



DEEP LEARNING CNN & RNN

Lecture-6 , Day-3

STTP on “Deep Learning, Computer Vision and
Speech Processing”

By: Suprava, Patnaik, Professor, ExTC, XIE, Mumbai

PART I:

DEEPLARNING :

WEIGHT SHARING(CNN)

Why Convolution?

- Translation invariance



- Object recognition:
 - Feed-forward NN has different weights everywhere and for every input location.
 - If we train it on left image, it wouldn't recognize the right image

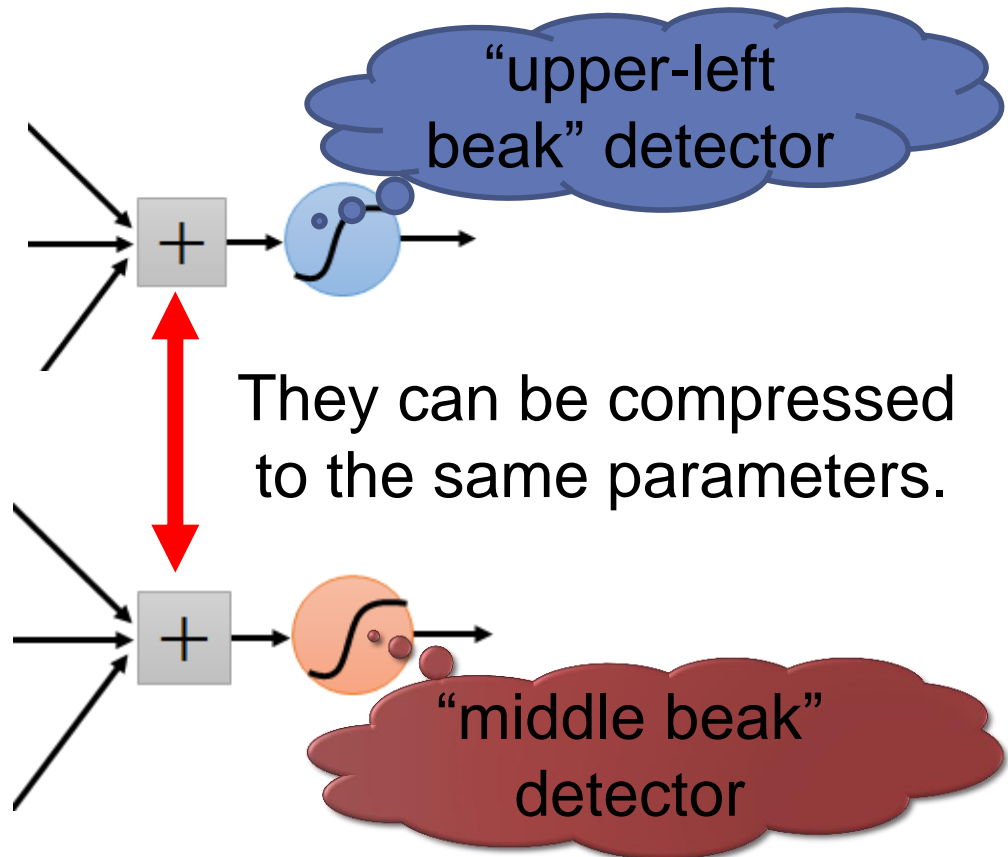
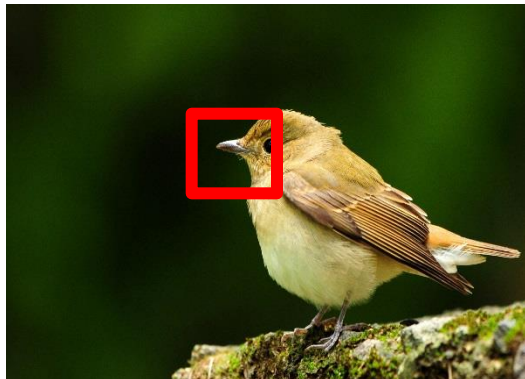
Consider learning an image:

- Some patterns are much smaller than the whole image

Can represent a small region with fewer parameters

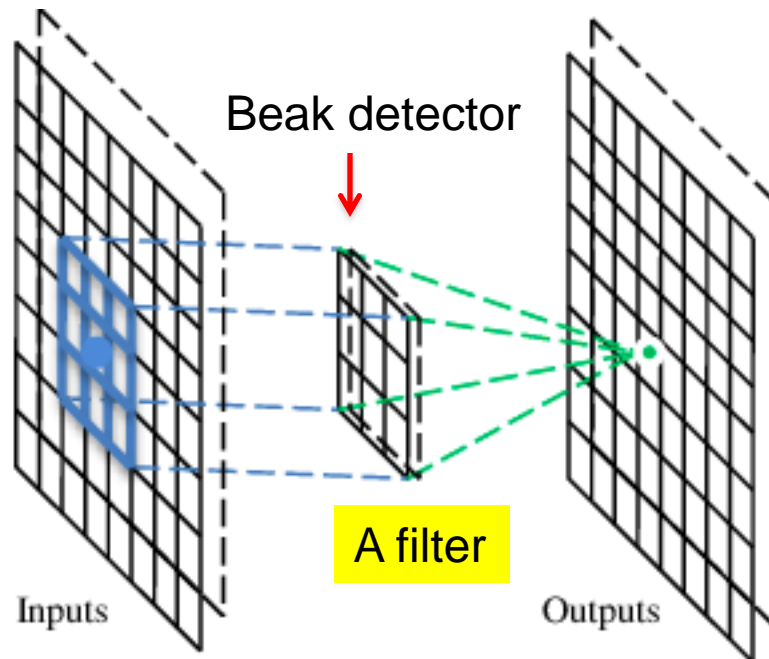


What about training a lot of such “small” detectors and each detector must “move around”.



A convolutional layer

A CNN is a neural network with some convolutional layers (and some other layers). A convolutional layer has a number of filters that does convolutional operation.



Convolution

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

These are the network parameters to be learned.

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

-1	1	-1
-1	1	-1
-1	1	-1

Filter 2

⋮ ⋮

Each filter detects a small pattern (3 x 3).

Convolution

stride=1

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

Dot
product



3

-1

Convolution

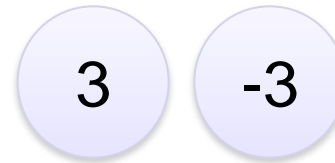
If stride=2

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1



Convolution

stride=1

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

3	-1	-3	-1
-3	1	0	-3
-3	-3	0	1
3	-2	-2	-1

Convolution

stride=1

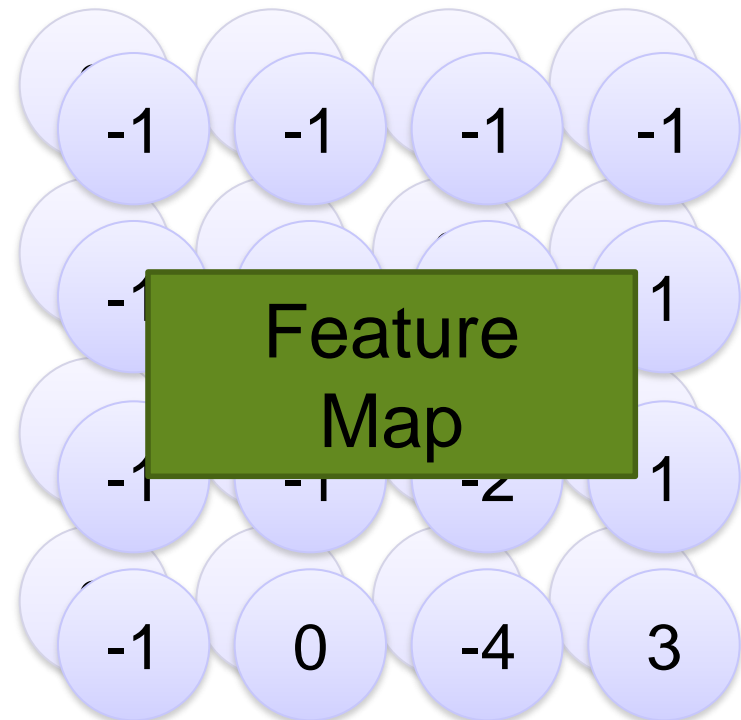
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

-1	1	-1
-1	1	-1
-1	1	-1

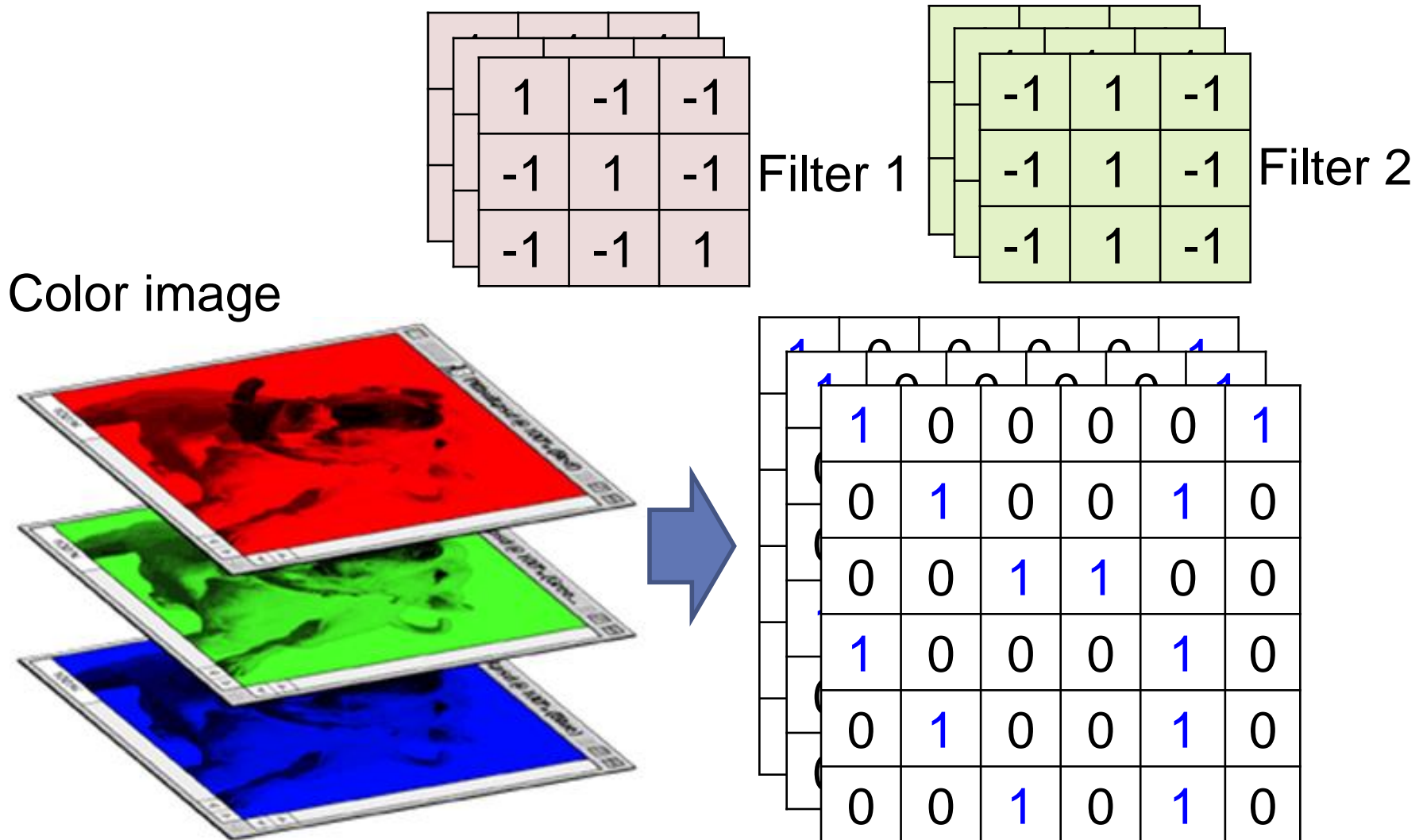
Filter 2

Repeat this for each filter

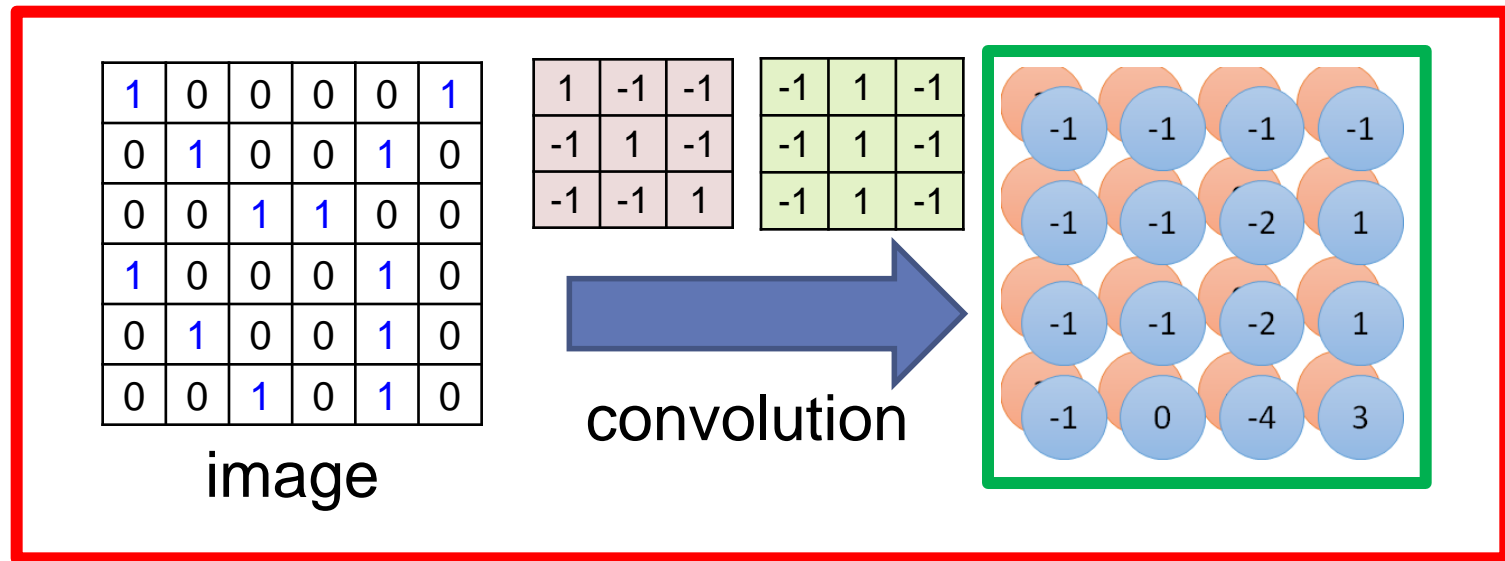


Two 4 x 4 images
Forming 2 x 4 x 4 matrix

Color image: RGB 3 channels

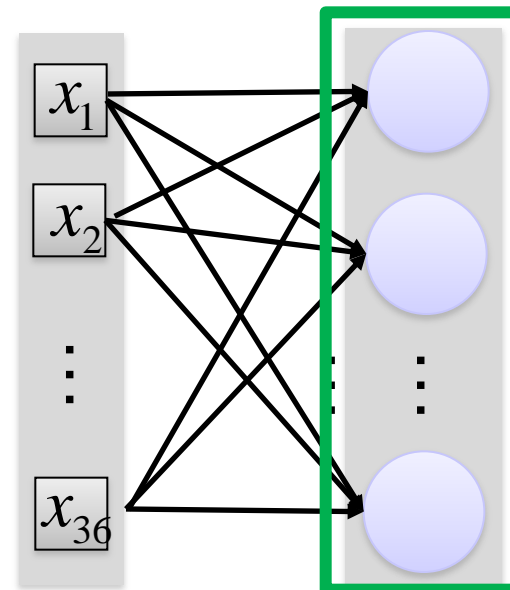


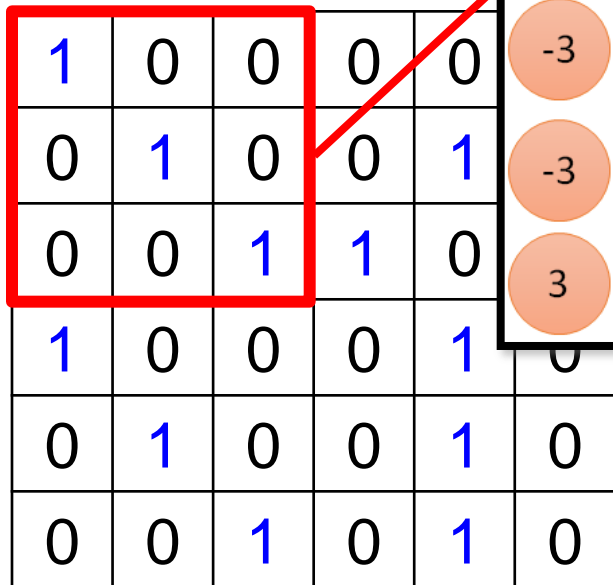
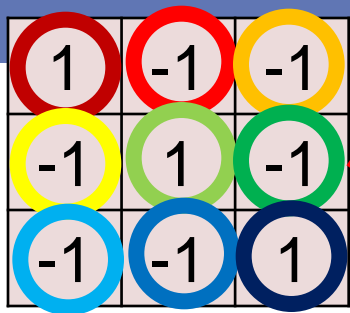
Convolution v.s. Fully Connected



Fully-
connected

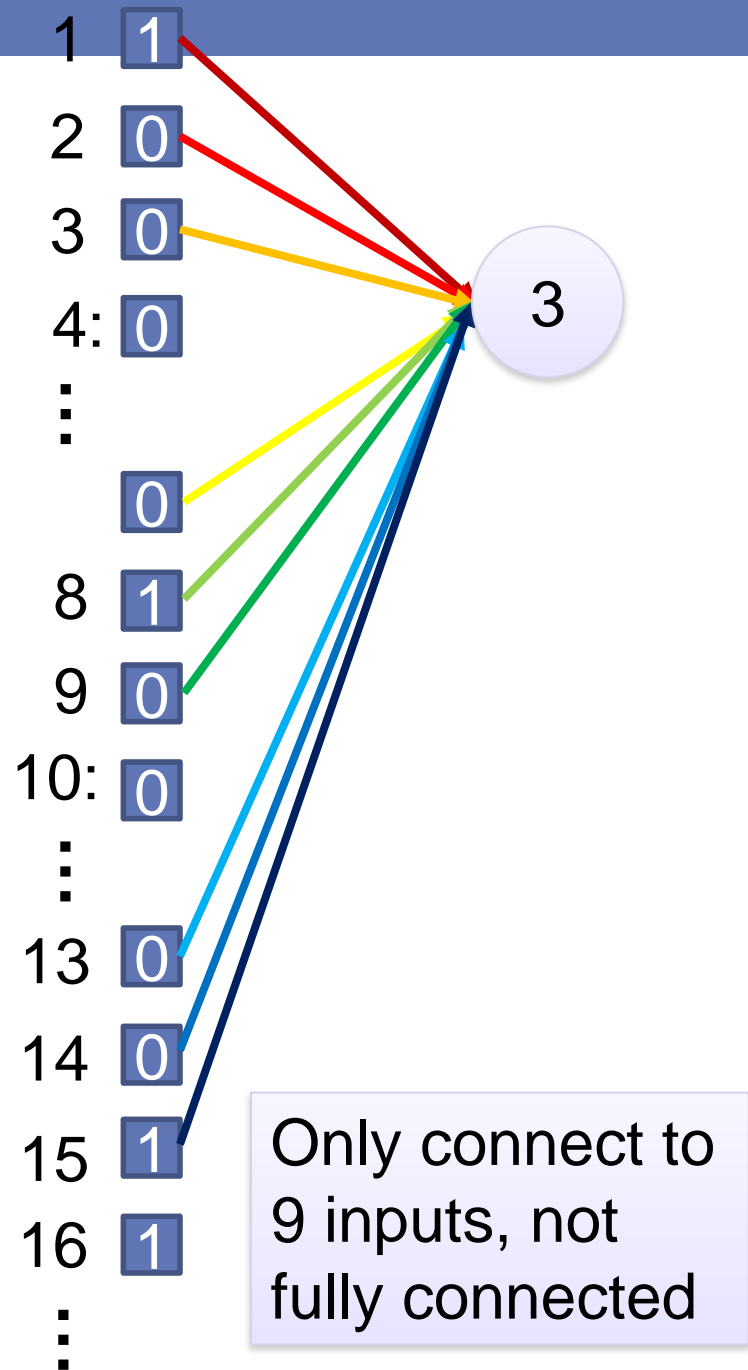
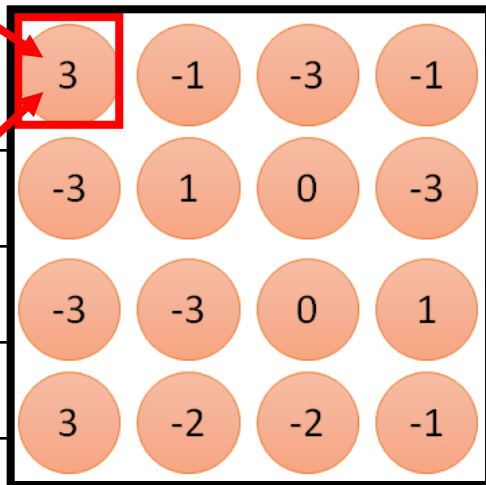
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

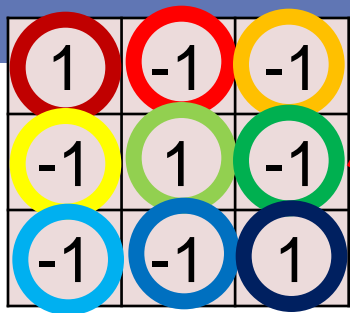




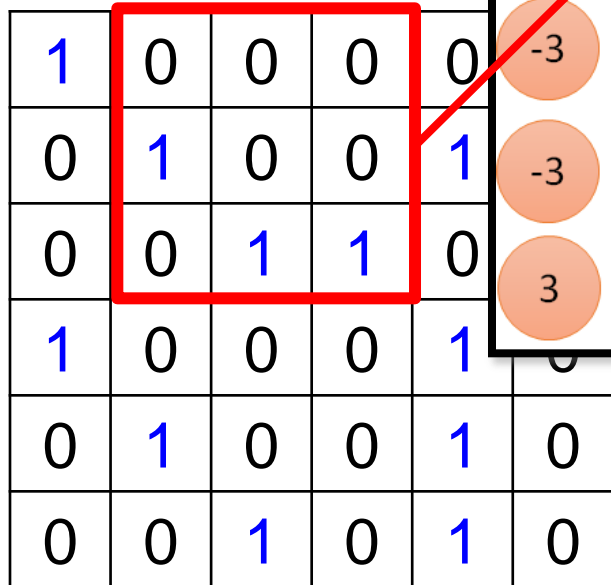
6 x 6 image

fewer parameters!

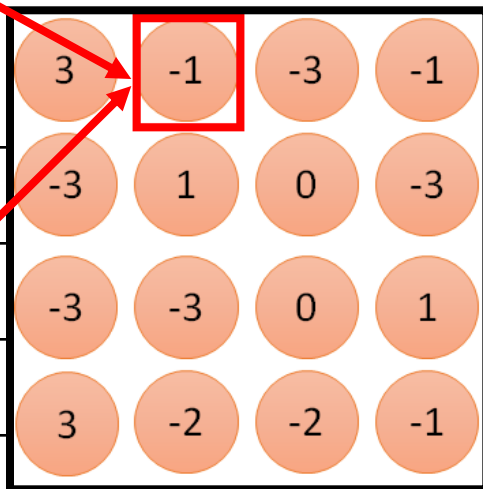




Filter 1

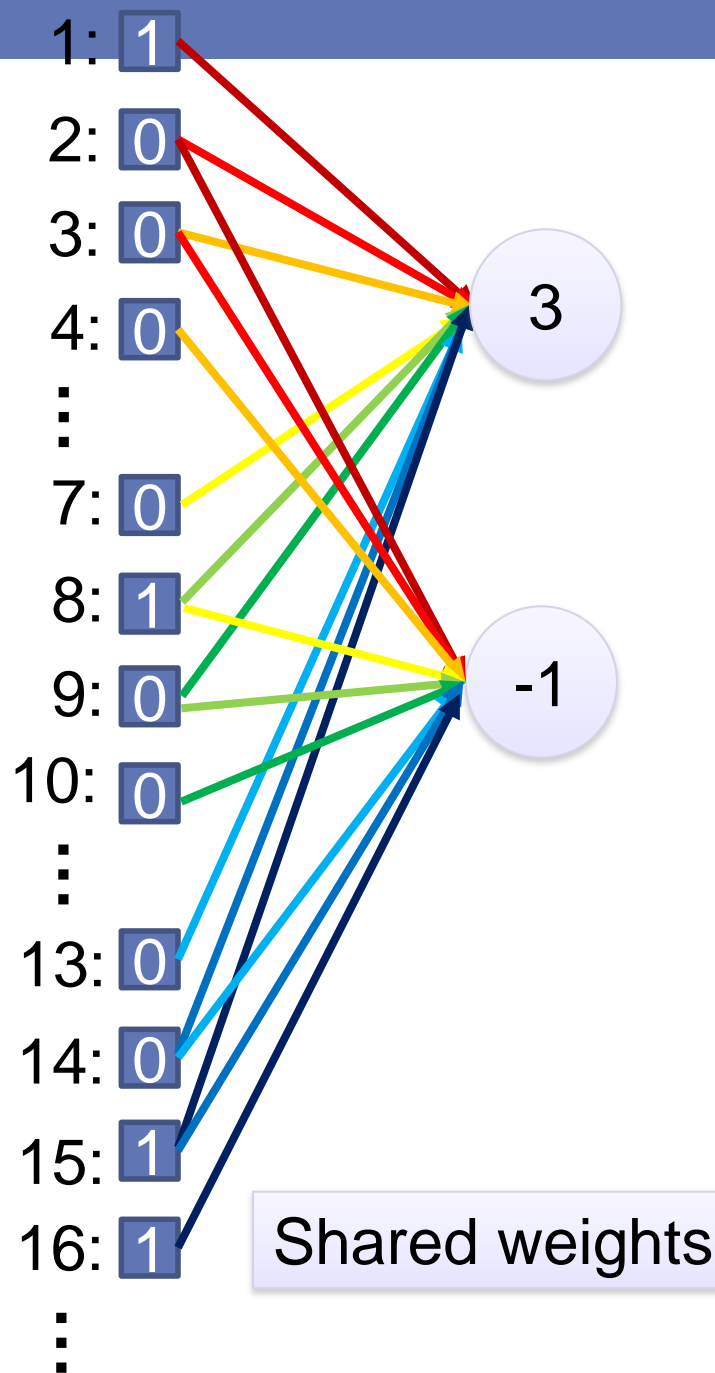


6 x 6 image



Fewer parameters

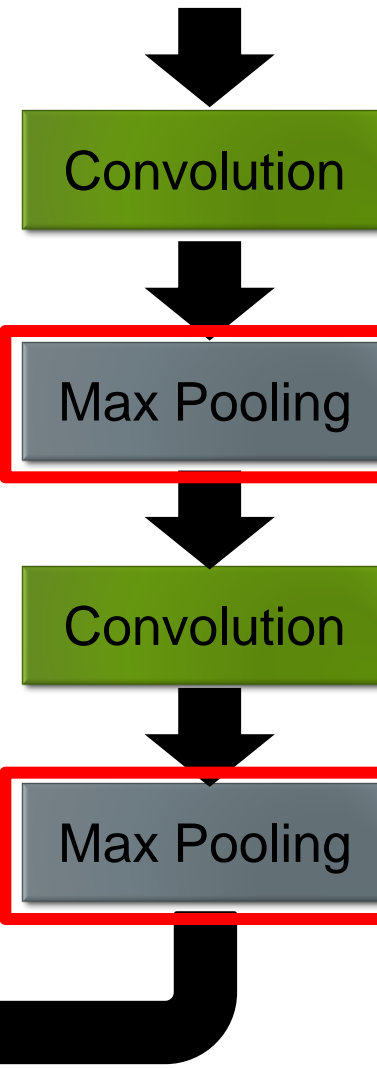
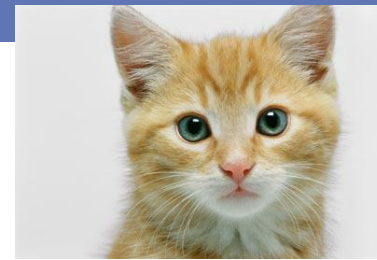
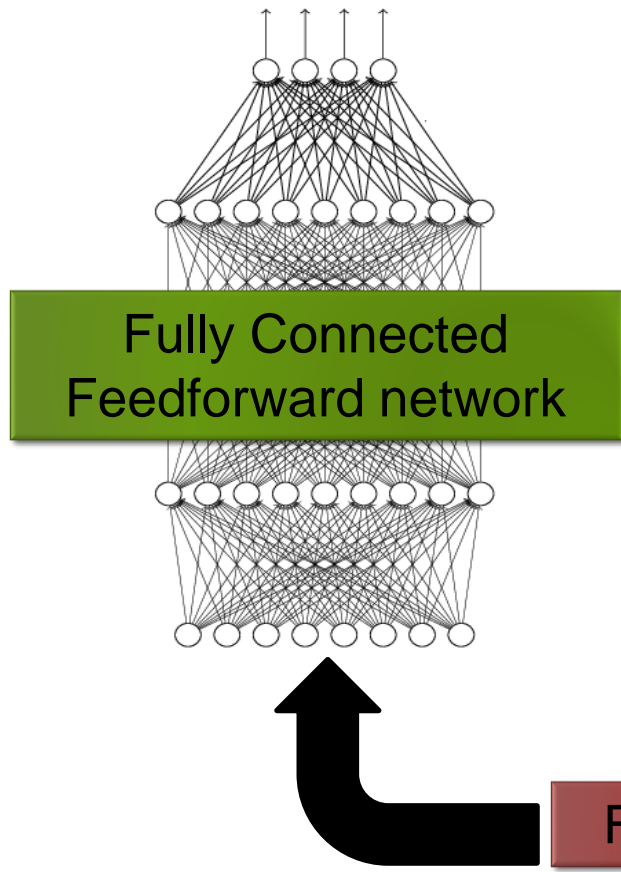
Even fewer parameters



Shared weights

The whole CNN

cat dog



Can
repeat
many
times

Max Pooling

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

-1	1	-1
-1	1	-1
-1	1	-1

Filter 2

3	-1	-3	-1
-3	1	0	-3
-3	-3	0	1
3	-2	-2	-1

-1	-1	-1	-1
-1	-1	-2	1
-1	-1	-2	1
-1	0	-4	3

Why Pooling

- Subsampling pixels will not change the object

bird



Subsampling

bird



We can subsample the pixels to make image smaller



fewer parameters to characterize the image

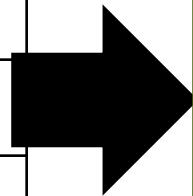
A CNN compresses a fully connected network in two ways:

- Reducing number of connections
- Shared weights on the edges to recognize objects independent of location
- Max pooling further reduces the complexity

Max Pooling

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

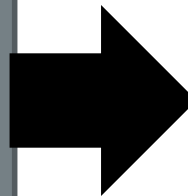
6 x 6 image



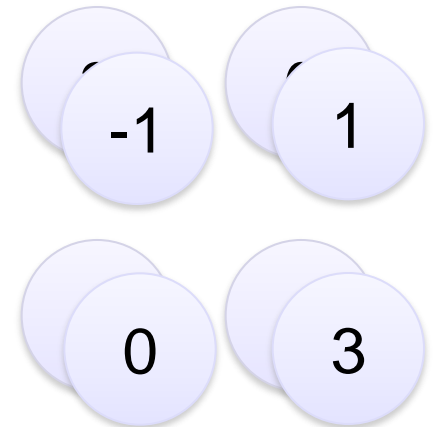
Conv



Max
Pooling



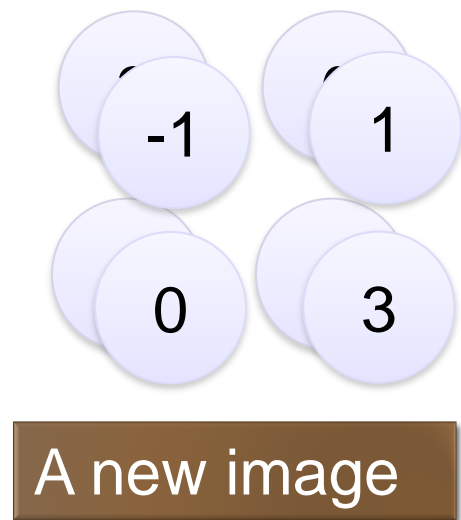
New image
but smaller



2 x 2 image

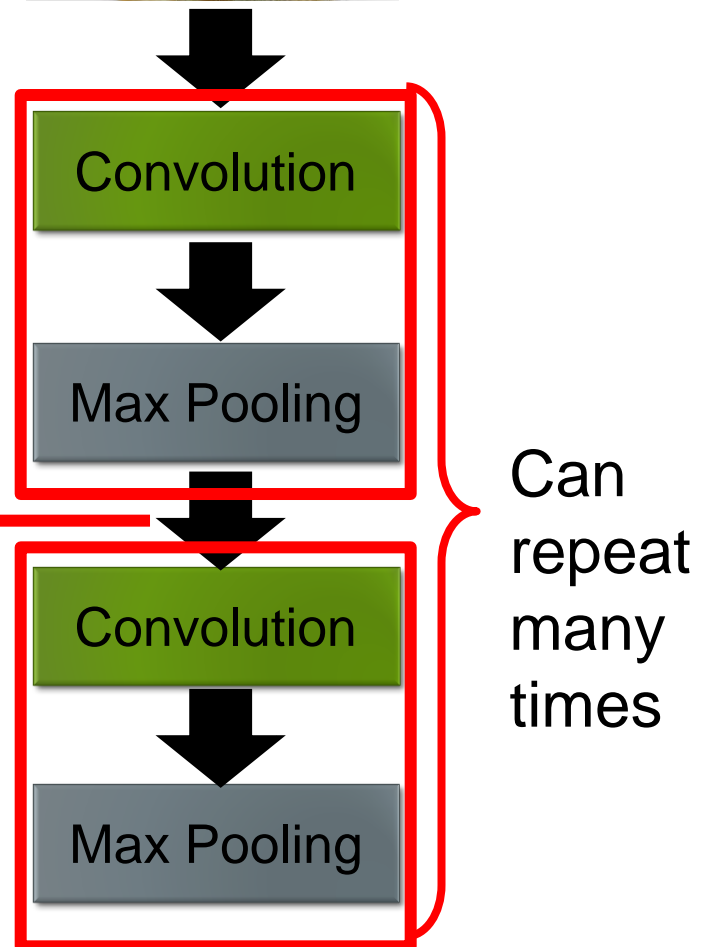
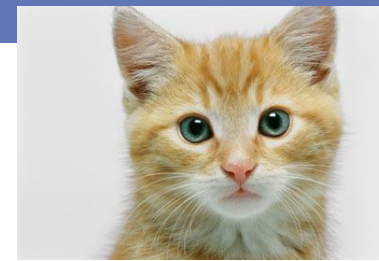
Each filter
is a channel

The whole CNN



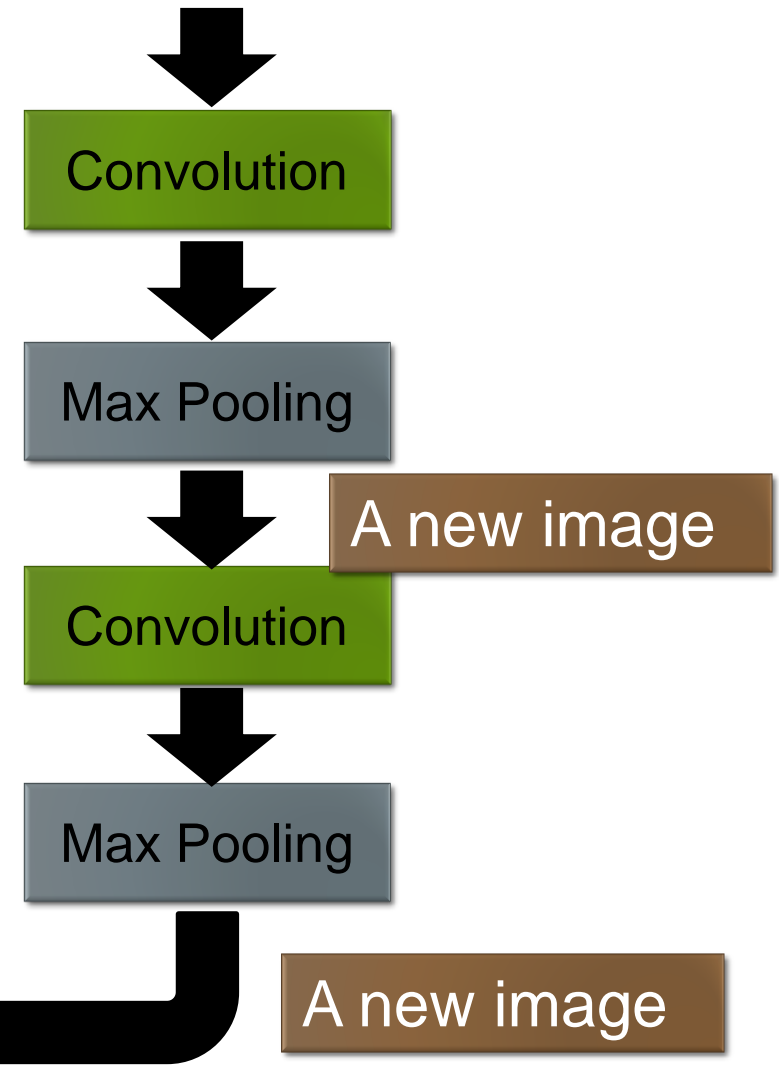
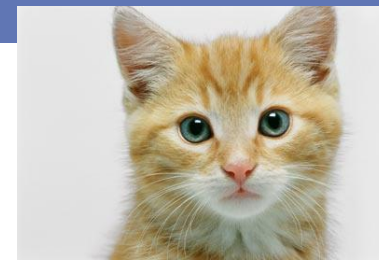
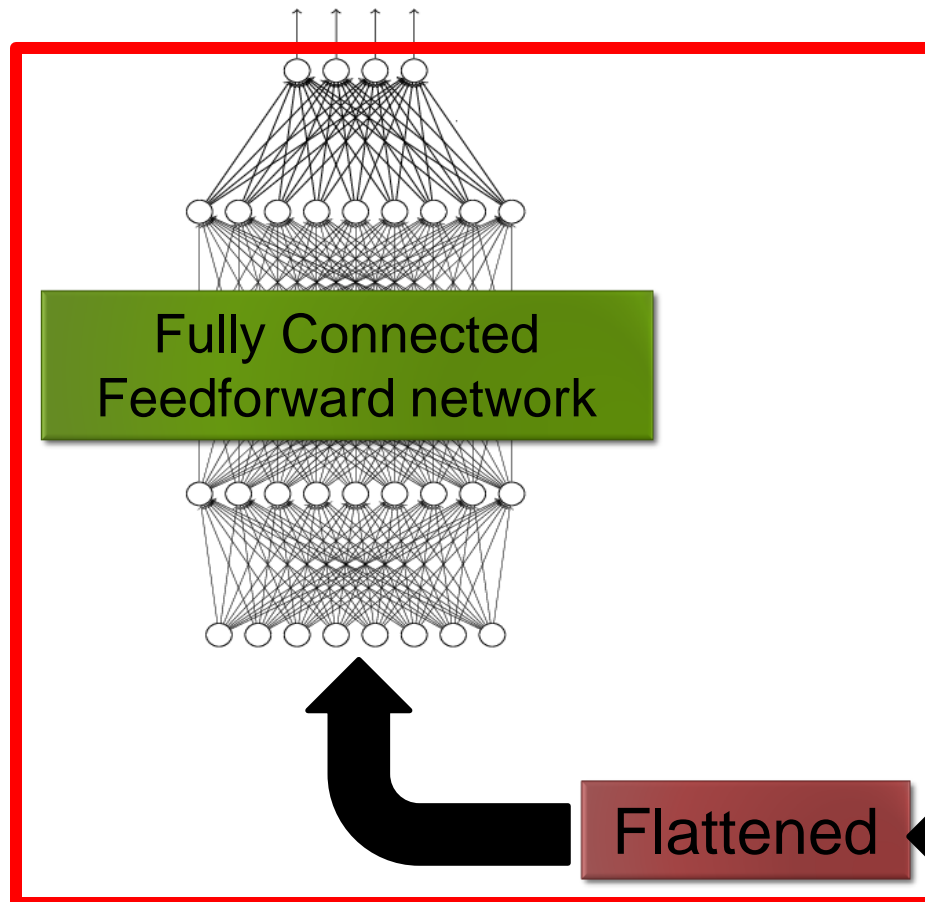
Smaller than the original image

The number of channels is the number of filters

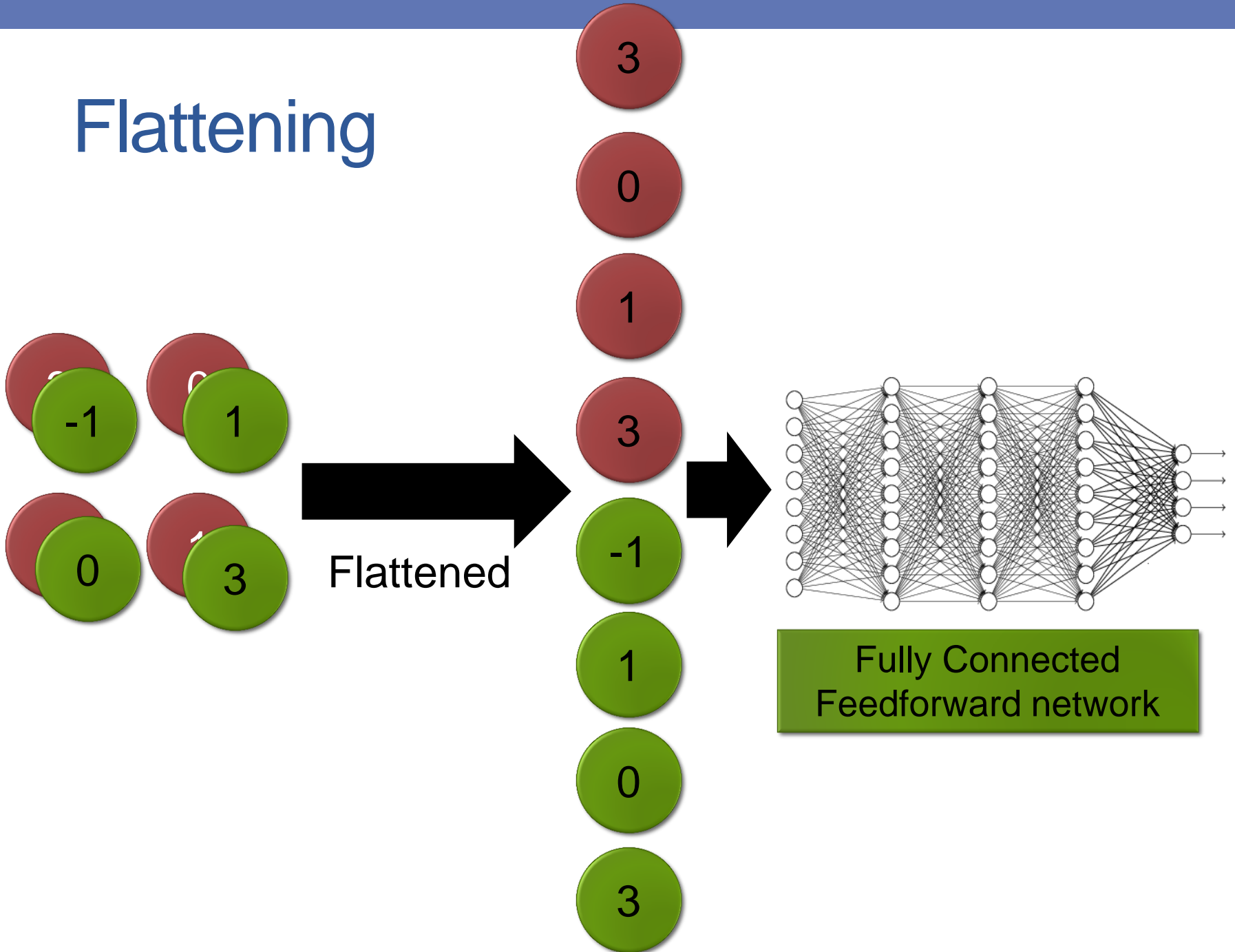


The whole CNN

cat dog



Flattening



CNN in Keras

Only modified the *network structure* and *input format* (vector -> 3-D tensor)

```
model2.add( Convolution2D( 25, 3, 3,  
                           input_shape=(28, 28, 1)) )
```

1	-1	-1
-1	1	-1
-1	-1	-1

-1	1	-1
-1	1	-1
-1	1	-1

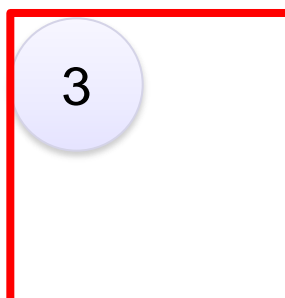
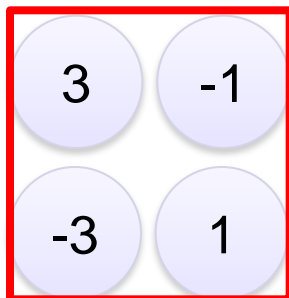
There are
25 3x3
filters.

Input_shape = (28 , 28 , 1)

28 x 28 pixels

1: black/white, 3: RGB

```
model2.add(MaxPooling2D( (2, 2) ))
```



input

Convolution

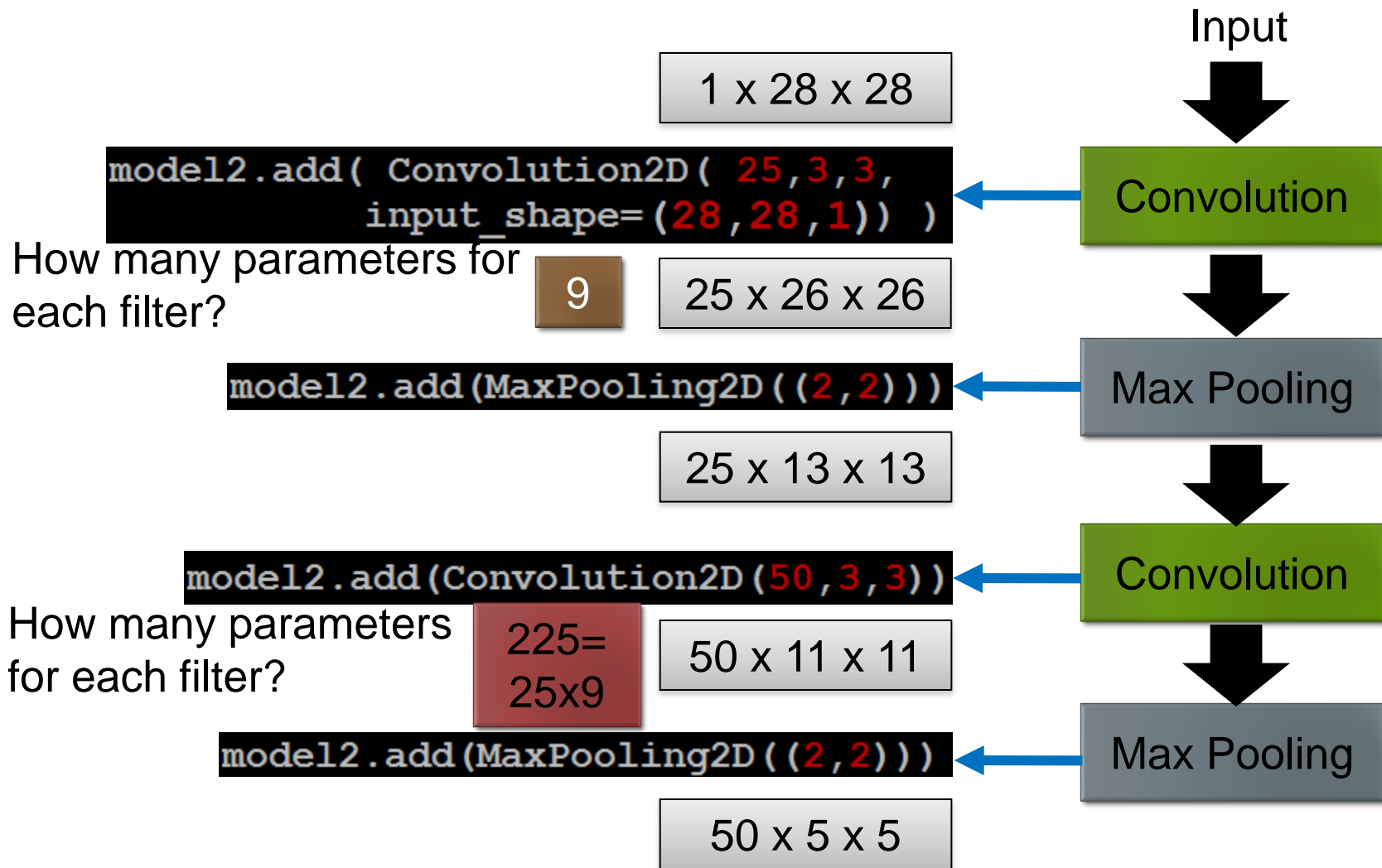
Max Pooling

Convolution

Max Pooling

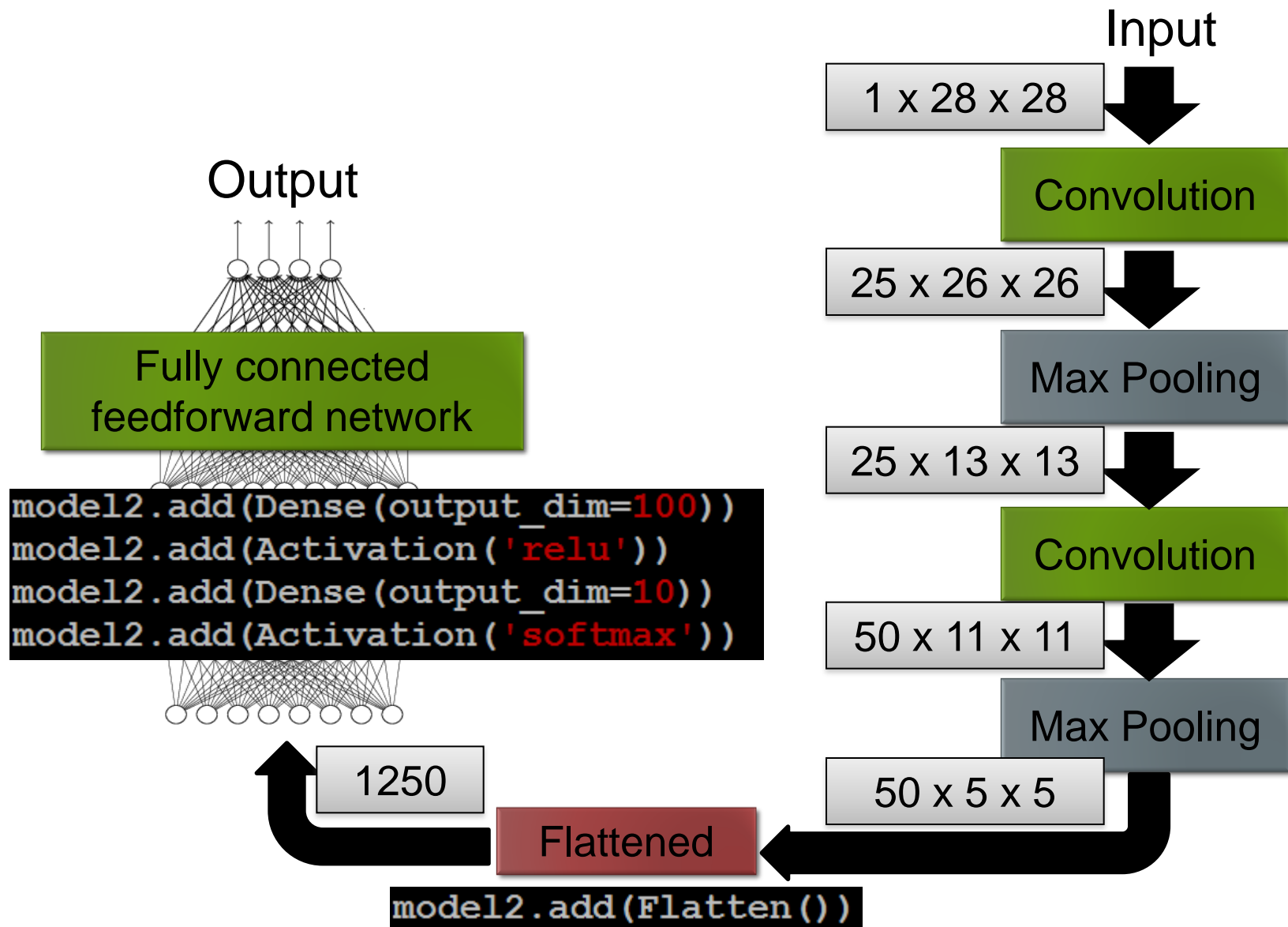
CNN in Keras

Only modified the *network structure* and *input format* (vector -> 3-D array)



CNN in Keras

Only modified the *network structure* and *input format (vector -> 3-D array)*



Advantage of convolution

- Multiplies the same weight everywhere on the image
- Suppose the filter is a feature finder that identifies the car, ie. outputs a comparatively high +ve number, when input is a car. It will find a car no matter where it is on the image
- Layers of convolution behave the same. First layer will find small, simple features anywhere on the image(eg. A line or edge). Next layer will find more complex features, and so on.
- Weight sharing leads to better generalization

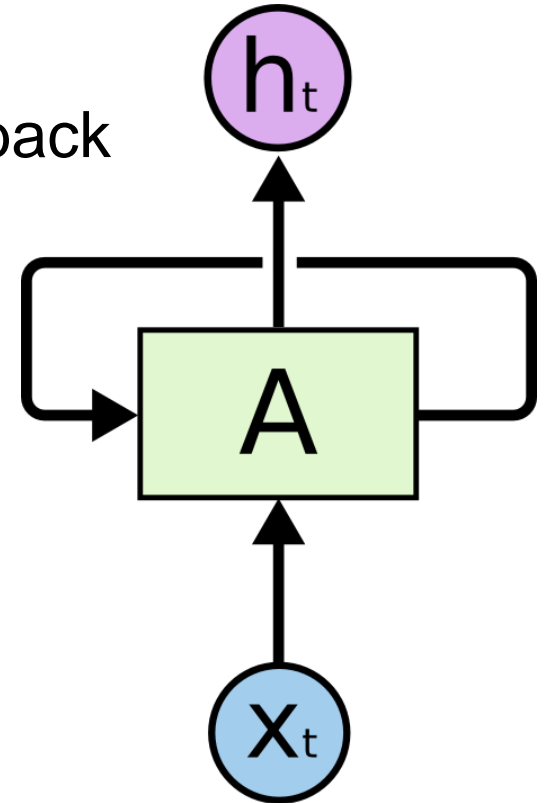
PART II:

NEURAL NETWORK

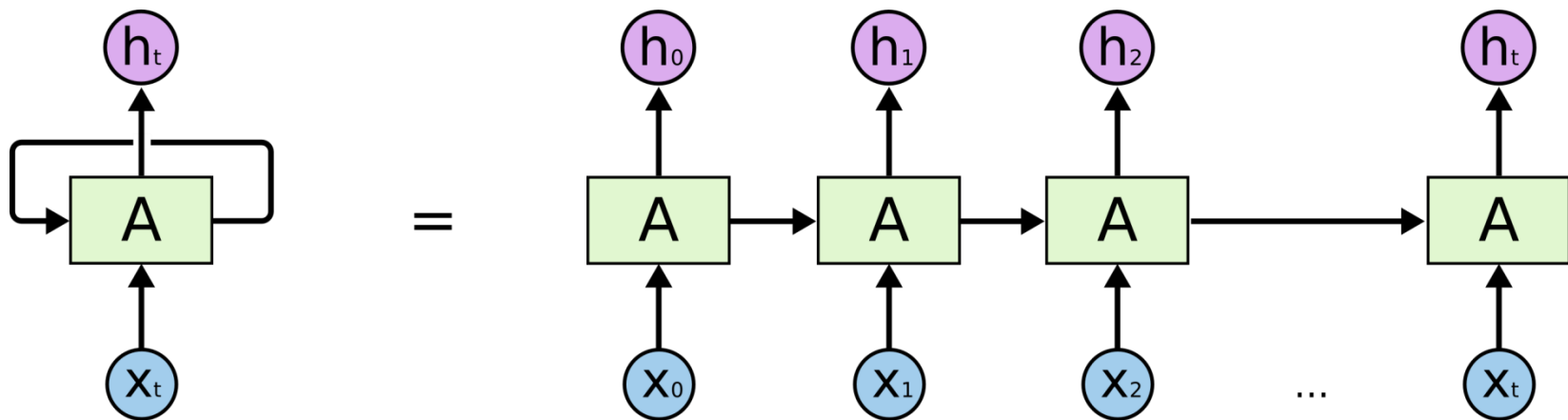
WITH MEMORY (RNN)

Construing from Memory

- RNN are the NN architectures with feedback



- An unrolled recurrent neural network

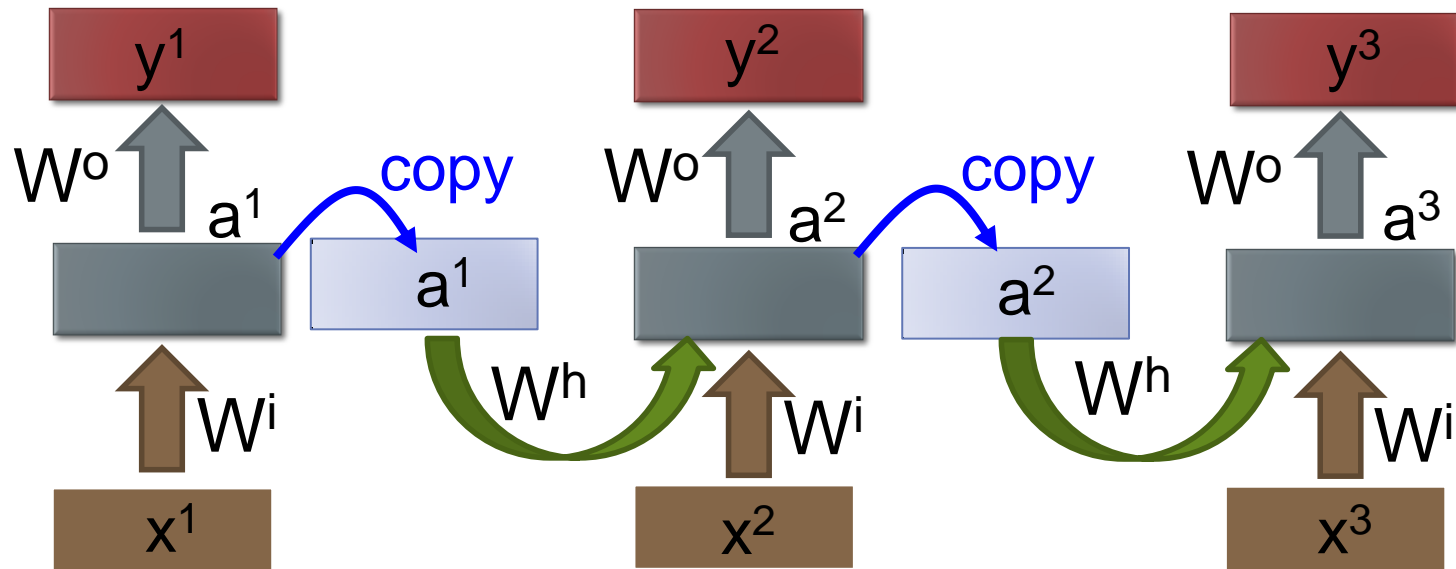


- Long memory vs Short memory
- Examples:

Clouds are in sky.

She grew up in France.....She speaks fluent French.

RNN

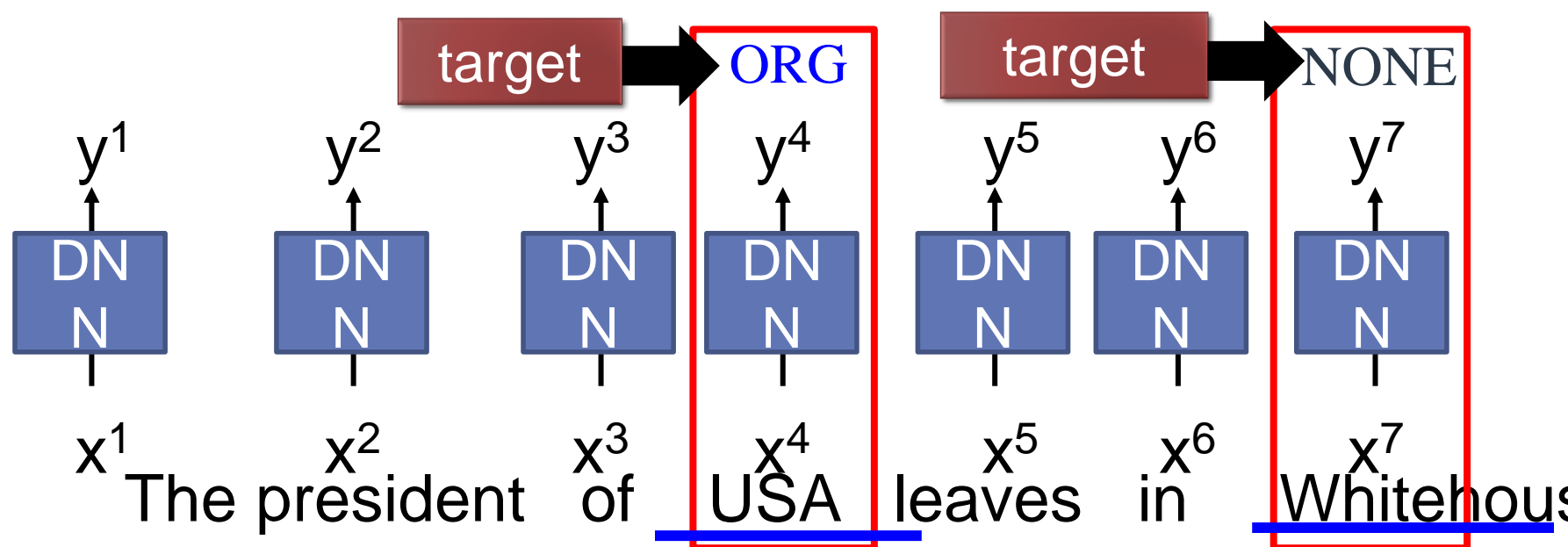


The same network is used again and again.

Output y^i depends on x^1, x^2, \dots, x^i

Neural Network needs Memory

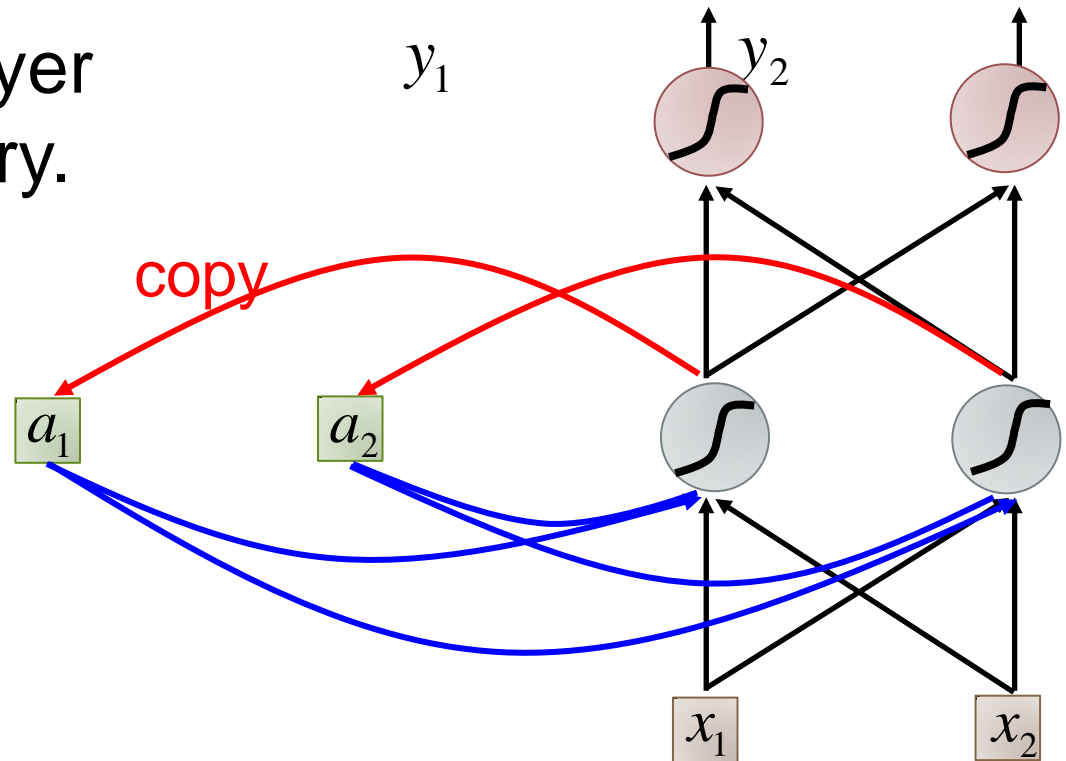
- Name Entity Recognition
 - Detecting entities like name of people, locations, organization, etc. in a sentence.



DNN needs memory!

Recurrent Neural Network (RNN)

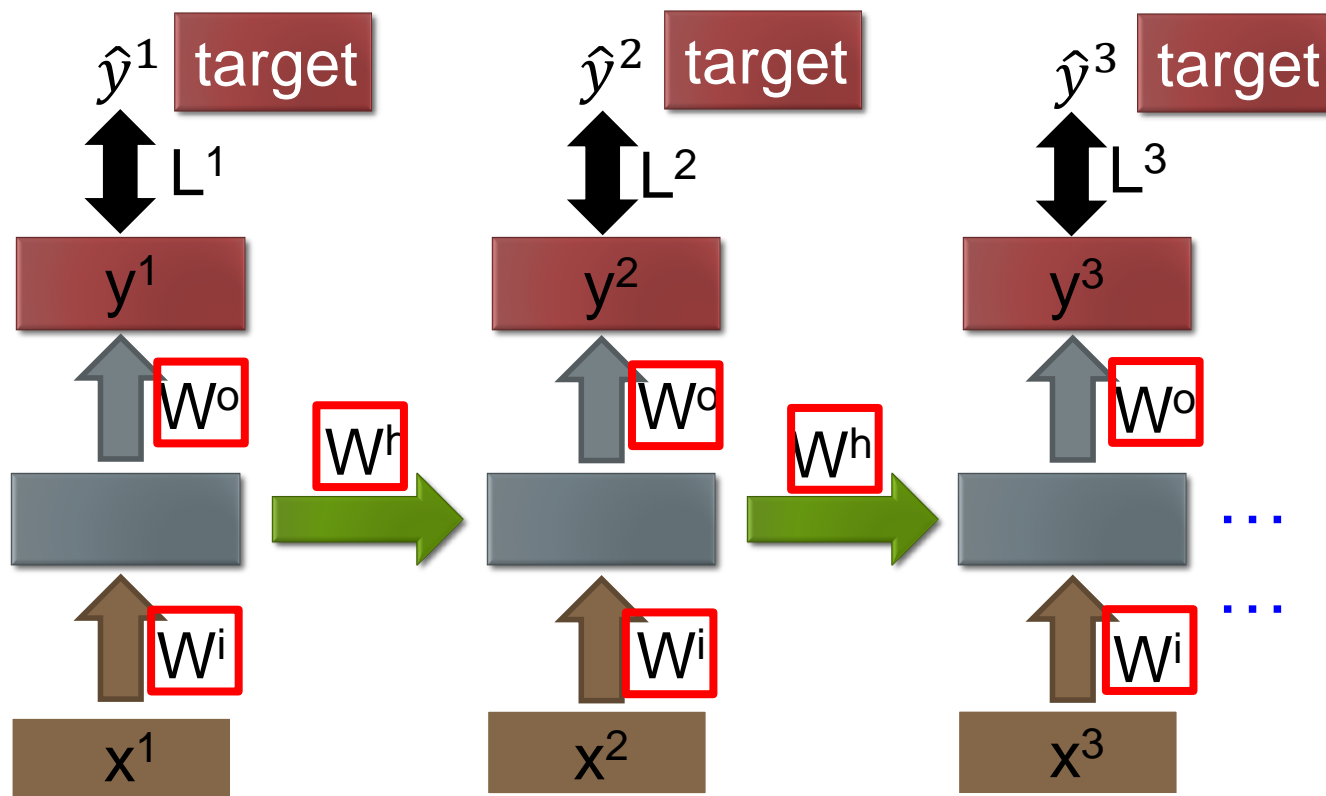
The output of hidden layer are stored in the memory.



- Memory can be considered as another input.
- Unfortunately, as that gap grows, RNNs become unable to learn to connect the information.

RNN

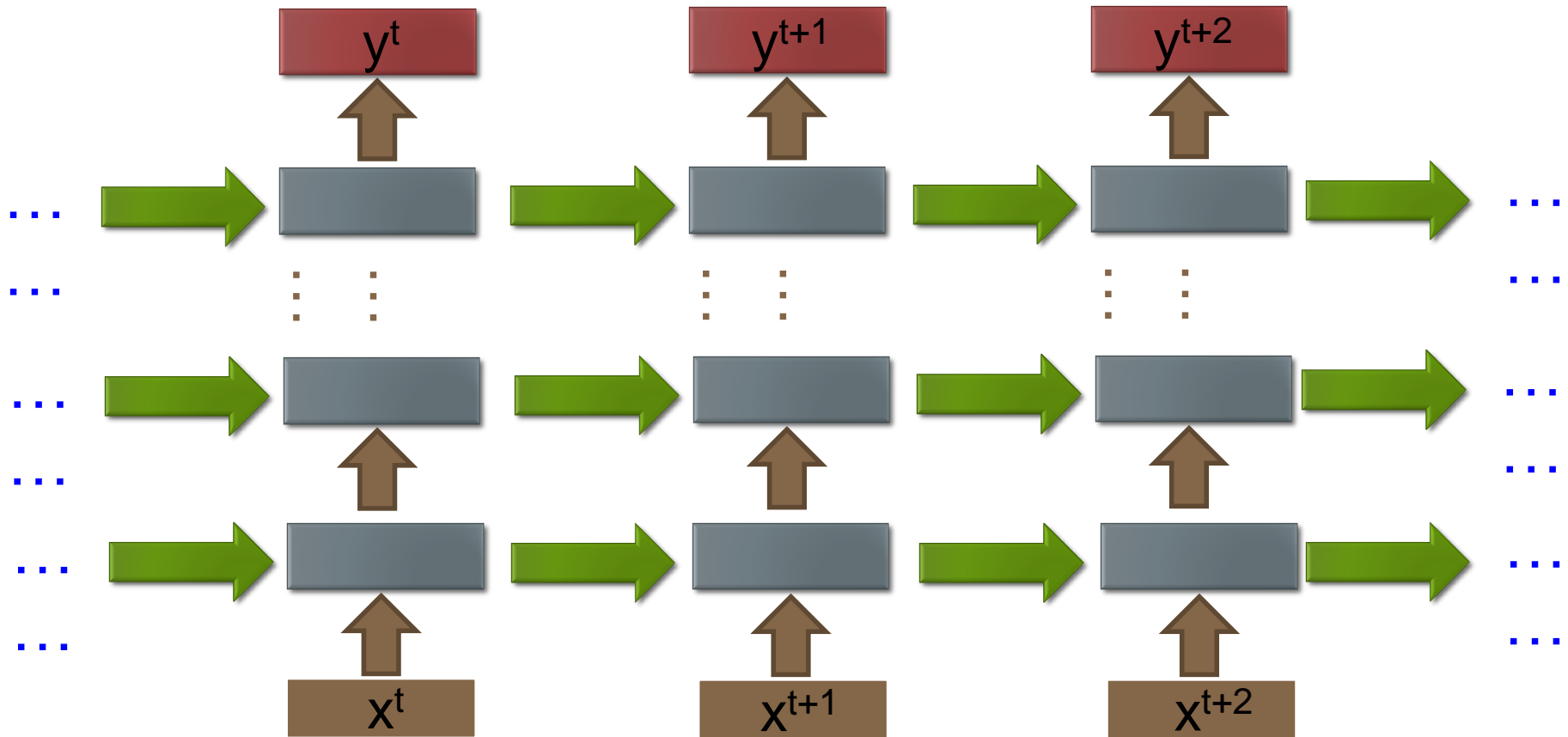
How to train?



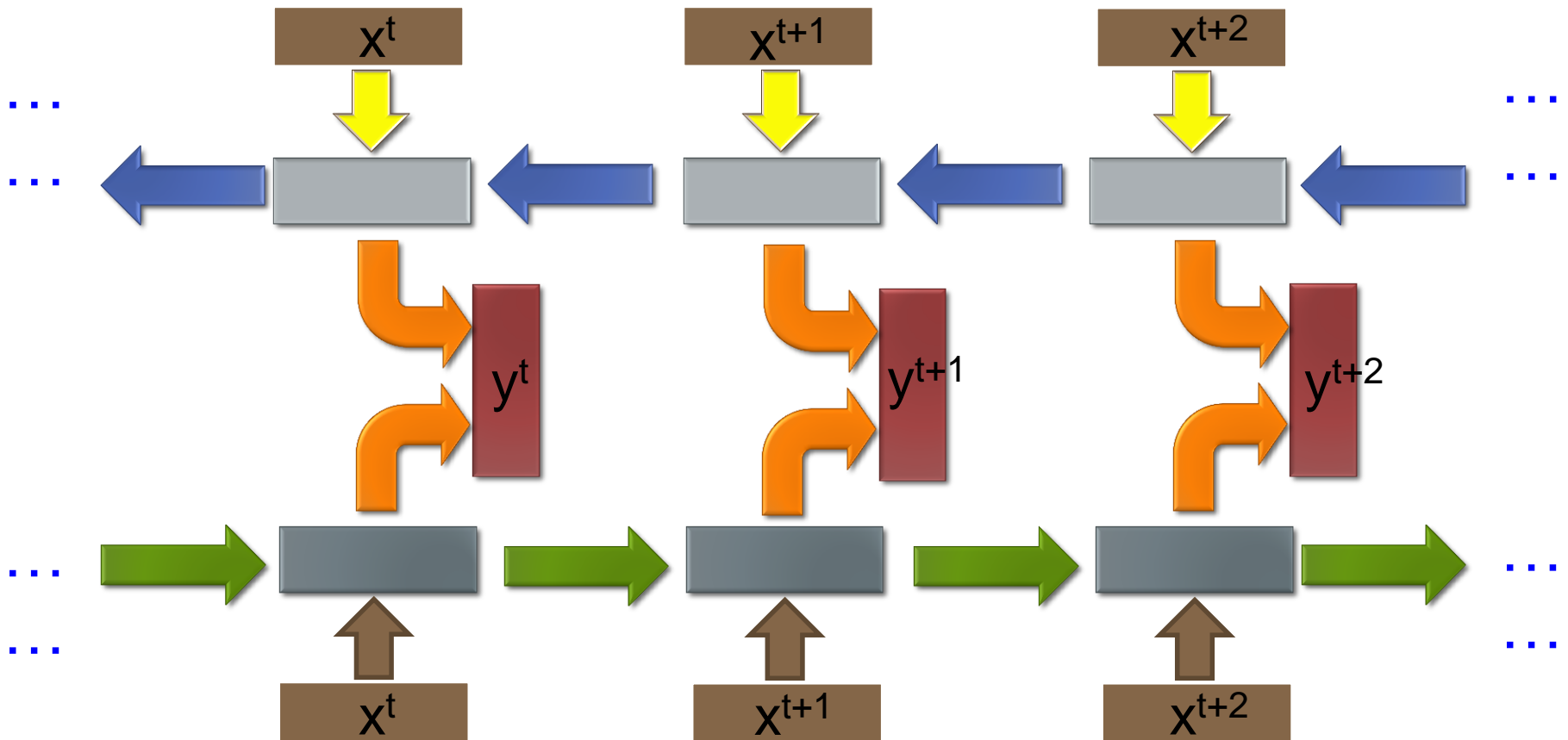
Find the network parameters to minimize the total cost:

Backpropagation through time (BPTT)

Of course it can be deep ...



Bidirectional RNN

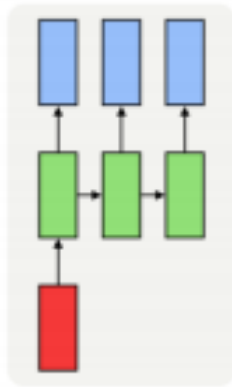


No size limitation

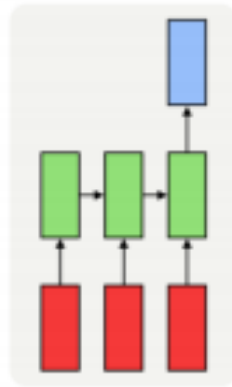
one to one



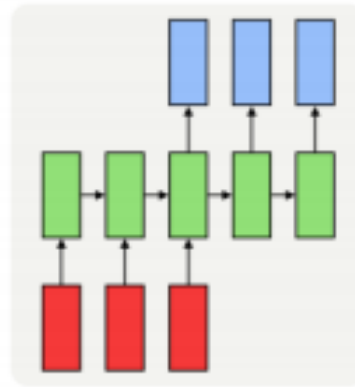
one to many



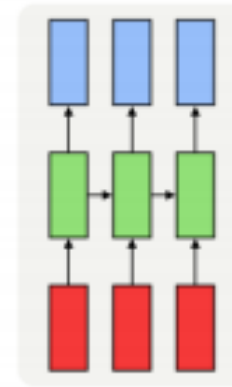
many to one



many to many



many to many



Vanishing or exploding gradient

$$\begin{array}{ll} w = 1 & \longrightarrow y^{1000} = 1 \\ w = 1.01 & \longrightarrow y^{1000} \approx 20000 \end{array}$$

Large
gradient

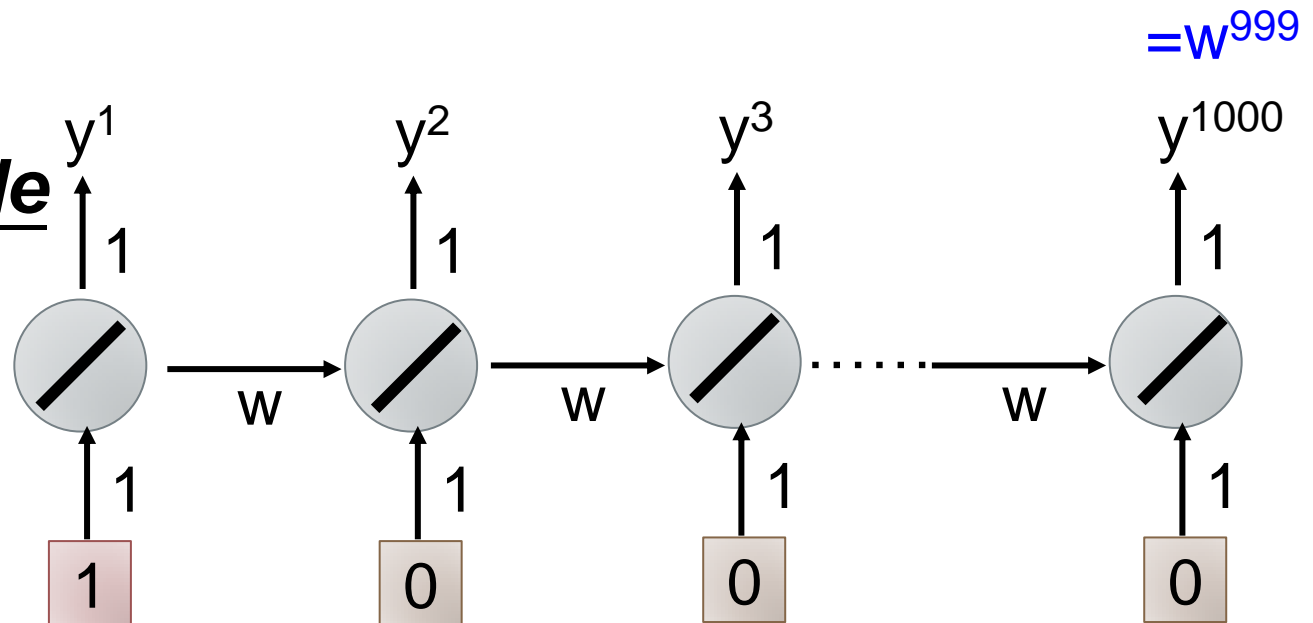
Small
Learning rate?

$$\begin{array}{ll} w = 0.99 & \longrightarrow y^{1000} \approx 0 \\ w = 0.01 & \longrightarrow y^{1000} \approx 0 \end{array}$$

small
gradient

Large
Learning rate?

Toy Example

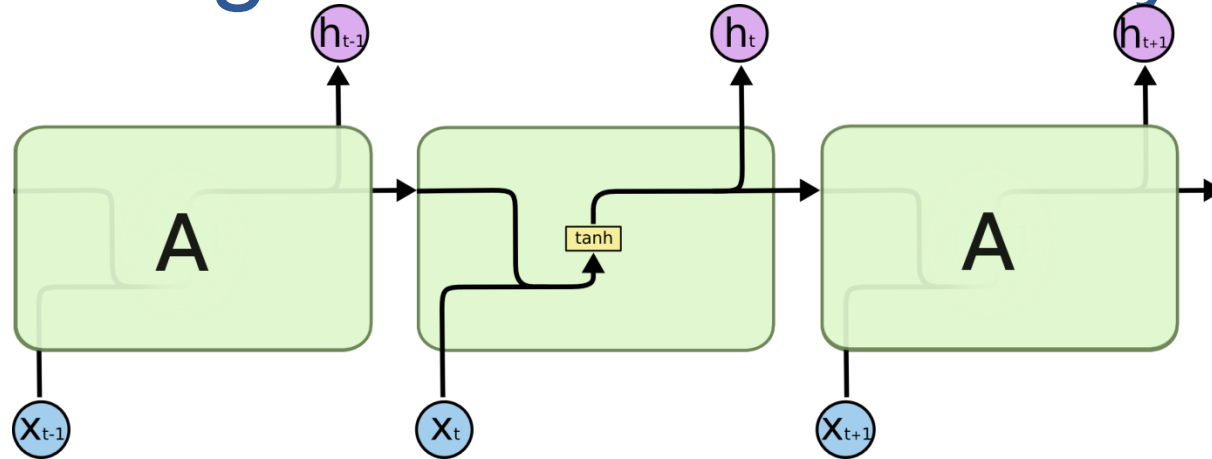


Helpful Techniques

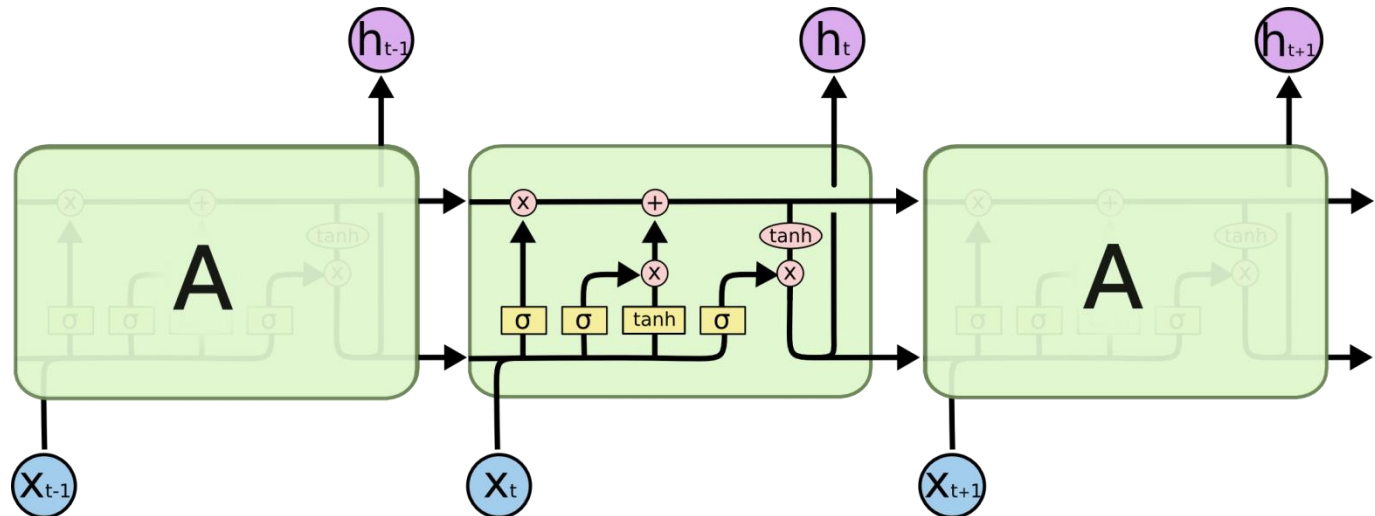
- Nesterov's Accelerated Gradient (NAG):
 - Advance momentum method
- RMS Prop
 - Advanced approach to give each parameter different learning rates
 - Considering the change of Second derivatives
- Long Short-term Memory (LSTM)
 - Can deal with gradient vanishing (not gradient explode)

LSTM: Long Short Term Memory

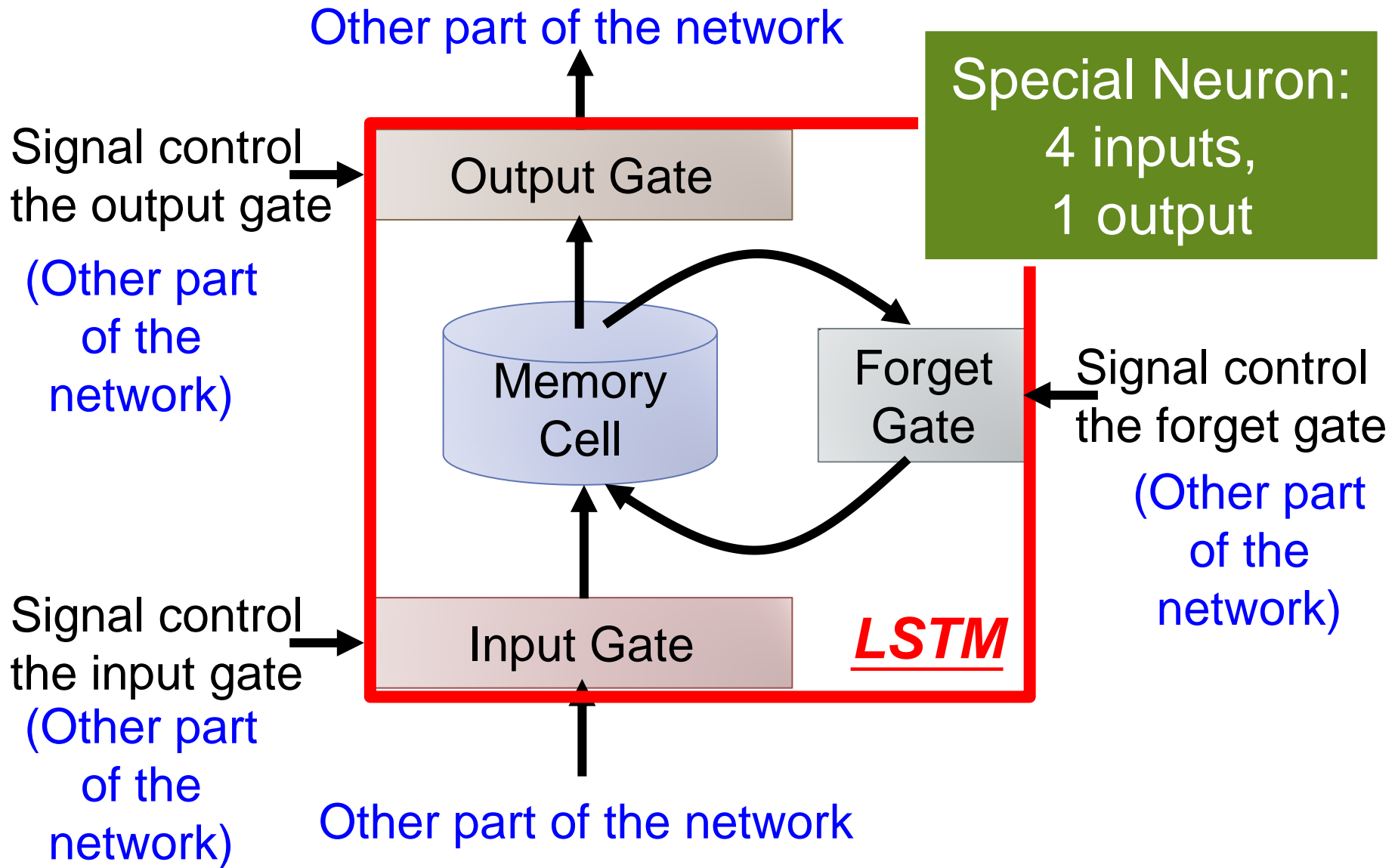
- RNN



- LSTM

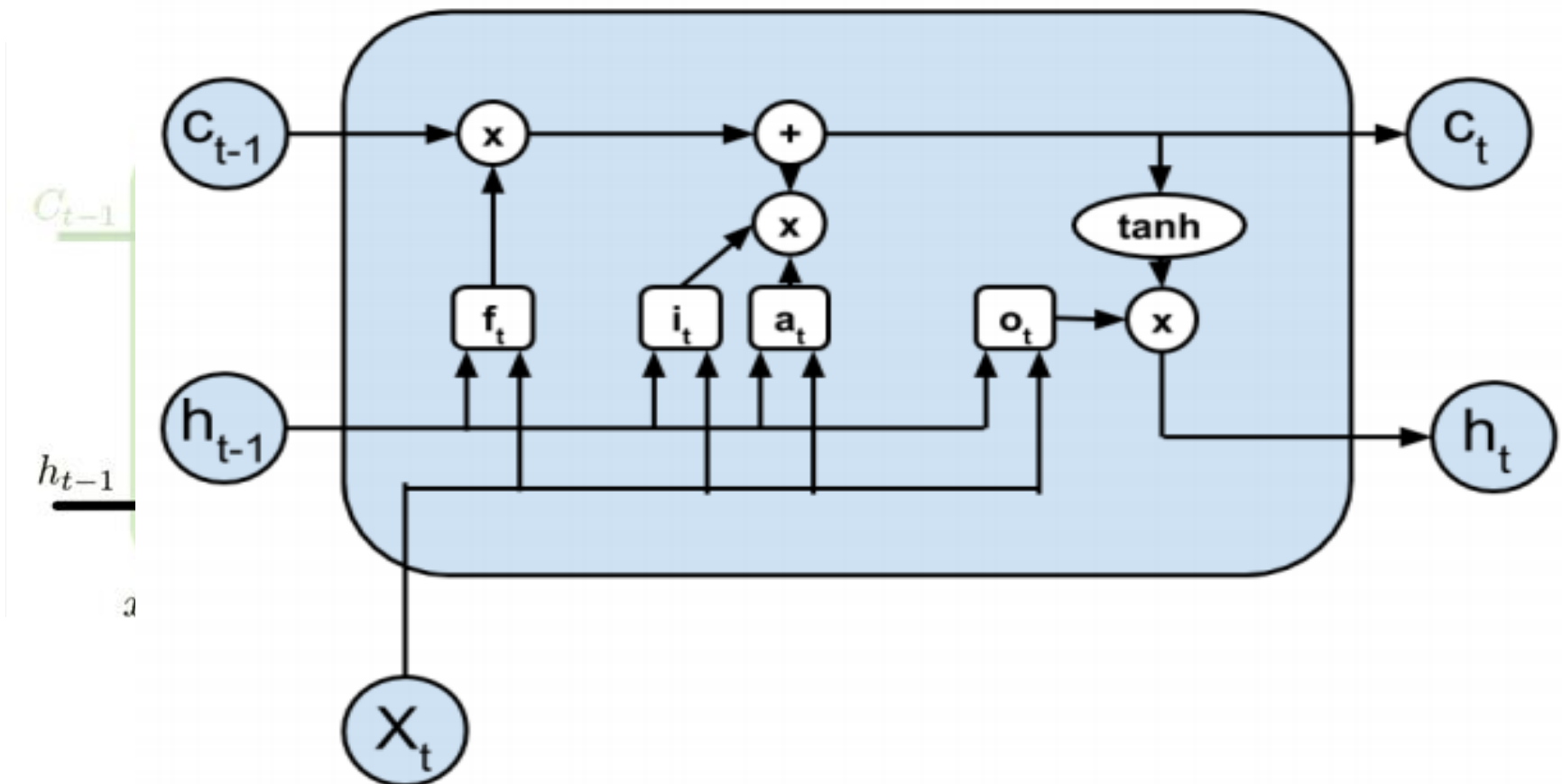


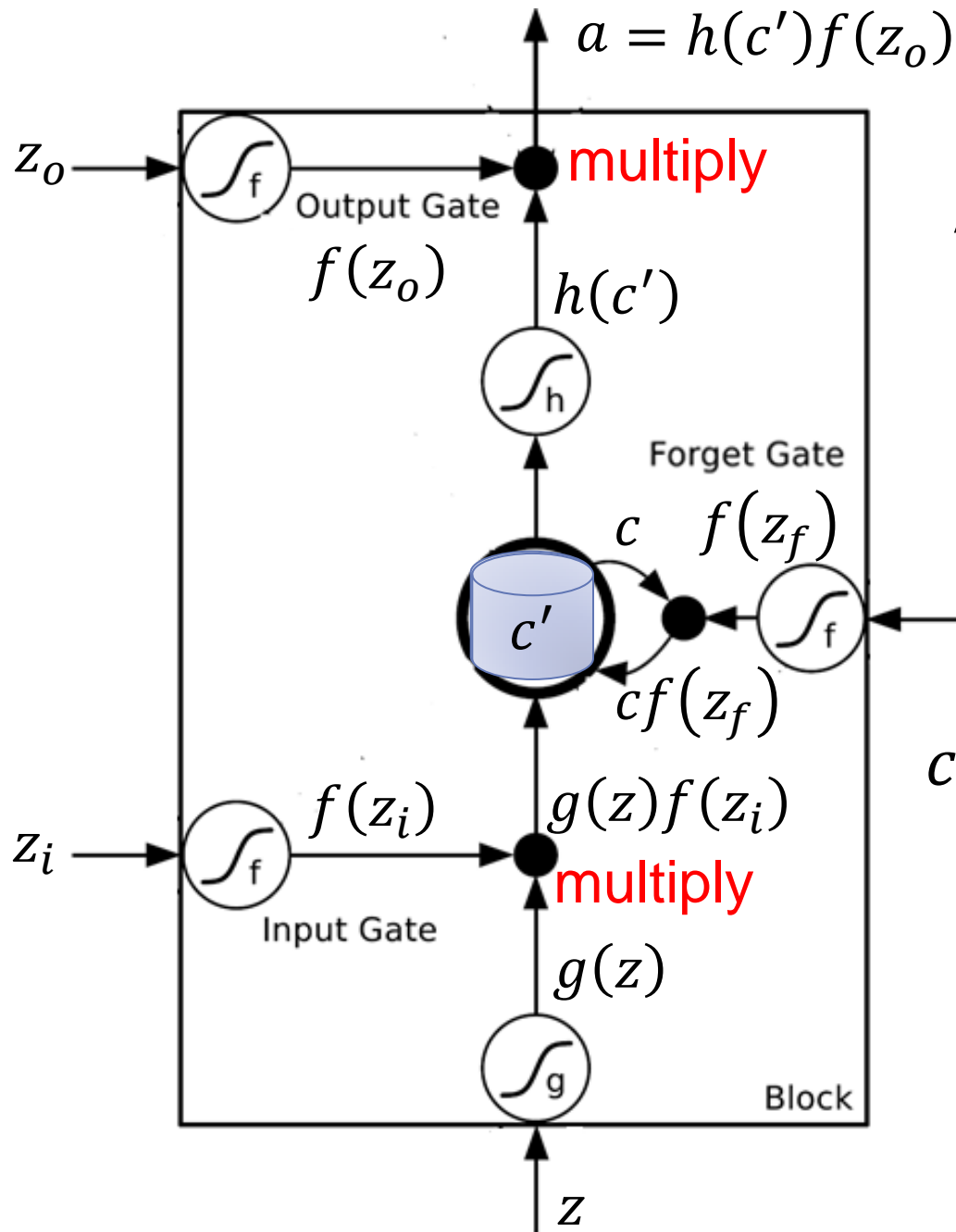
Long Short-term Memory (LSTM)



Core concepts behind LSTM

- Cell state (like a conveyor belt)
- Ability to remove or add information to the cell state





Activation function f is usually a sigmoid function

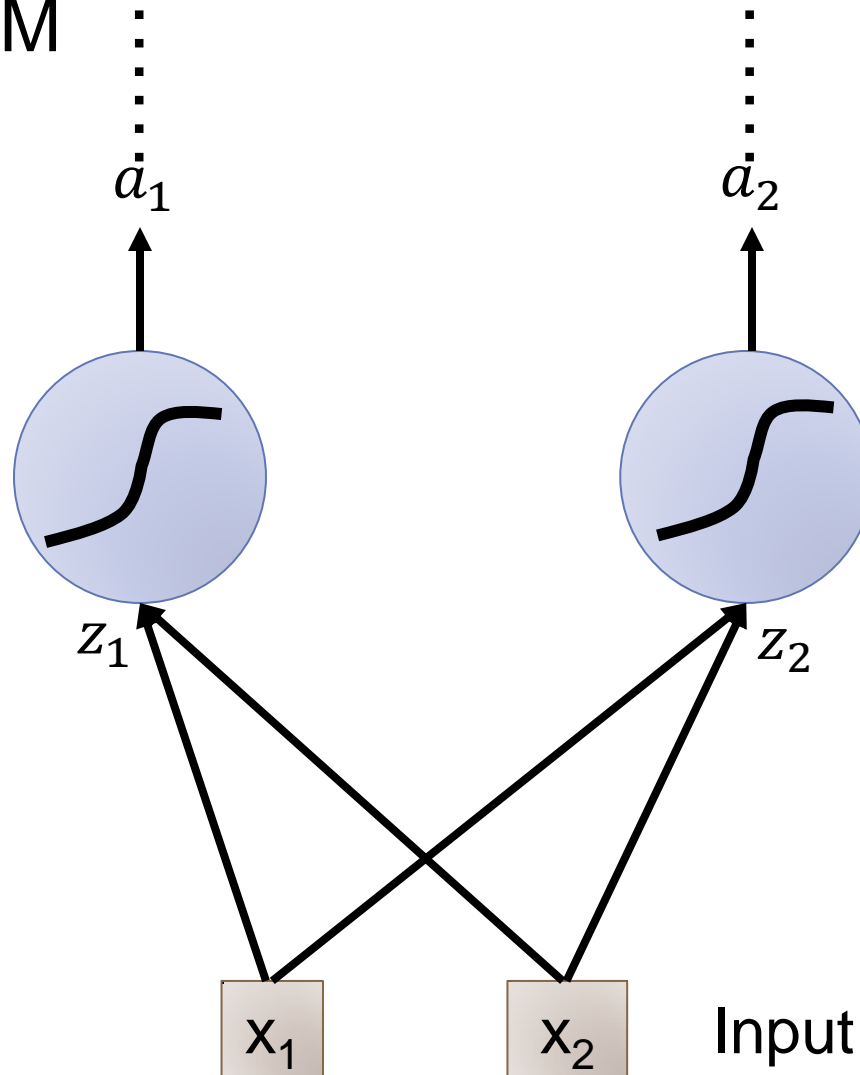
Between 0 and 1

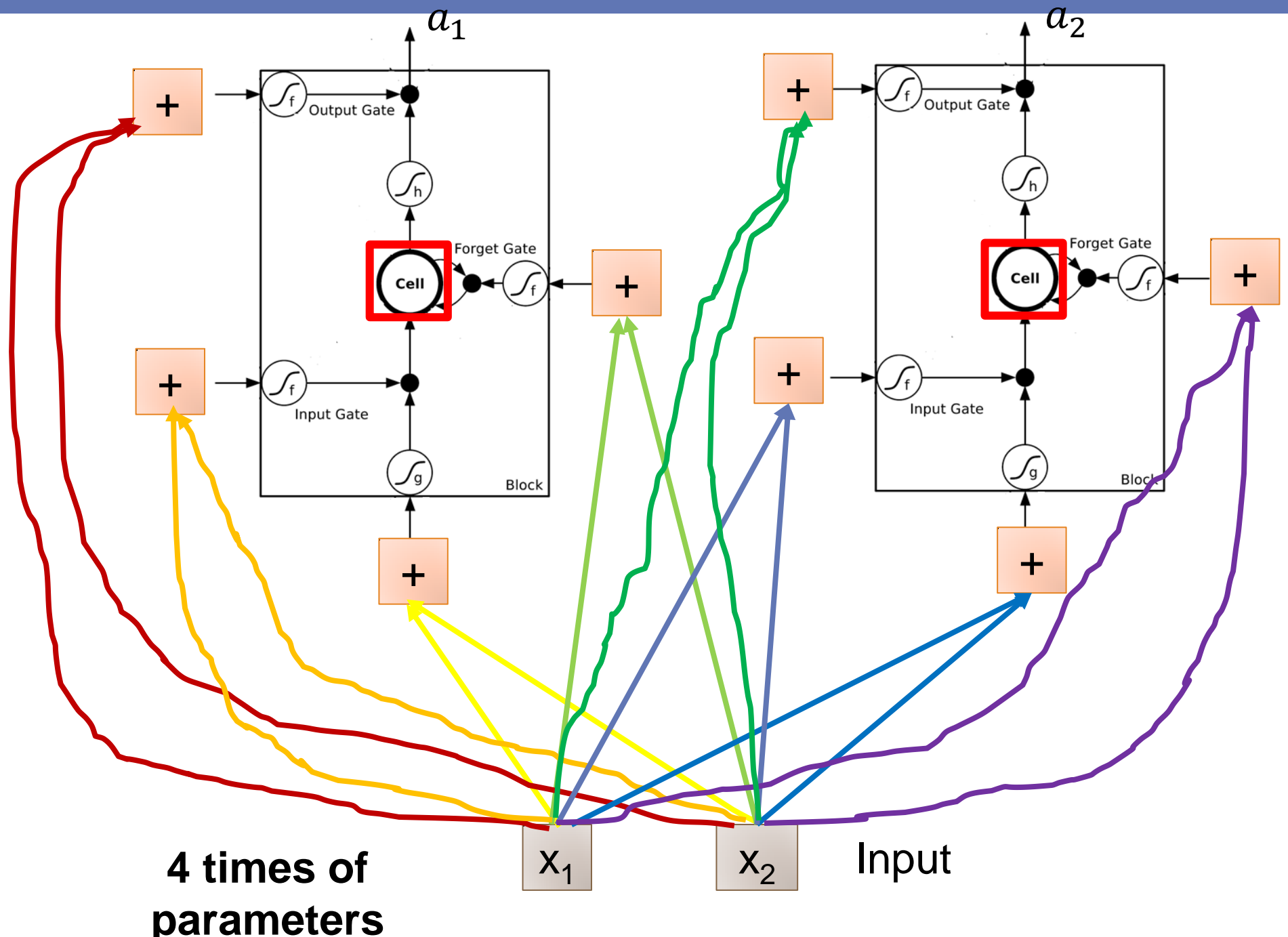
Mimic open and close gate

$$c' = g(z)f(z_i) + cf(z_f)$$

Original Network:

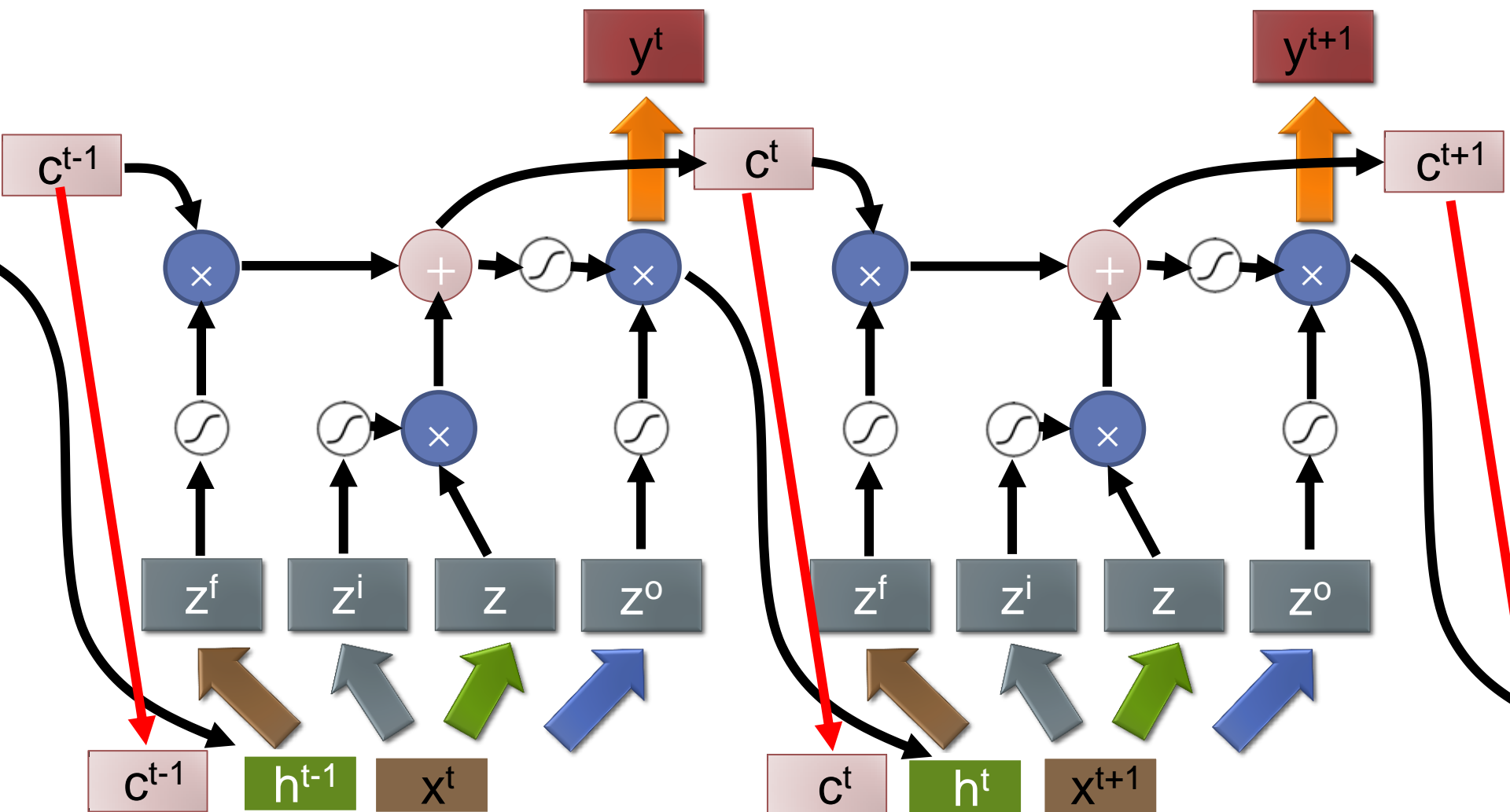
- Simply replace the neurons with LSTM





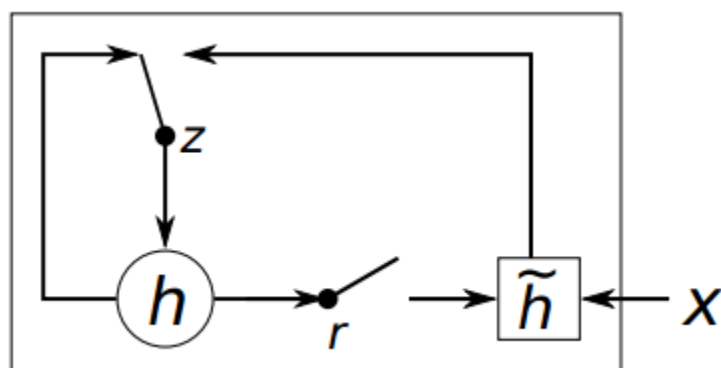
LSTM

Extension: "peephole"



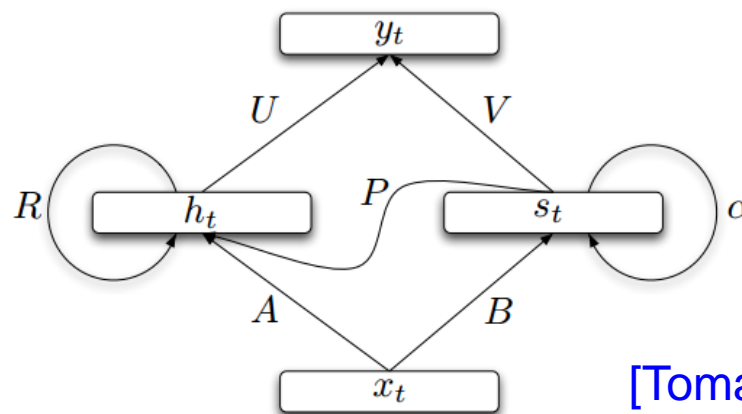
Other Simpler Alternatives

Gated Recurrent Unit (GRU)



[Cho, EMNLP'14]

Structurally Constrained Recurrent Network (SCRN)



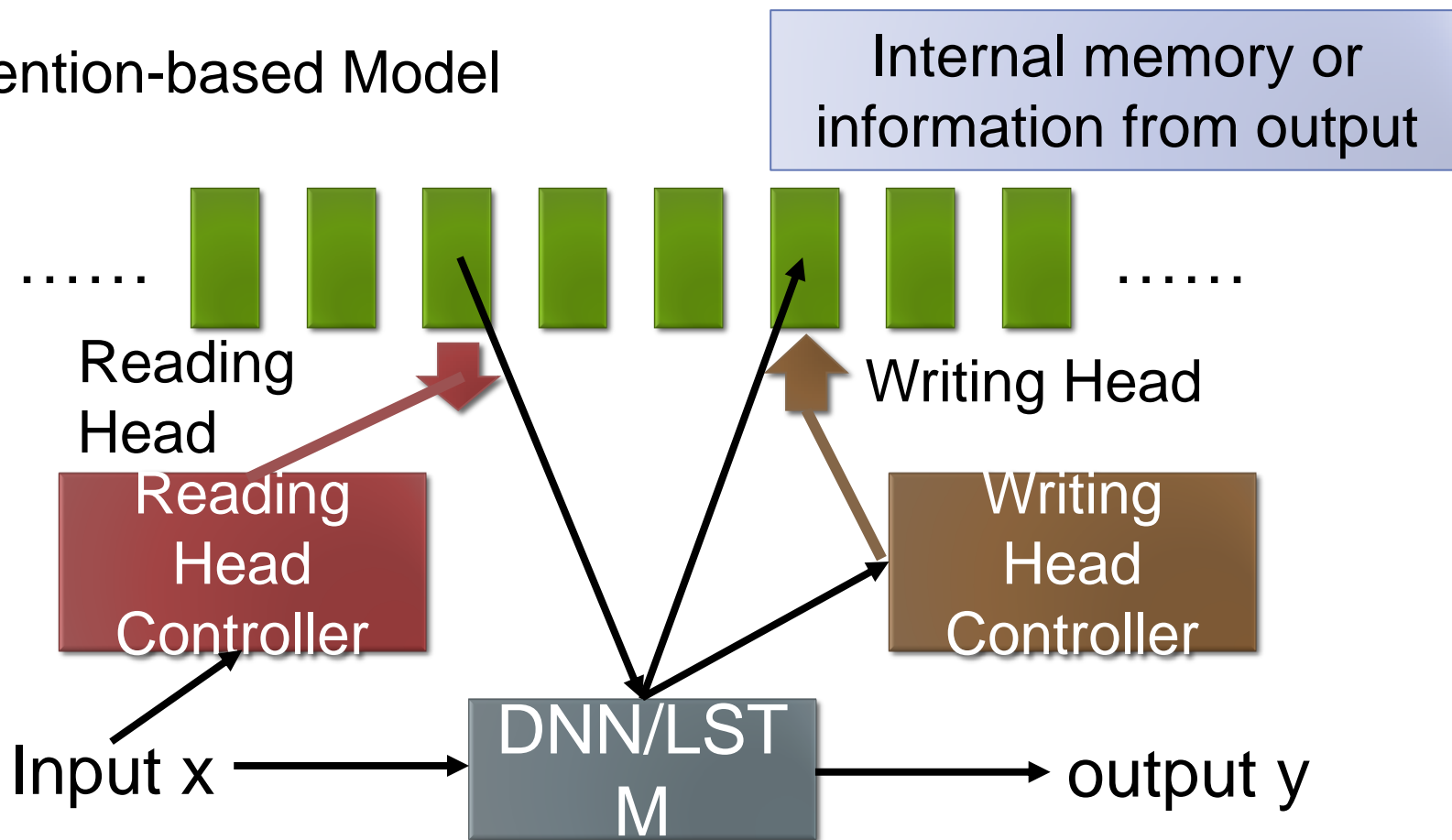
[Tomas
Mikolov,
ICLR'15]

Vanilla RNN Initialized with Identity matrix + ReLU activation function [Quoc V. Le, arXiv'15]

- Outperform or be comparable with LSTM in 4 different tasks

What is the next wave?

- Attention-based Model

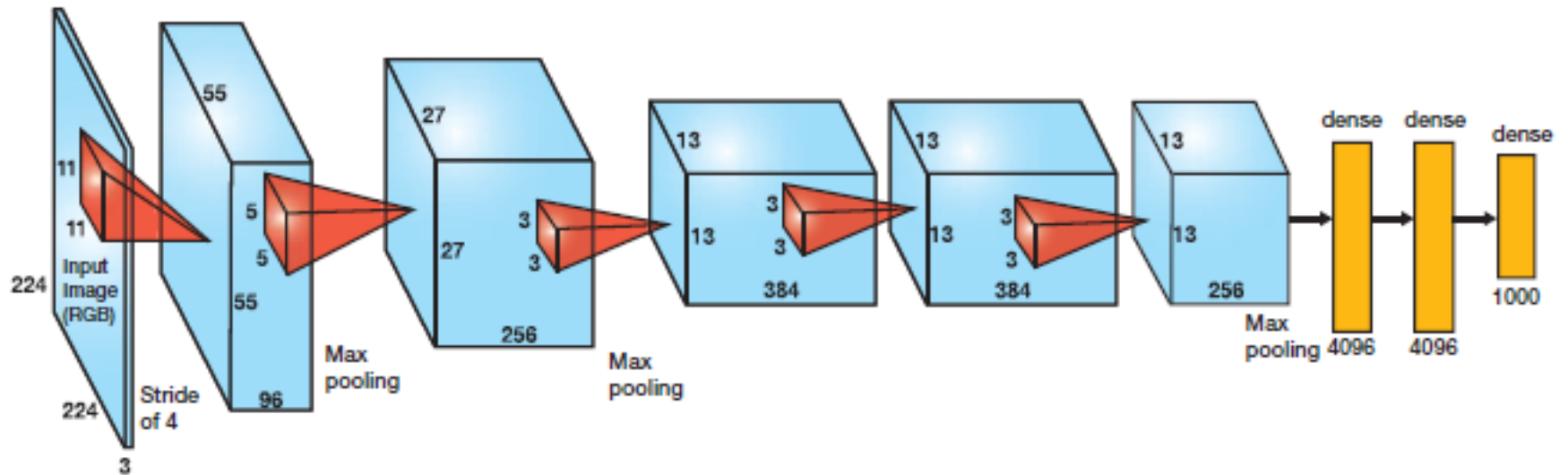


Already applied on speech recognition,
caption generation, QA, visual QA

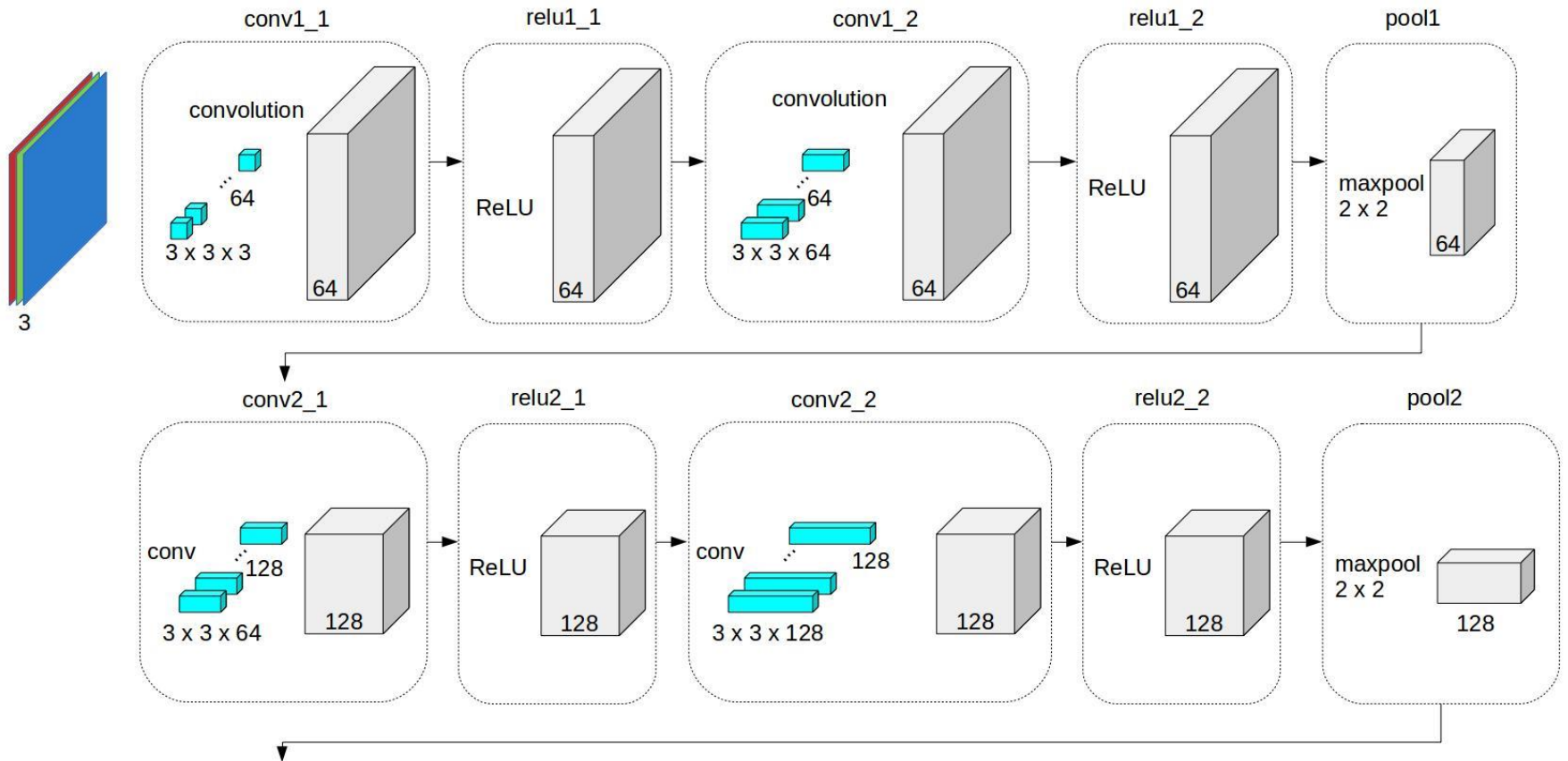
What is the next wave?

- **Attention-based Model**
- End-To-End Memory Networks. S. Sukhbaatar, A. Szlam, J. Weston, R. Fergus. arXiv Pre-Print, 2015.
- Neural Turing Machines. Alex Graves, Greg Wayne, Ivo Danihelka. arXiv Pre-Print, 2014
- Ask Me Anything: Dynamic Memory Networks for Natural Language Processing. Kumar et al. arXiv Pre-Print, 2015
- Neural Machine Translation by Jointly Learning to Align and Translate. D. Bahdanau, K. Cho, Y. Bengio; International Conference on Representation Learning 2015.
- Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. Kelvin Xu et. al.. arXiv Pre-Print, 2015.
- Attention-Based Models for Speech Recognition. Jan Chorowski, Dzmitry Bahdanau, Dmitriy Serdyuk, Kyunghyun Cho, Yoshua Bengio. arXiv Pre-Print, 2015.
- Recurrent models of visual attention. V. Mnih, N. Hees, A. Graves and K. Kavukcuoglu. In NIPS, 2014.
- A Neural Attention Model for Abstractive Sentence Summarization. A. M. Rush, S. Chopra and J. Weston. EMNLP 2015.

AlexNet

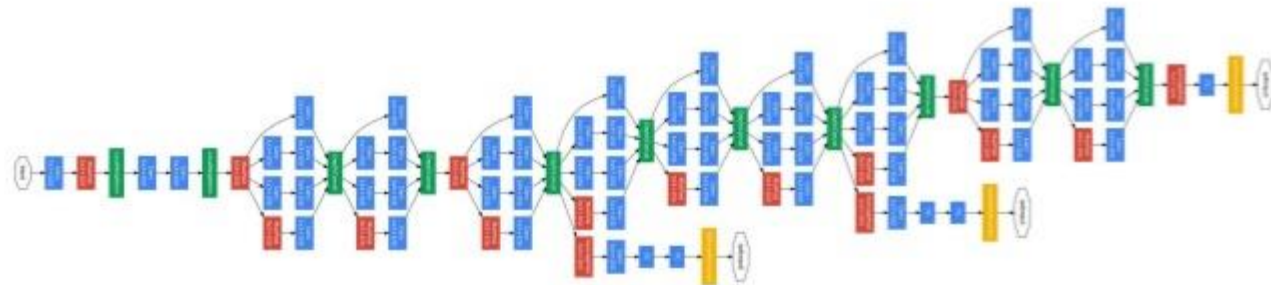


VGG



GoogleNet

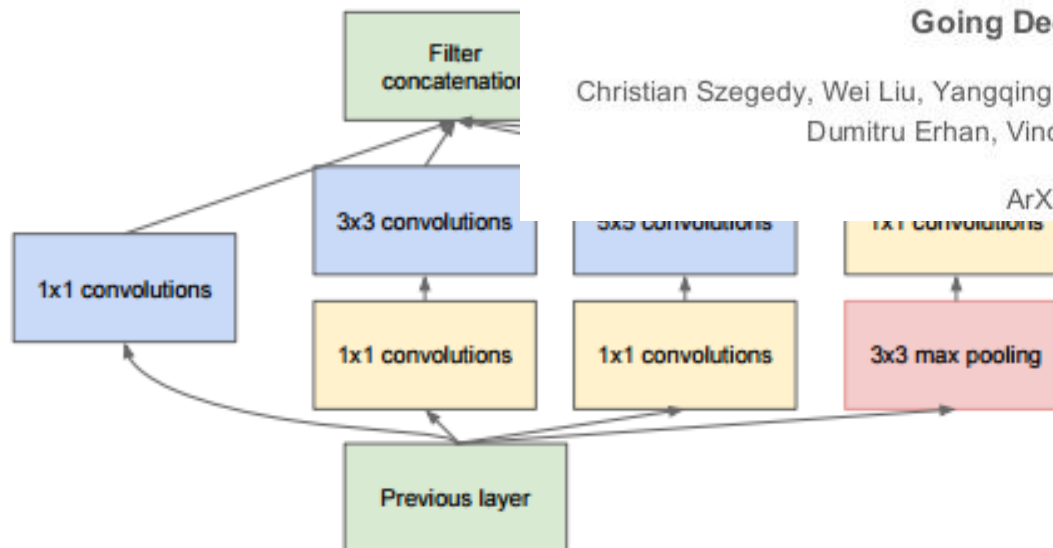
The Inception Architecture (GoogLeNet, 2014)



Going Deeper with Convolutions

Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, Andrew Rabinovich

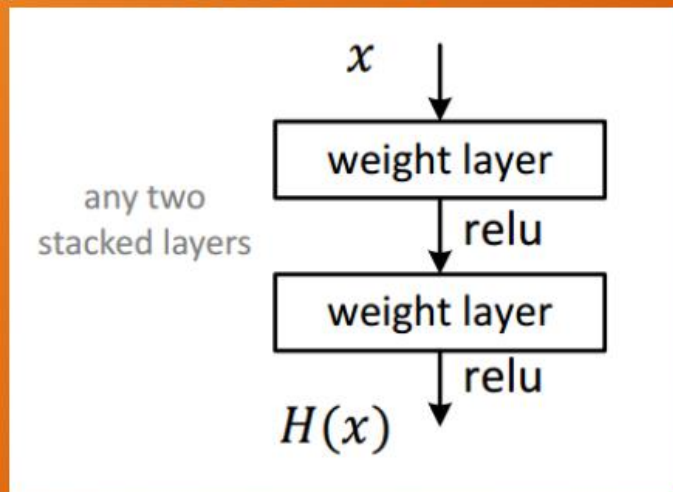
ArXiv 2014, CVPR 2015



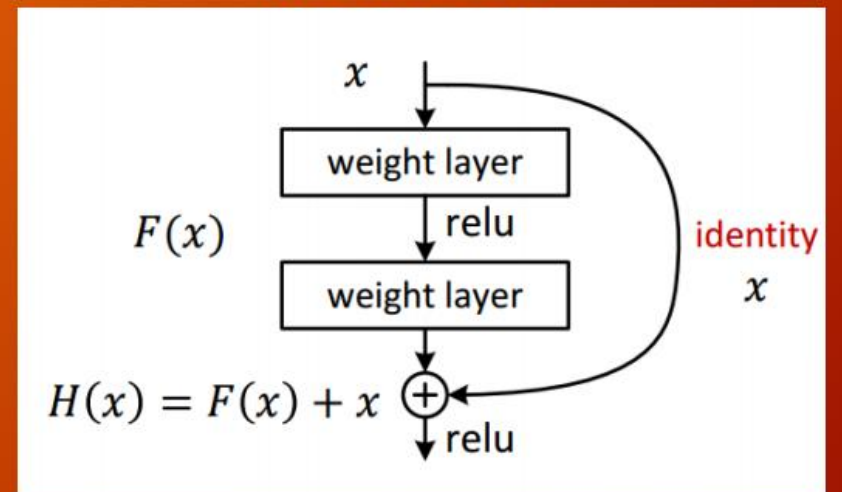
(b) Inception module with dimension reductions

ResNet

- Not just stacking layers on top of each other...
„Normal”



- „Residual”

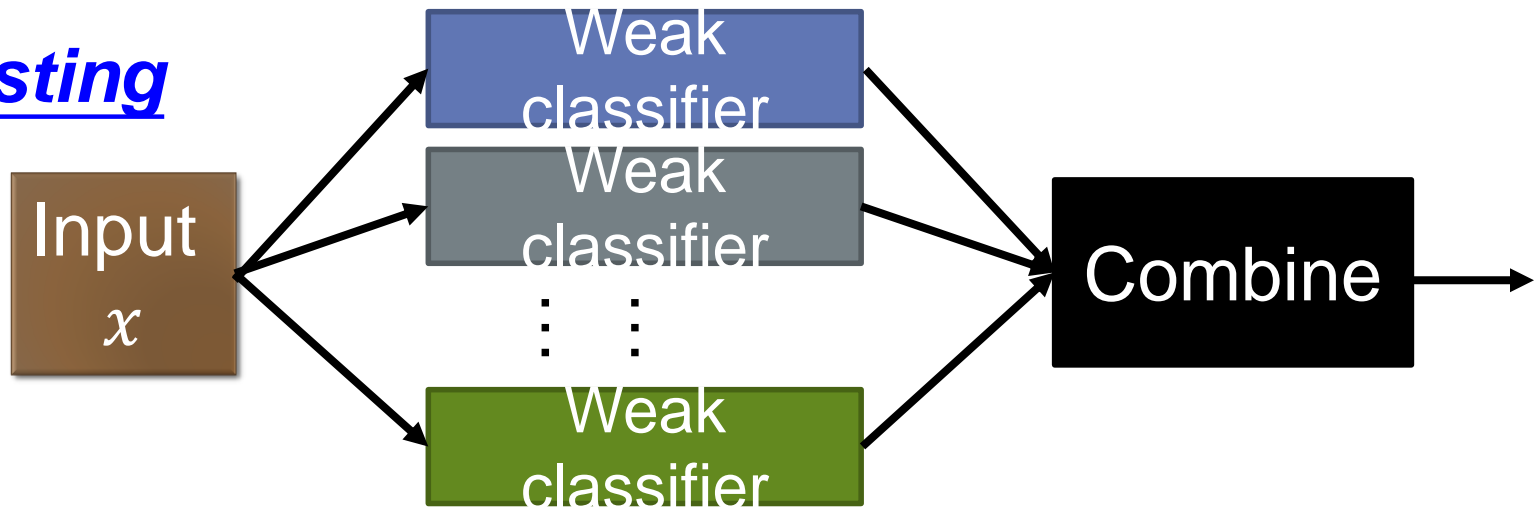


R-FCN

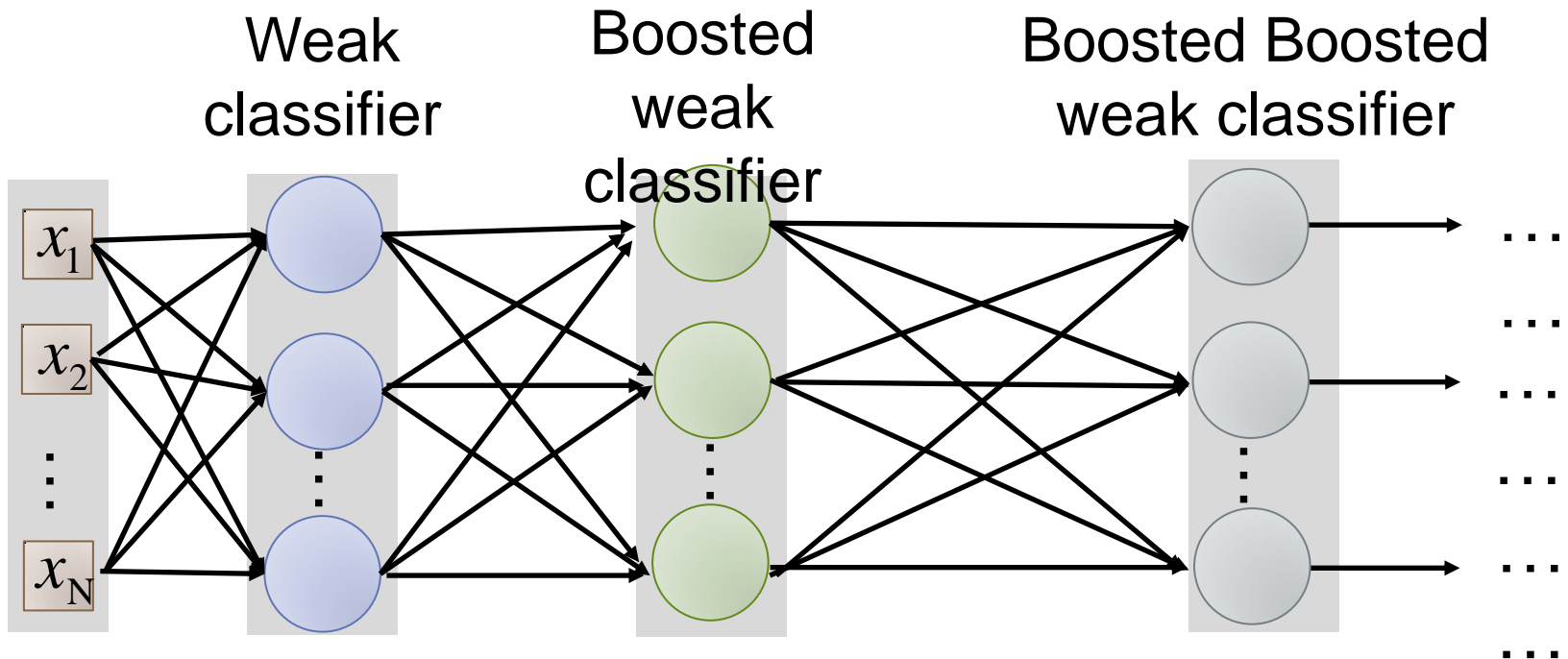
- Object Detection by Region based fully connected CNN
 - Regression to boundary box coordinates + object classification inside the bounding box
- Networks are rarely trained from scratch. Pre-trained weights are often available for popular networks. Proceed with fine tuning (transform learning)

THANK YOU
~~FOR YOUR ATTENTION!~~

Boosting

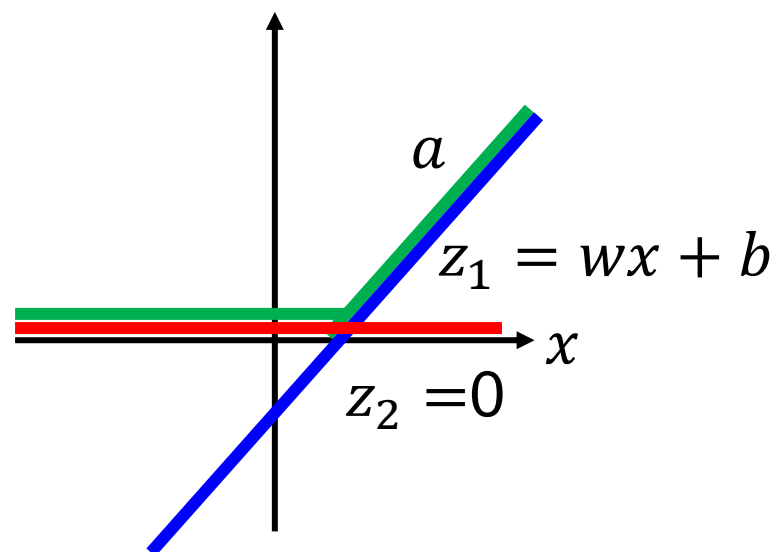
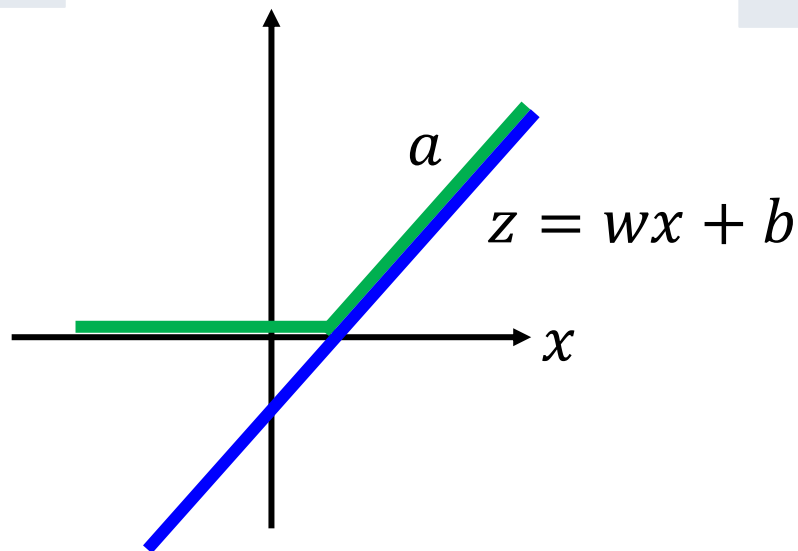
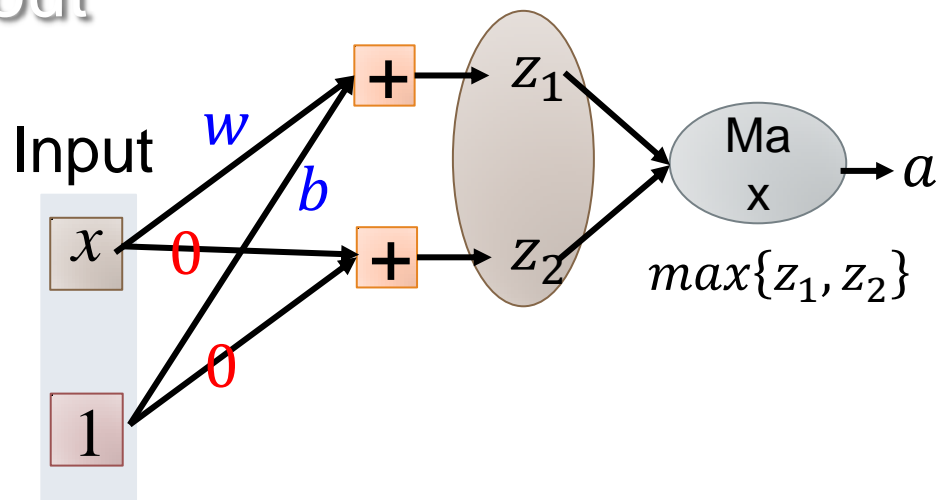
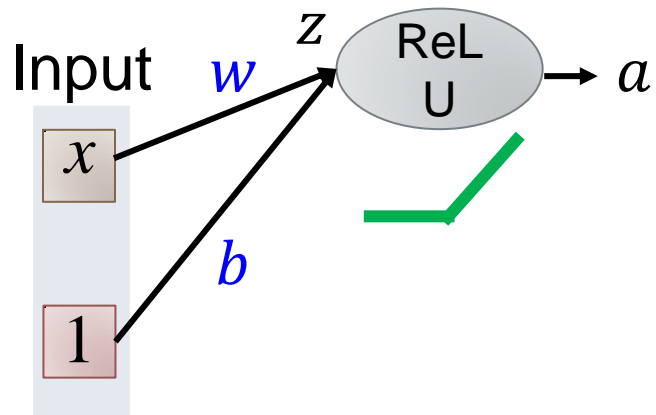


Deep Learning



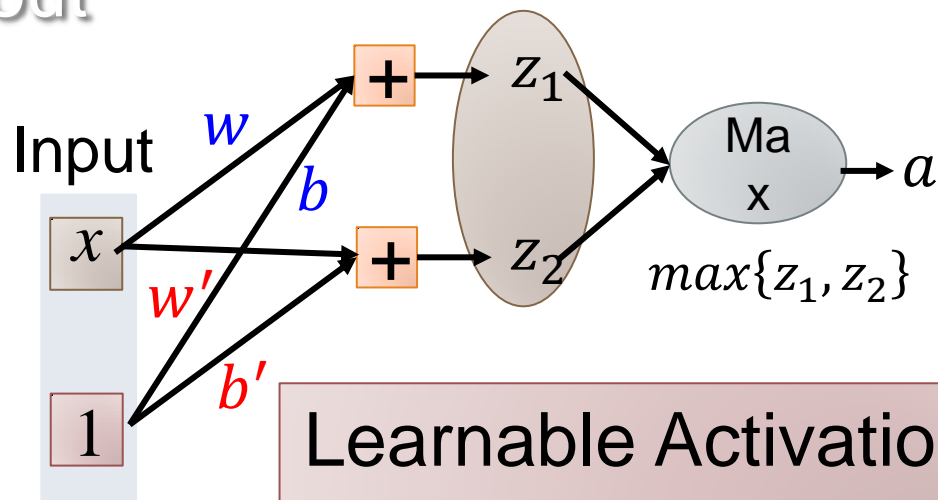
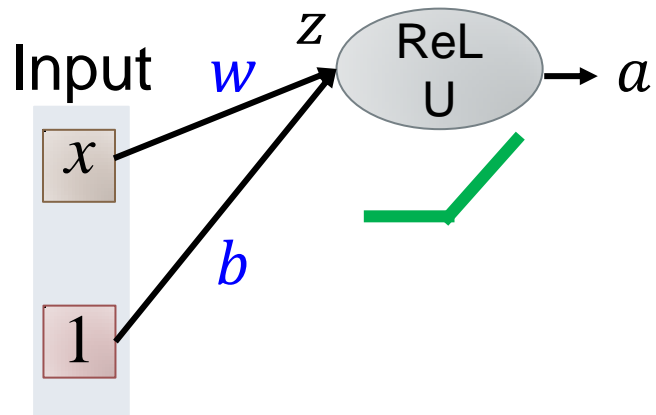
Maxout

ReLU is a special cases of
Maxout



Maxout

ReLU is a special cases of
Maxout



Learnable Activation
Function

