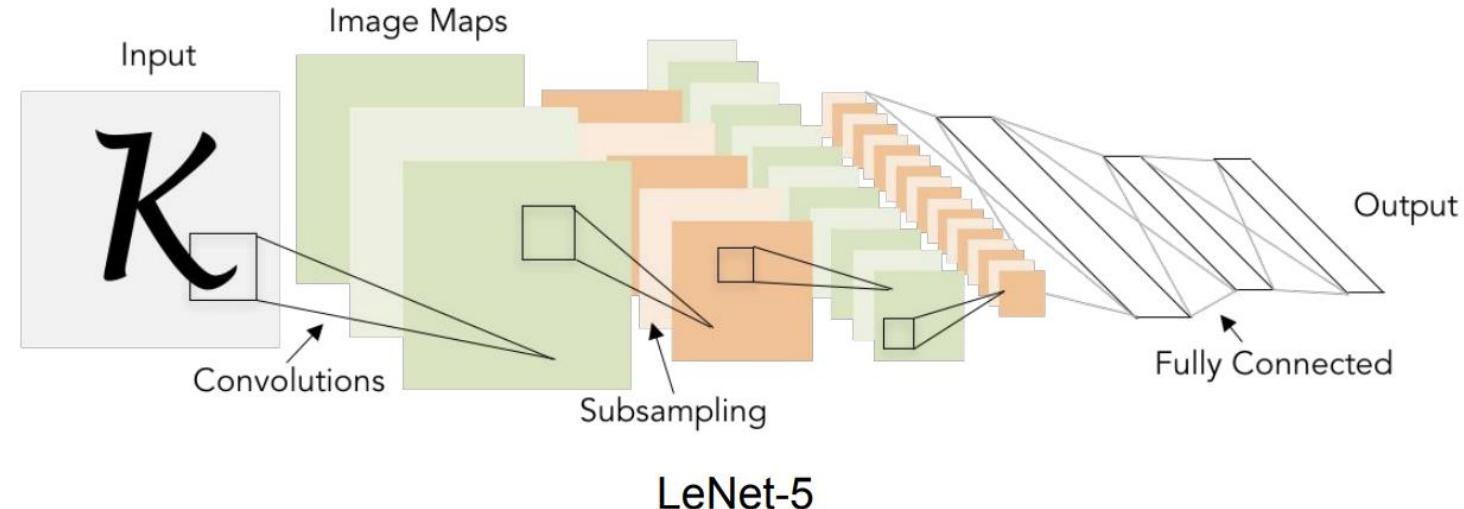
Complex cells:
Response to light orientation and movement

Hypercomplex cells:
response to movement with an end point



A bit of history: **Gradient-based learning applied to document recognition**

[LeCun, Bottou, Bengio, Haffner 1998]



A bit of history: **ImageNet Classification with Deep Convolutional Neural Networks**

[Krizhevsky, Sutskever, Hinton, 2012]

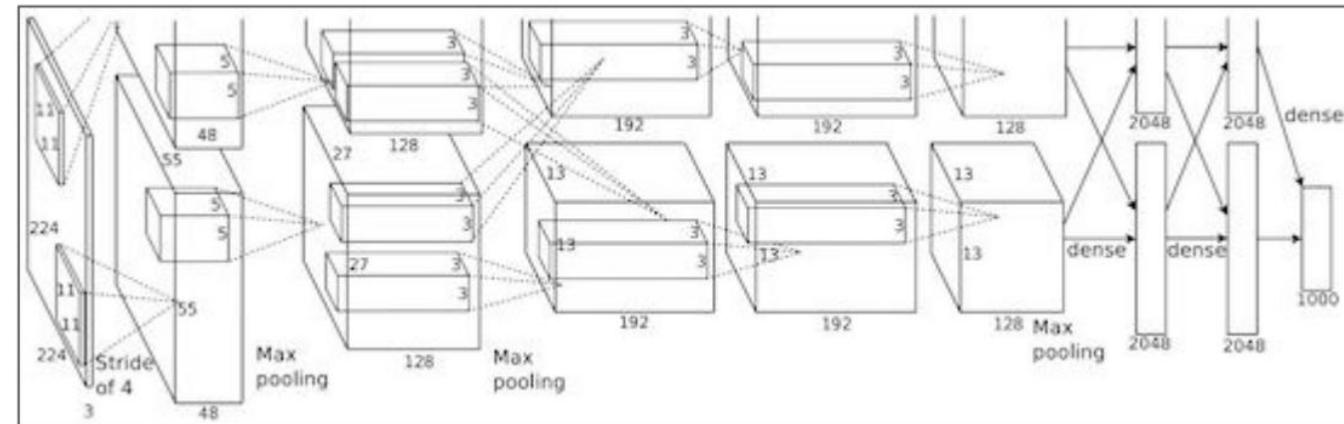


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

“AlexNet”

Fast-forward to today: ConvNets are everywhere

Classification



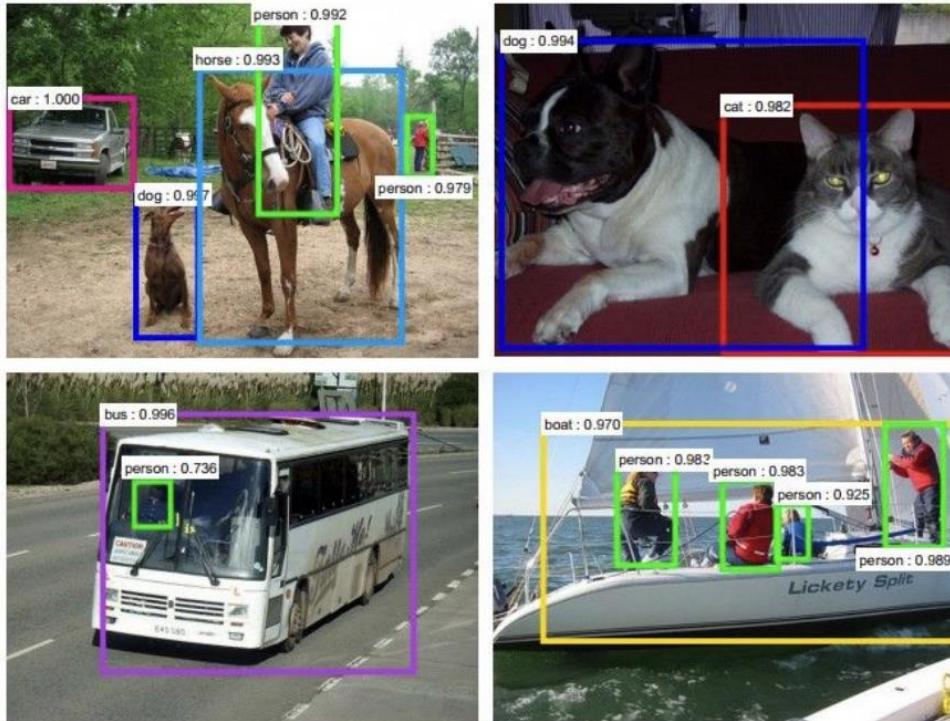
Retrieval



Figures copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

Fast-forward to today: ConvNets are everywhere

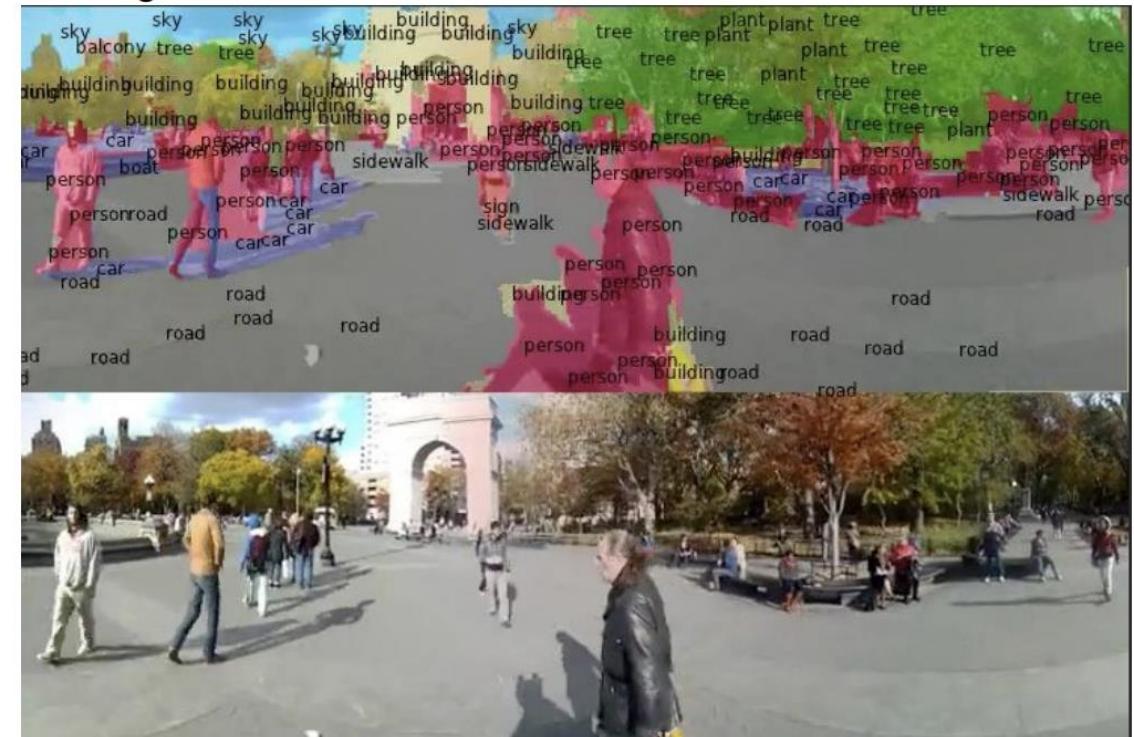
Detection



Figures copyright Shaoqing Ren, Kaiming He, Ross Girshick, Jian Sun, 2015. Reproduced with permission.

[Faster R-CNN: Ren, He, Girshick, Sun 2015]

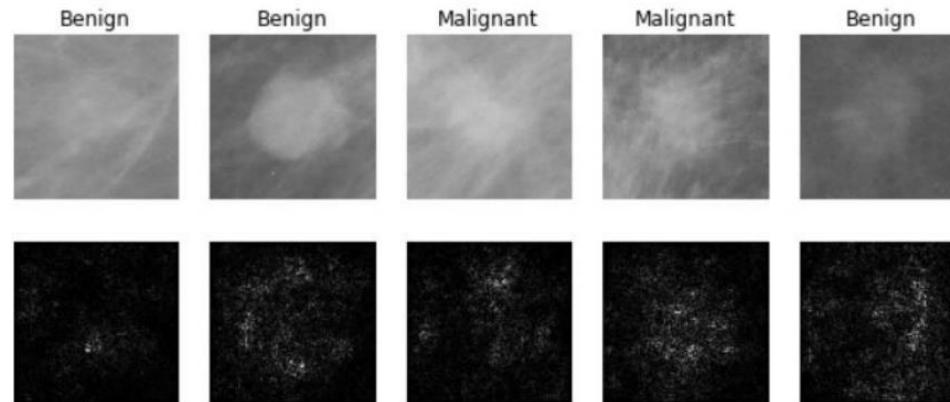
Segmentation



Figures copyright Clement Farabet, 2012.
Reproduced with permission.

[Farabet et al., 2012]

Fast-forward to today: ConvNets are everywhere



[Levy et al. 2016]

Figure copyright Levy et al. 2016.
Reproduced with permission.



[Dieleman et al. 2014]

From left to right: [public domain by NASA](#), usage [permitted](#) by
ESA/Hubble, [public domain by NASA](#), and [public domain](#).



[Sermanet et al. 2011]
[Ciresan et al.]

Photos by Lane McIntosh.
Copyright CS231n 2017.

Fast-forward to today: ConvNets are everywhere

No errors



A white teddy bear sitting in the grass



A man riding a wave on top of a surfboard

Minor errors



A man in a baseball uniform throwing a ball



A cat sitting on a suitcase on the floor

Somewhat related



A woman is holding a cat in her hand



A woman standing on a beach holding a surfboard

Image Captioning

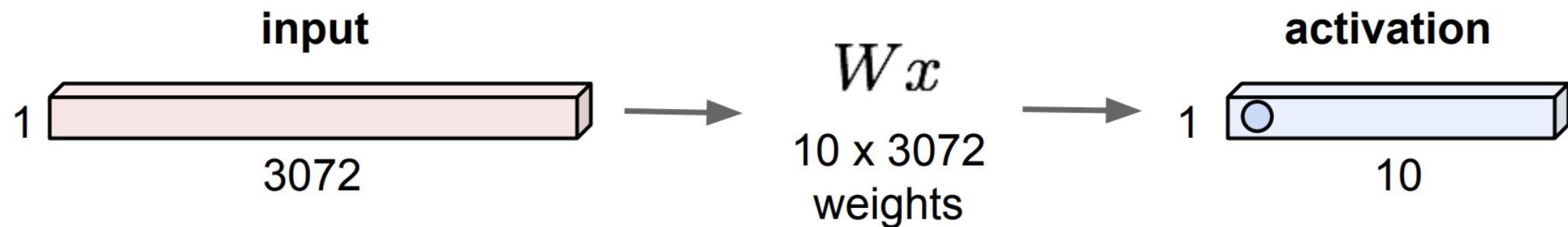
[Vinyals et al., 2015]
[Karpathy and Fei-Fei, 2015]

All images are CC0 Public domain:
<https://pixabay.com/en/luggage-antique-cat-1643010/>
<https://pixabay.com/en/teddy-plush-bears-cute-teddy-bear-1623436/>
<https://pixabay.com/en/surf-wave-summer-sport-litoral-1668716/>
<https://pixabay.com/en/woman-female-model-portrait-adult-983967/>
<https://pixabay.com/en/handstand-lake-meditation-496008/>
<https://pixabay.com/en/baseball-player-shortstop-infield-1045263/>

Captions generated by Justin Johnson using [Neuraltalk2](#)

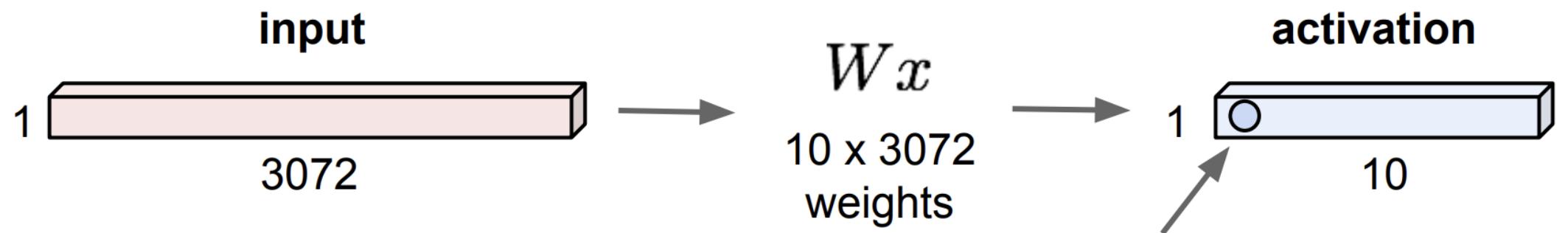
FULLY CONNECTED LAYER

32x32x3 image -> stretch to 3072 x 1



FULLY CONNECTED LAYER

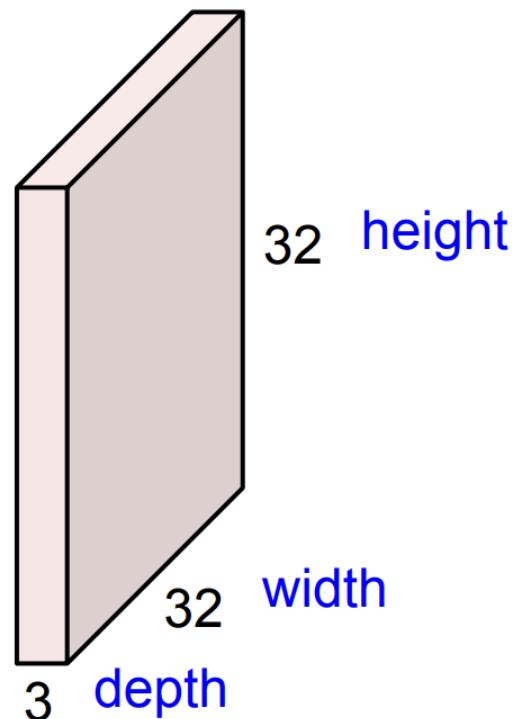
32x32x3 image -> stretch to 3072 x 1



1 number:
the result of taking a dot product
between a row of W and the input
(a 3072-dimensional dot product)

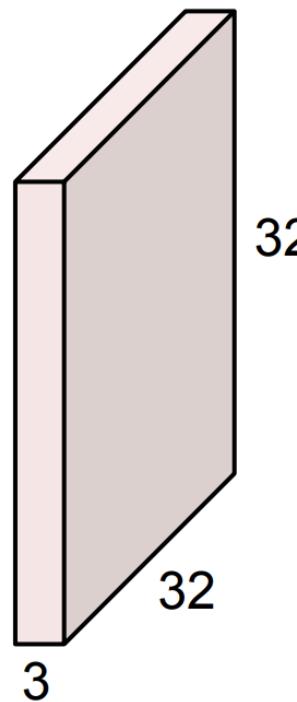
CONVOLUTION LAYER

32x32x3 image -> preserve spatial structure



CONVOLUTION LAYER

32x32x3 image

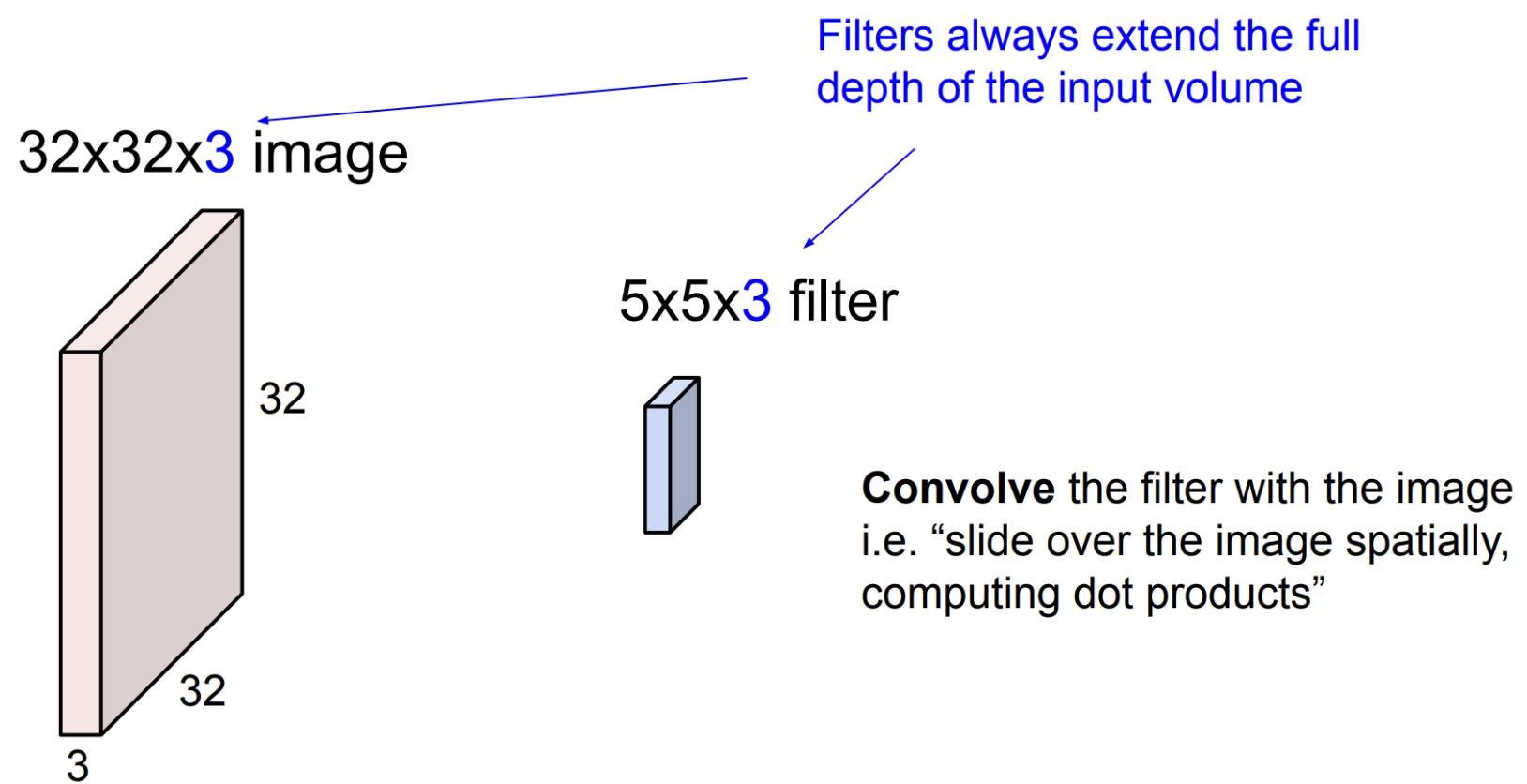


5x5x3 filter

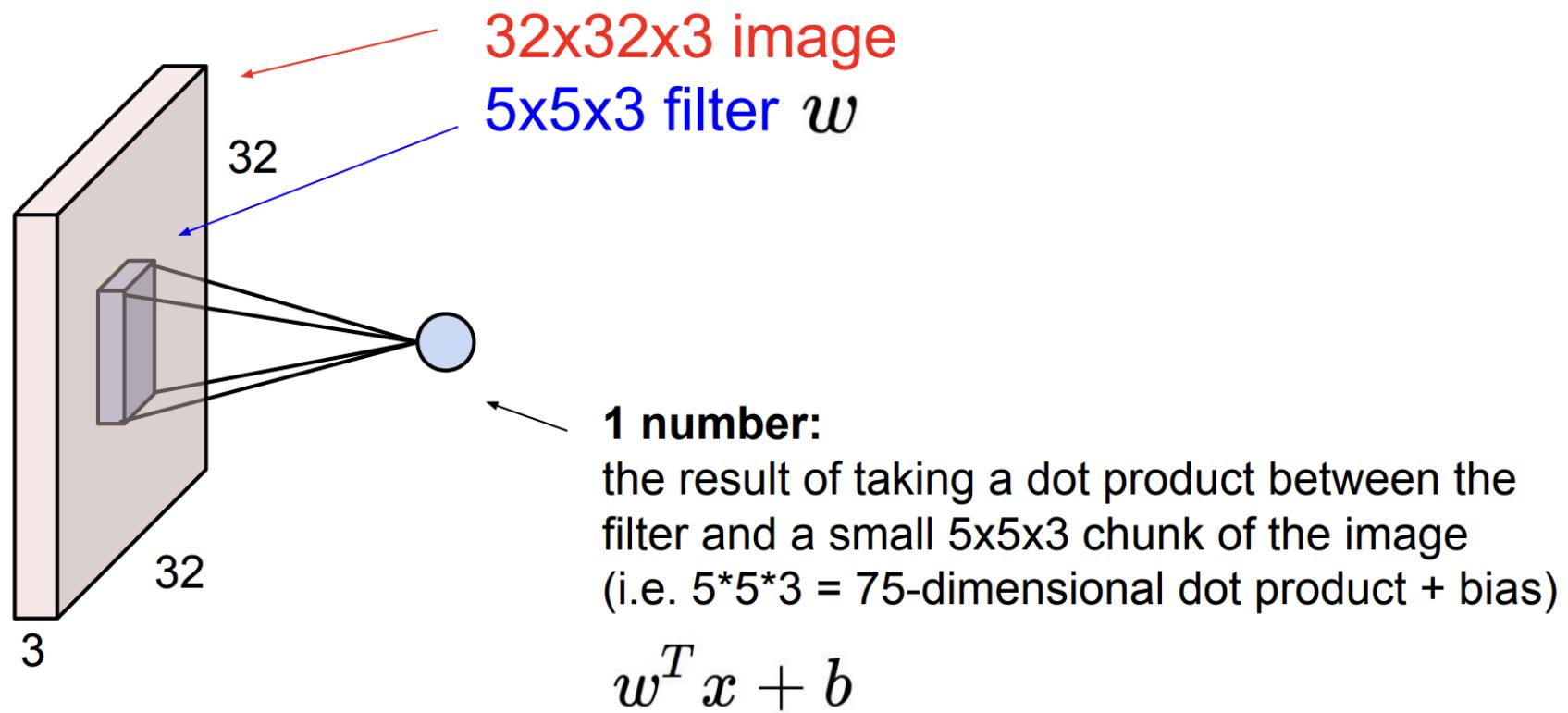


Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

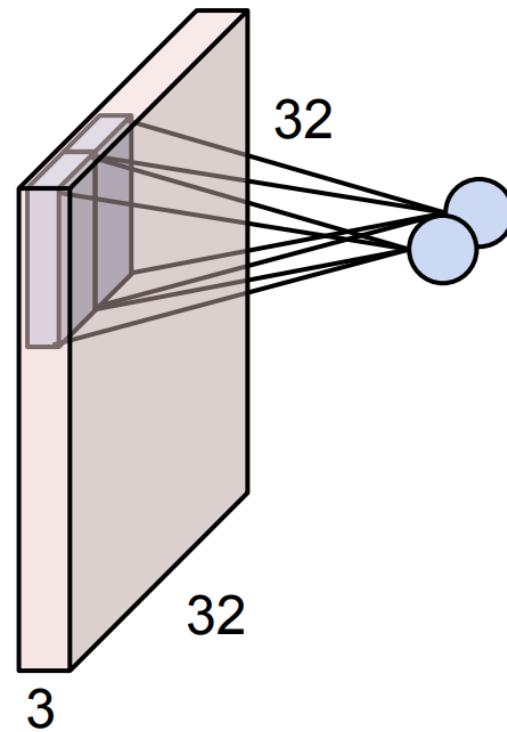
CONVOLUTION LAYER



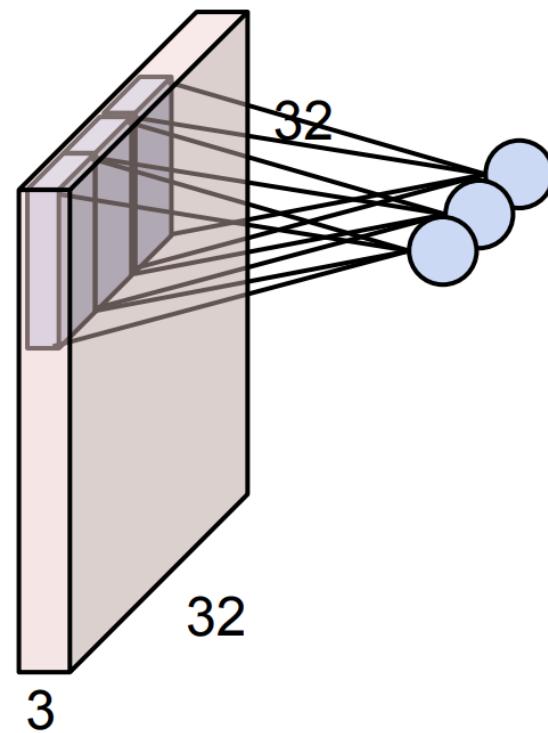
CONVOLUTION LAYER



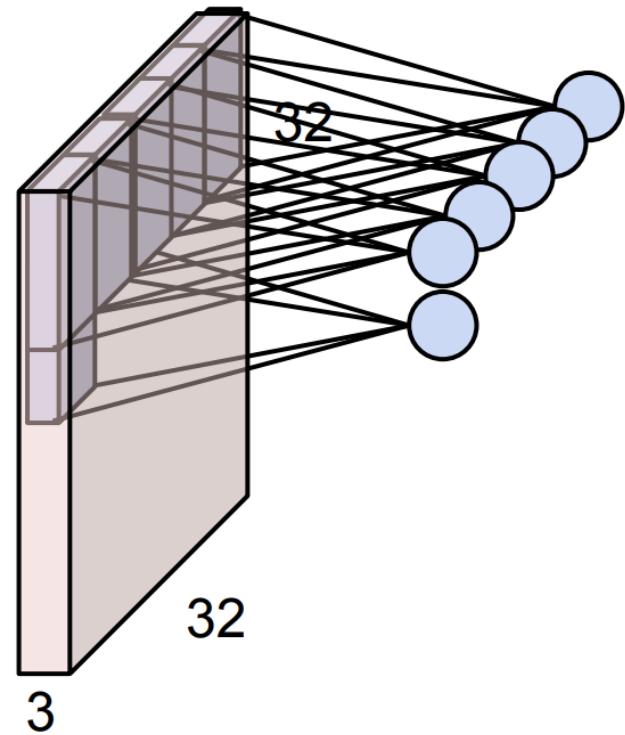
CONVOLUTION LAYER



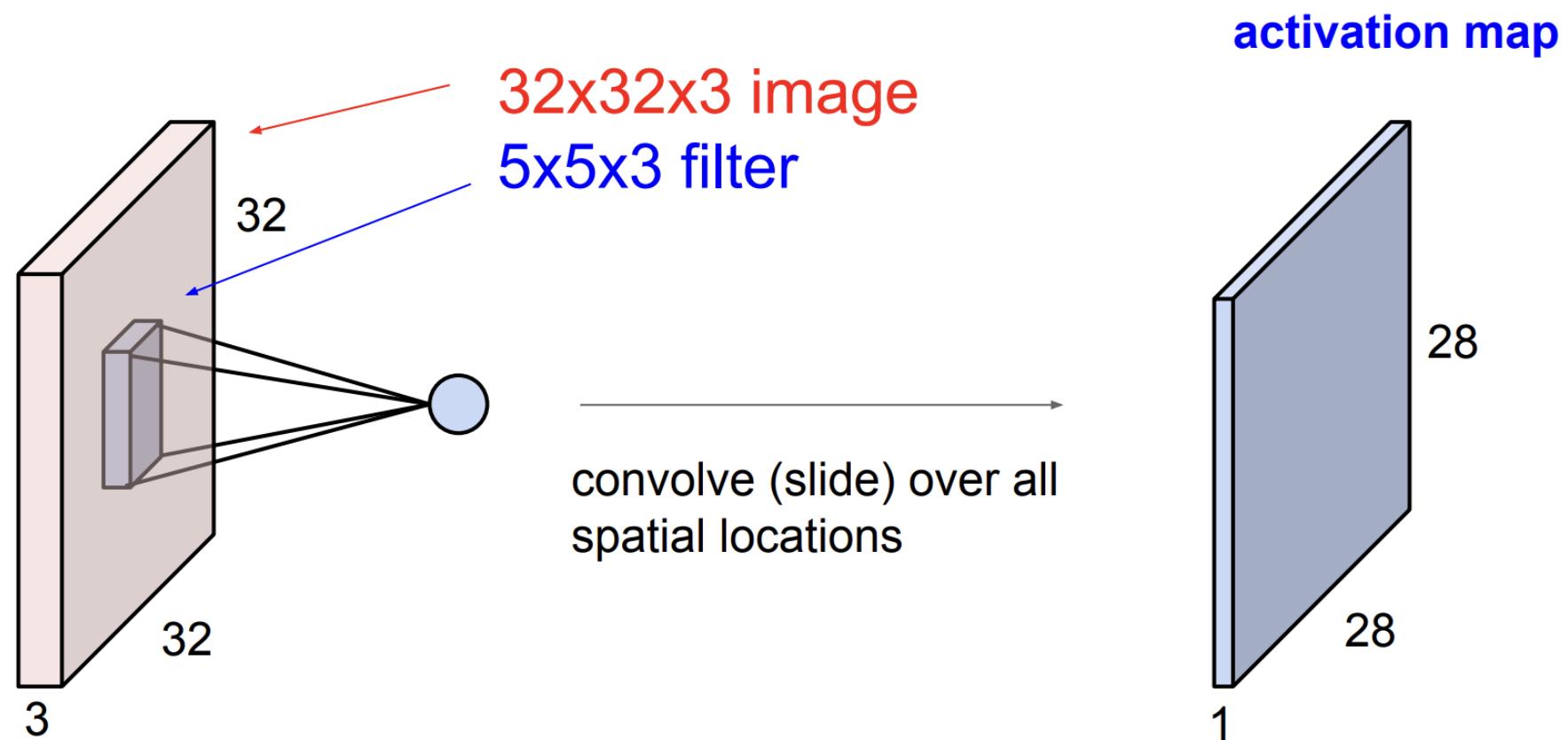
CONVOLUTION LAYER



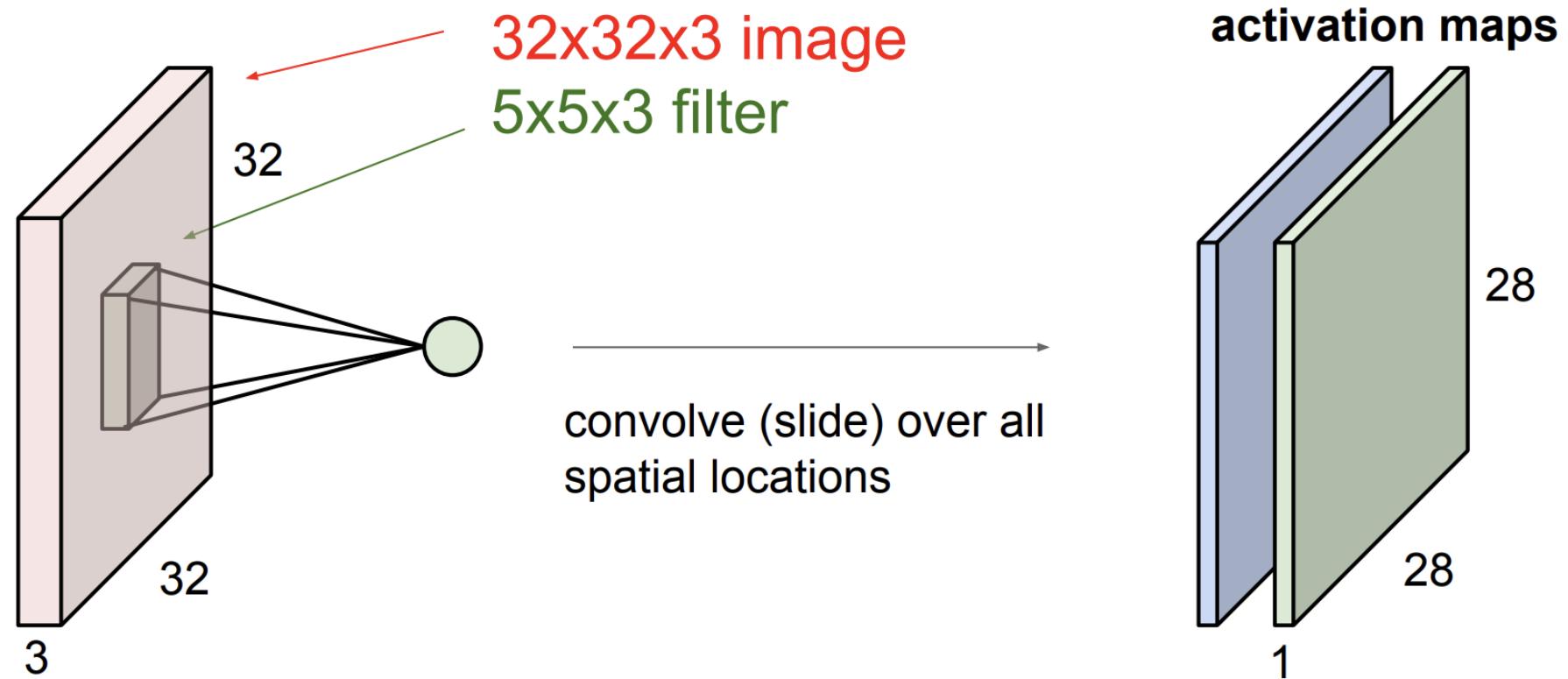
CONVOLUTION LAYER



CONVOLUTION LAYER

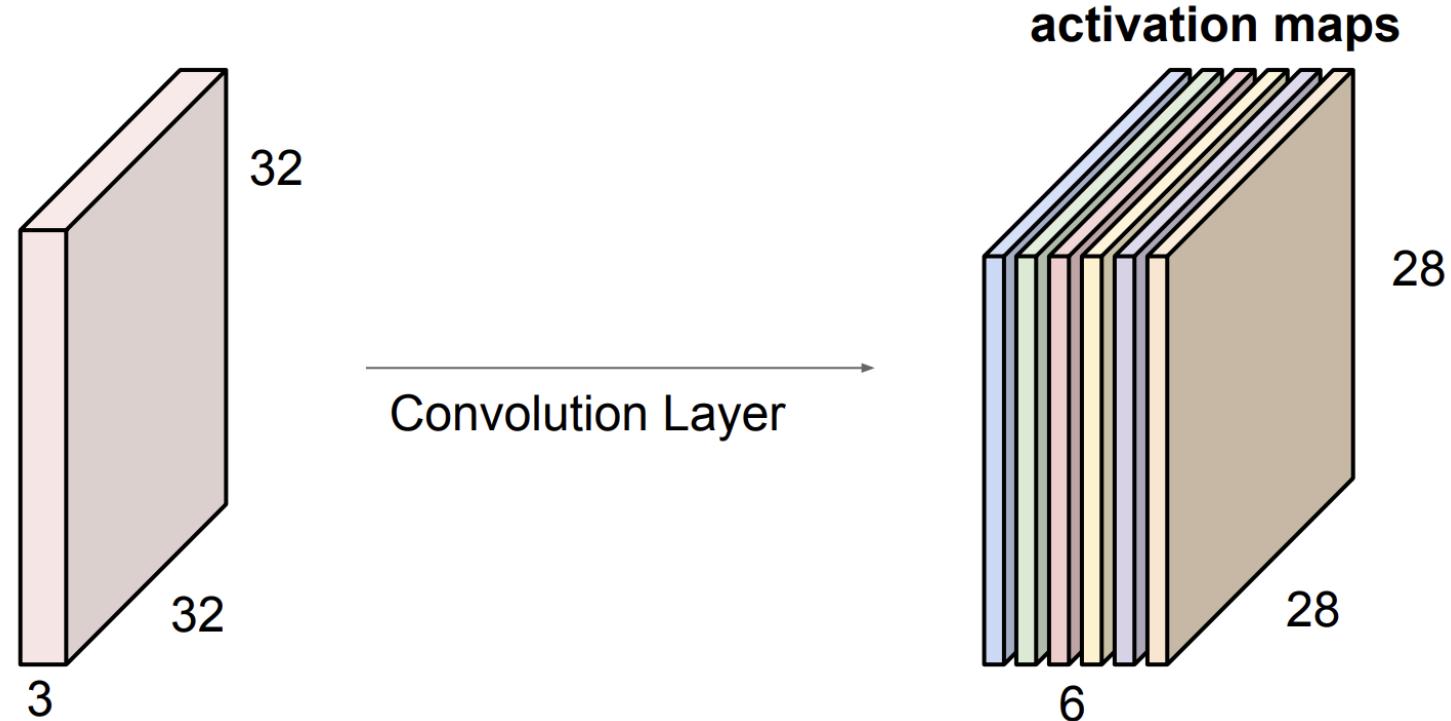


CONVOLUTION LAYER



CONVOLUTION LAYER

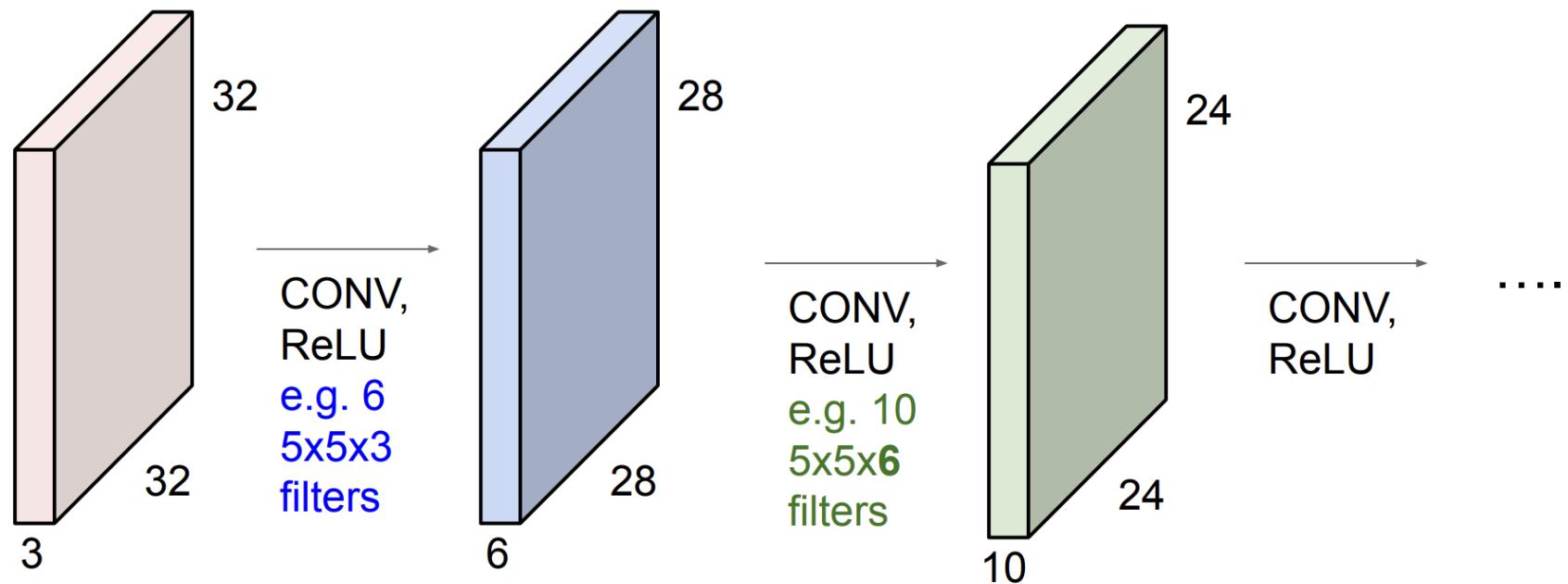
- For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



We stack these up to get a “new image” of size 28x28x6!

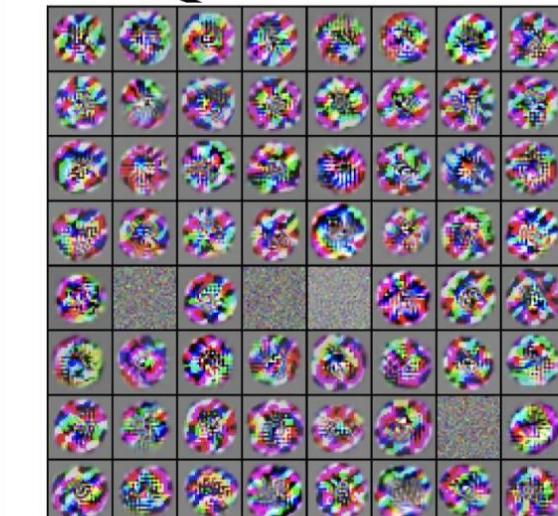
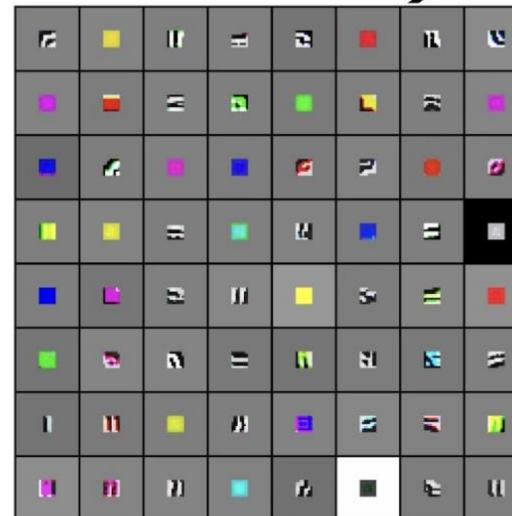
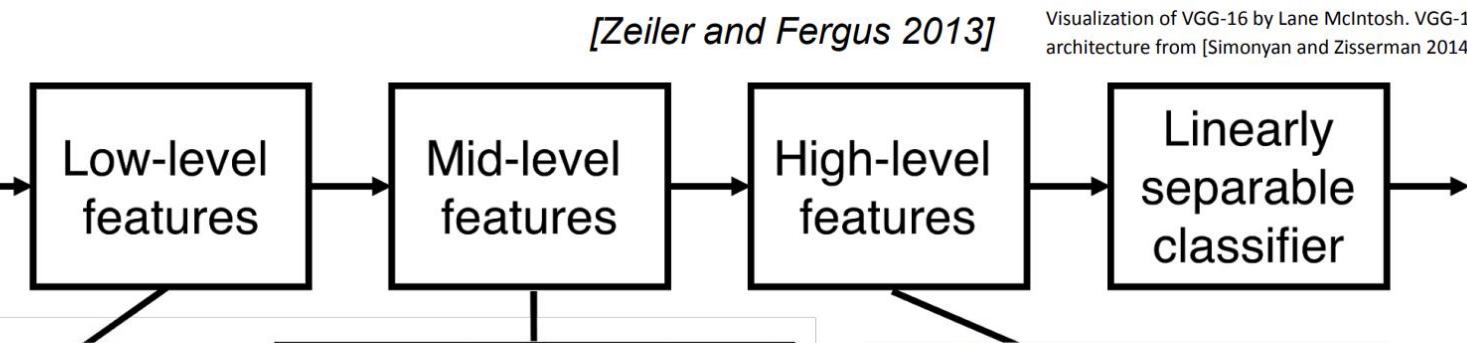
CNN: SEQUENCE OF CONVOLUTION LAYERS

Preview: ConvNet is a sequence of Convolution Layers, interspersed with activation functions



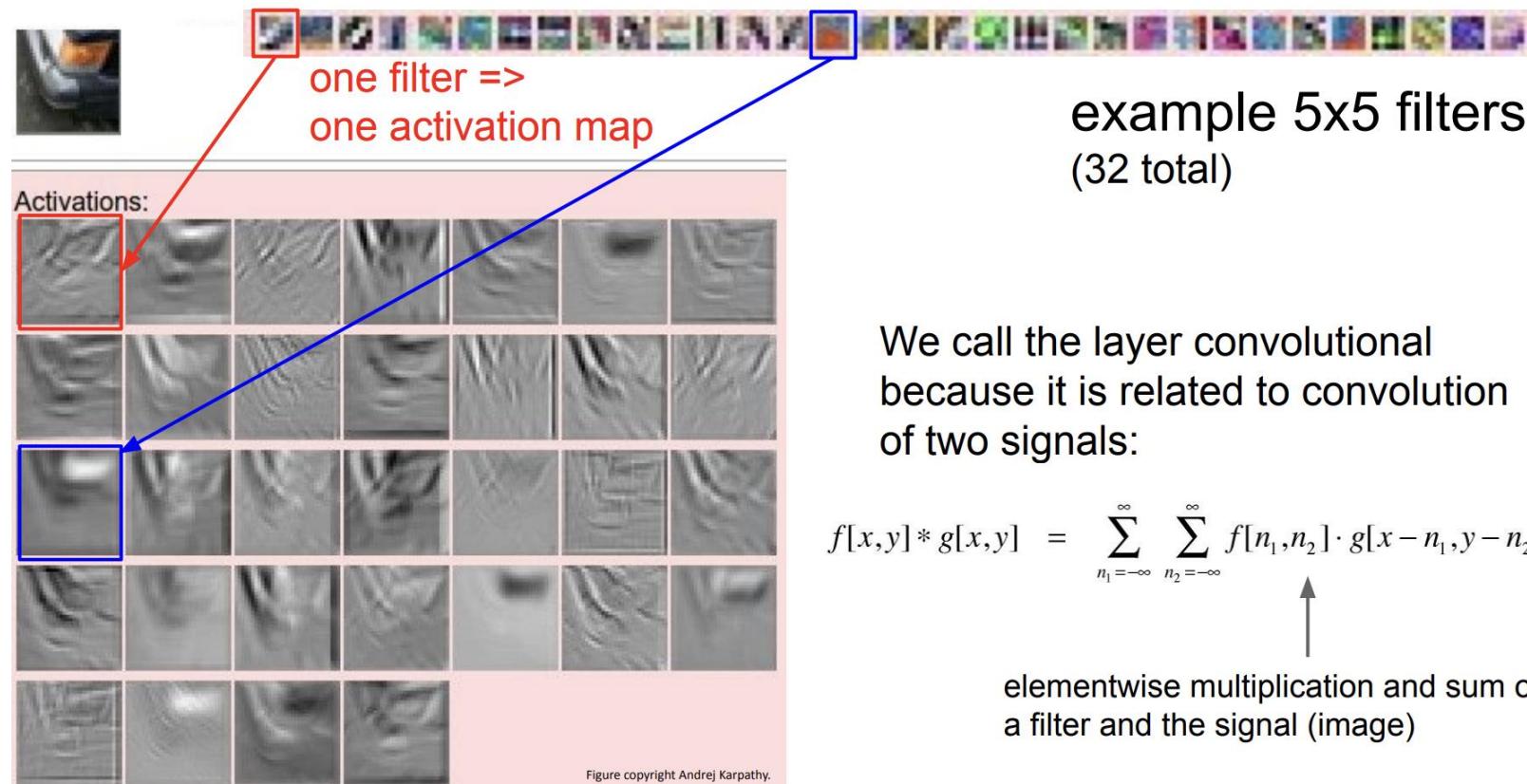
HIERARCHICAL FEATURES

Preview



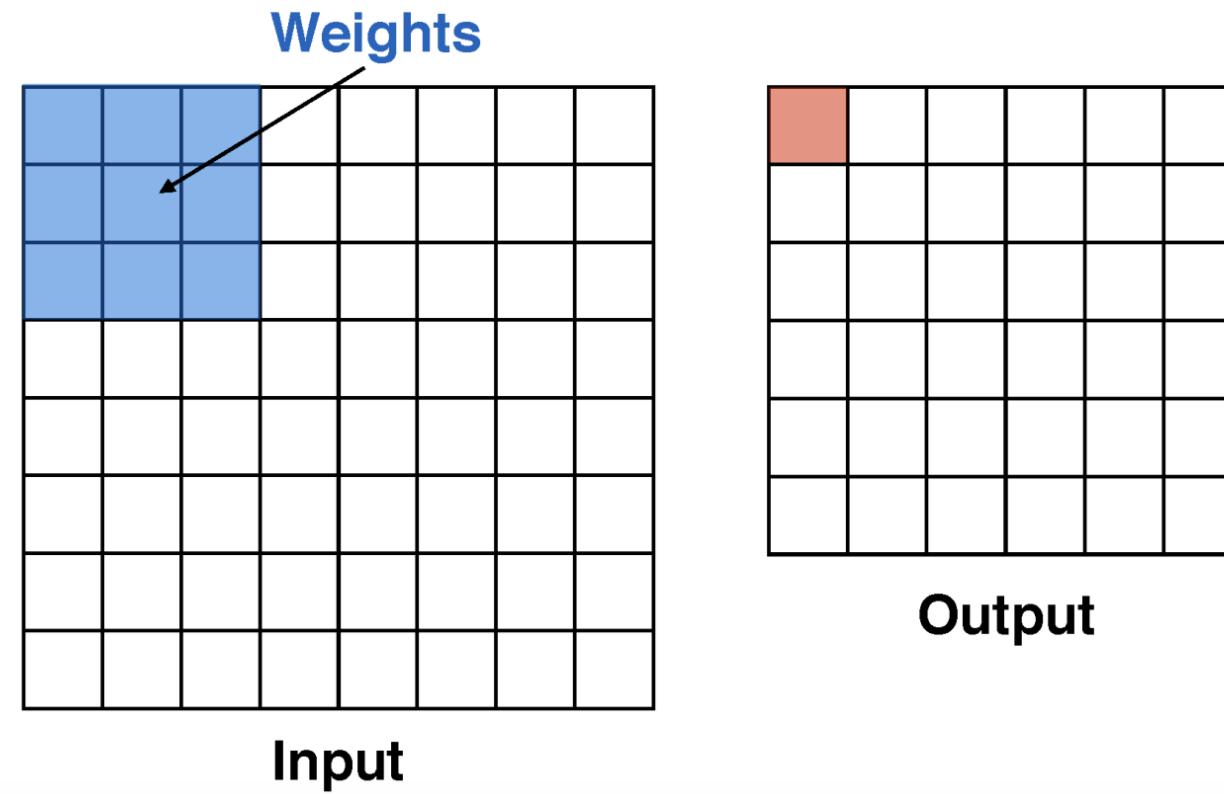
Visualization of VGG-16 by Lane McIntosh. VGG-16 architecture from [Simonyan and Zisserman 2014].

EXAMPLE



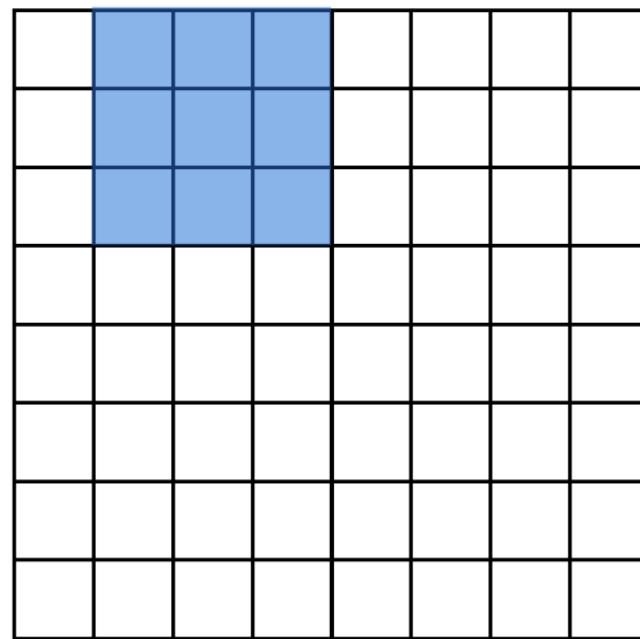
Convolution: Stride

During convolution, the weights “slide” along the input to generate each output

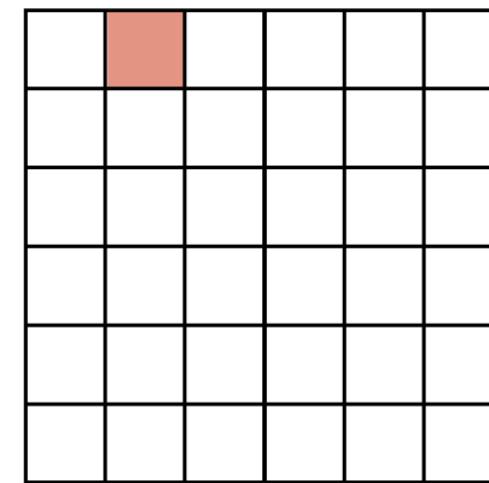


Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



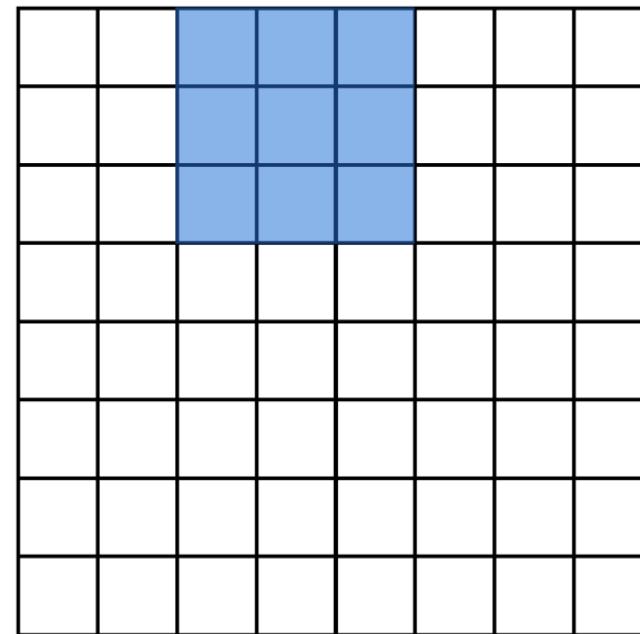
Input



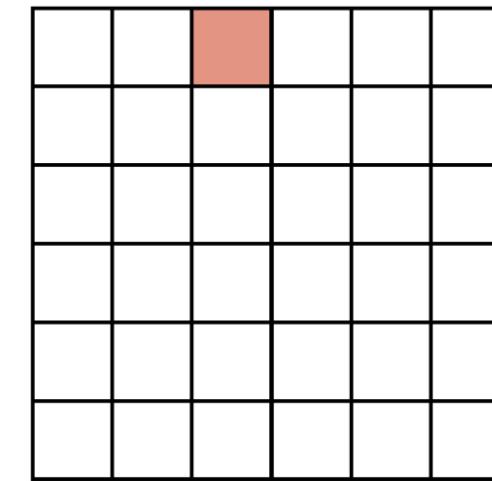
Output

Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



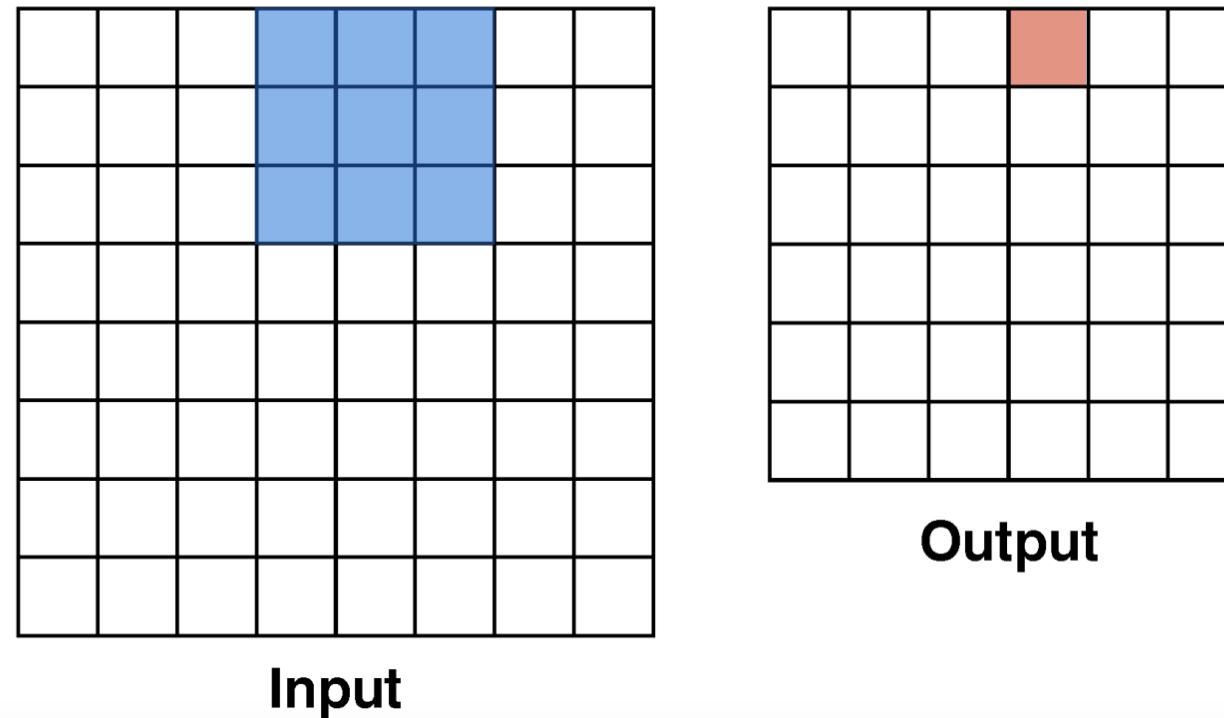
Input



Output

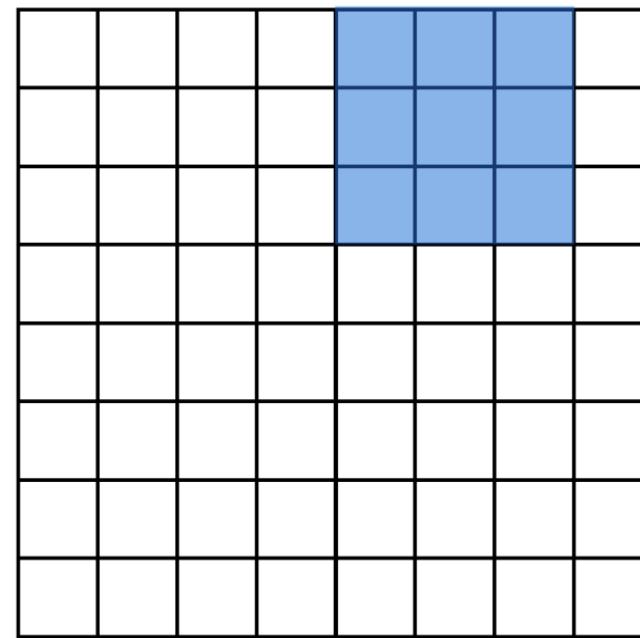
Convolution: Stride

During convolution, the weights “slide” along the input to generate each output

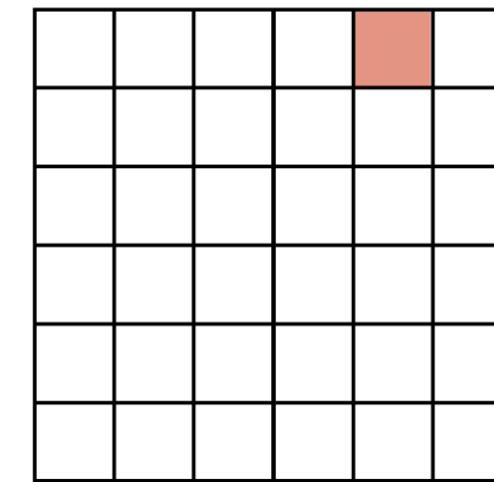


Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



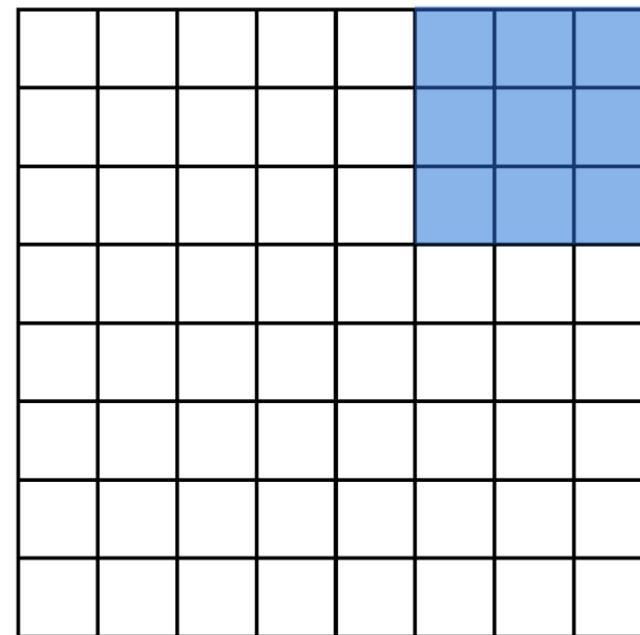
Input



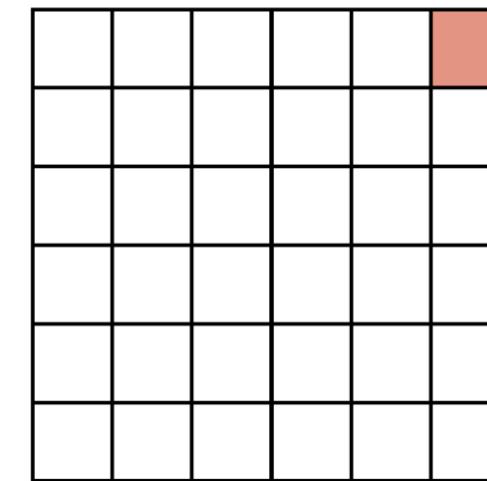
Output

Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



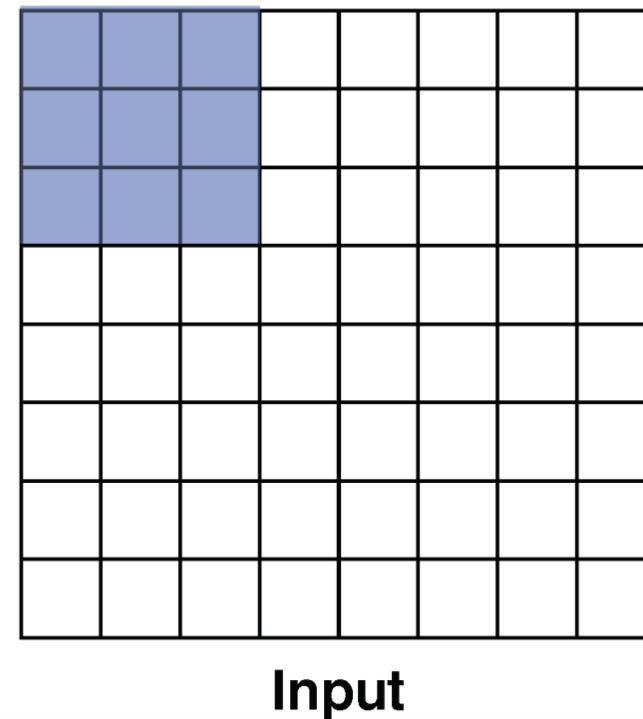
Input



Output

Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



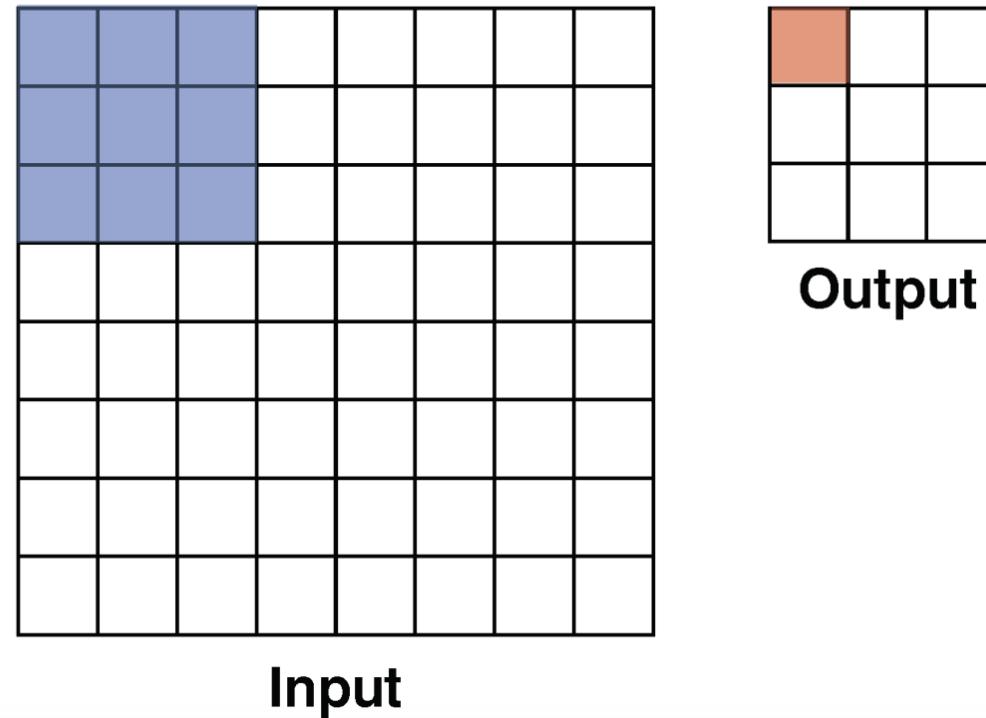
Recall that at each position we are doing a **3D** sum:

$$h^r = \sum_{ijk} x^r_{ijk} W_{ijk} + b$$

(channel, row, column)

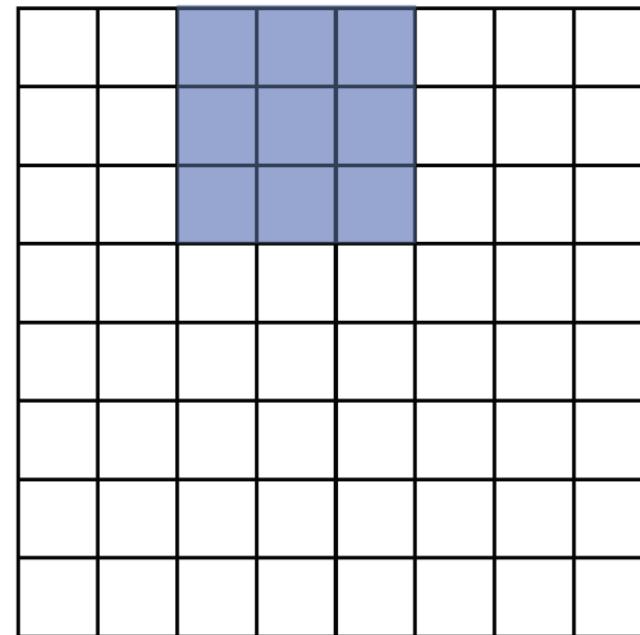
Convolution: Stride

But we can also convolve with a **stride**, e.g. stride = 2

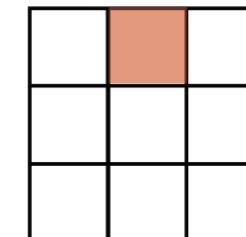


Convolution: Stride

But we can also convolve with a **stride**, e.g. stride = 2



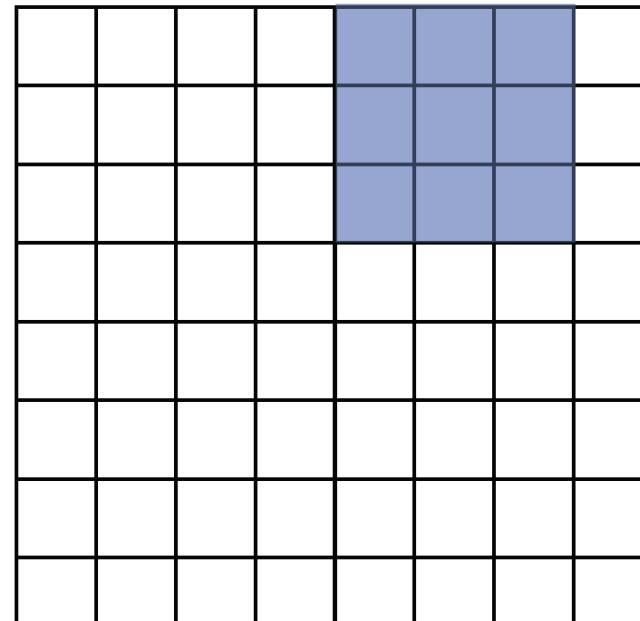
Input



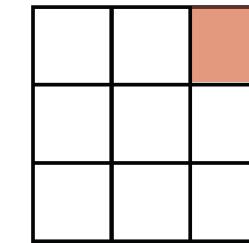
Output

Convolution: Stride

But we can also convolve with a **stride**, e.g. stride = 2



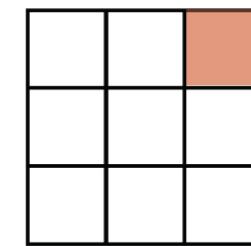
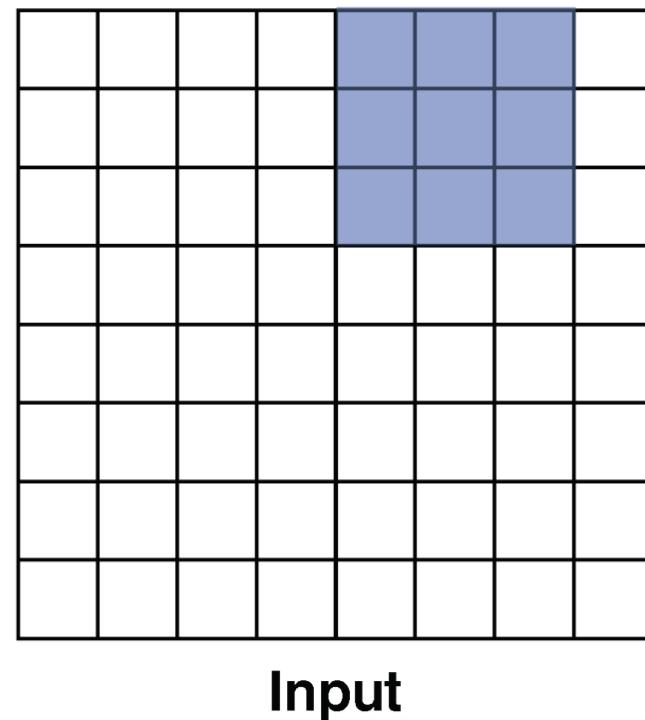
Input



Output

Convolution: Stride

But we can also convolve with a **stride**, e.g. stride = 2



- Notice that with certain strides, we may not be able to cover all of the input
- The output is also half the size of the input

Convolution: Padding

We can also pad the input with zeros.

Here, **pad = 1, stride = 2**

0	0	0	0	0	0	0	0	0	0
0									0
0									0
0									0
0									0
0									0
0									0
0									0
0	0	0	0	0	0	0	0	0	0

Input

0			

Output

Convolution: Padding

We can also pad the input with zeros.

Here, **pad = 1, stride = 2**

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input

Output

Convolution: Padding

We can also pad the input with zeros.

Here, **pad = 1, stride = 2**

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input

Output

Convolution: Padding

We can also pad the input with zeros.

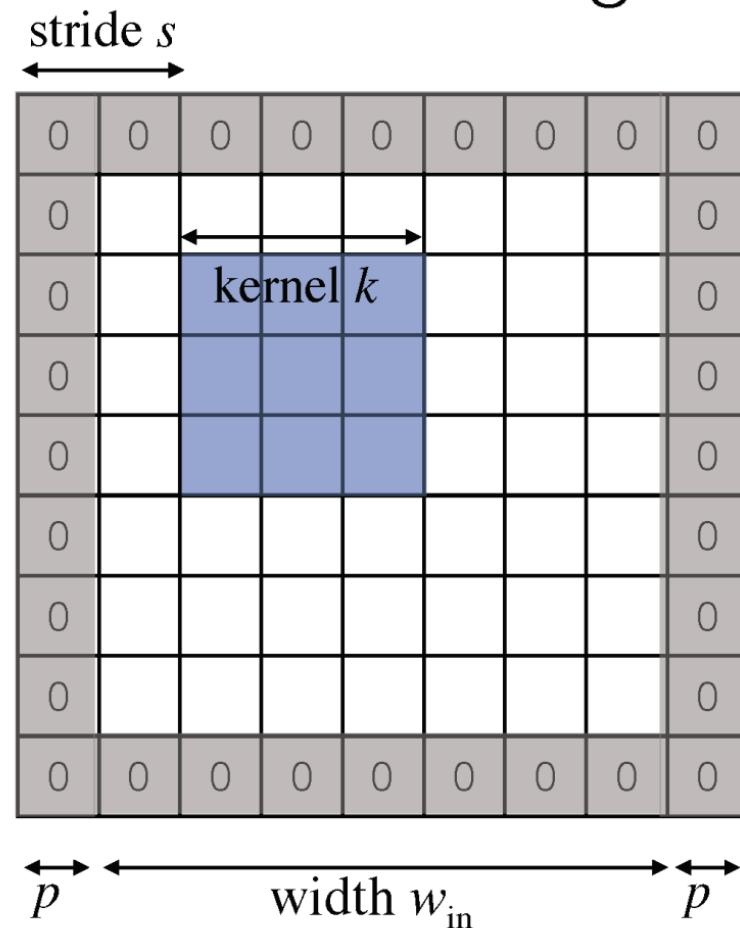
Here, **pad = 1, stride = 2**

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input

Output

Convolution: How big is the output?

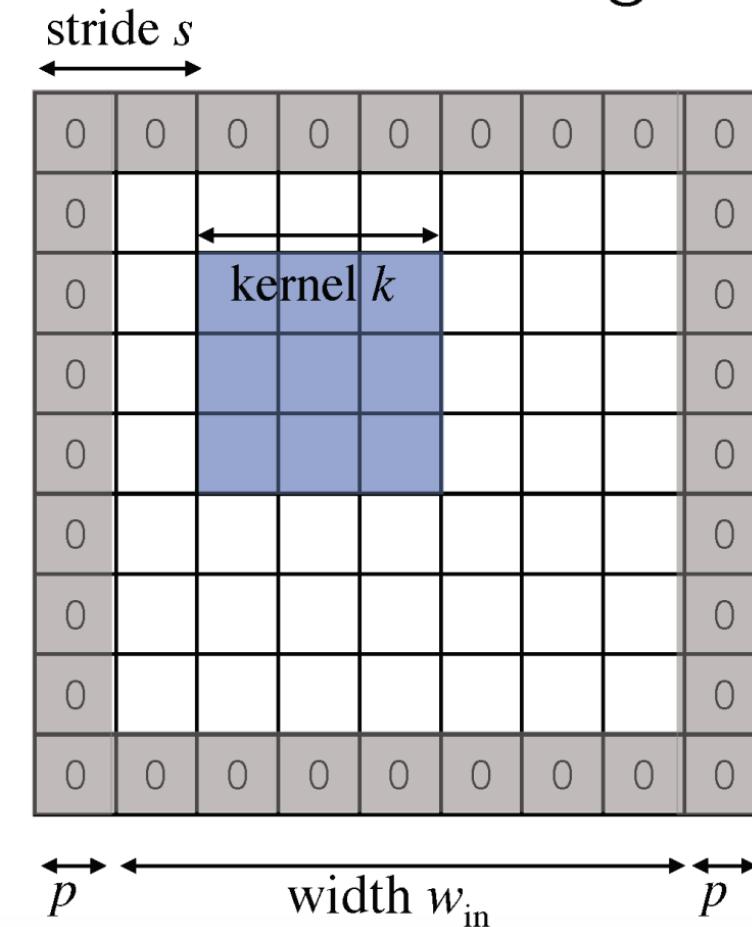


In general, the output has size:

$$w_{out} = \left\lfloor \frac{w_{in} + 2p - k}{s} \right\rfloor + 1$$

Convolution:

How big is the output?



Example: $k=3$, $s=1$, $p=1$

$$\begin{aligned}w_{\text{out}} &= \left\lfloor \frac{w_{\text{in}} + 2p - k}{s} \right\rfloor + 1 \\&= \left\lfloor \frac{w_{\text{in}} + 2 - 3}{1} \right\rfloor + 1 \\&= w_{\text{in}}\end{aligned}$$

VGGNet [Simonyan 2014]
uses filters of this shape

Pooling

For most ConvNets, **convolution** is often followed by **pooling**:

- Creates a smaller representation while retaining the most important information
- The “max” operation is the most common
- Why might “avg” be a poor choice?

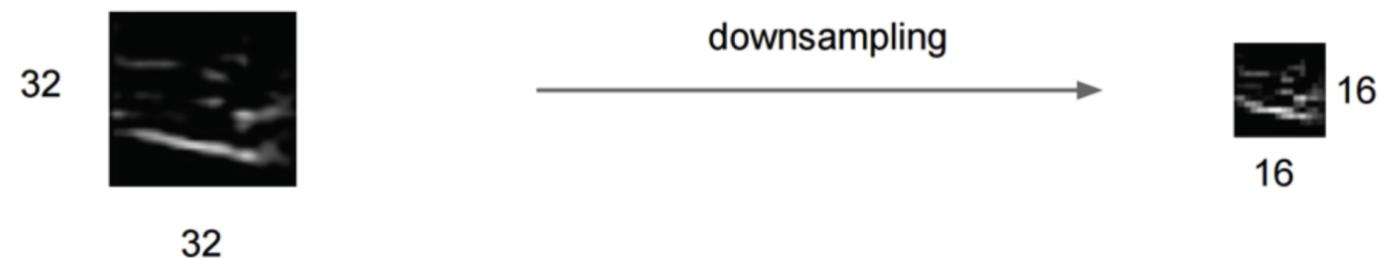
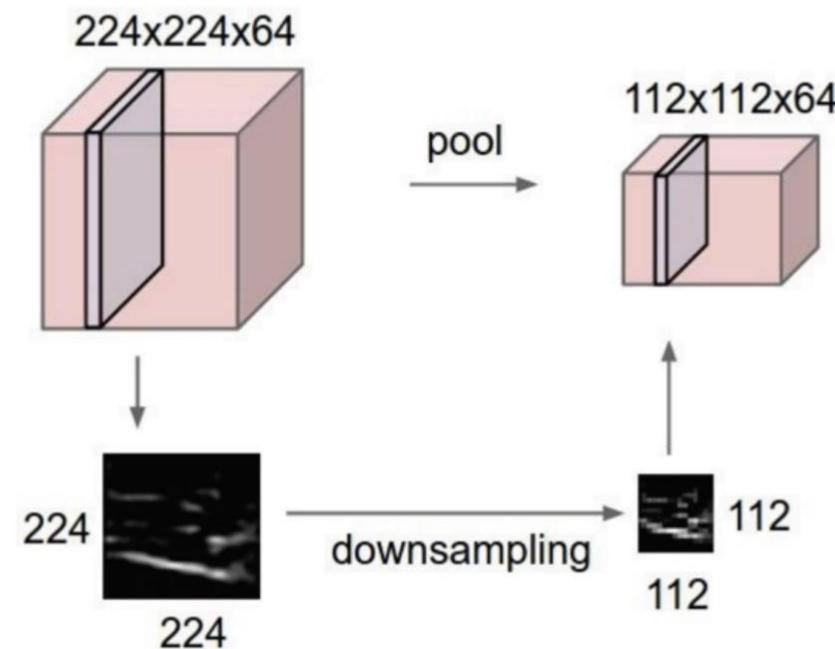


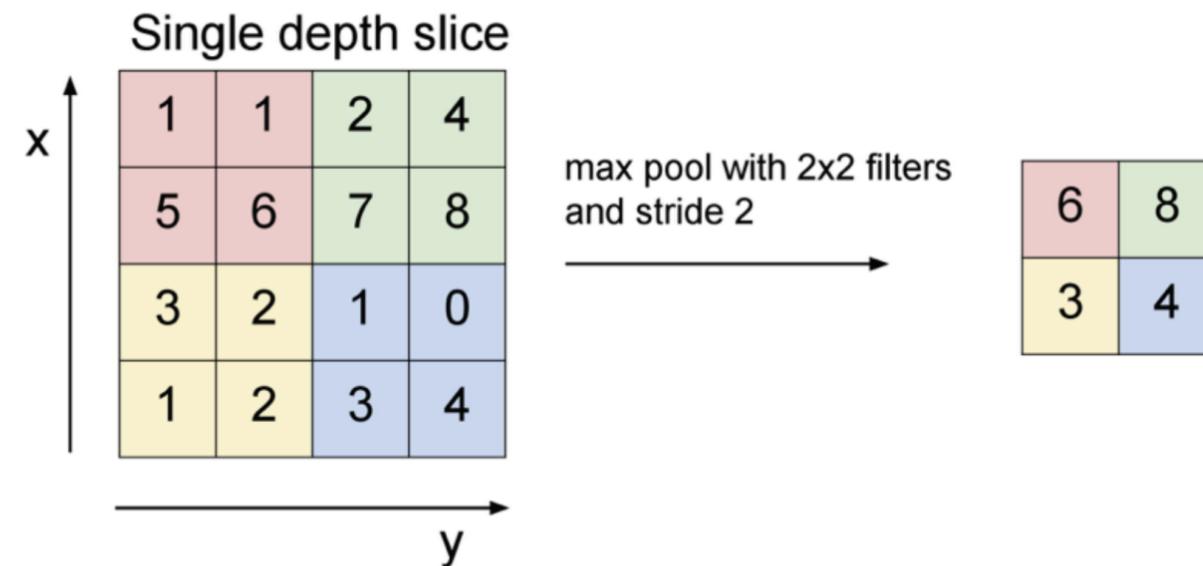
Figure: Andrey Karpathy

Pooling

- makes the representations smaller and more manageable
- operates over each activation map independently:



Max Pooling



What's the backprop rule for max pooling?

- In the forward pass, store the index that took the max
- The backprop gradient is the input gradient at that index

Figure: Andrej Karpathy

Example ConvNet

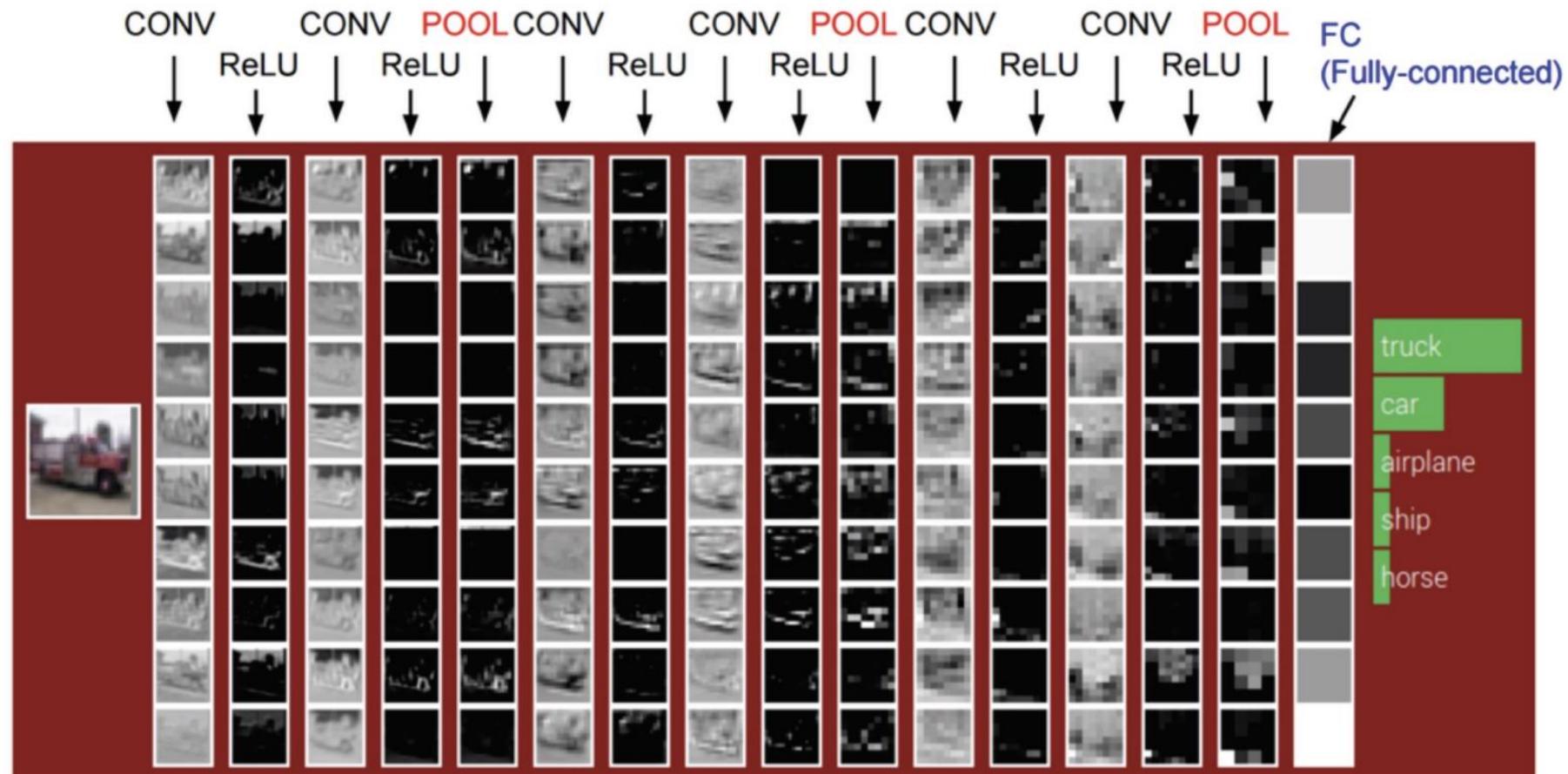


Figure: Andrej Karpathy

RESOURCES

- **CIFAR-10**
- <https://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html>
- **MNIST**
- <https://cs.stanford.edu/people/karpathy/convnetjs/demo/mnist.html>

IMAGE CAPTIONING WITH VISUAL ATTENTION

- We can similarly use attention for image captioning (image → text)
- Builds directly on previous work

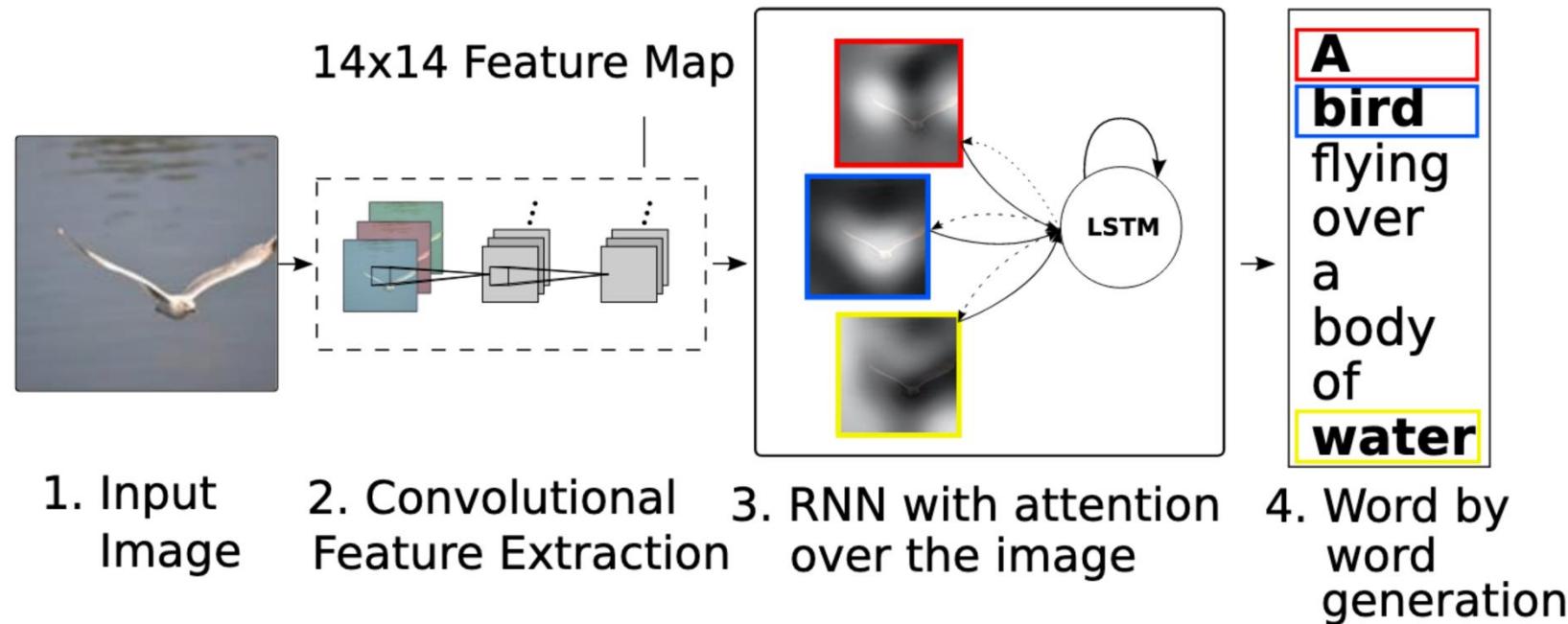


IMAGE CAPTIONING WITH VISUAL ATTENTION

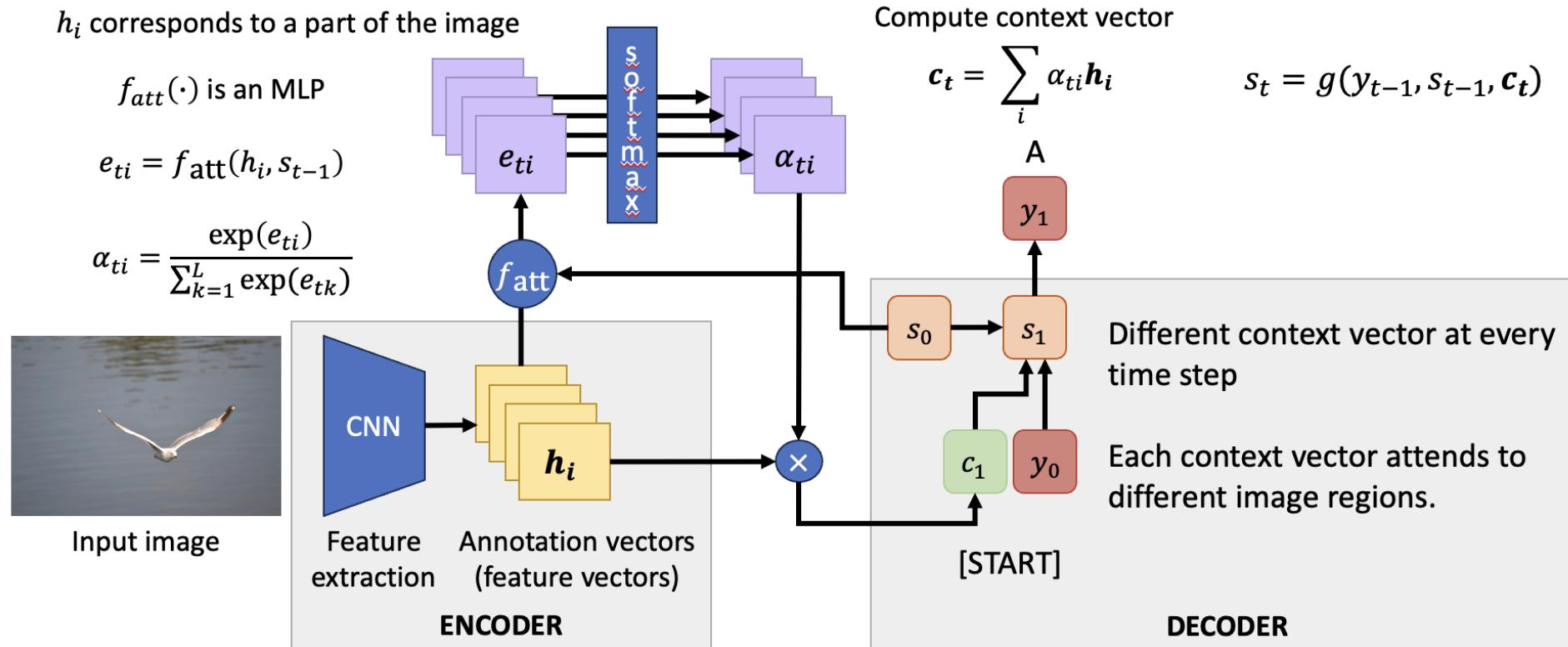


IMAGE CAPTIONING WITH VISUAL ATTENTION

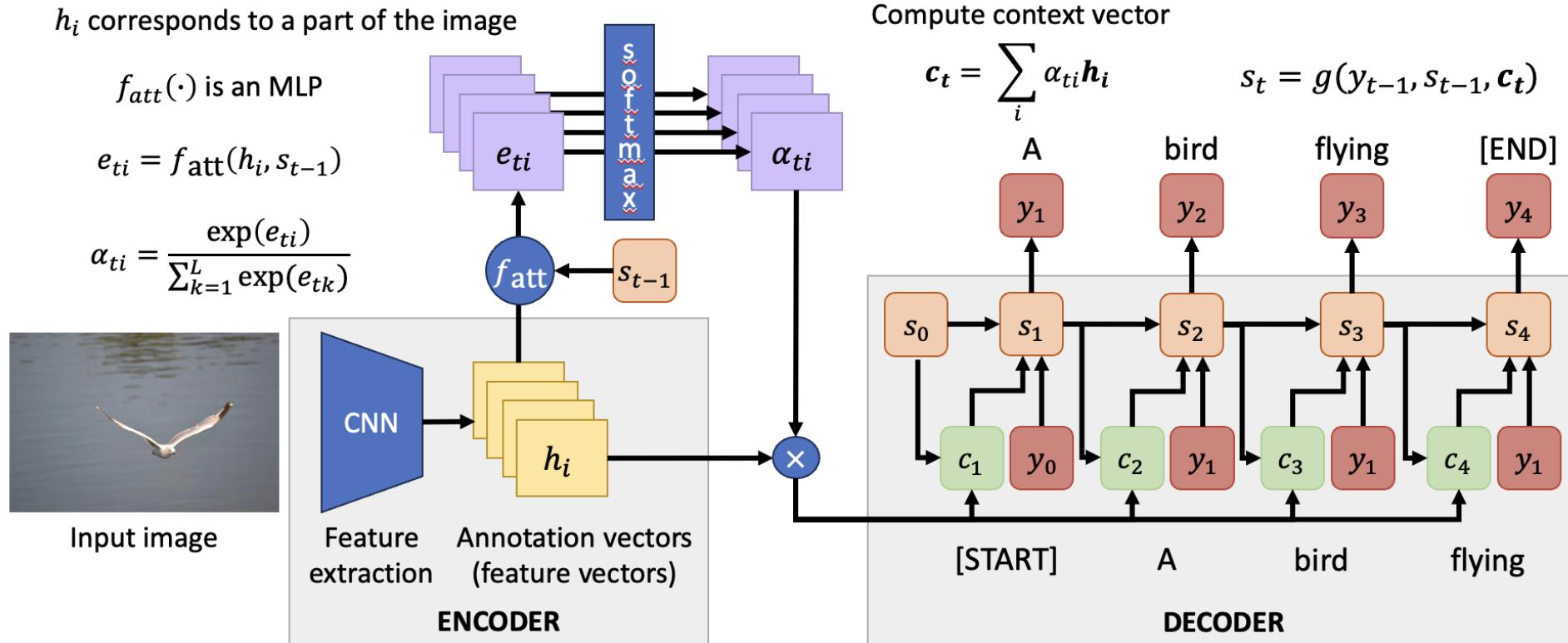


IMAGE CAPTIONING WITH VISUAL ATTENTION

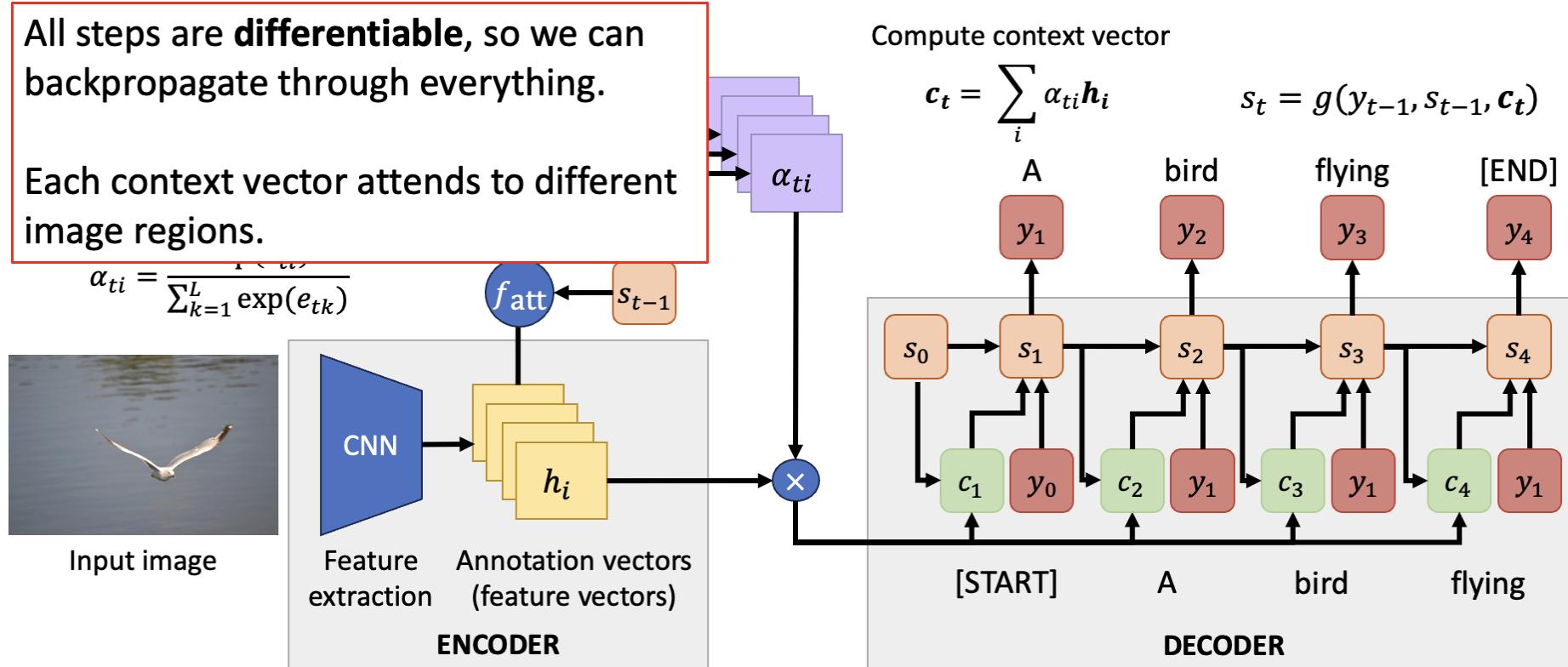


IMAGE CAPTIONING WITH VISUAL ATTENTION

- Visualization of the attention for each generated word
 - Gives insight to “where” and “what” the attention focused on when generating each word

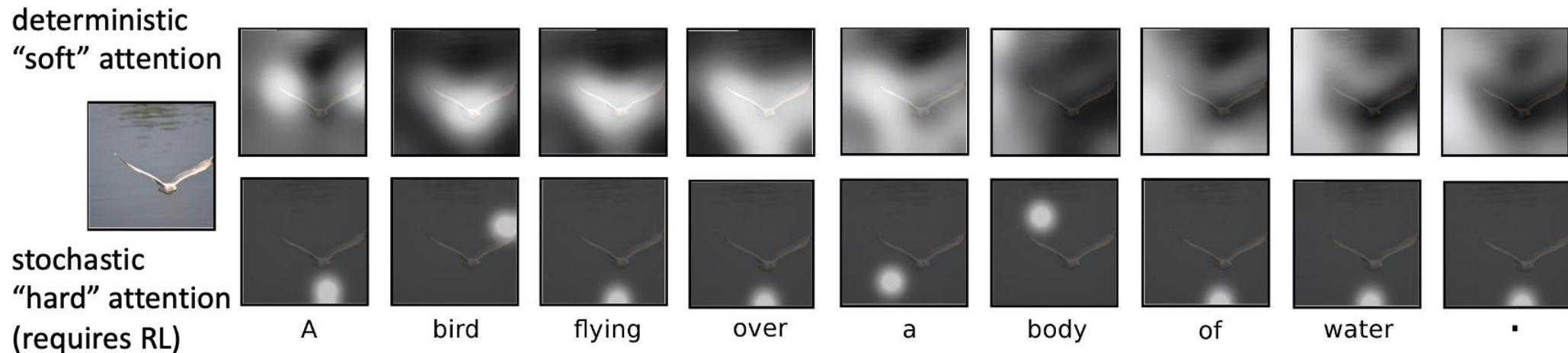


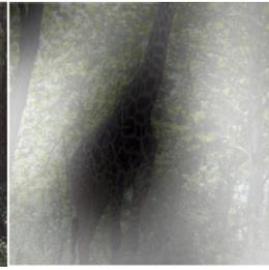
IMAGE CAPTIONING WITH VISUAL ATTENTION



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.

RESOURCES

- <https://medium.com/@deepeshrishu09/automatic-image-captioning-with-pytorch-cf576c98d319>
- <https://www.kaggle.com/code/mdteach/image-captioning-with-attention-pytorch>