



# BIG DATA and AI for business

## Deep Learning (3)

**Decisions, Operations & Information Technologies**  
**Robert H. Smith School of Business**  
**Fall, 2020**



# Variants of Neural Networks

---

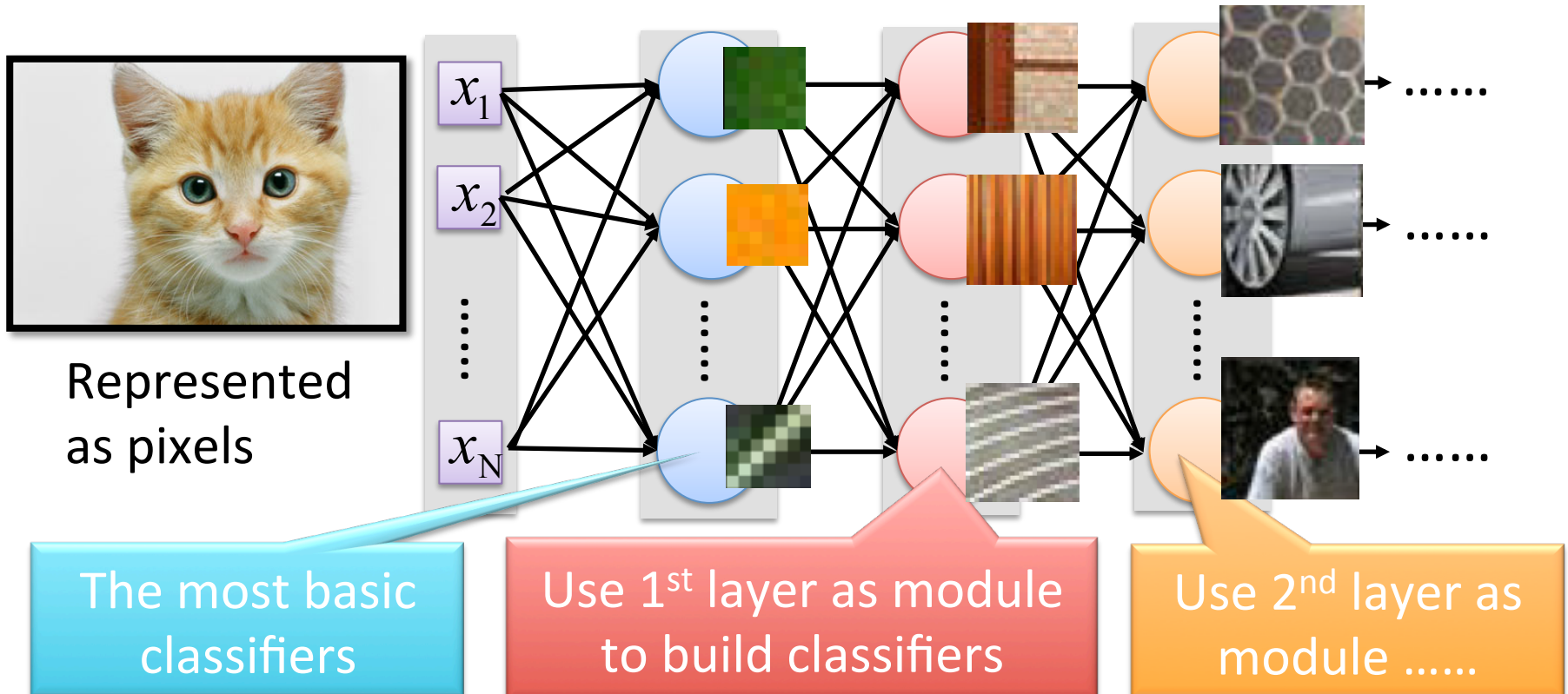
Convolutional Neural  
Network (CNN)

Widely used in  
image processing

Recurrent Neural Network  
(RNN)

# Why CNN for Image?

[Zeiler, M. D., *ECCV 2014*]



Can the network be simplified by considering the properties of images?

# Why CNN for Image

- Some patterns are much smaller than the whole image

A neuron does not have to see the whole image to discover the pattern.

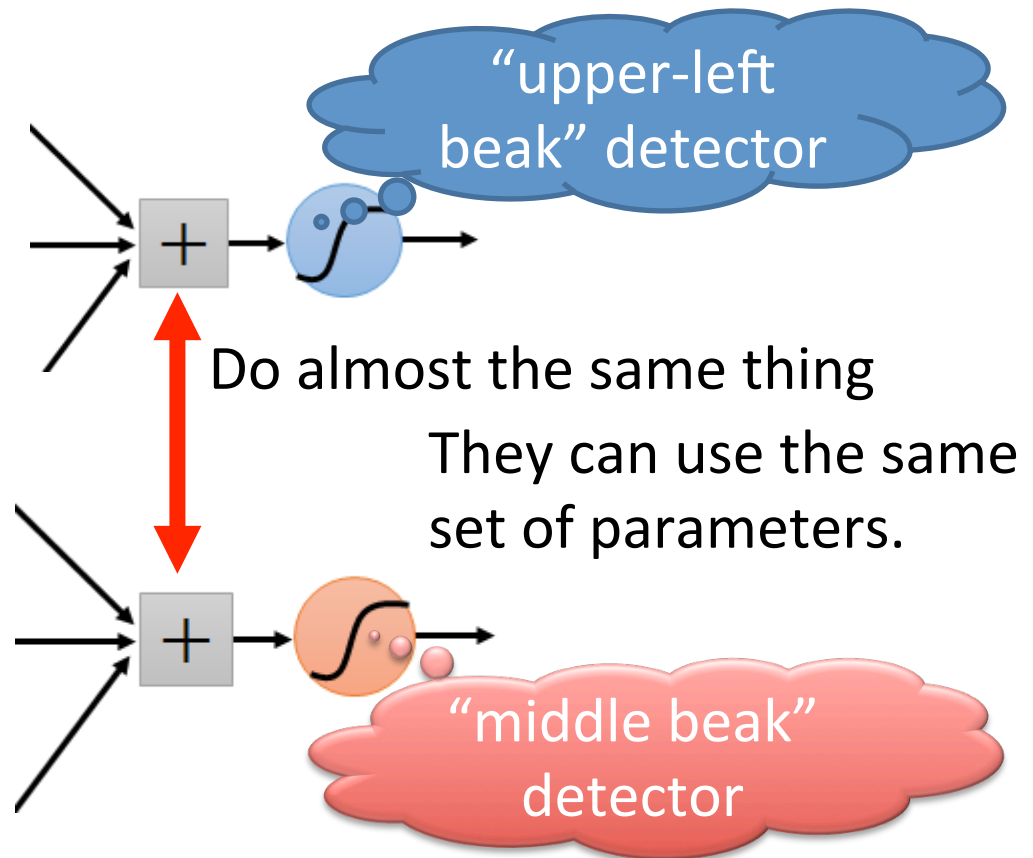
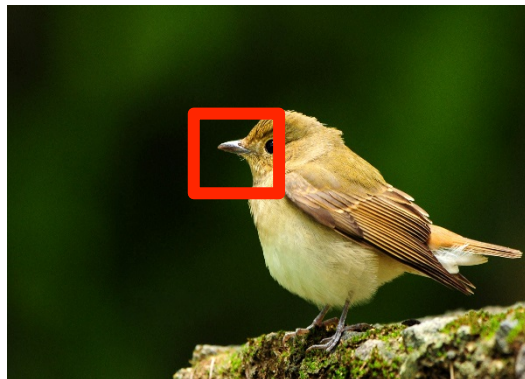
Connecting to small region with less parameters





# Why CNN for Image

- The same patterns appear in different regions.



# Why CNN for Image

- Subsampling the pixels will not change the object bird

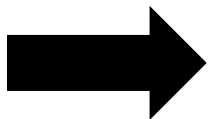


subsampling

bird



We can subsample the pixels to make image smaller

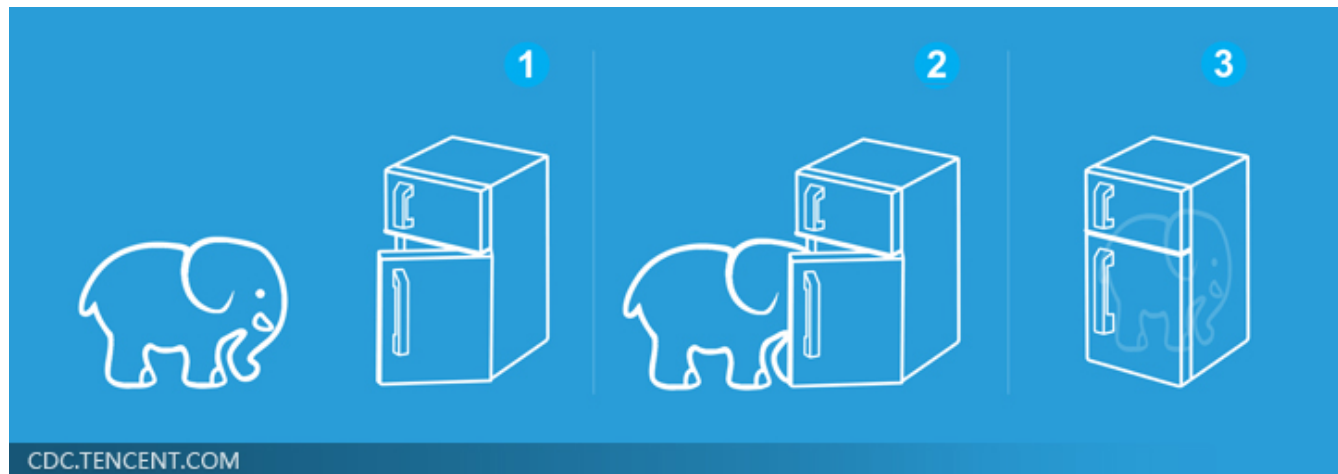


Less parameters for the network to process the image

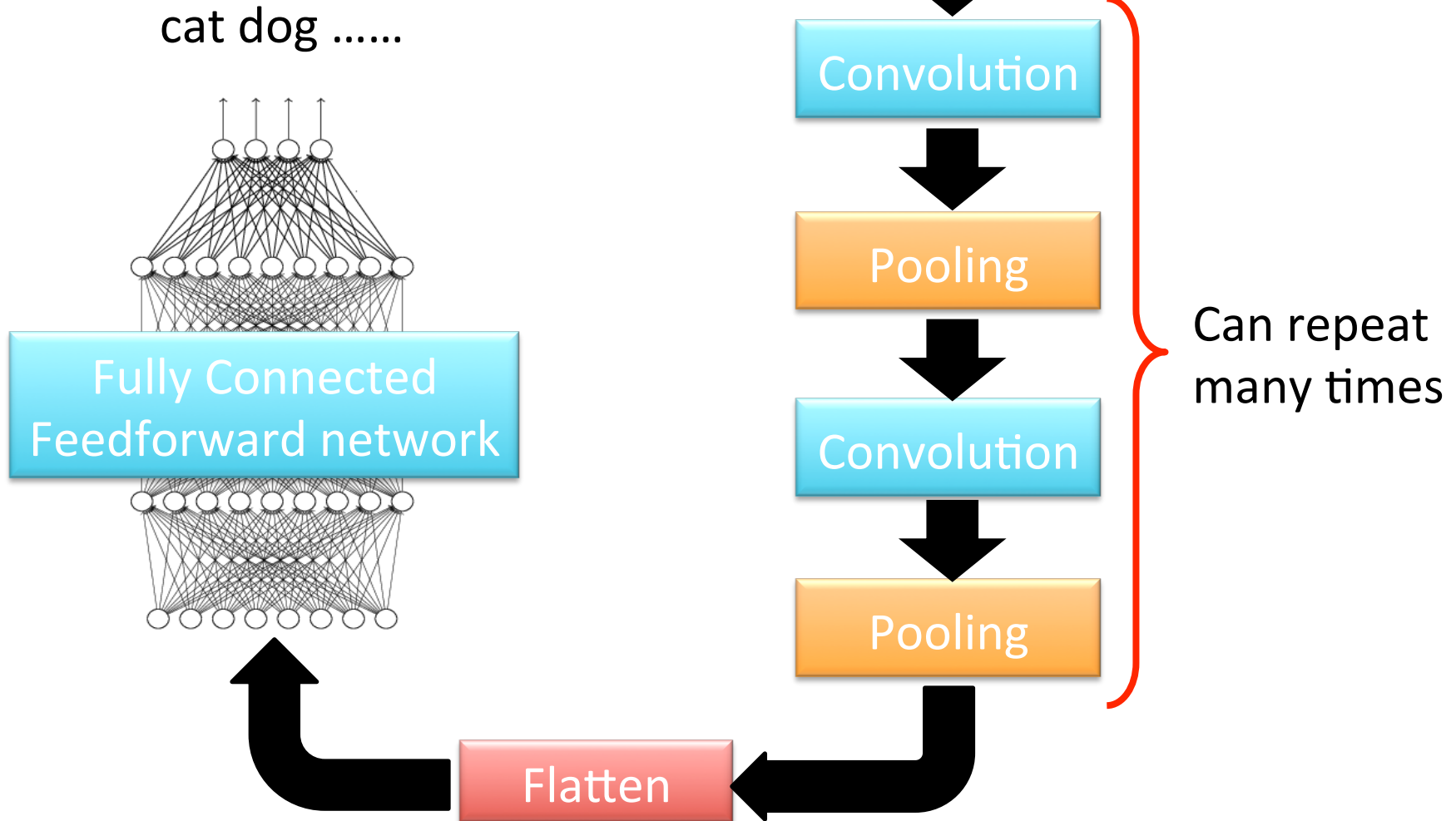
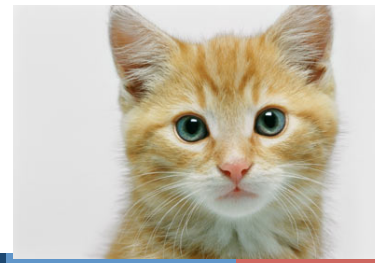
# Three Steps for Deep Learning



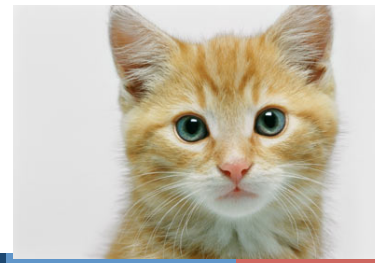
Deep Learning is so simple .....



# The whole CNN



# The whole CNN



## Property 1

- Some patterns are much smaller than the whole image

## Property 2

- The same patterns appear in different regions.

## Property 3

- Subsampling the pixels will not change the object

Convolution

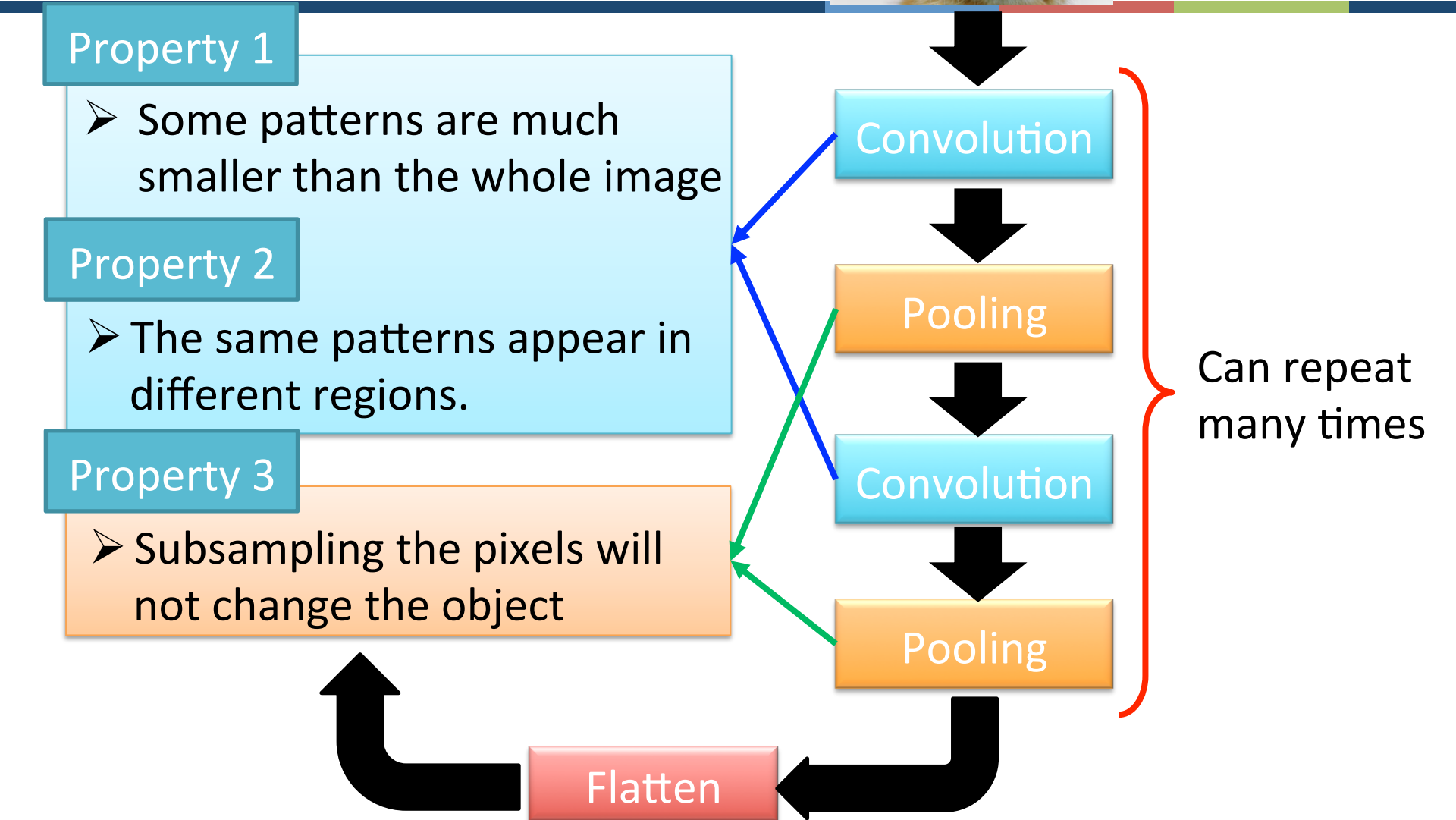
Pooling

Convolution

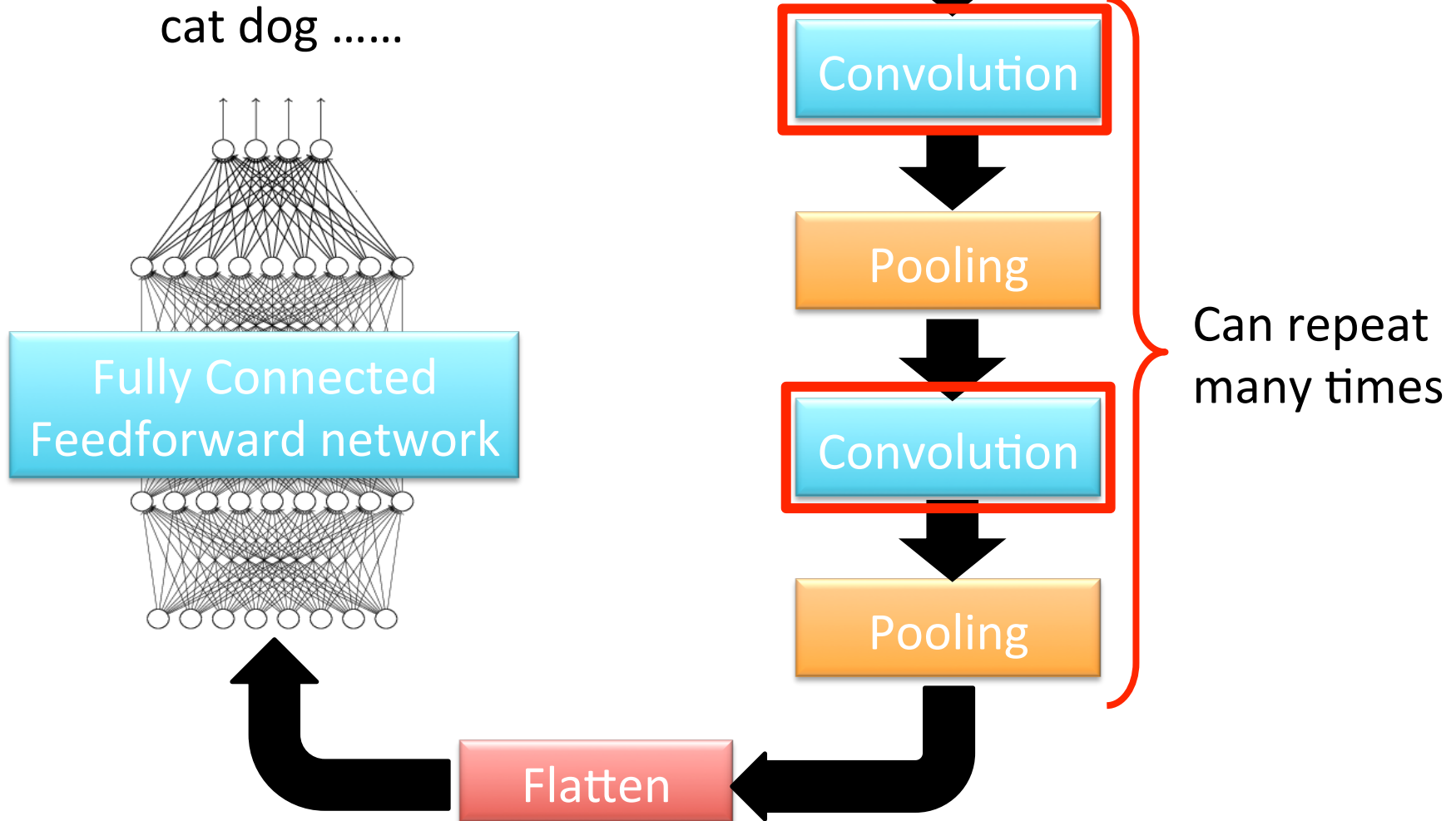
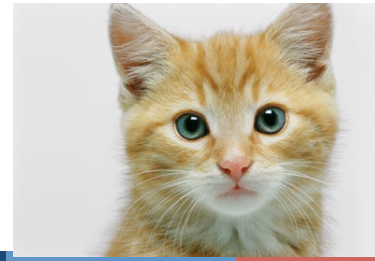
Pooling

Can repeat many times

Flatten



# The whole CNN



# CNN – 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

Those are the network parameters to be learned.

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

Filter 1  
Matrix

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

Filter 2  
Matrix

⋮

Property 1

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

# CNN – Convolution

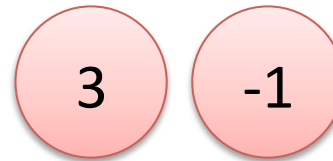
1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

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





# CNN – Convolution

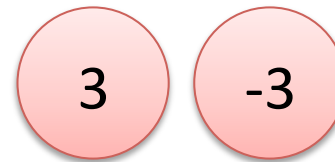
1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

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



We set stride=1 below

# CNN – Convolution

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

Filter 1

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

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

Property 2

# CNN – Convolution

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

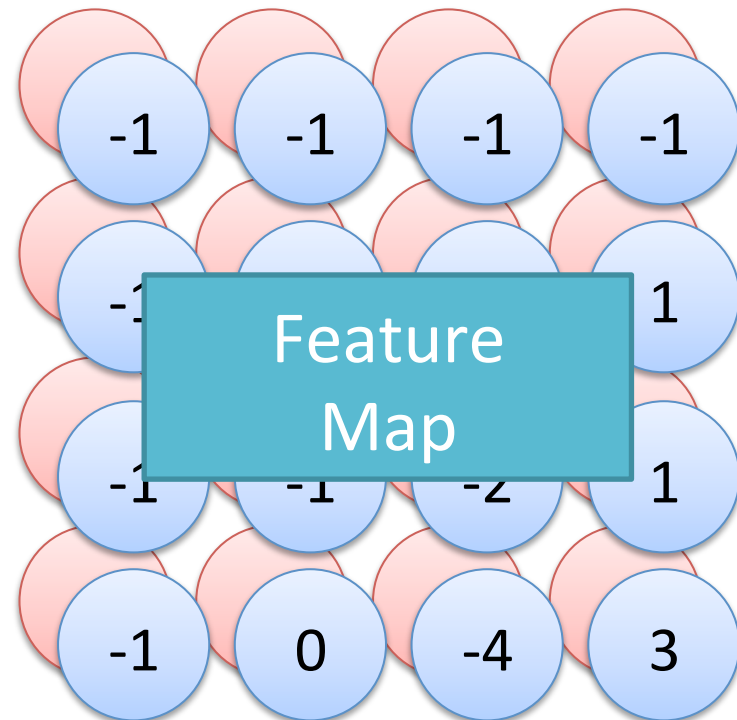
Filter 2

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

Do the same process for every filter



4 x 4 image

# CNN – Zero Padding

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

Filter 1

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

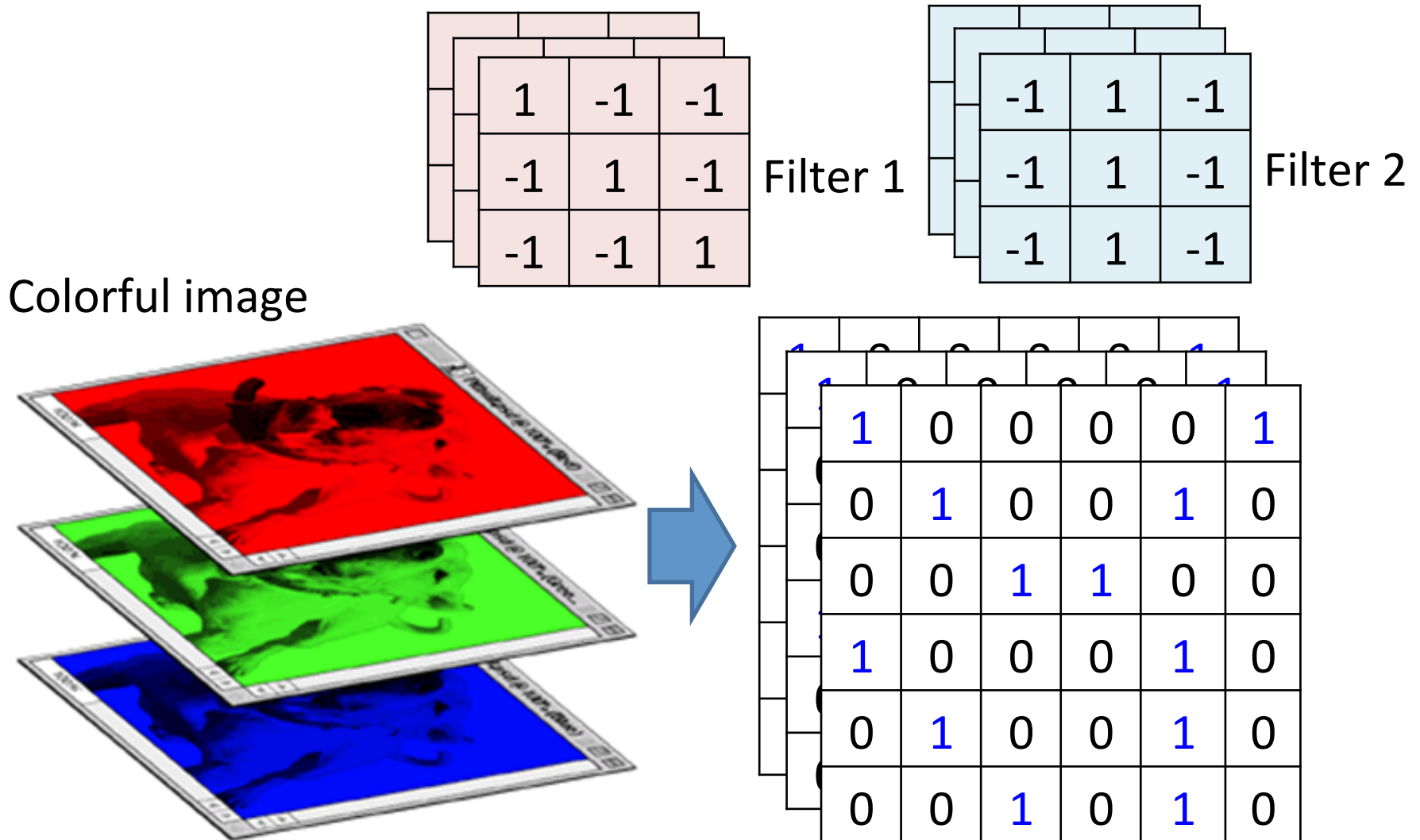
6 x 6 image

You will get another 6 x 6 images in this way

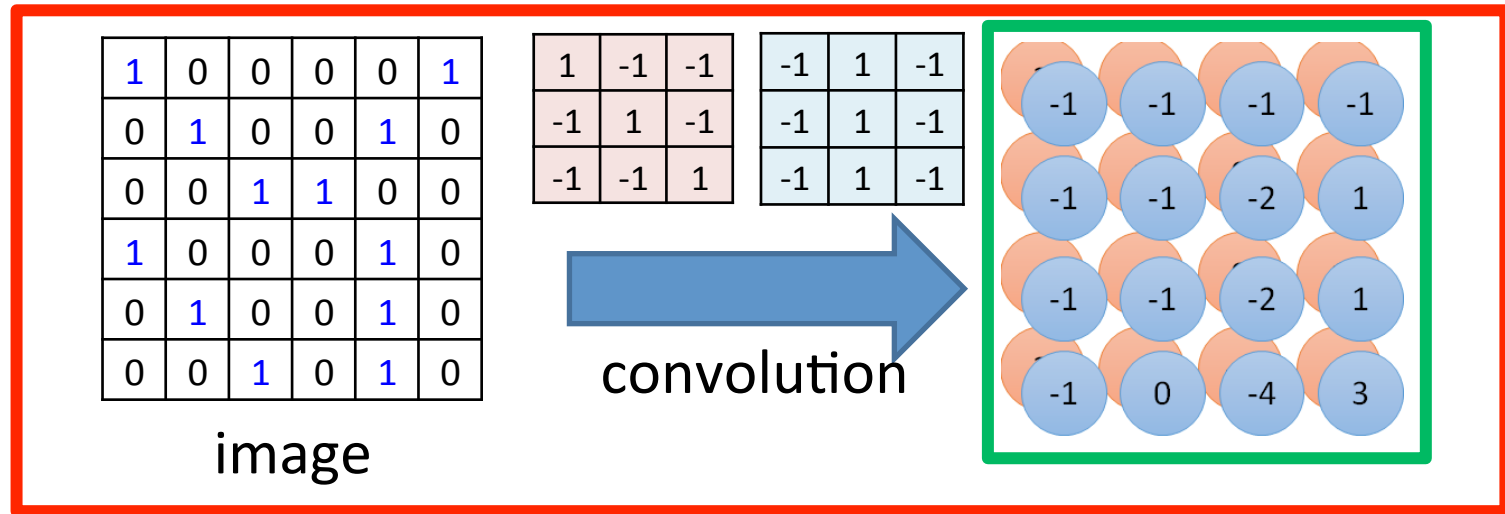


Zero padding

# CNN – Colorful image

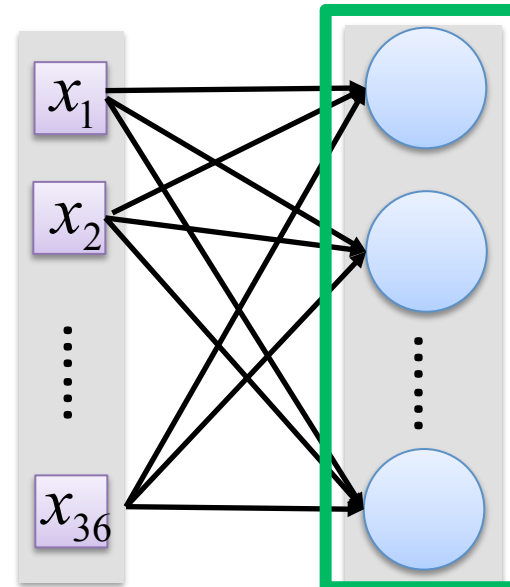


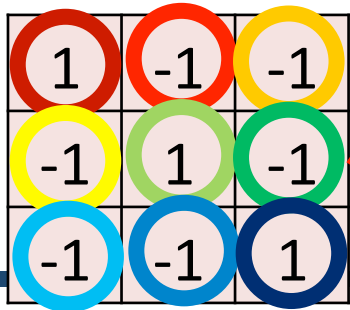
# Convolution v.s. Fully Connected



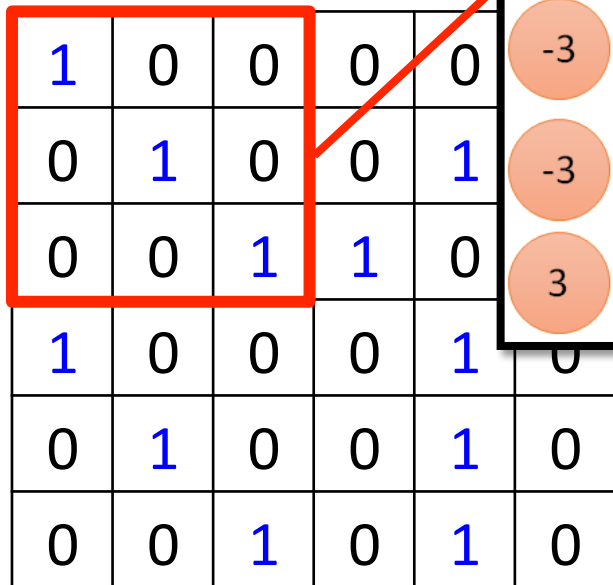
Fully-  
connected

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



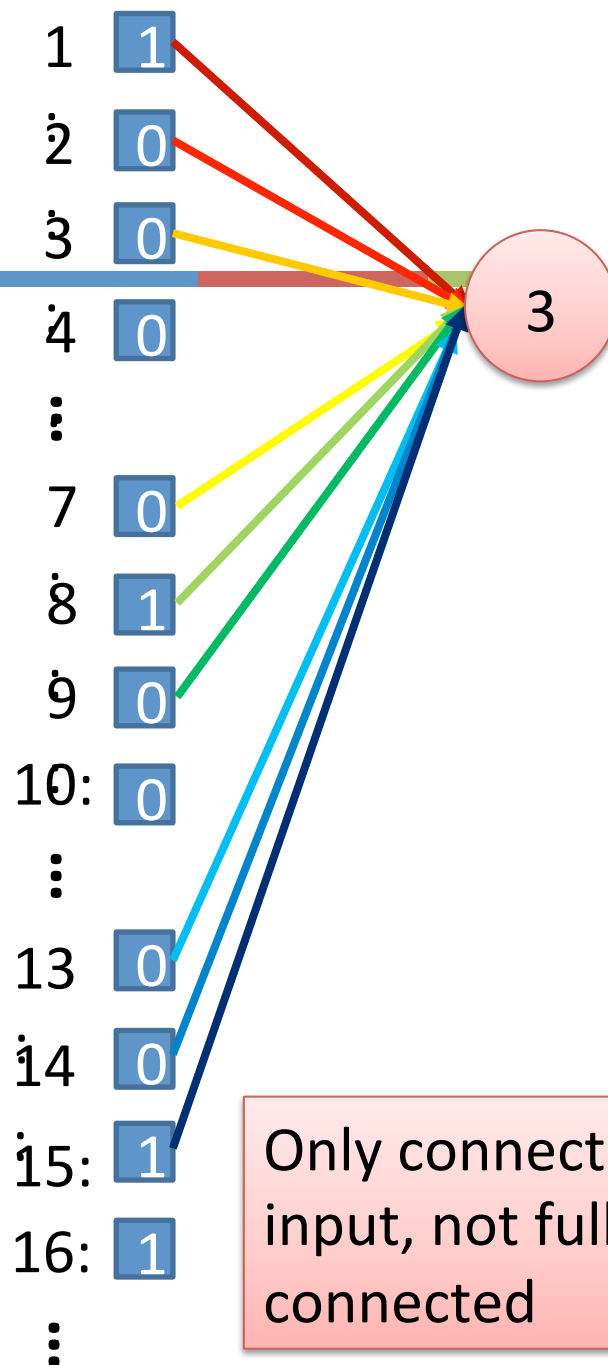
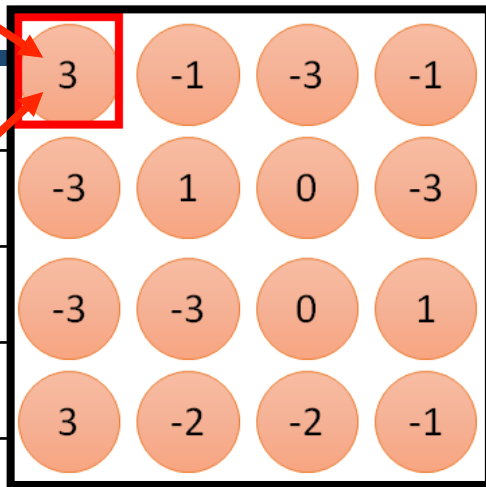


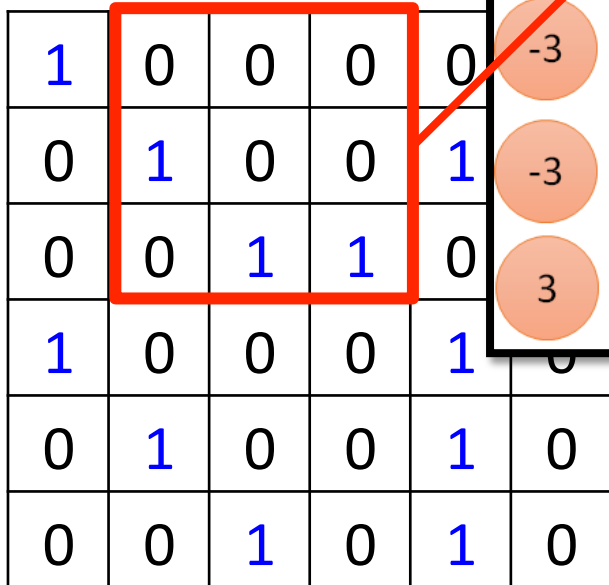
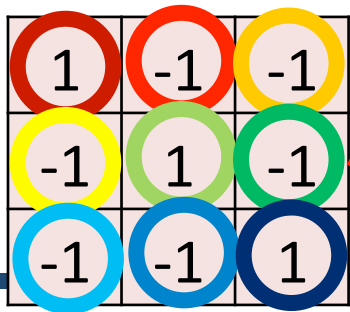
Filter 1



6 x 6 image

Less parameters!

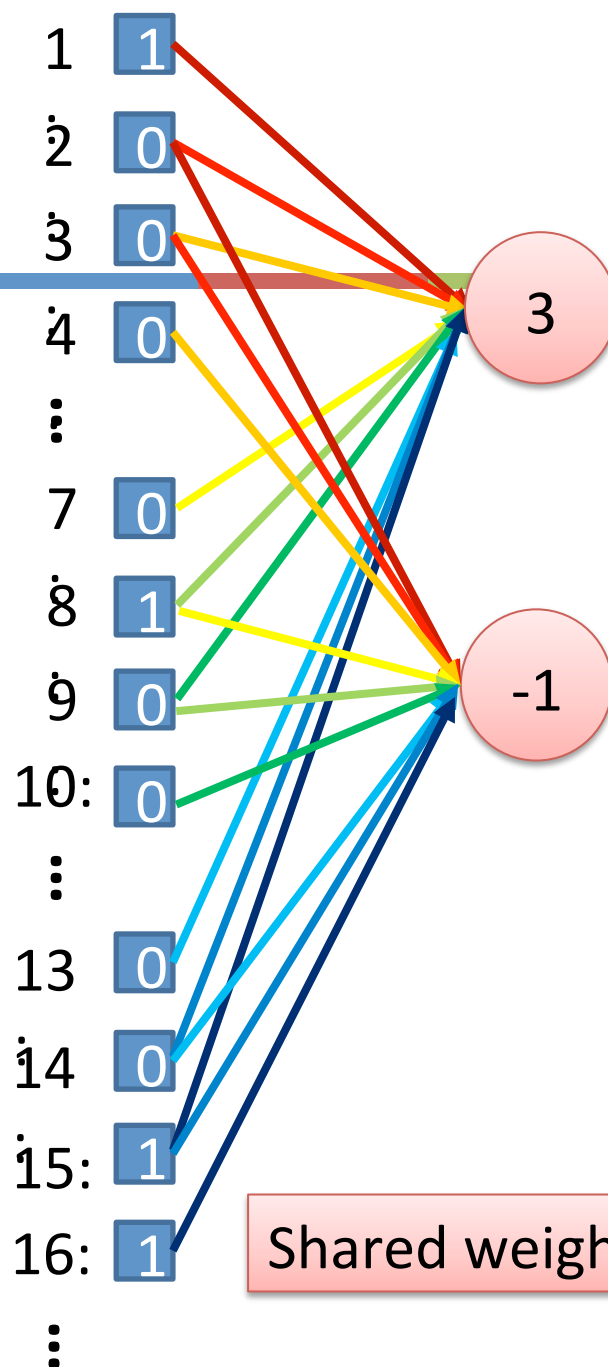
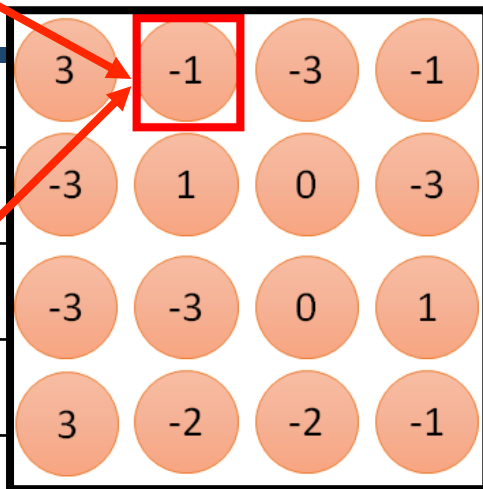




6 x 6 image

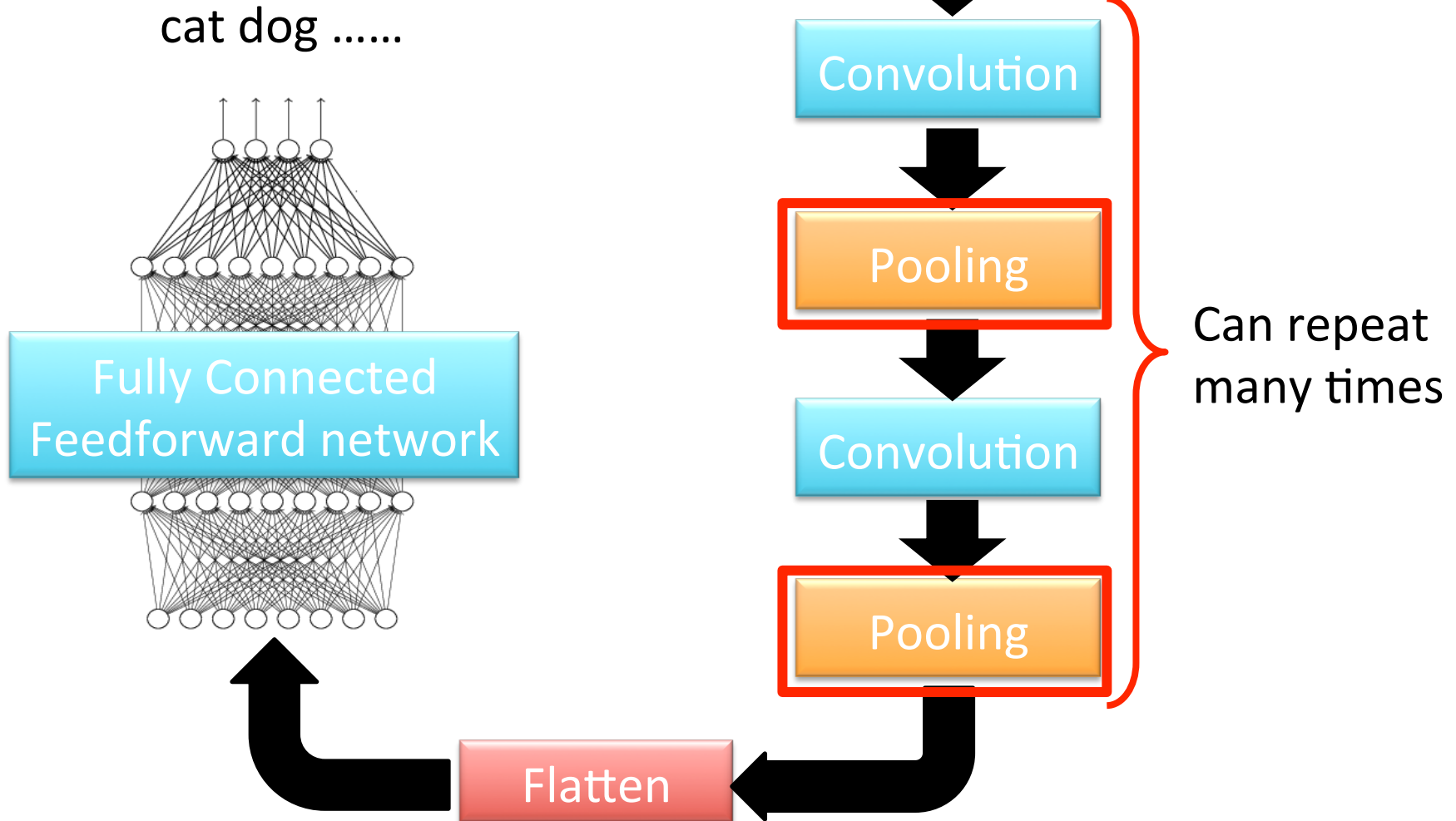
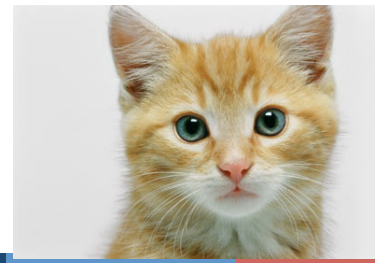
Less parameters!

Even less parameters!





# The whole CNN



# CNN – Max Pooling

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

Filter 1

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

Filter 2

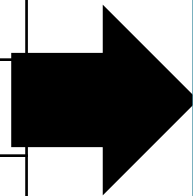
3	-1	-3	-1
-3	1	0	-3
-3	-3	0	1
3	-2	-2	-1

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

# CNN – 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

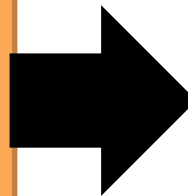
6 x 6 image



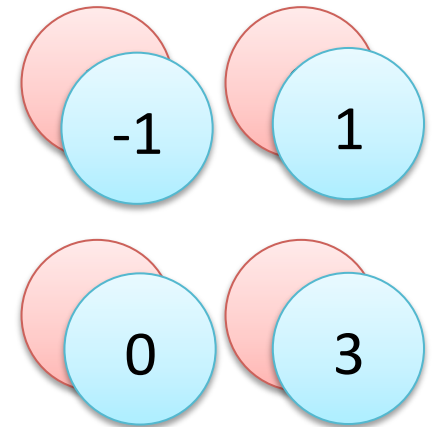
Conv



Max  
Pooling



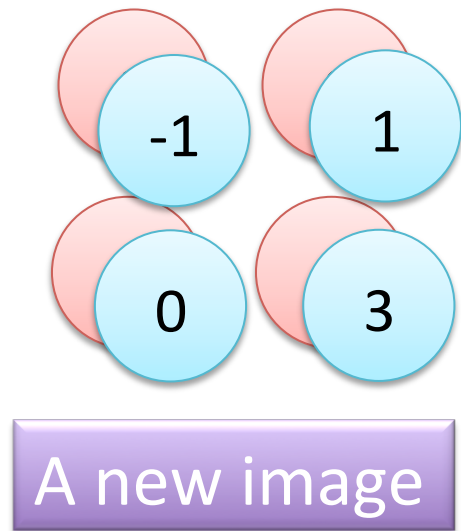
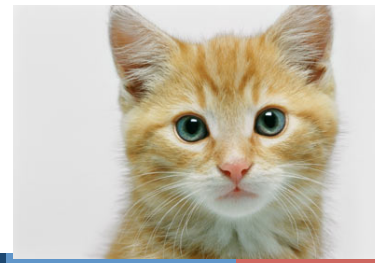
New image  
but smaller



2 x 2 image

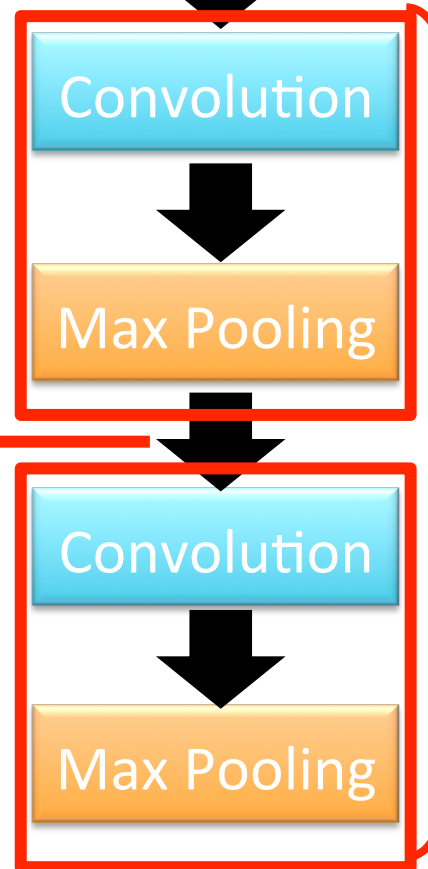
Each filter  
is a channel

# The whole CNN



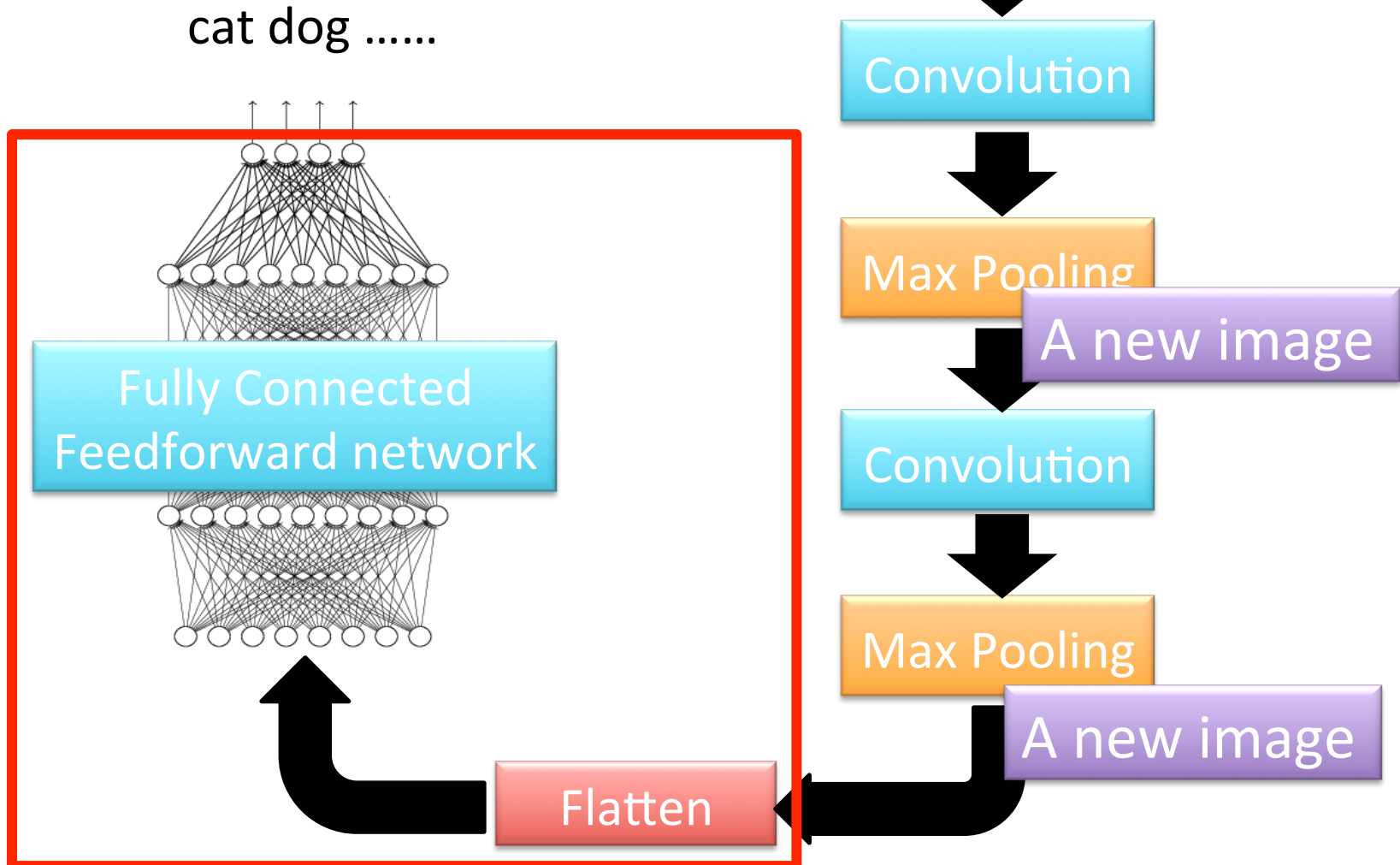
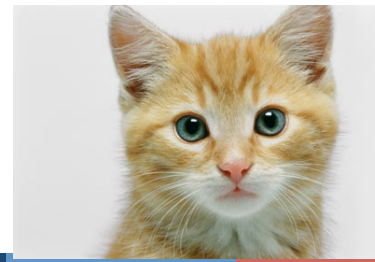
Smaller than the original image

The number of the channel is the number of filters

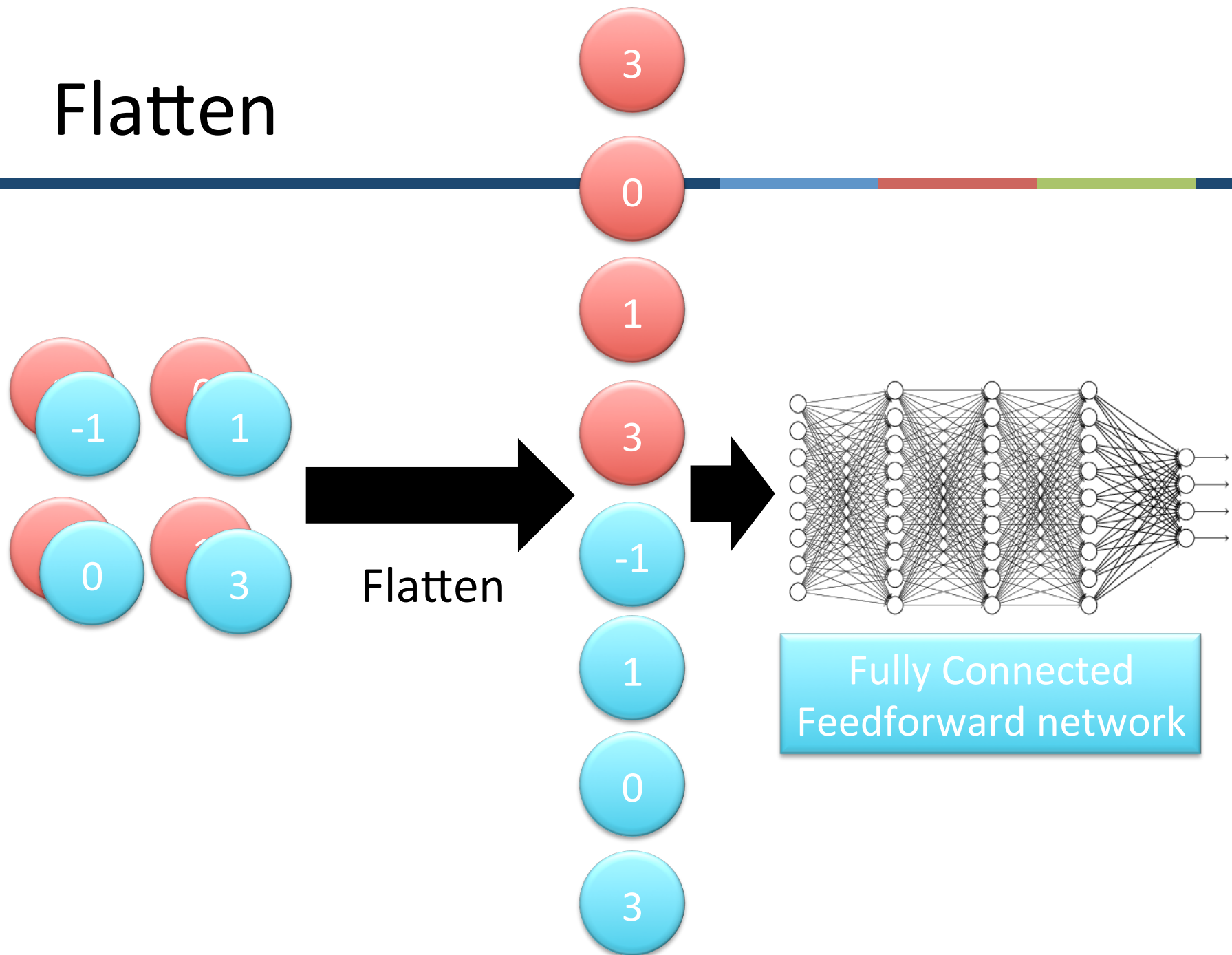


Can repeat many times

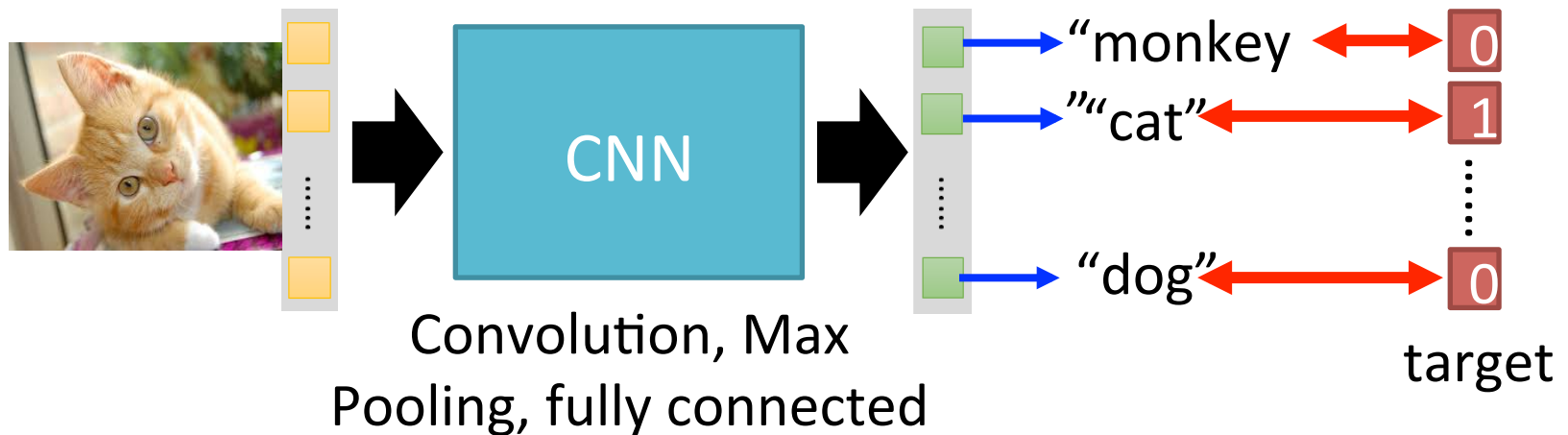
# The whole CNN



# Flatten



# Convolutional Neural Network



Learning: Nothing special, just gradient descent .....

# CNN in Keras

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

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

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

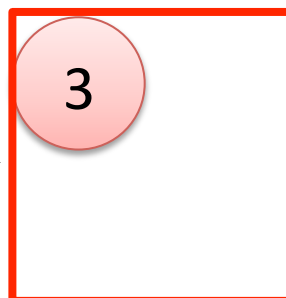
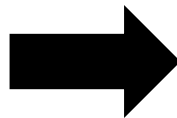
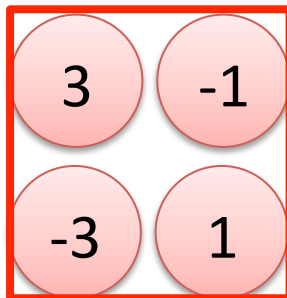
-1	1	-1
-1	1	-1
-1	1	-1

..... There are 25  
3x3 filters.

Input\_shape = ( 1, 28, 28 )

1: black/weight, 3: RGB 28 x 28 pixels

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



input

Convolution

Max Pooling

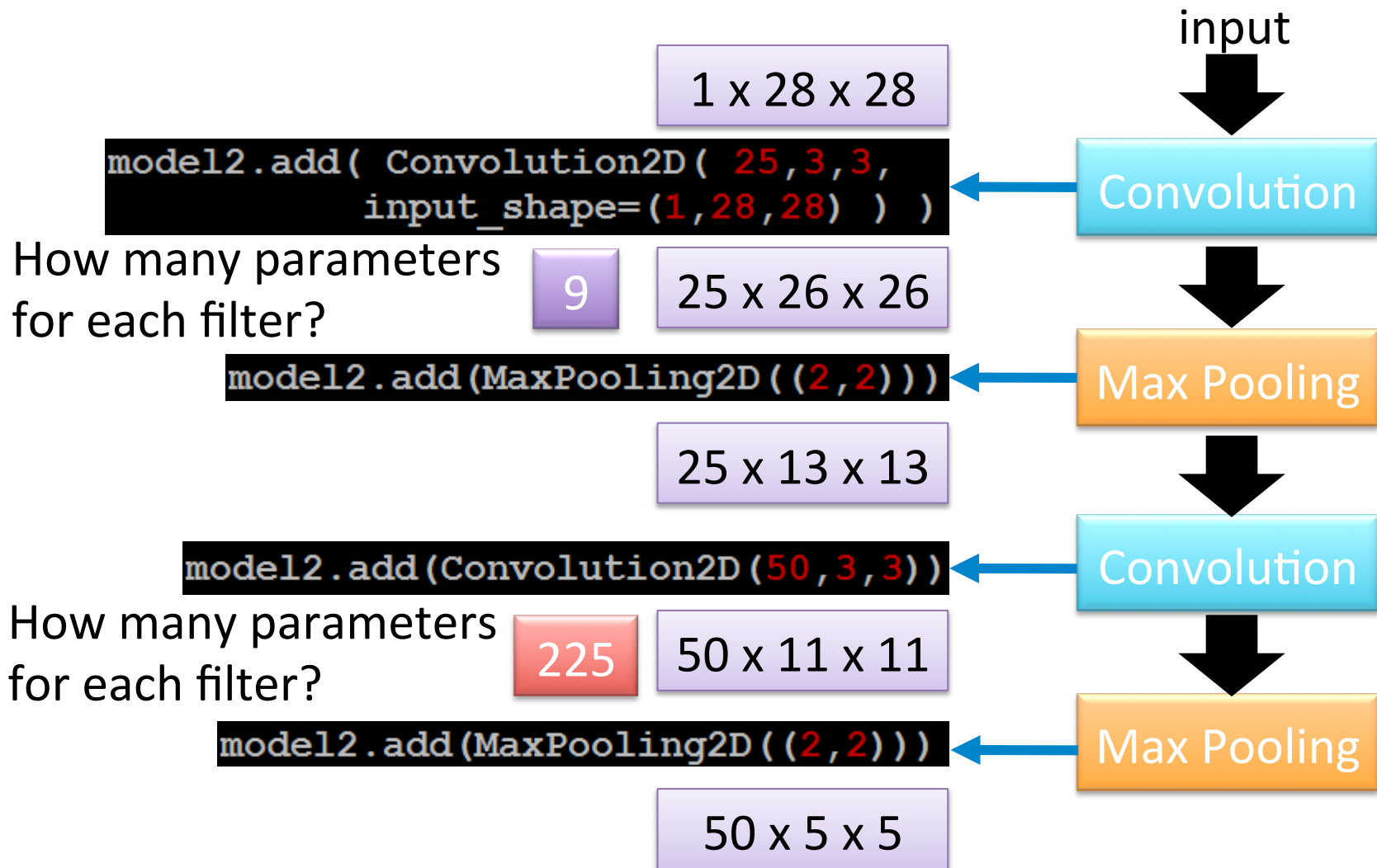
Convolution

Max Pooling



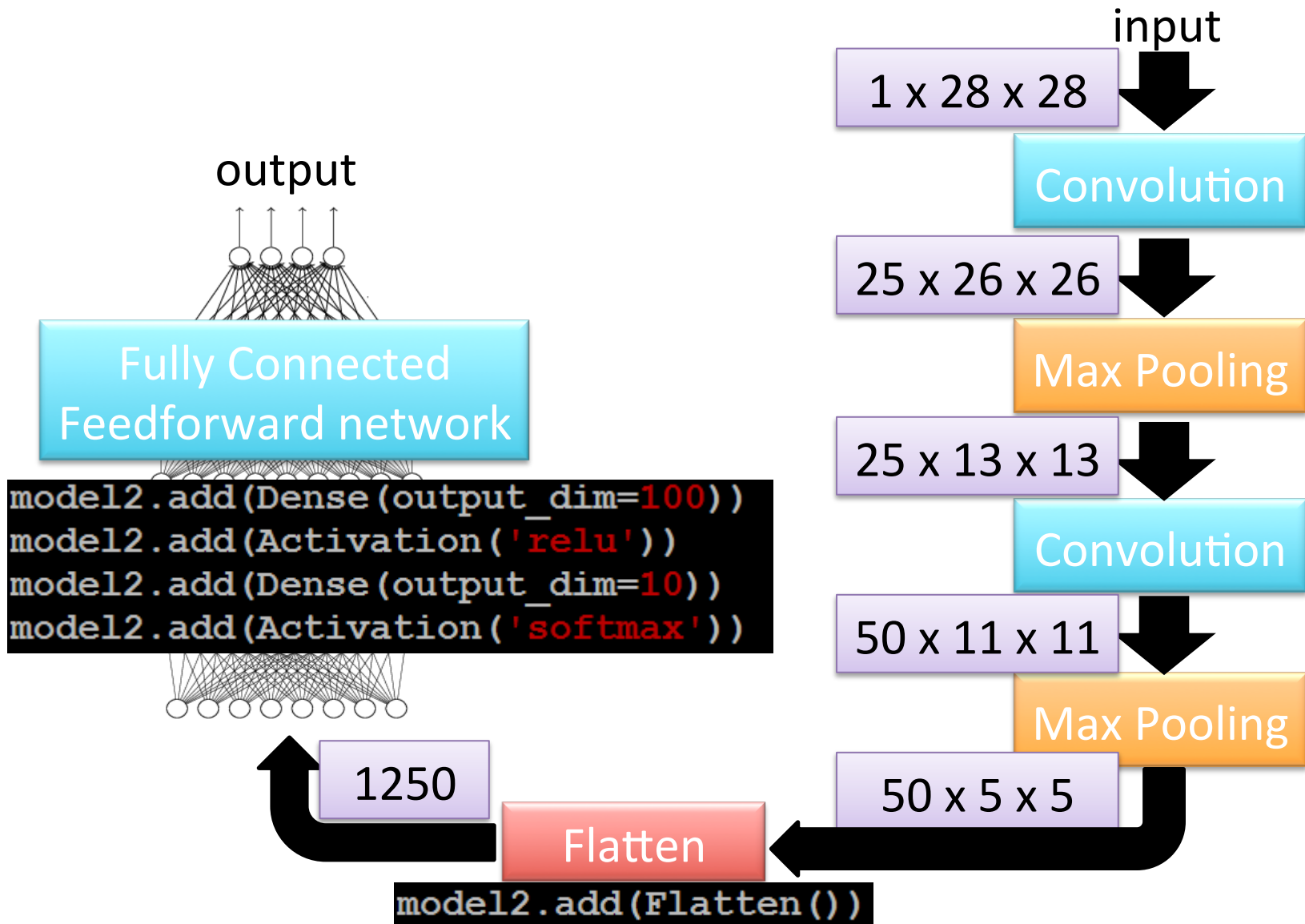
## CNN in Keras

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



# CNN in Keras

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



```
class Net(nn.Module):
```

```
    def __init__(self):
```

```
        super().__init__()
```

```
        self.conv1 = nn.Conv2d(1, 25, 3)
```

```
        self.pool = nn.MaxPool2d(2, 2)
```

```
        self.conv2 = nn.Conv2d(25, 50, 3)
```

```
        self.pool = nn.MaxPool2d(2, 2)
```

```
        self.fc1 = nn.Linear(50*5*5, 100)
```

```
        self.fc2 = nn.Linear(100, 10)
```

```
    def forward(self, x):
```

```
        x = self.pool(F.relu(self.conv1(x)))
```

```
        x = self.pool(F.relu(self.conv2(x)))
```

```
        x = x.view(-1, self.num_flat_features(x))
```

```
        x = F.relu(self.fc1(x))
```

```
        x = F.log_softmax(self.fc2(x))
```

```
    return x
```

```
    def num_flat_features(self, x):
```

```
        size = x.size()[1:]
```

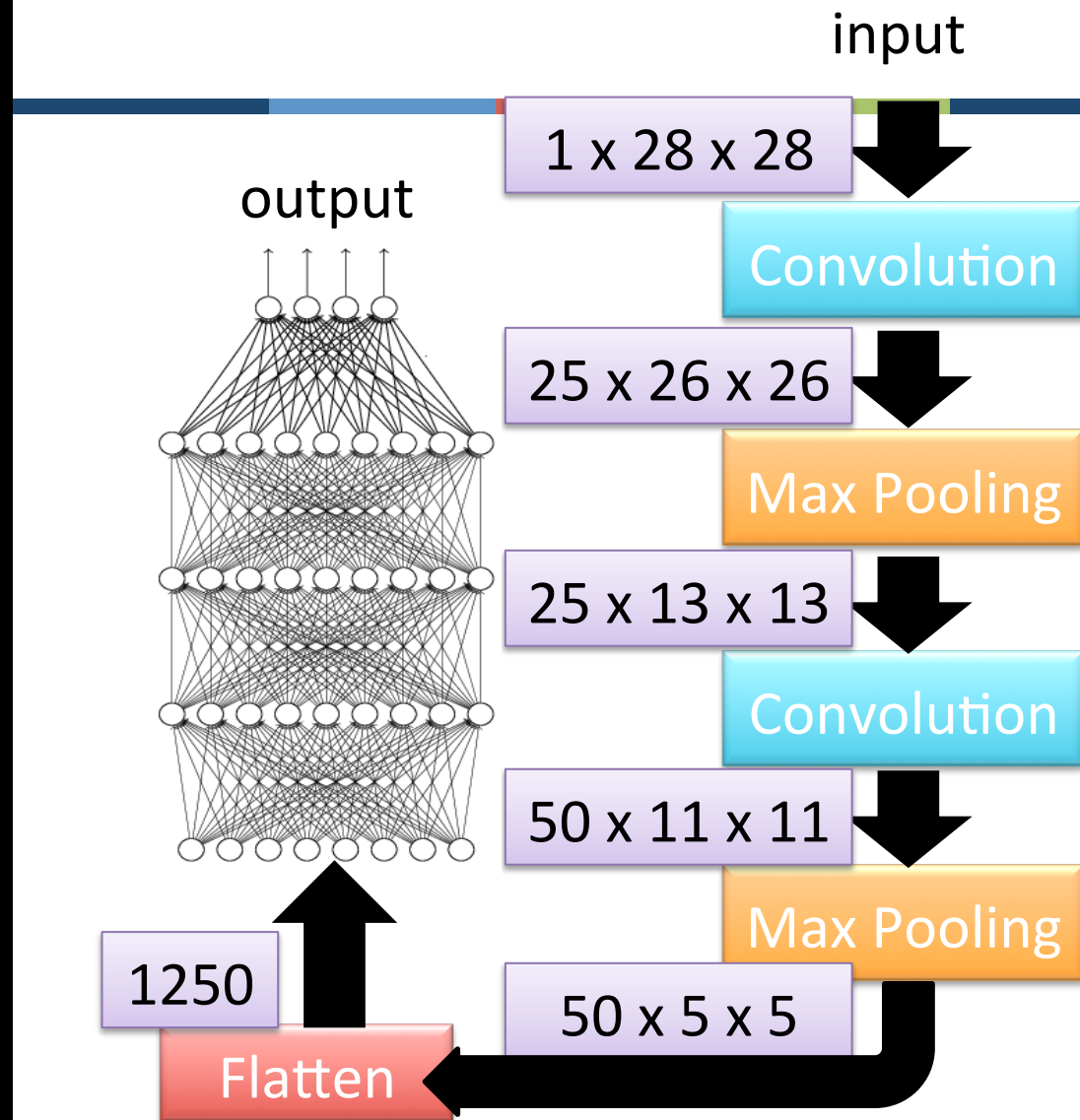
```
        num_features = 1
```

```
        for s in size:
```

```
            num_features *= s
```

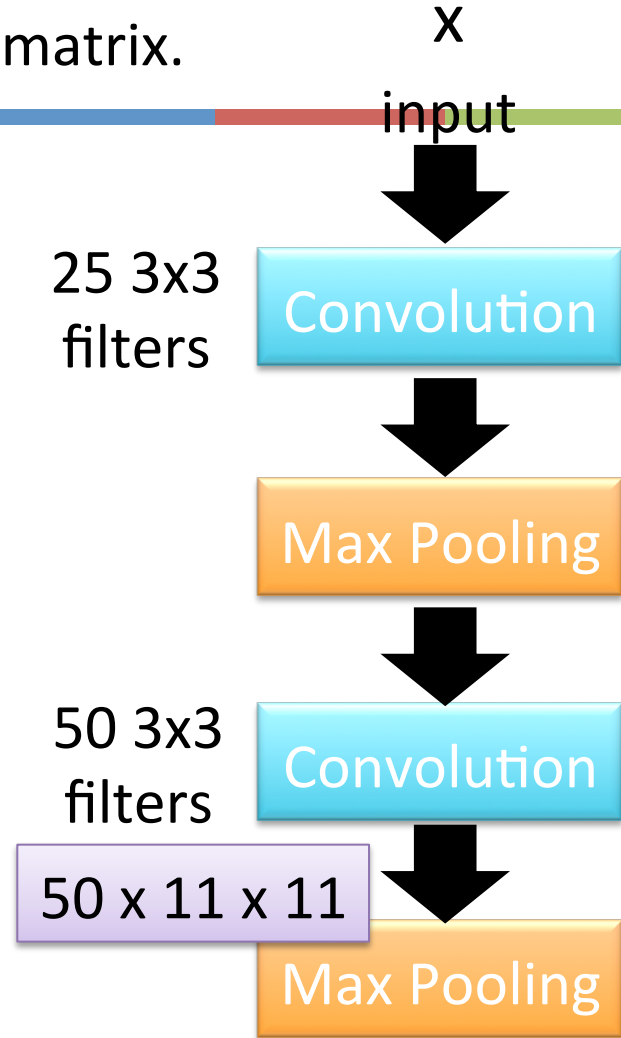
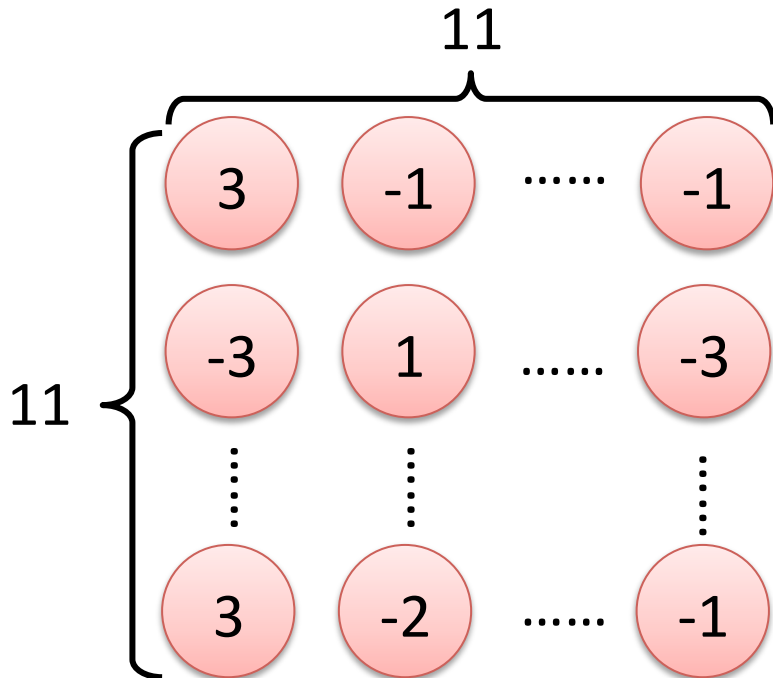
```
    return num_features
```

## CNN in PyTorch



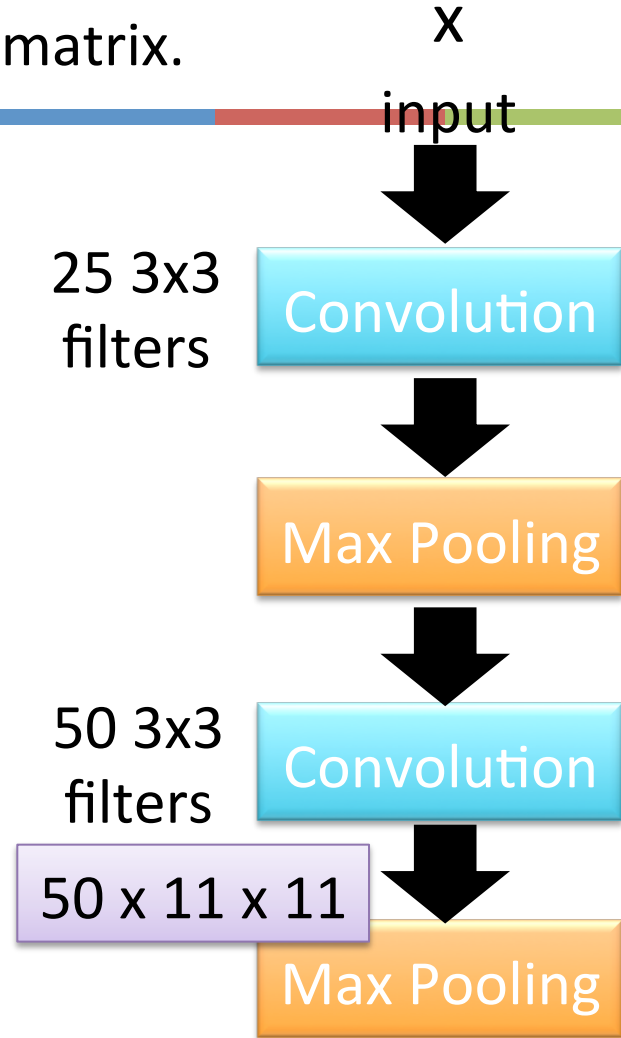
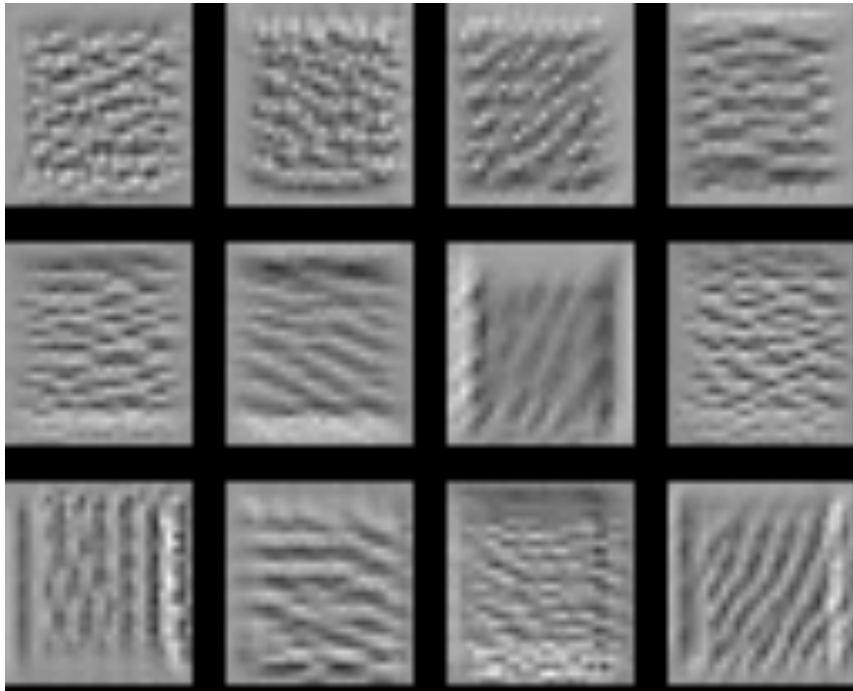
# What does CNN learn?

The output of the k-th filter is a 11 x 11 matrix.



# What does CNN learn?

The output of the k-th filter is a 11 x 11 matrix.



# More Application: Playing Go



Black: 1  
white: -1  
none: 0



Network



Next move  
(19 x 19  
positions)

19 x 19 vector

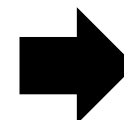
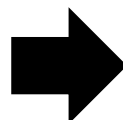
Fully-connected feedforward  
network can be used

But CNN performs much better.

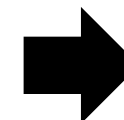
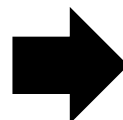
# More Application: Playing Go

Training: record of previous plays

black: 5之五 → white: 天元 → black: 五之5 ...



Target:  
“天元” = 1  
else = 0



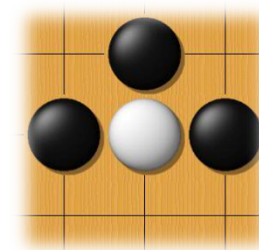
Target:  
“五之5” = 1  
else = 0



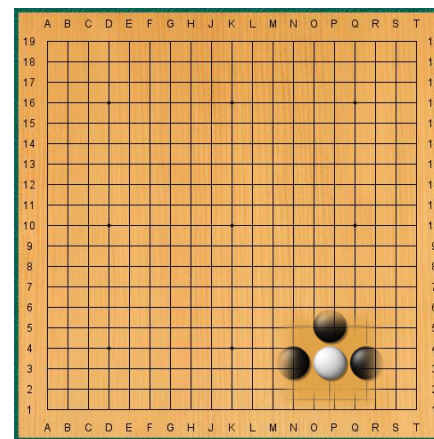
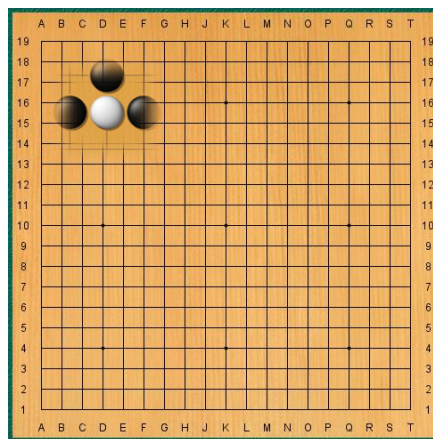
# Why CNN for playing Go?

- Some patterns are much smaller than the whole image

Alpha Go uses 5 x 5 for first layer



- The same patterns appear in different regions.





# Why CNN for playing Go?

- Subsampling the pixels will not change the object



Max Pooling

How to explain this???

**Neural network architecture.** The input to the policy network is a  $19 \times 19 \times 48$  image stack consisting of 48 feature planes. The first hidden layer zero pads the input into a  $23 \times 23$  image, then convolves  $k$  filters of kernel size  $5 \times 5$  with stride 1 with the input image and applies a rectifier nonlinearity. Each of the subsequent hidden layers 2 to 12 zero pads the respective previous hidden layer into a  $21 \times 21$  image, then convolves  $k$  filters of kernel size  $3 \times 3$  with stride 1, again followed by a rectifier nonlinearity. The final layer convolves 1 filter of kernel size  $1 \times 1$  with stride 1, with a different bias for each position, and applies a softmax function. The

Alpha Go does not use Max Pooling ..... Extended Data Table 3 additionally show the results of training with  $k = 128, 256$  and 384 filters.

# Variants of Neural Networks

---

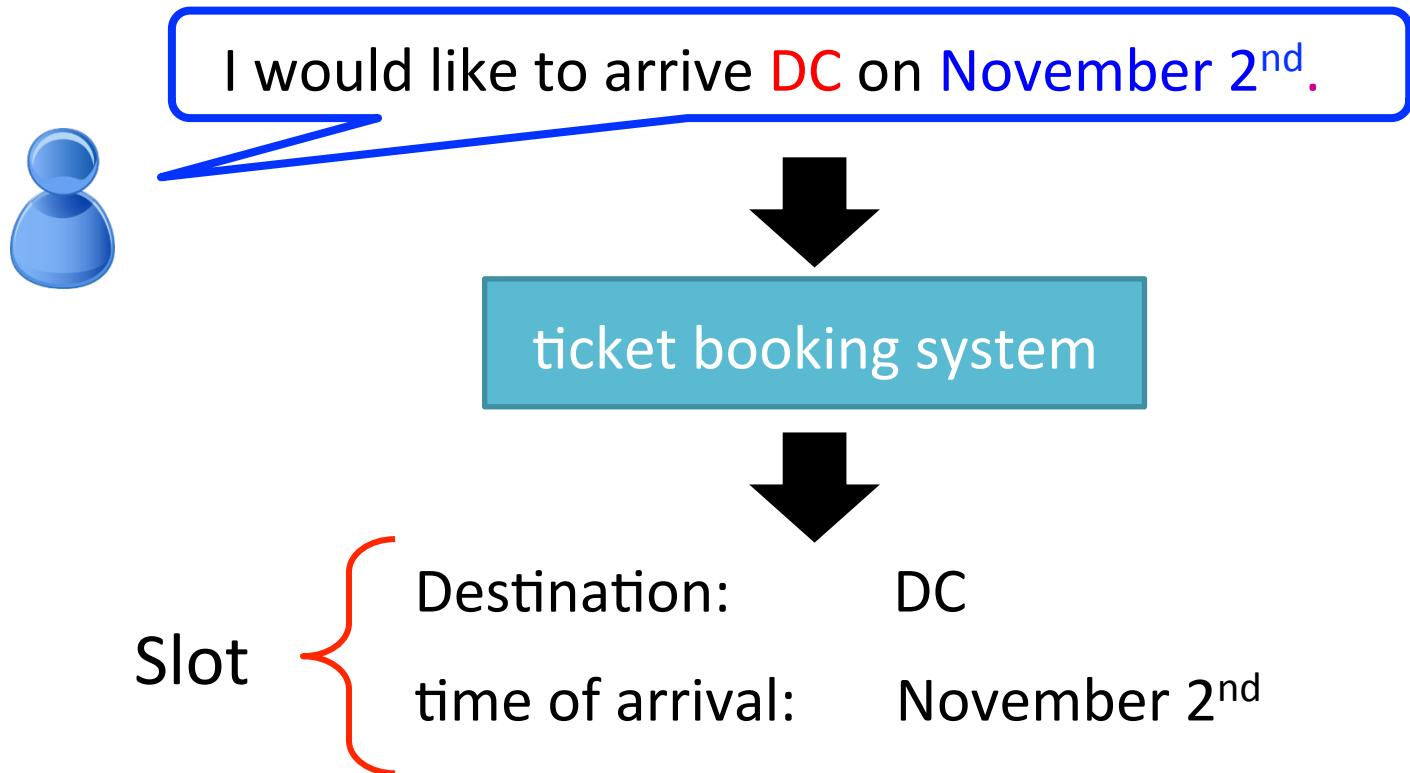
Convolutional Neural  
Network (CNN)

Recurrent Neural Network  
(RNN)

Neural Network with Memory

# Example Application

- Slot Filling

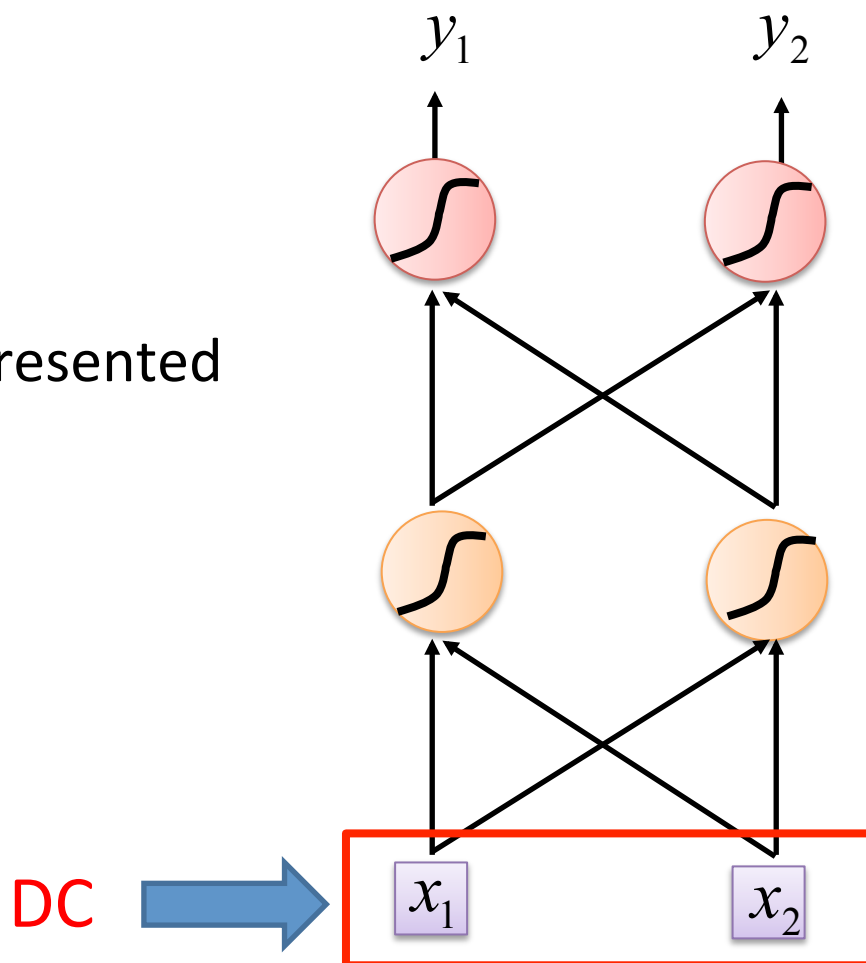


# Example Application

Solving slot filling by  
Feedforward network?

Input: a word

(Each word is represented  
as a vector)



# 1-of-N encoding

How to represent each word as a vector?

**1-of-N Encoding**    lexicon = {apple, bag, cat, dog, elephant}

The vector is lexicon size.

apple = [ 1   0   0   0   0 ]

Each dimension corresponds  
to a word in the lexicon

bag    = [ 0   1   0   0   0 ]

cat    = [ 0   0   1   0   0 ]

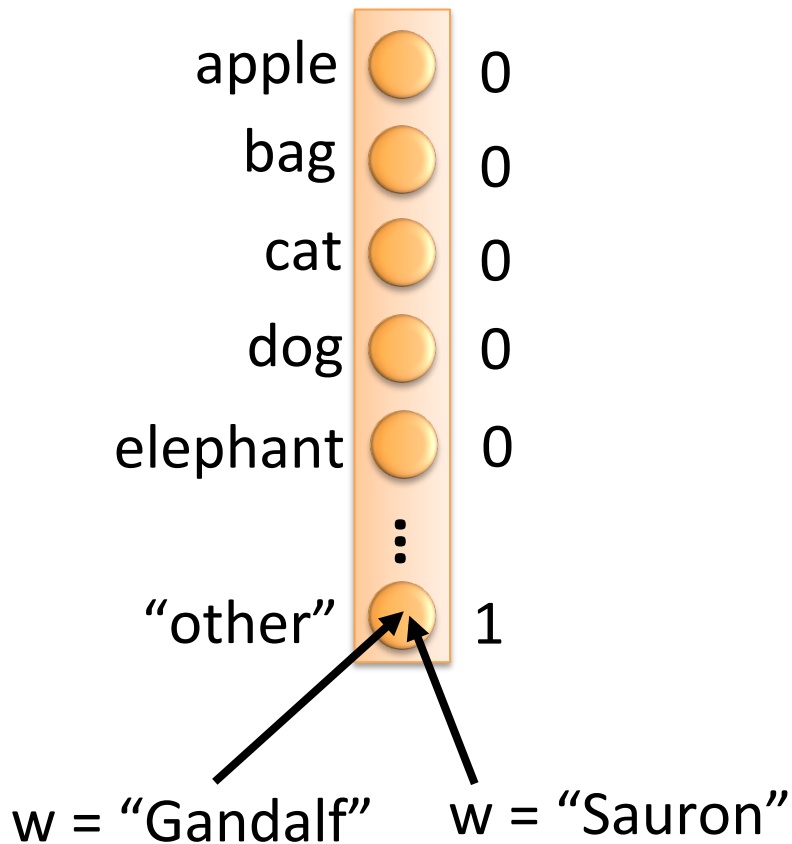
The dimension for the word  
is 1, and others are 0

dog    = [ 0   0   0   1   0 ]

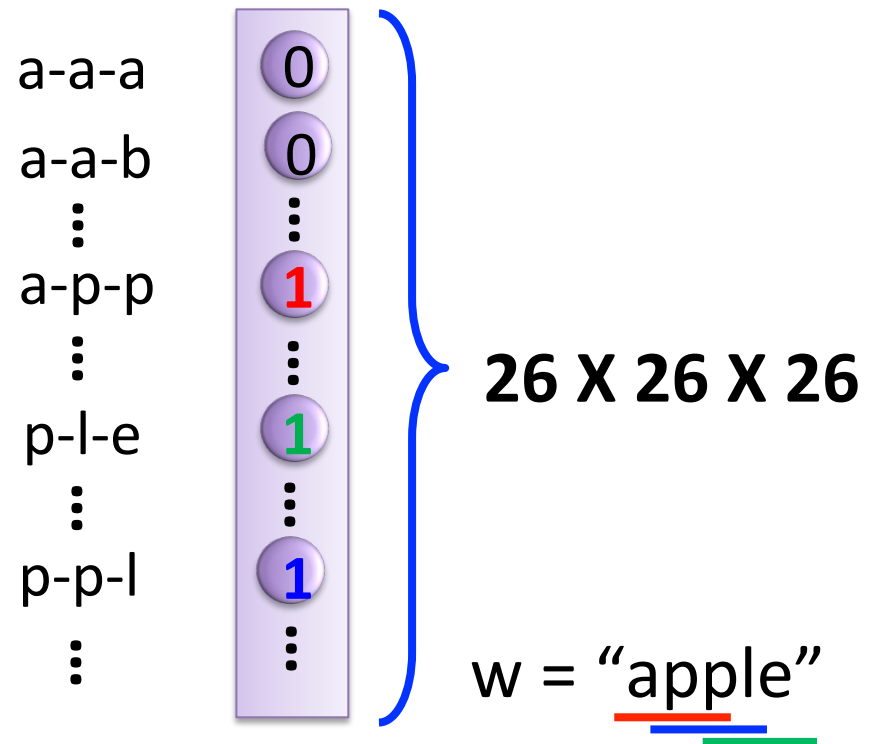
elephant = [ 0   0   0   0   1 ]

# Beyond 1-of-N encoding

## Dimension for "Other"



## Word hashing



# Example Application

time of  
departure

Solving slot filling by  
Feedforward network?

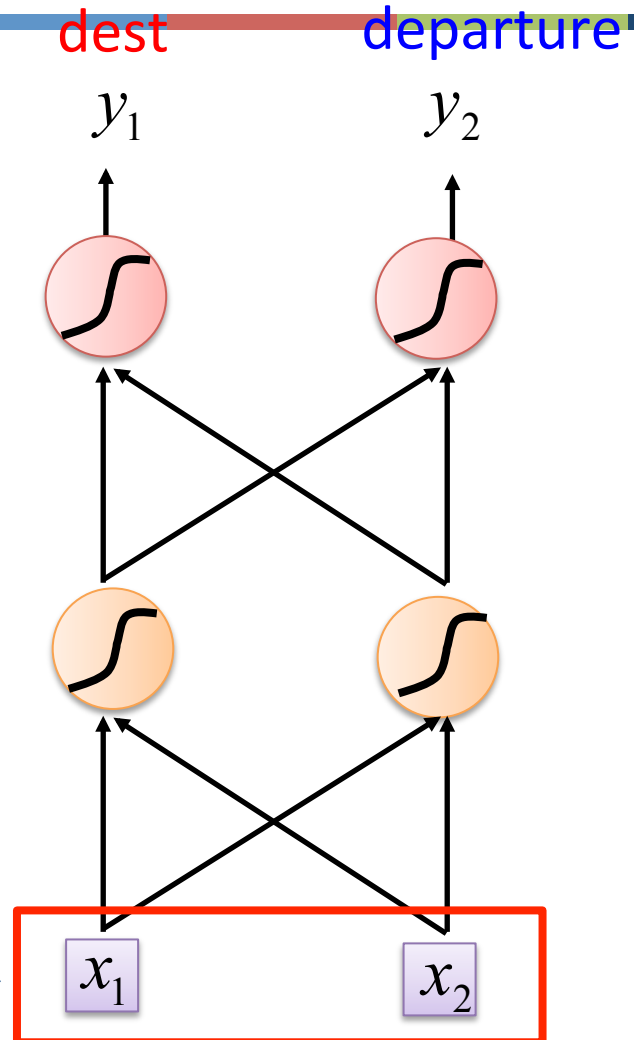
Input: a word

(Each word is represented  
as a vector)

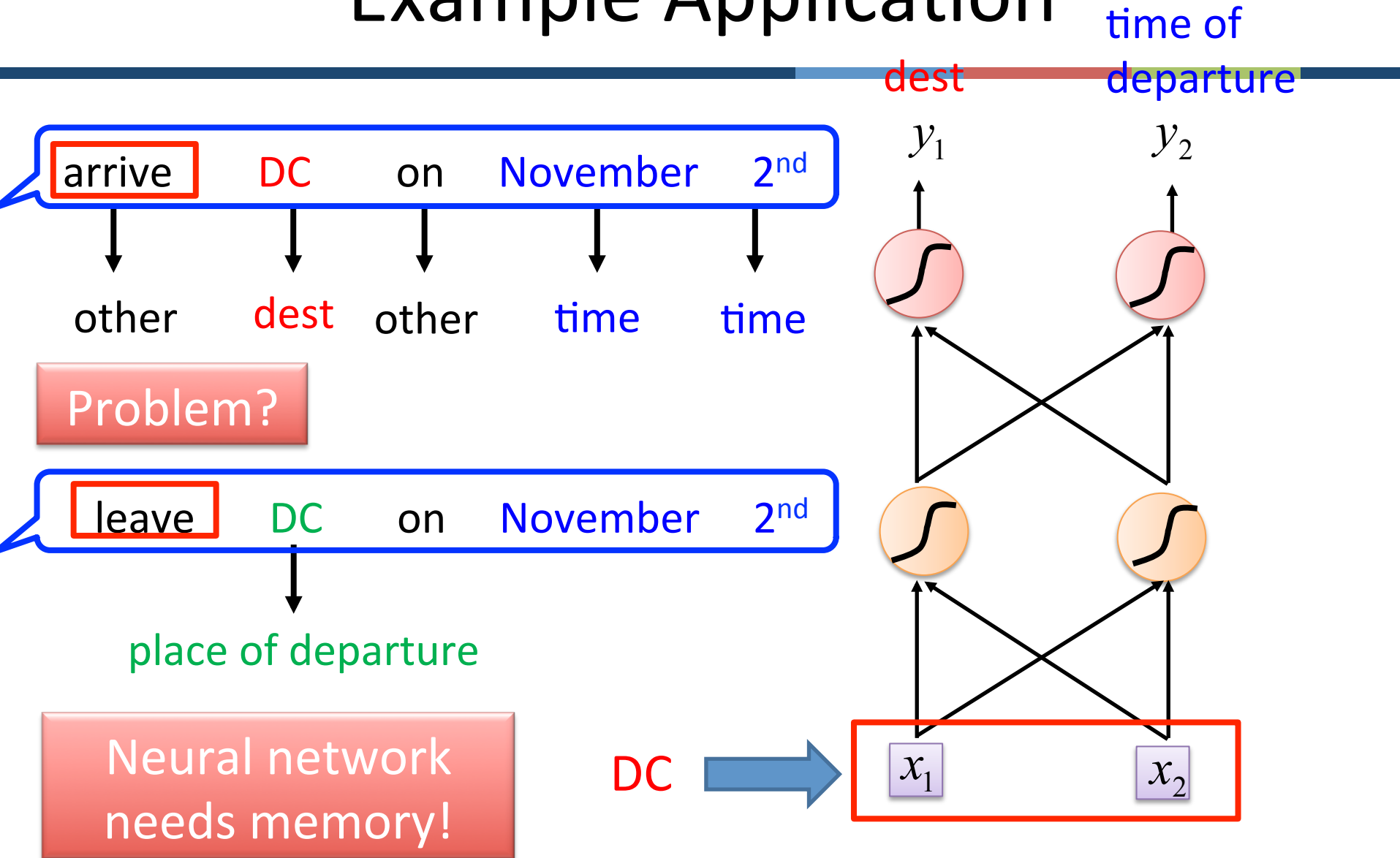
Output:

Probability distribution that  
the input word belonging to  
the slots

DC



# Example Application

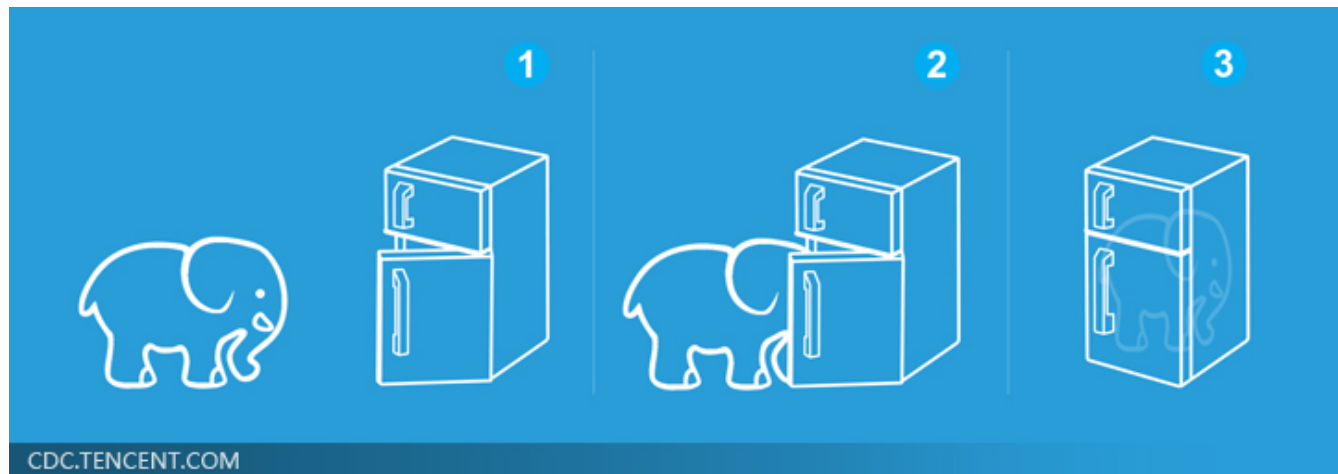




# Three Steps for Deep Learning

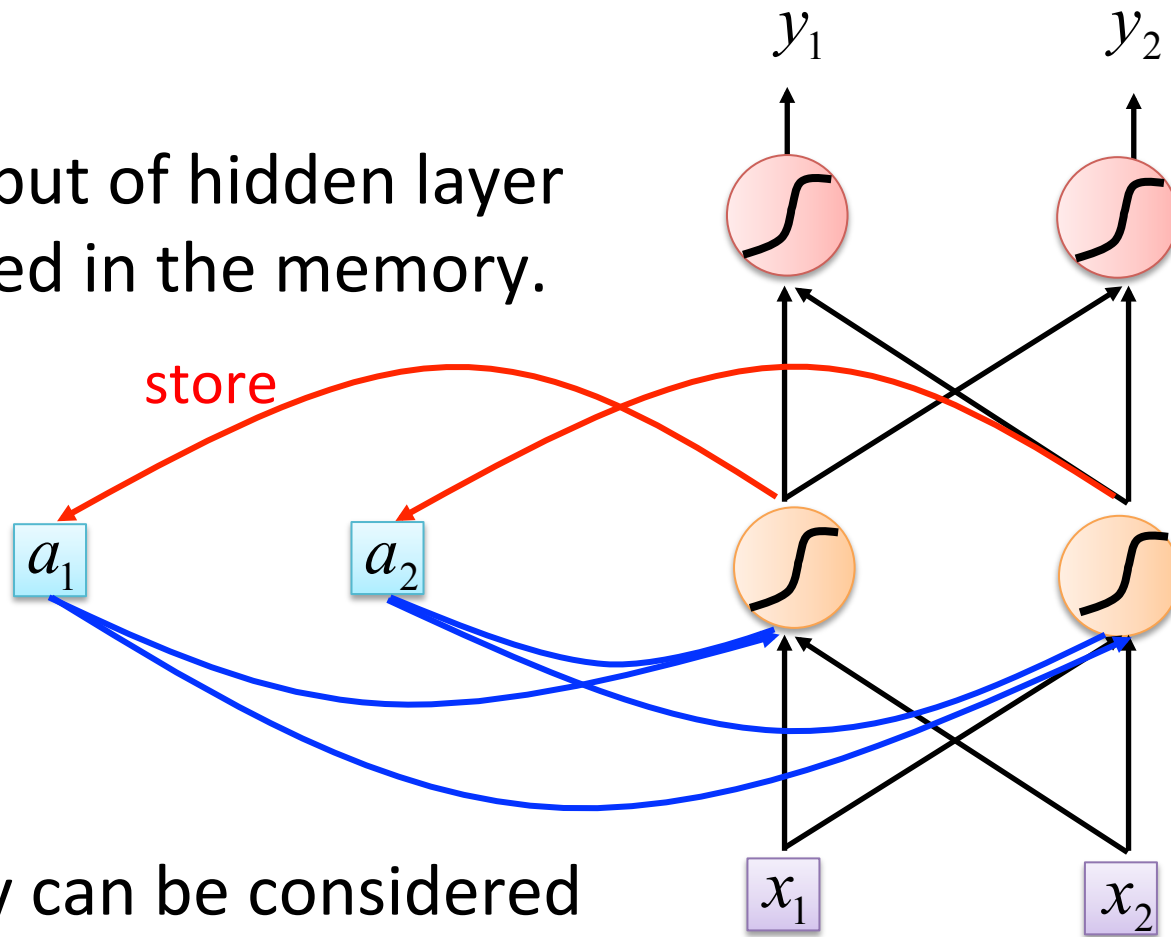


Deep Learning is so simple .....



# Recurrent Neural Network (RNN)

The output of hidden layer are stored in the memory.



Memory can be considered as another input.

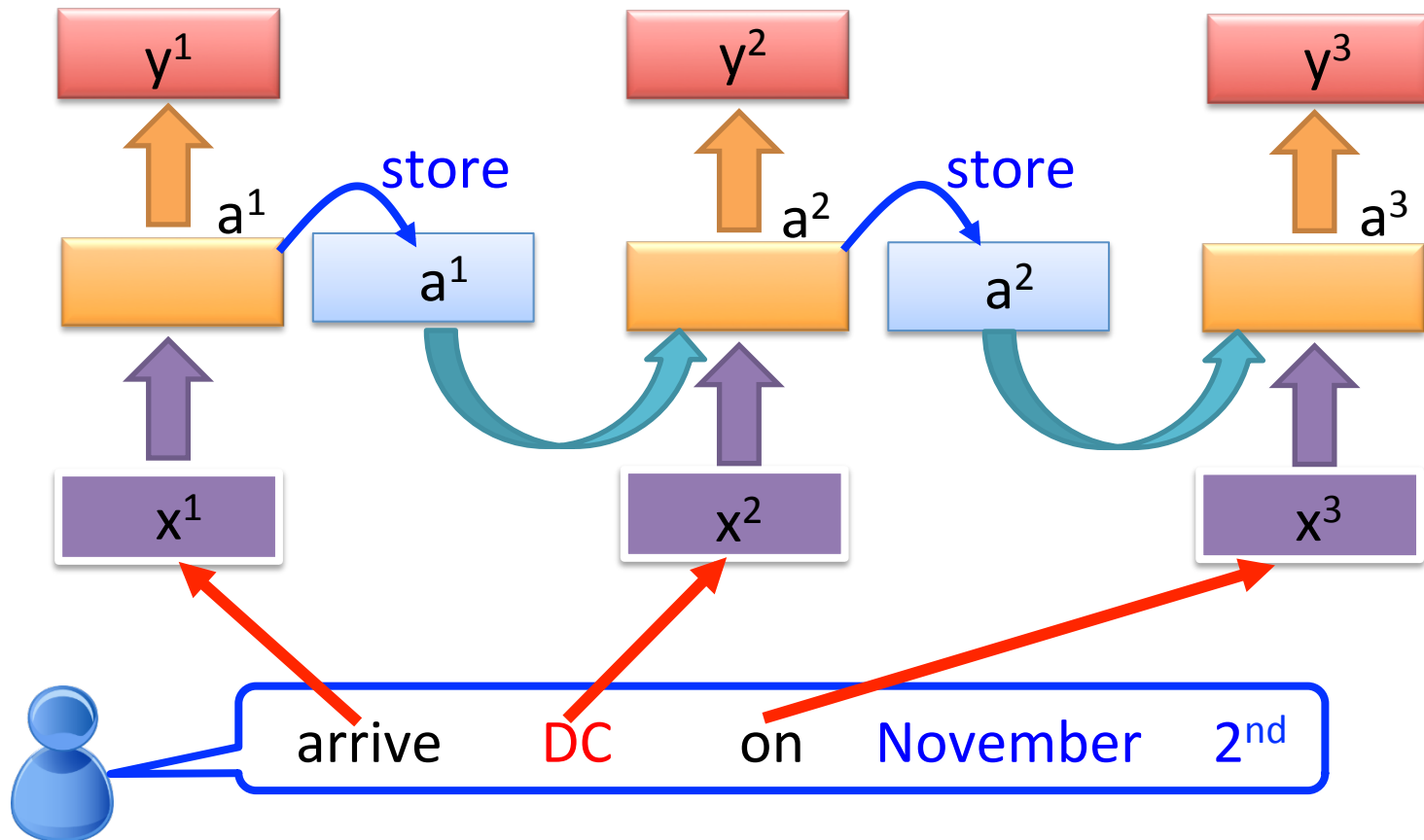
# RNN

The same network is used again and again.

Probability of  
“arrive” in each slot

Probability of “DC”  
in each slot

Probability of  
“on” in each slot



# RNN

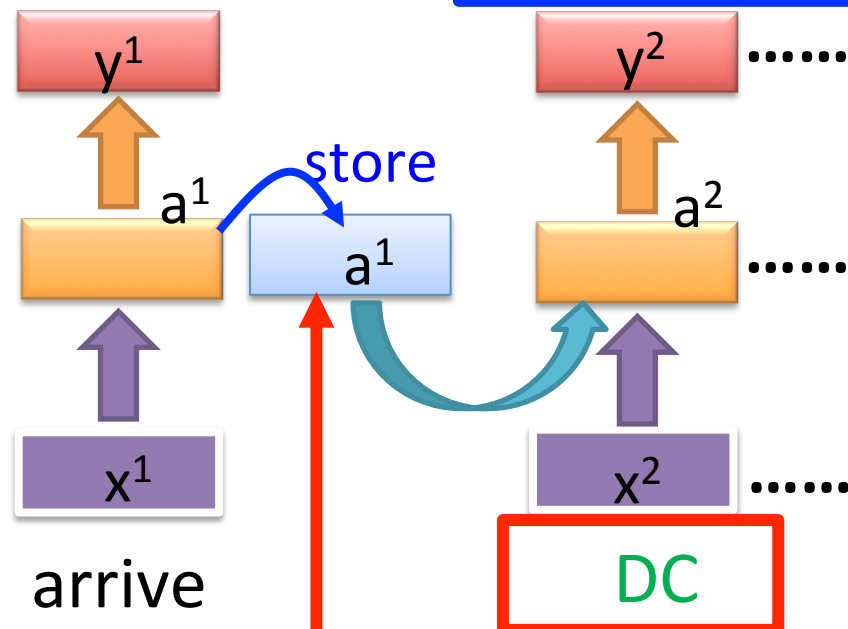
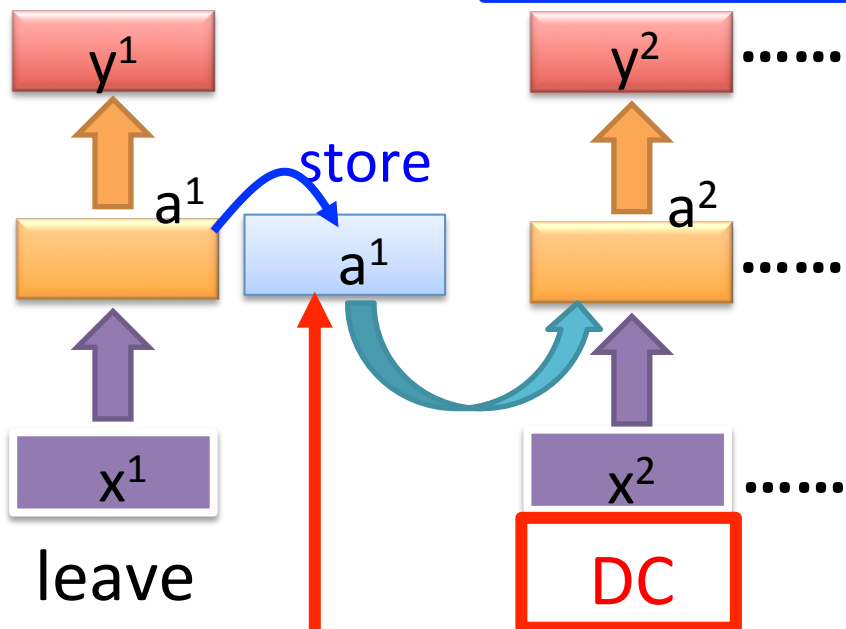
Different

Prob of "leave"  
in each slot

Prob of "DC" in  
each slot

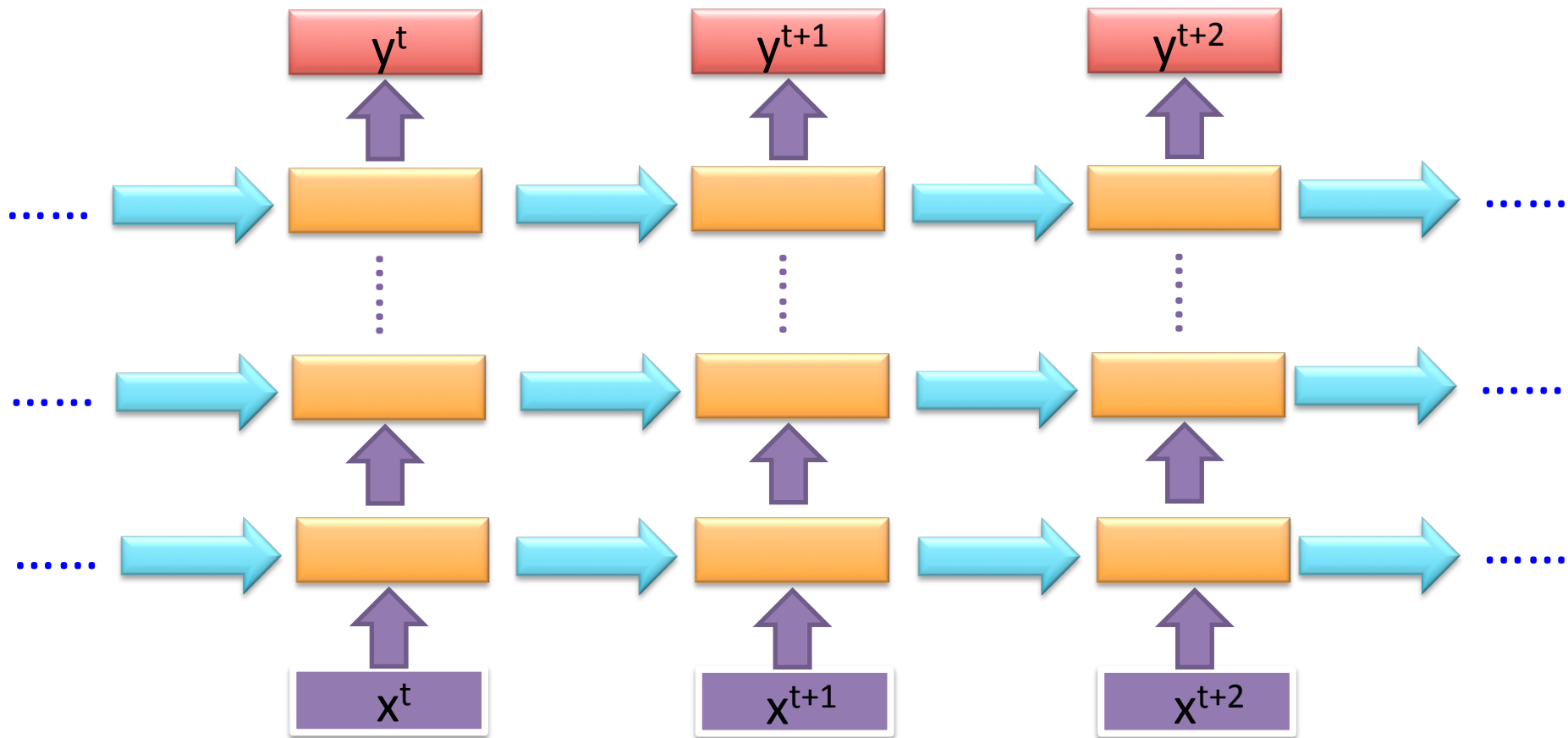
Prob of "arrive"  
in each slot

Prob of "DC" in  
each slot

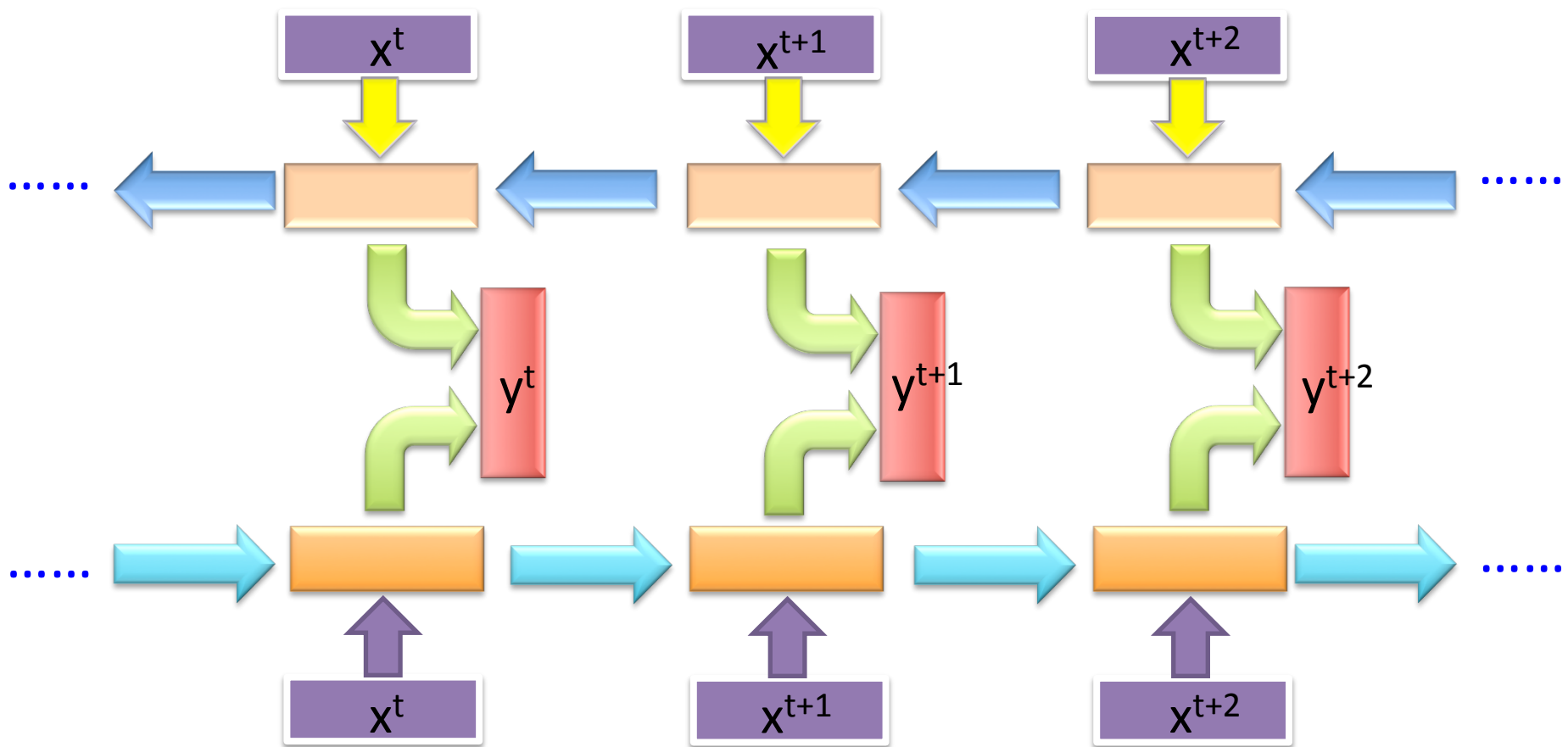


The values stored in the memory is different.

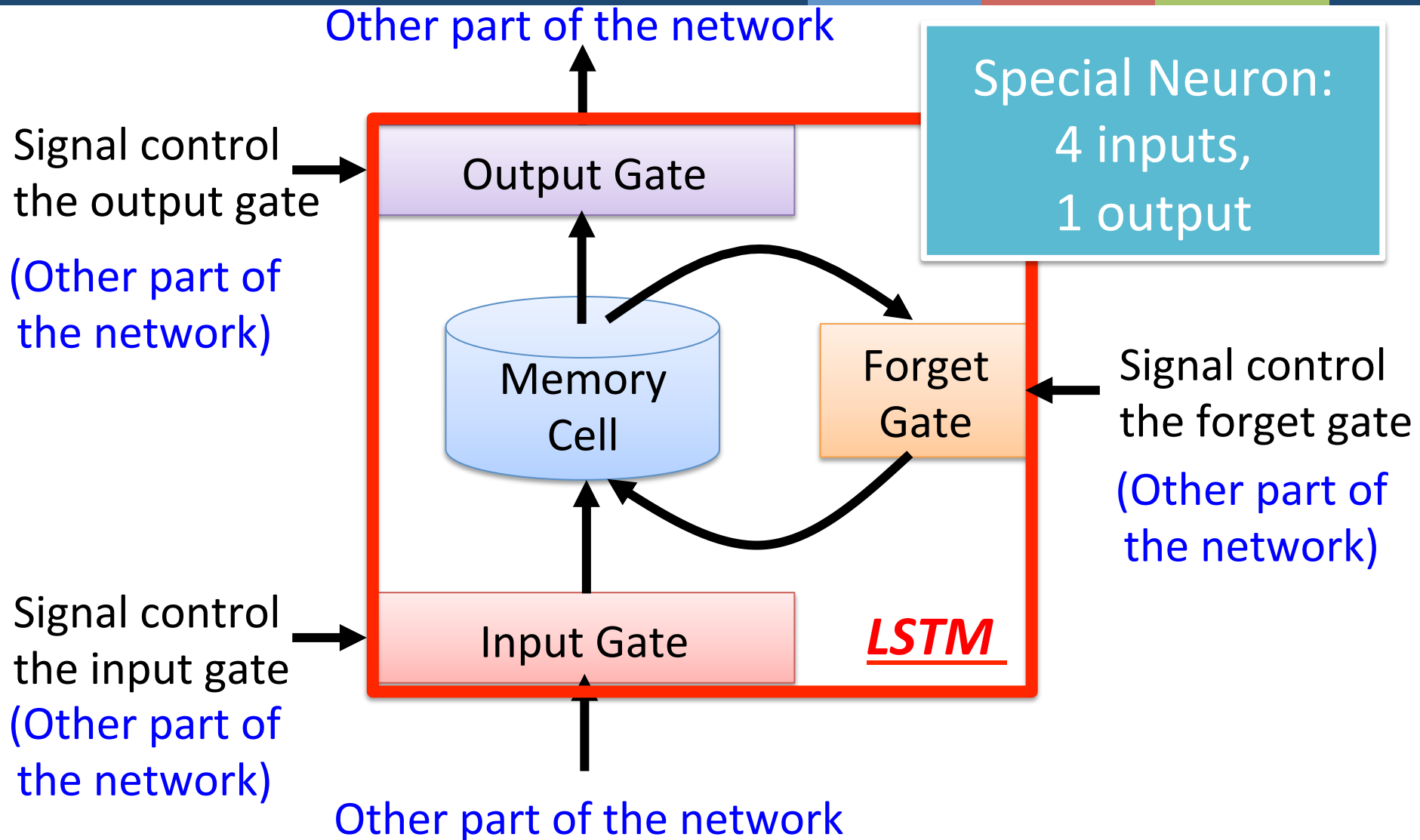
# Of course it can be deep ...

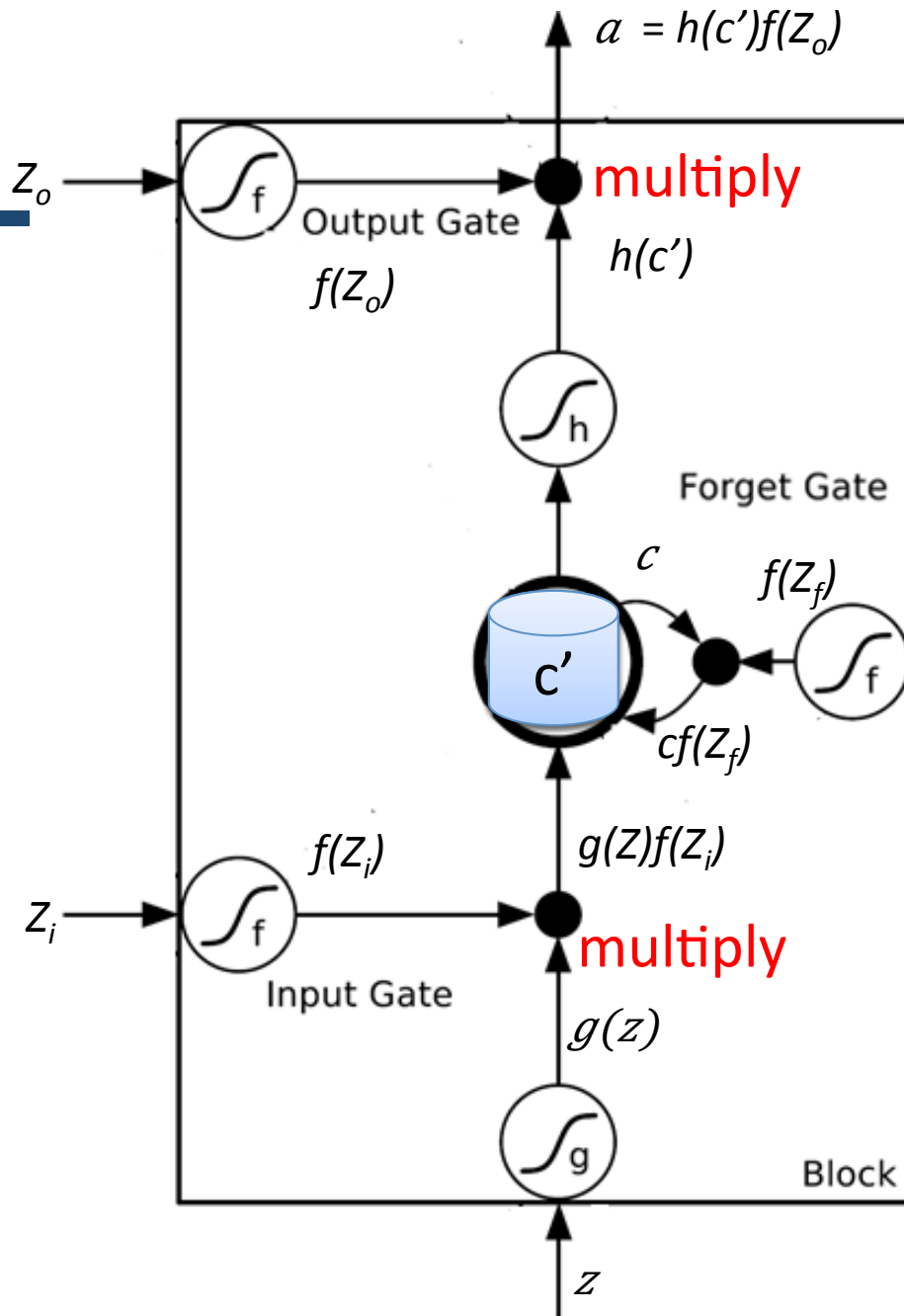


# Bidirectional RNN



# Long Short-term Memory (LSTM)





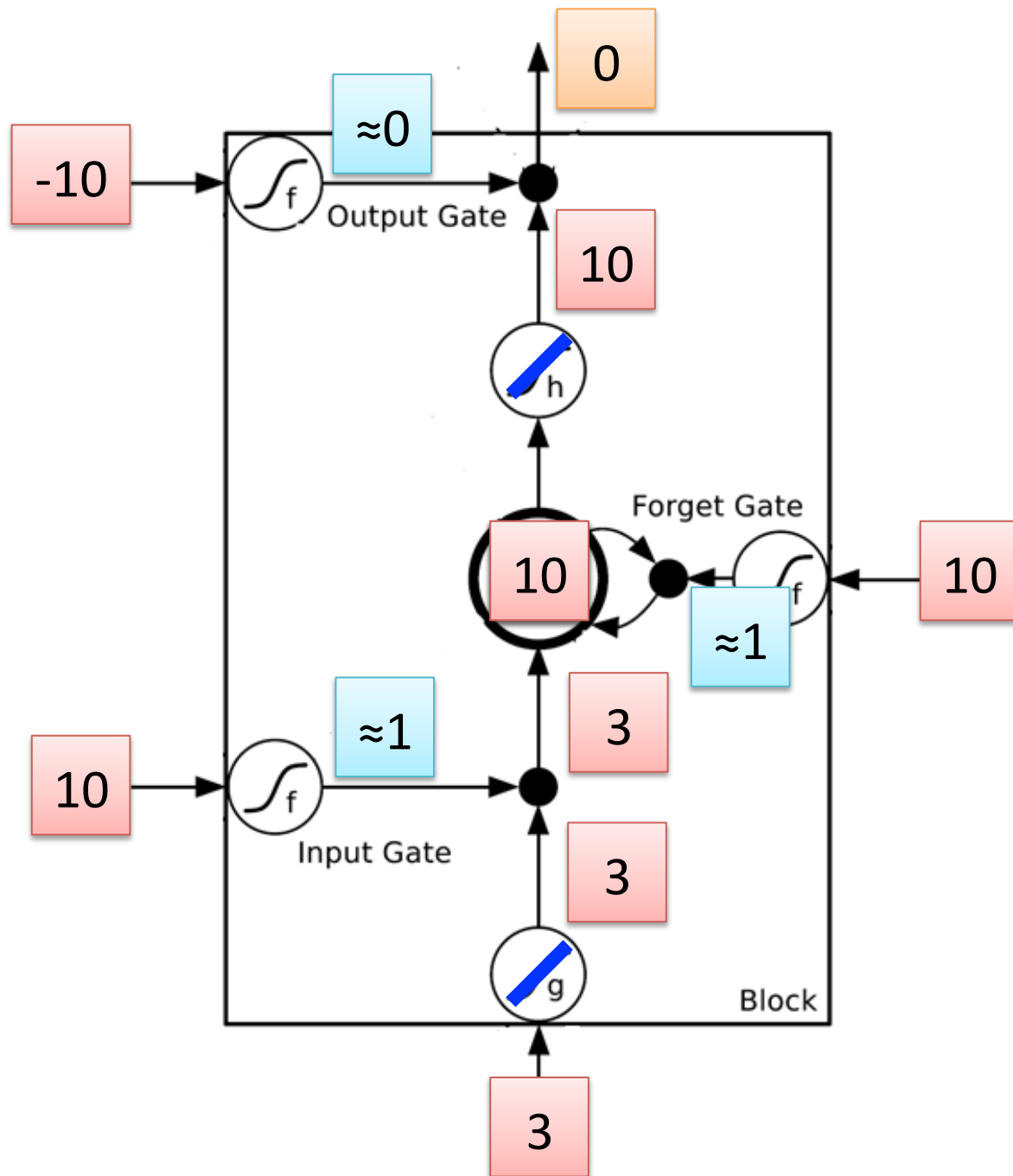
Activation function  $f$  is usually a sigmoid function

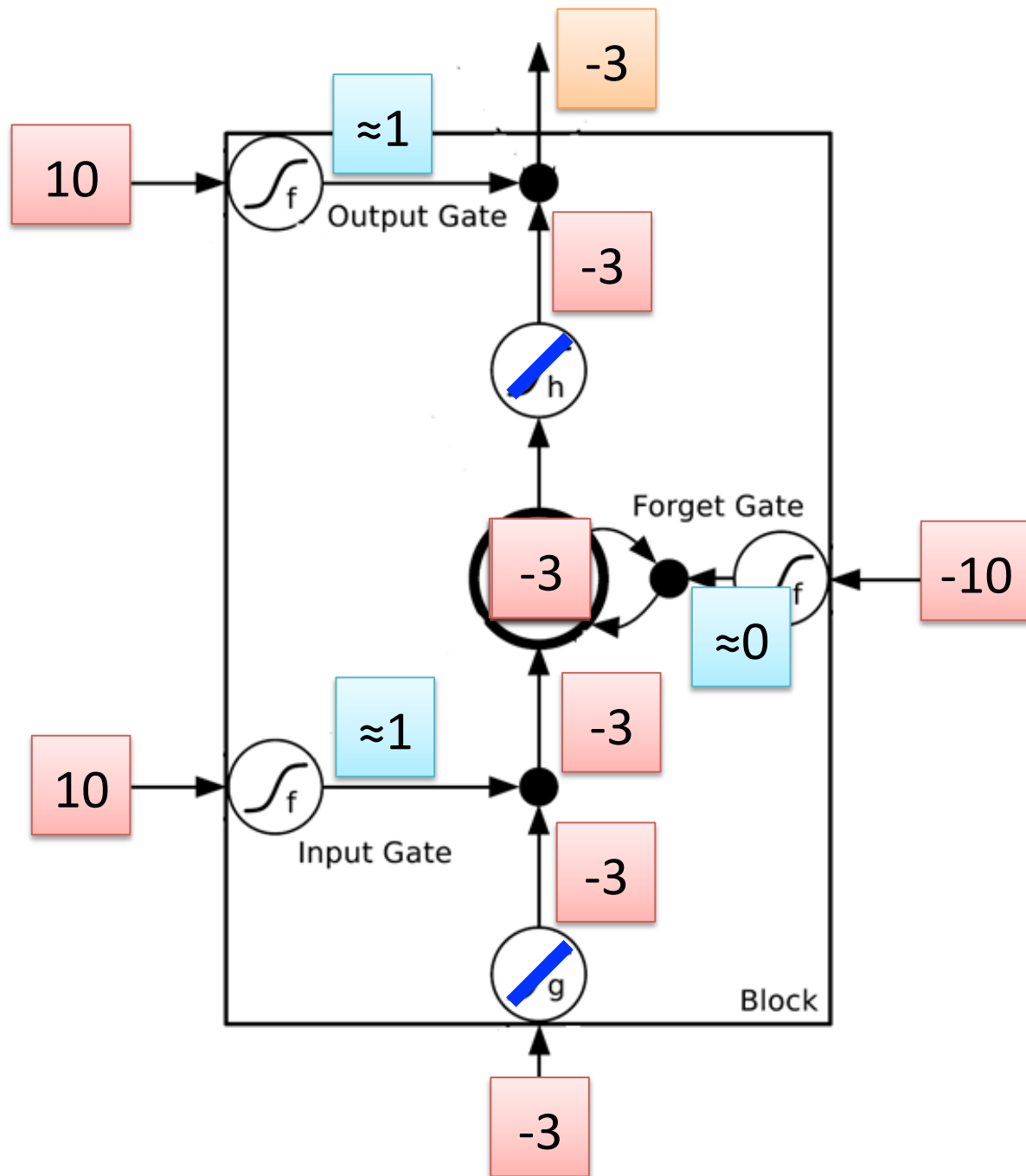
Between 0 and 1

Mimic open and close gate

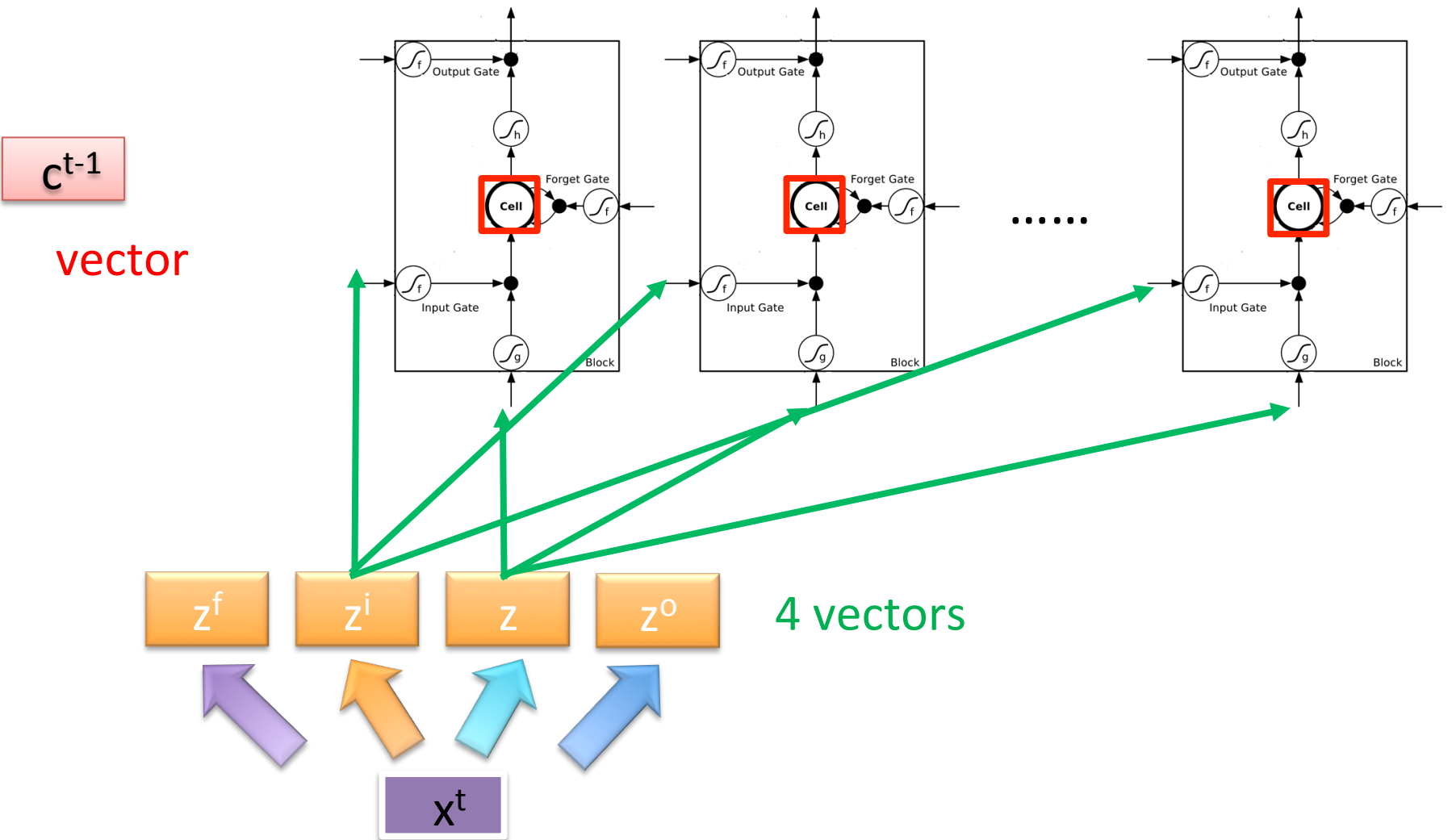
$$c' = g(z)f(Z_i) + cf(Z_f)$$



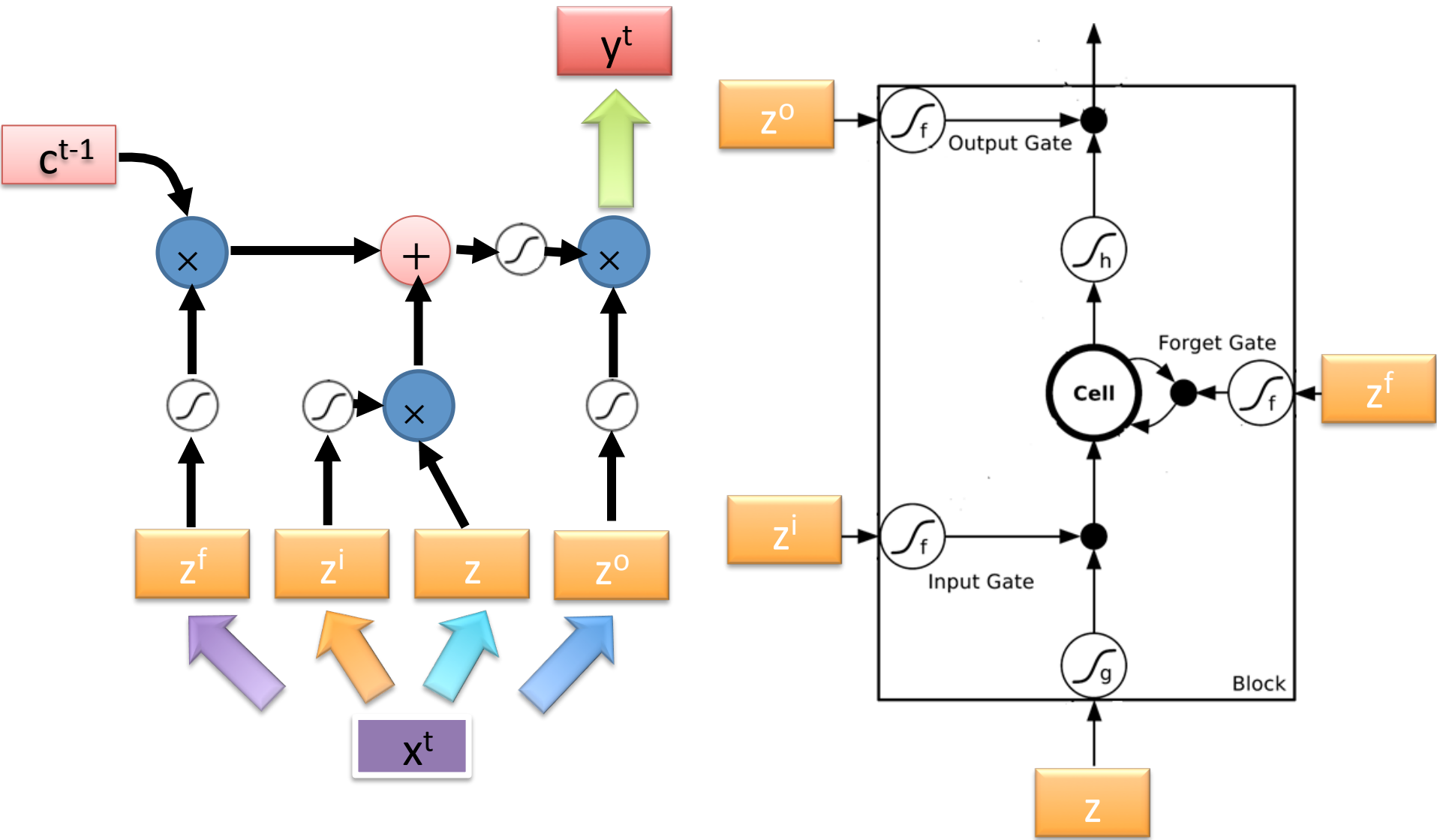




# LSTM

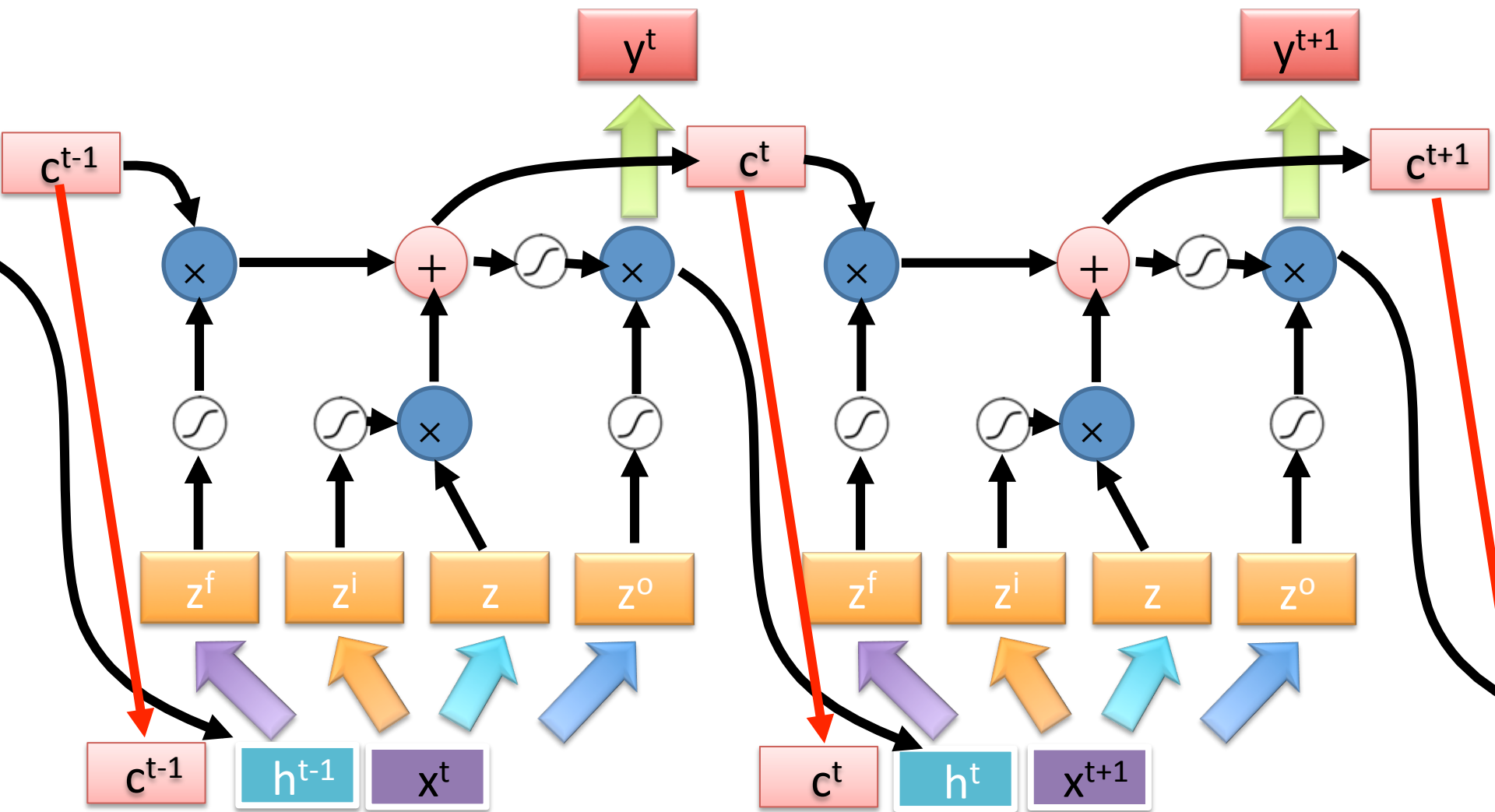


# LSTM

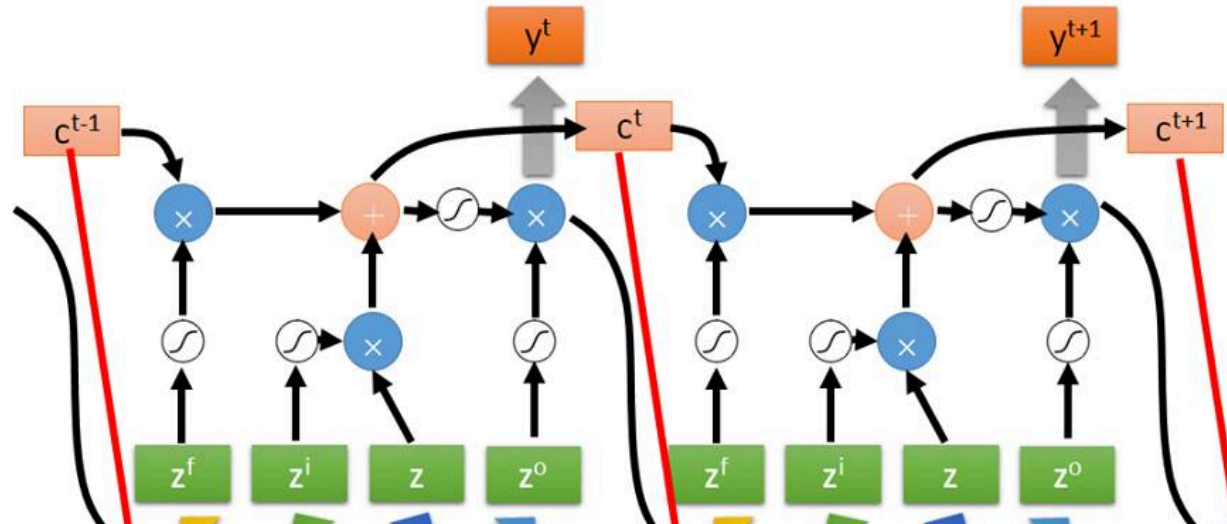


# LSTM

Extension: "peephole"



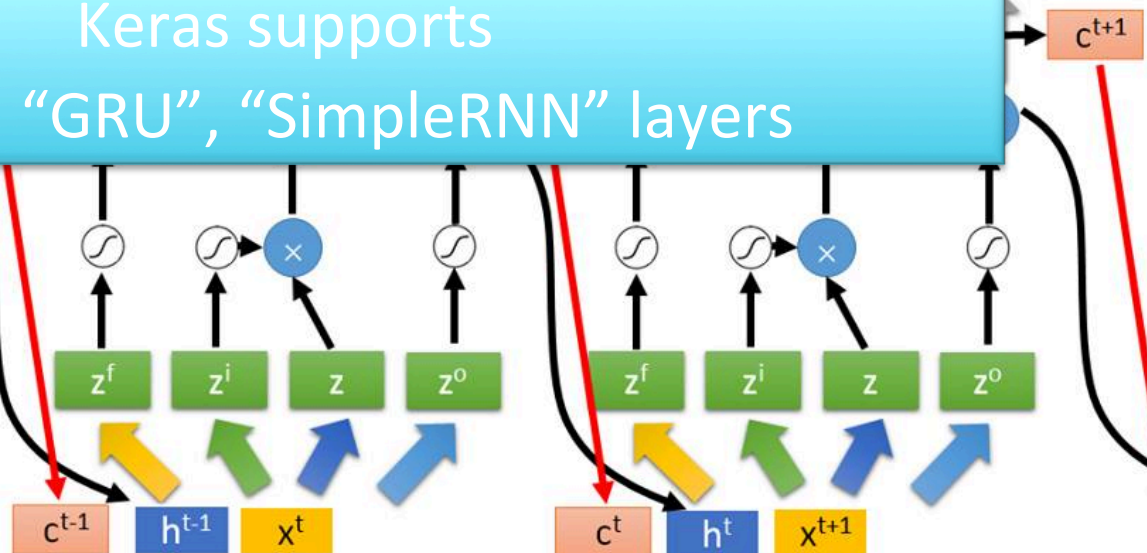
# Multiple-layer LSTM



Don't worry if you cannot understand this.  
Keras can handle it.

Keras supports  
“LSTM”, “GRU”, “SimpleRNN” layers

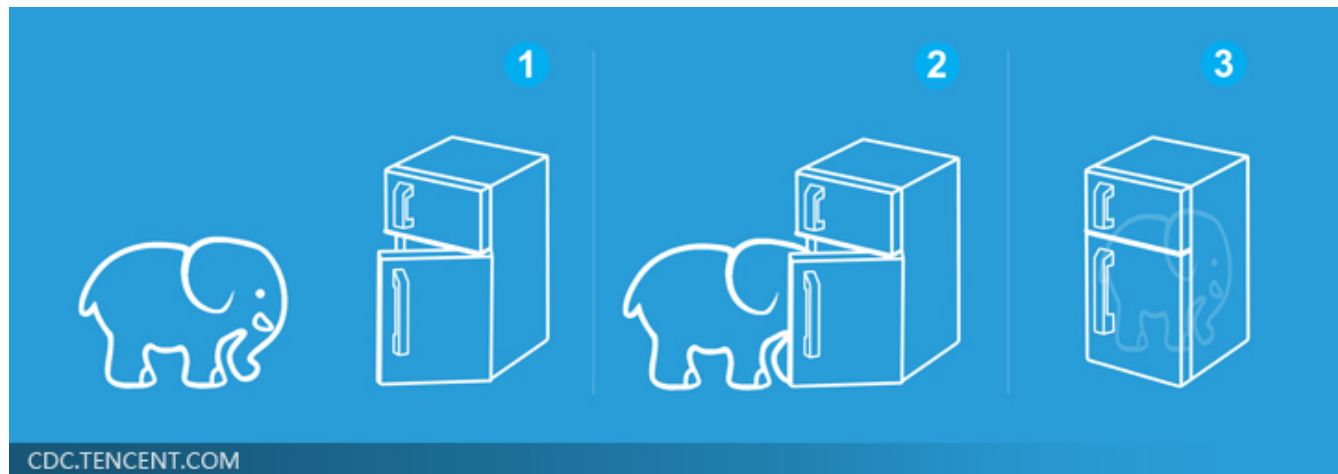
This is quite  
standard now.



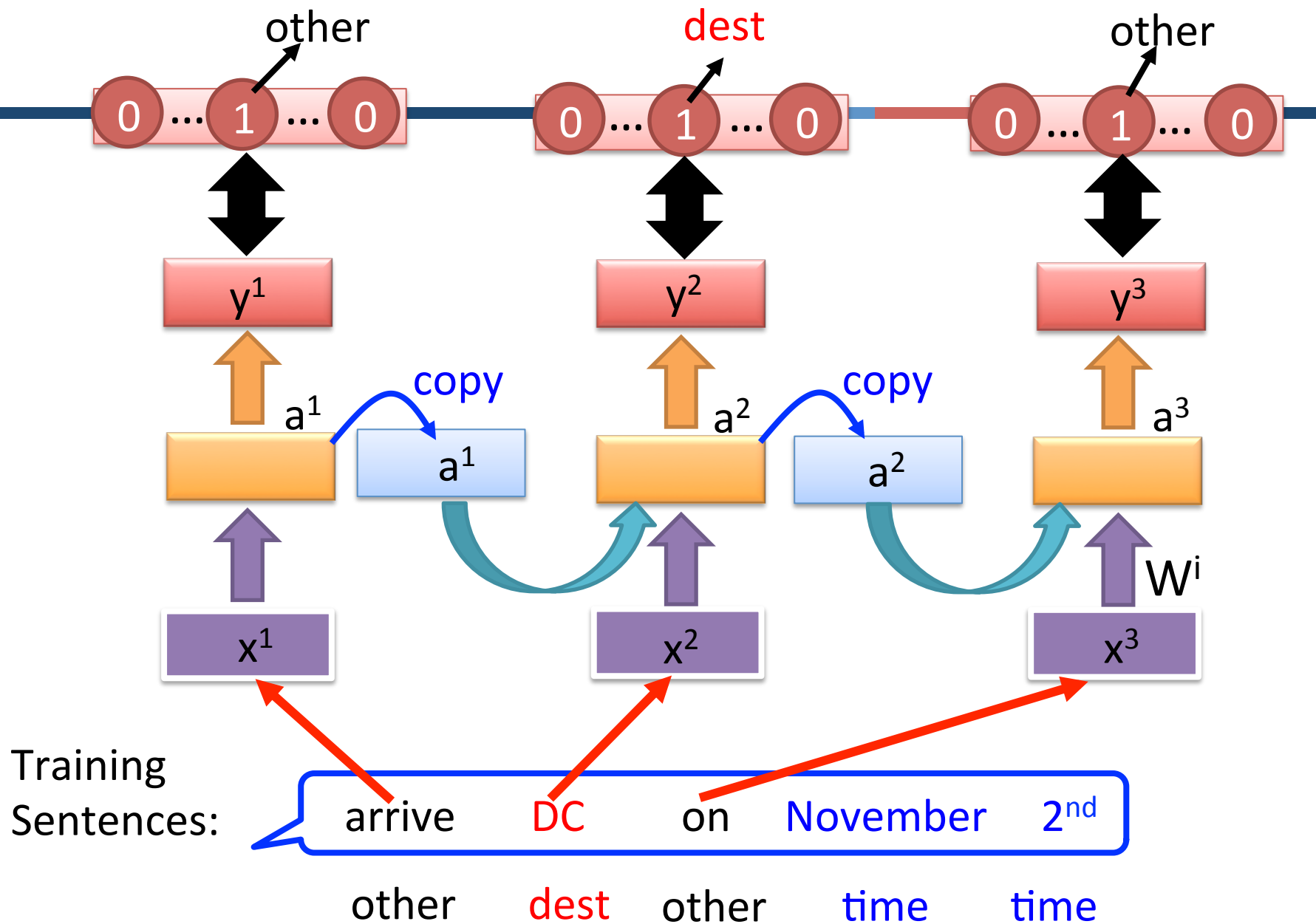
# Three Steps for Deep Learning



Deep Learning is so simple .....



# Learning Target

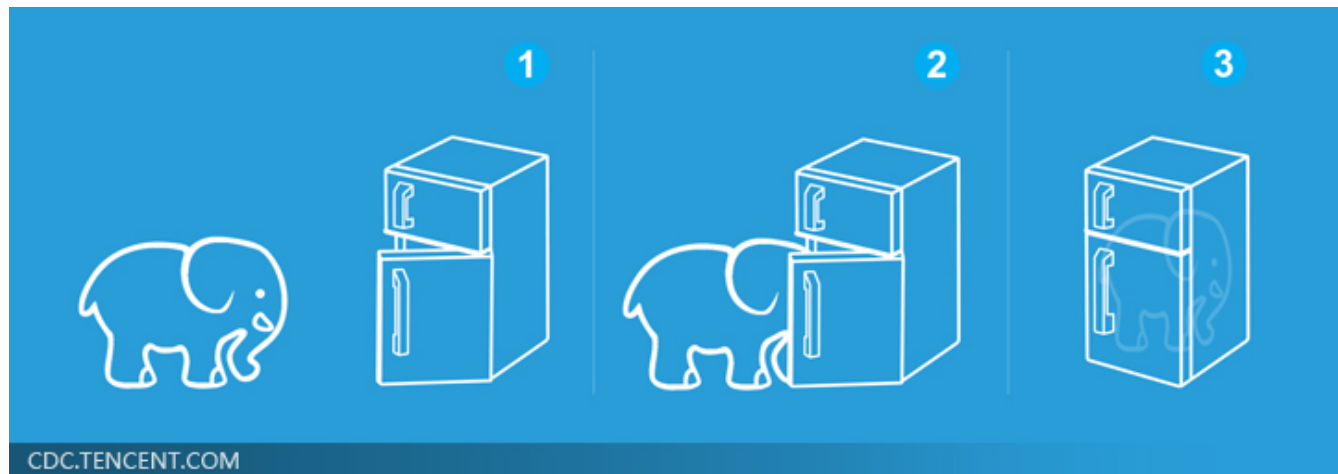




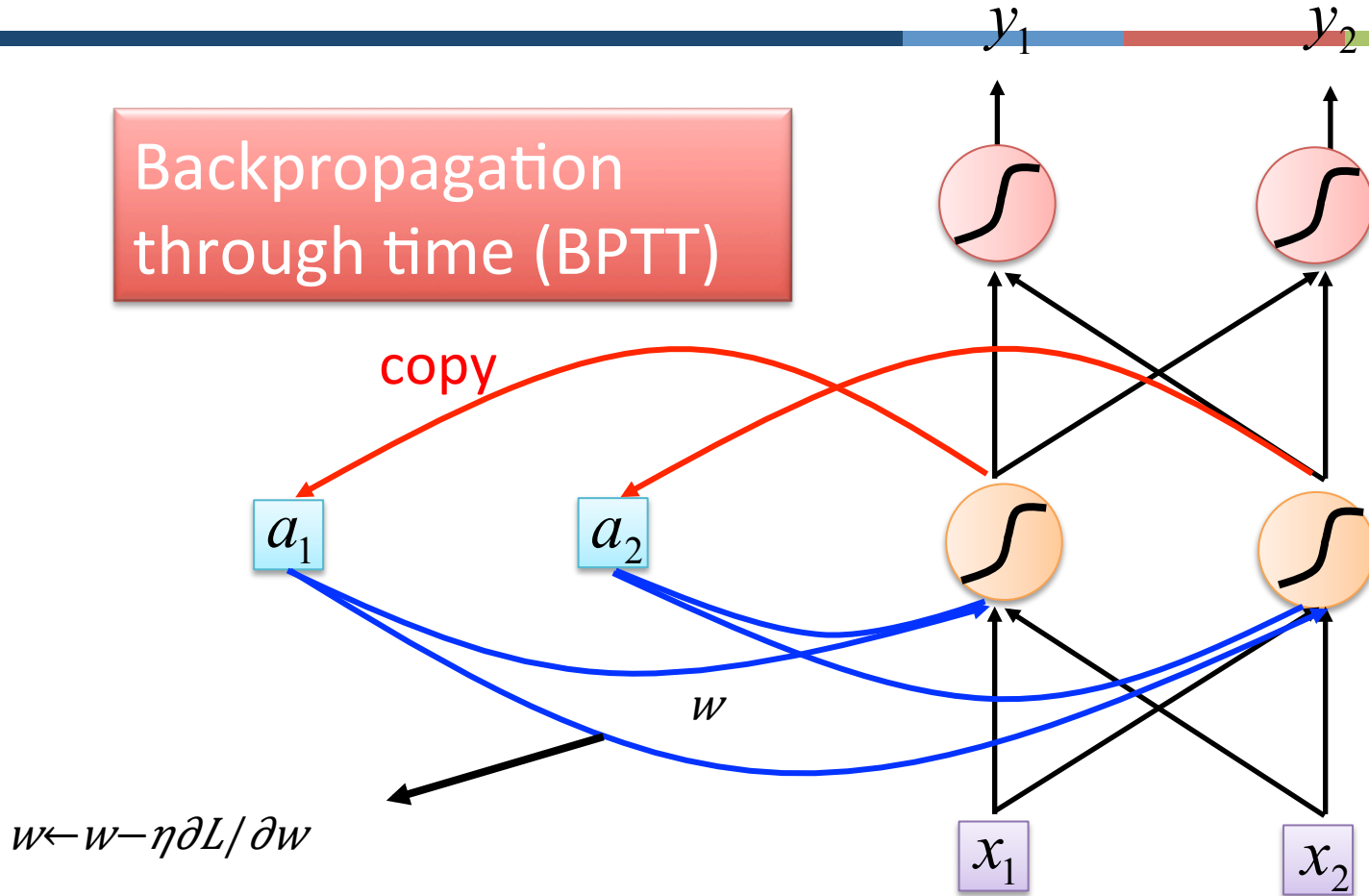
# Three Steps for Deep Learning



Deep Learning is so simple .....



# Learning

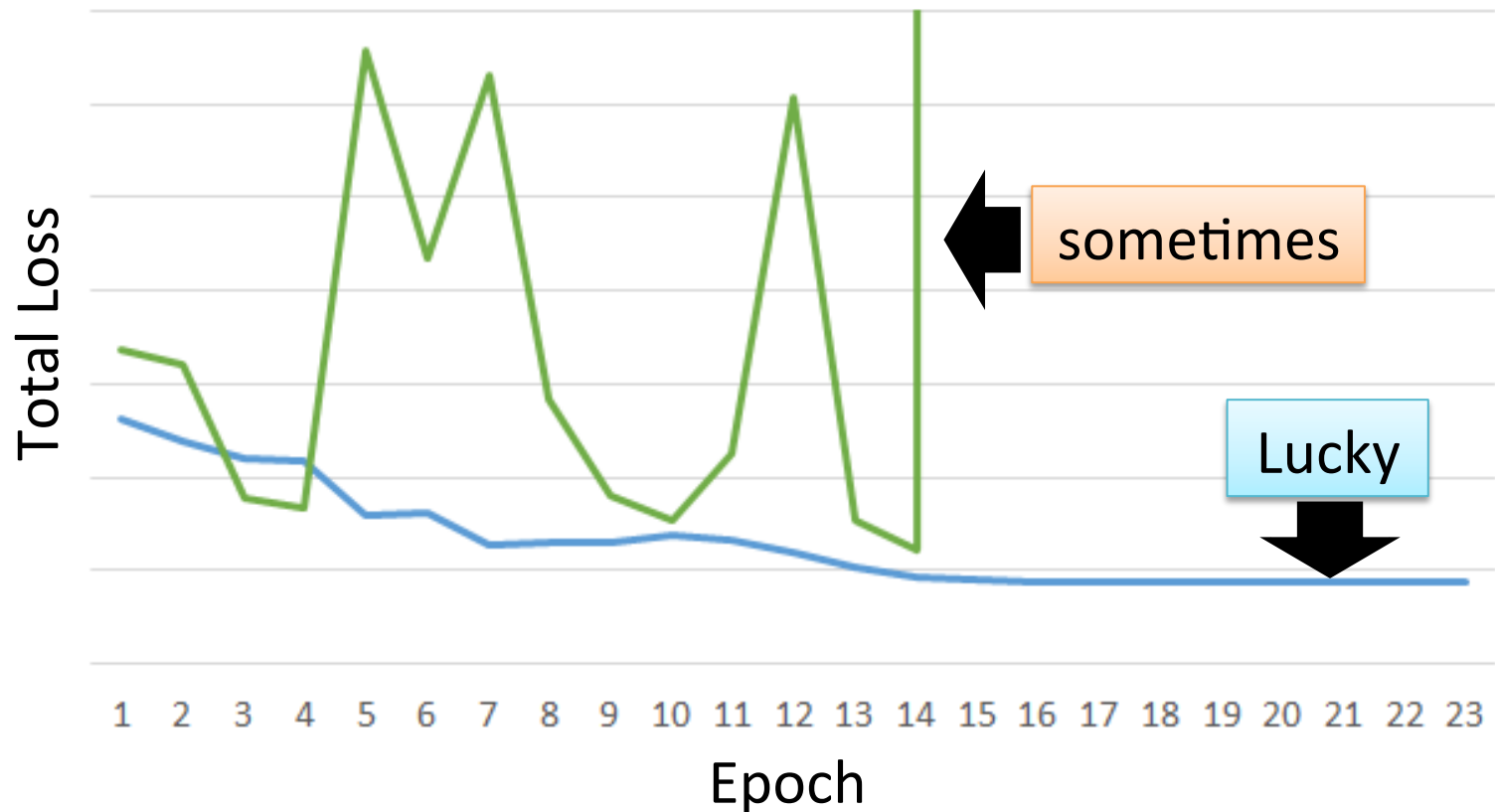


RNN Learning is very difficult in practice.

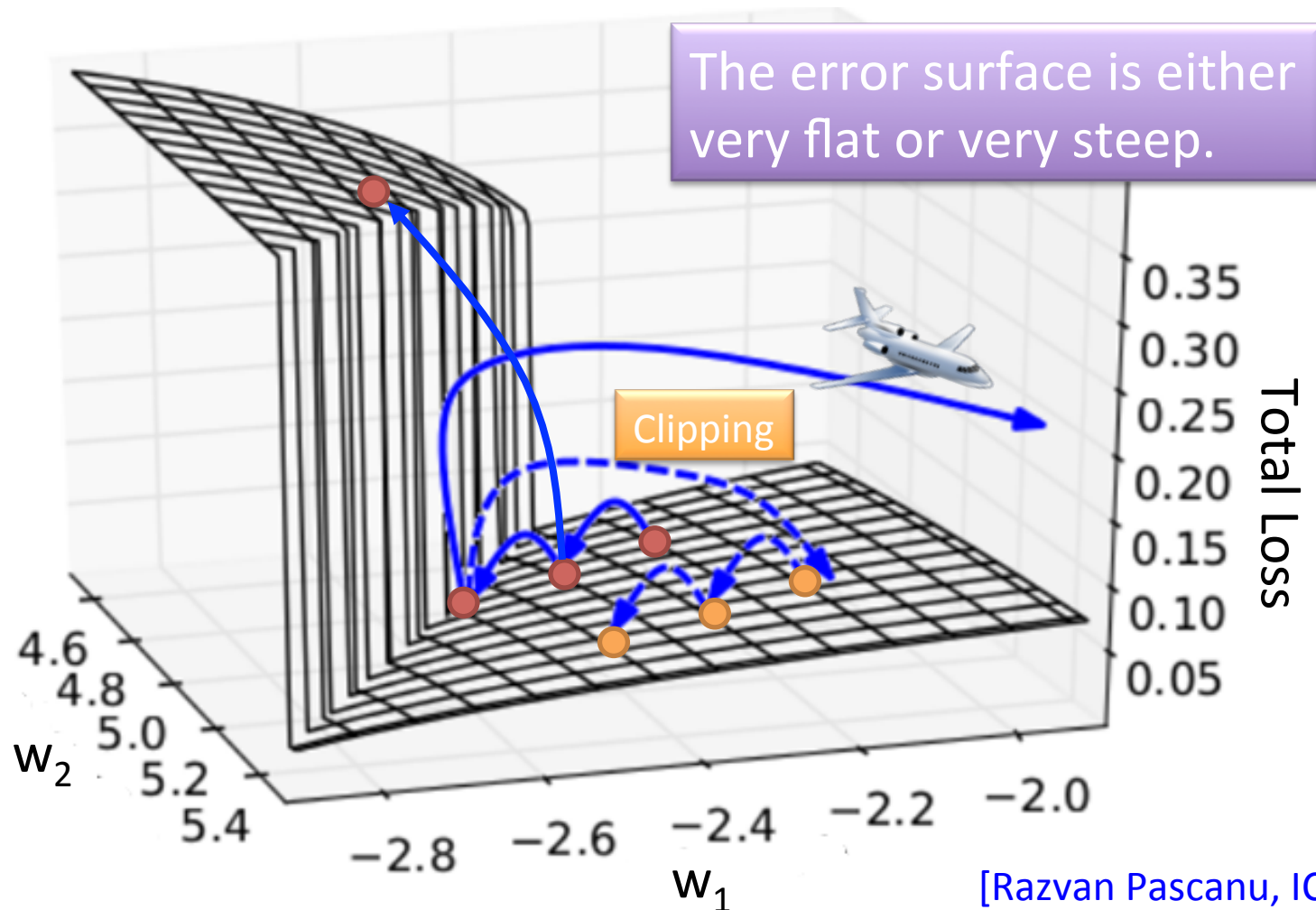
# Unfortunately .....

- RNN-based network is not always easy to learn

Real experiments on Language modeling

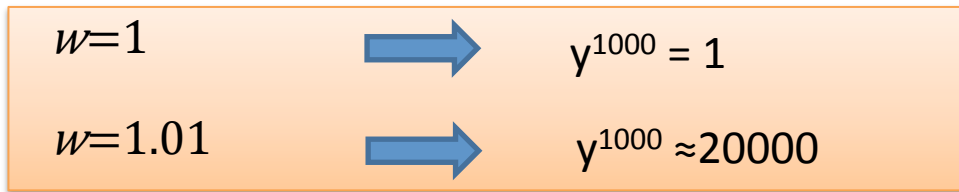


# The error surface is rough.



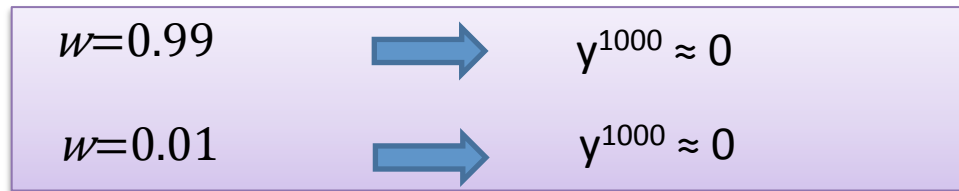
[Razvan Pascanu, ICML'13]

# Why?



Large  $\partial L / \partial w$

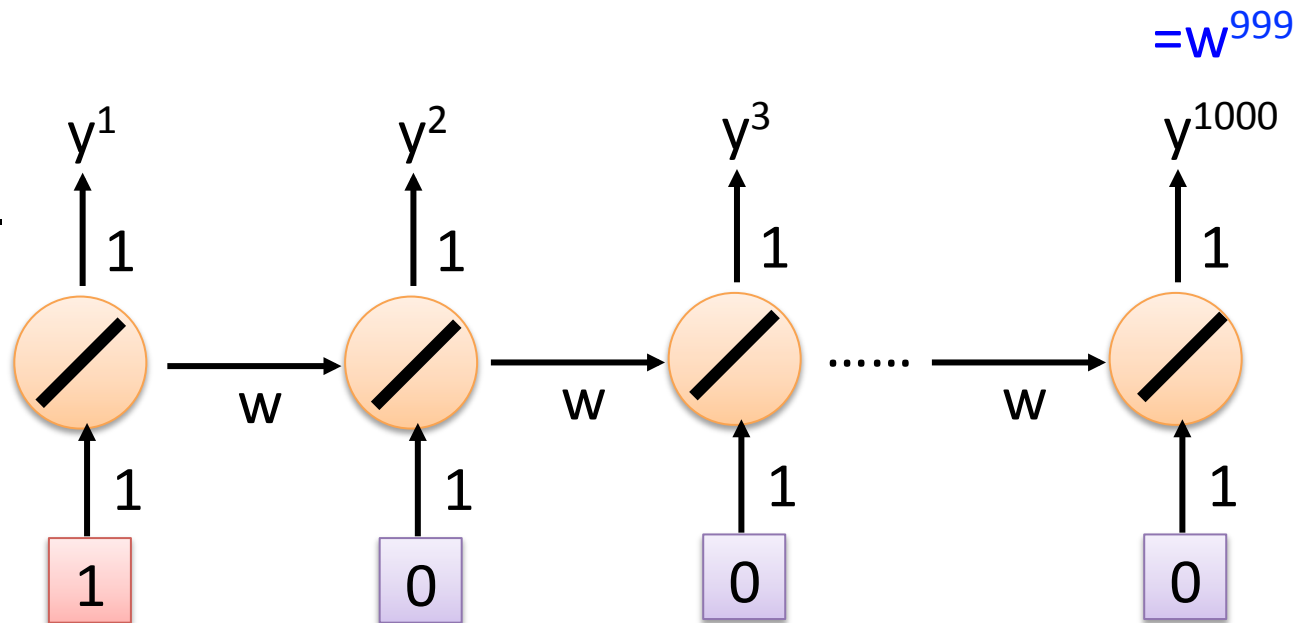
Small Learning rate?



small  $\partial L / \partial w$

Large Learning rate?

## Toy Example



# Helpful Techniques

- Long Short-term Memory (LSTM)

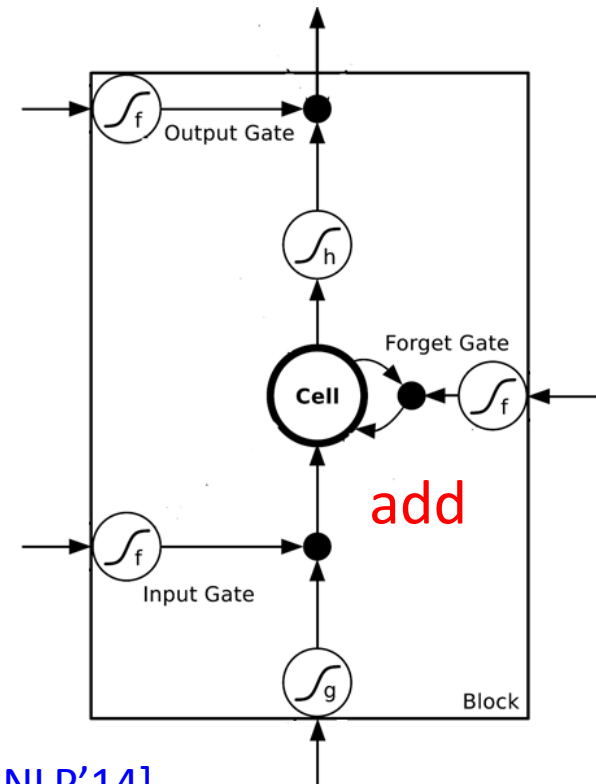
- Can deal with gradient vanishing (not gradient explode)

- Memory and input are **added**

- The influence never disappears unless forget gate is closed

➡ No Gradient vanishing  
(If forget gate is

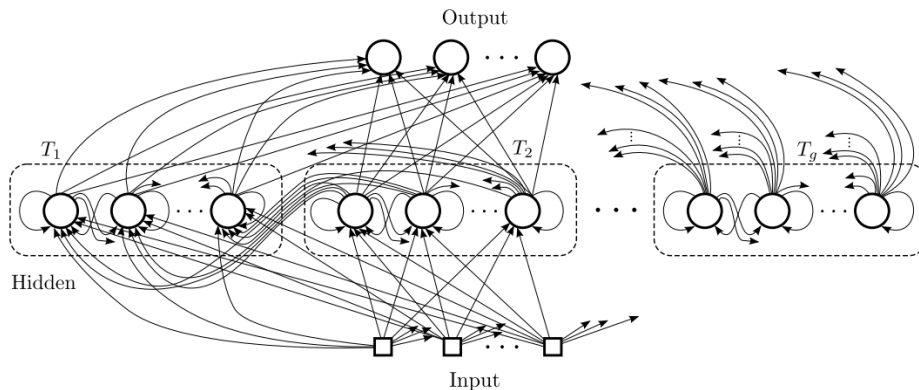
Gated Recurrent Unit (GRU):  
simpler than LSTM



[Cho, EMNLP'14]

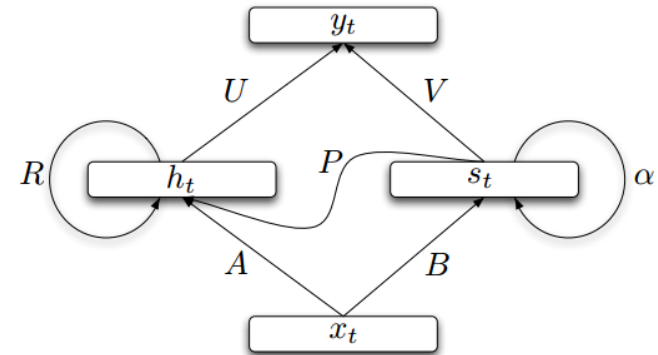
# Helpful Techniques

## Clockwise RNN



[Jan Koutnik, JMLR'14]

## Structurally Constrained Recurrent Network (SCRN)



[Tomas Mikolov, ICLR'15]

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

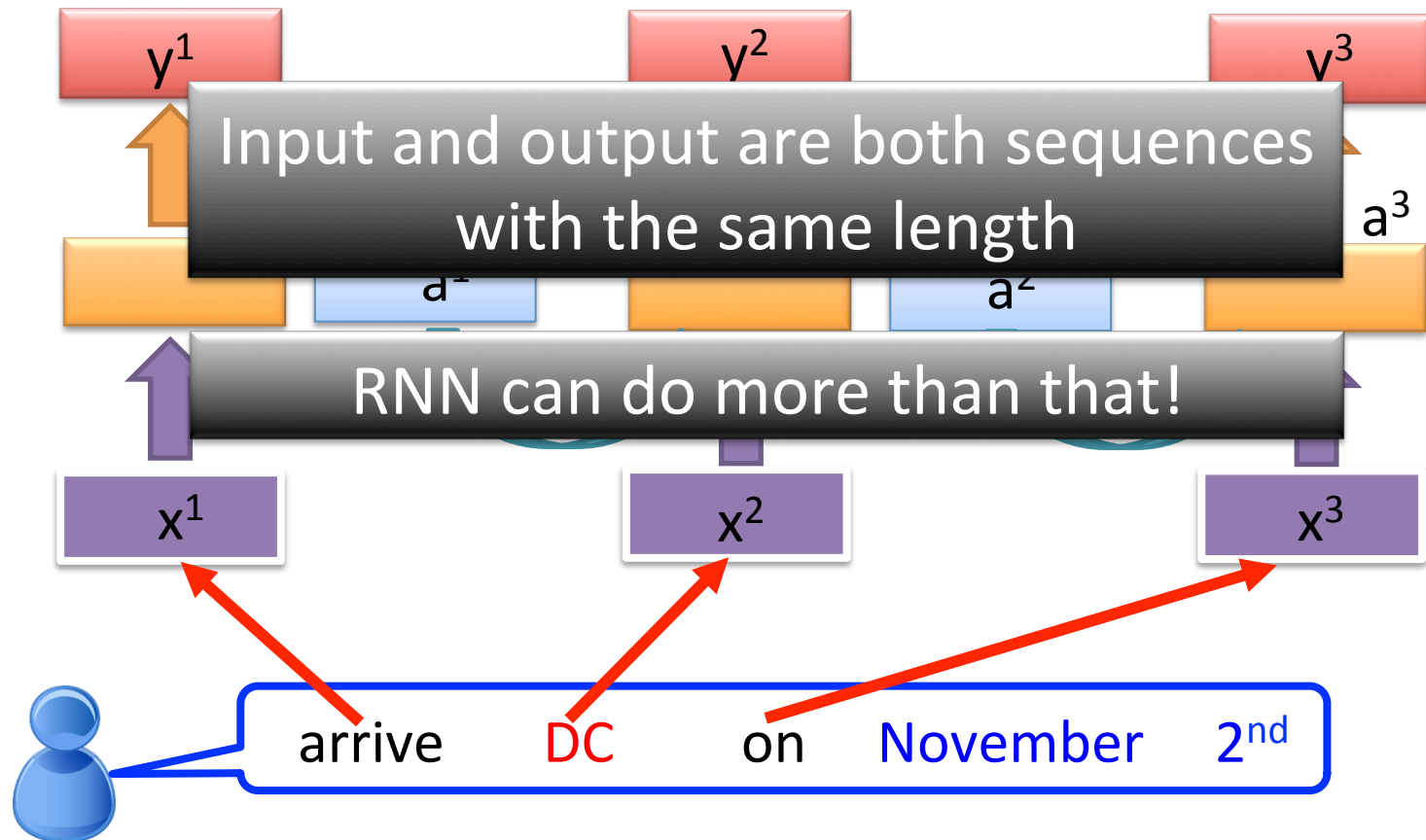
➤ Outperform or be comparable with LSTM in 4 different tasks

# More Applications .....

Probability of  
“arrive” in each slot

Probability of “**DC**”  
in each slot

Probability of  
“on” in each slot

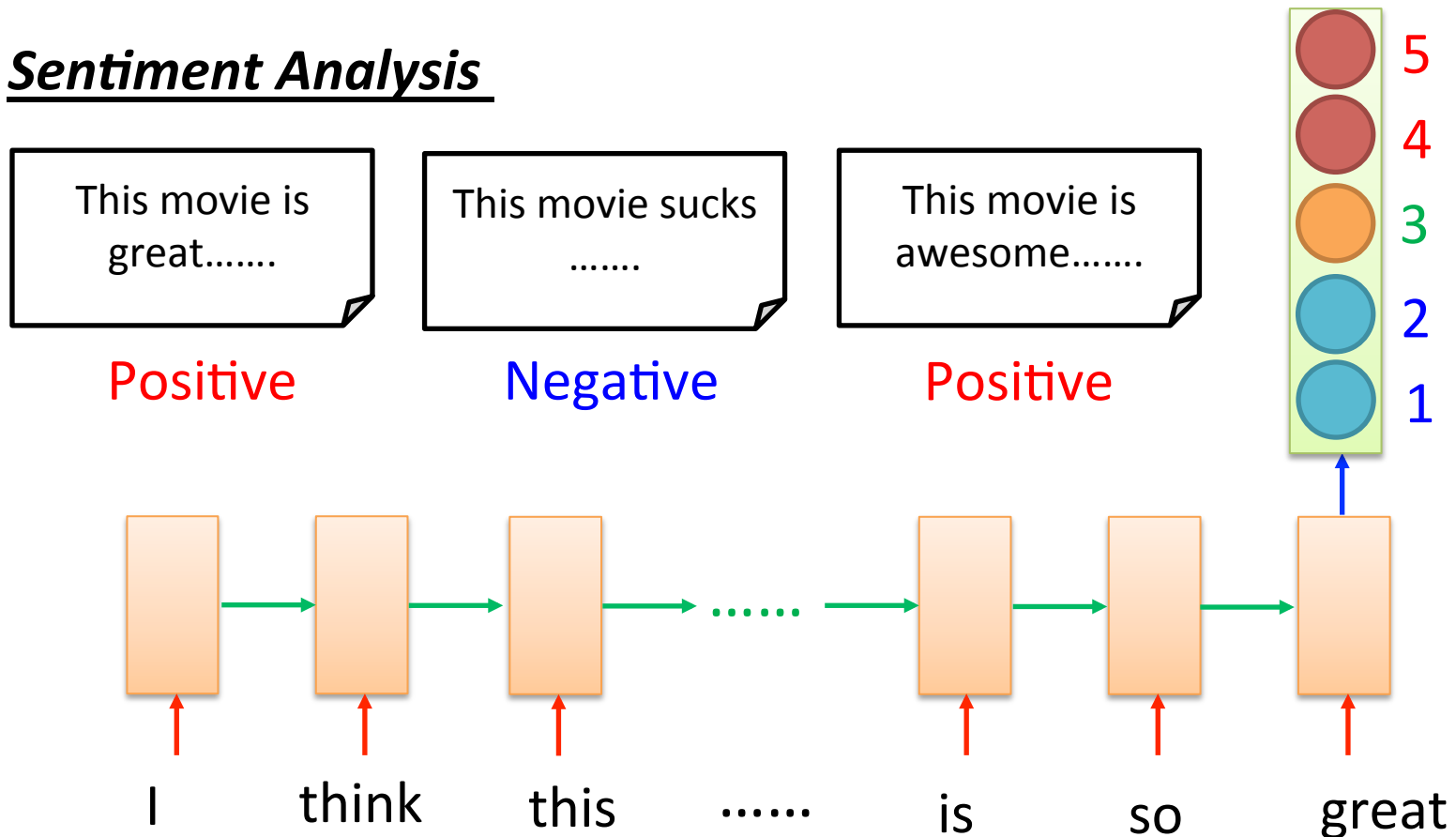




# Many to one

- Input is a vector sequence, but output is only one vector

## Sentiment Analysis

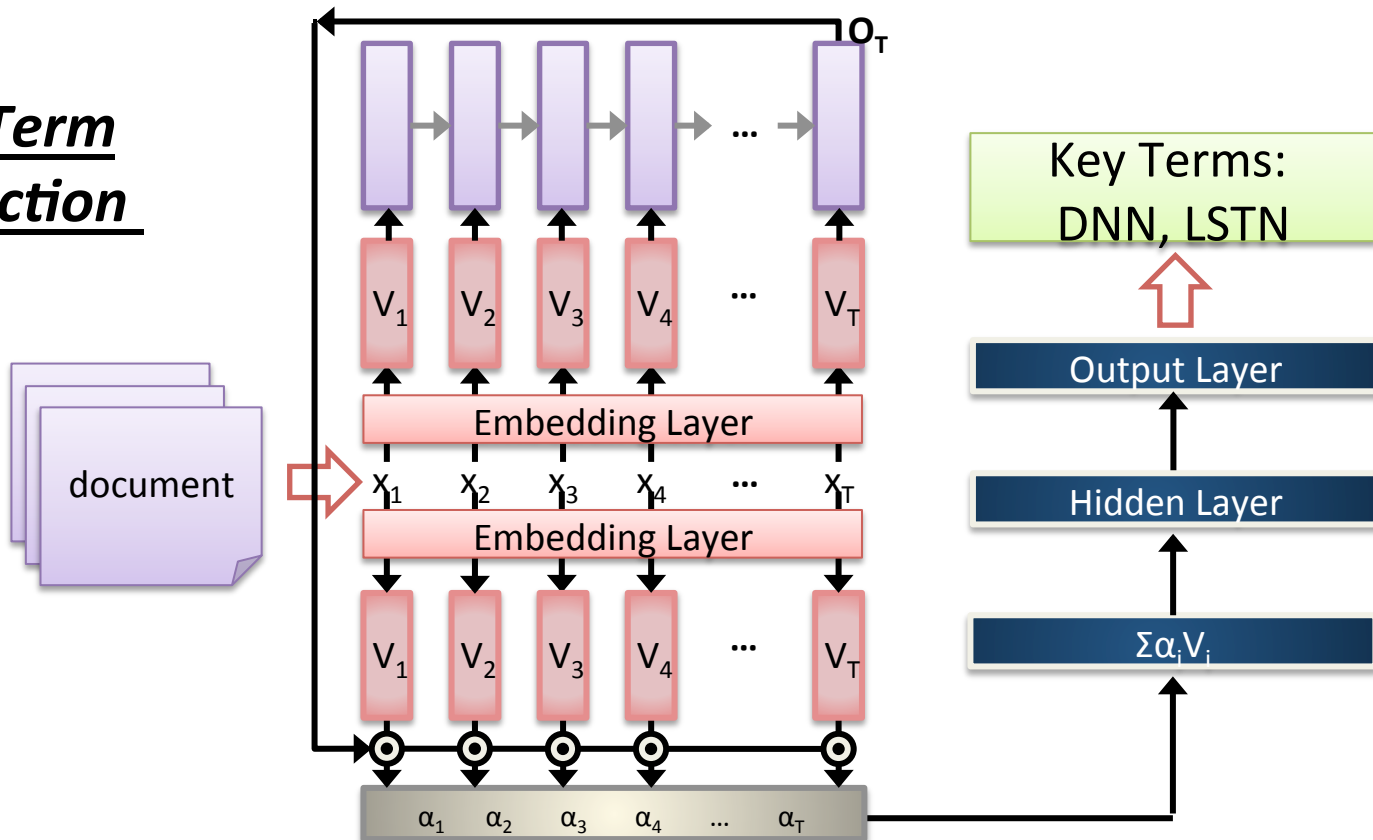


# Many to one

[Shen & Lee, Interspeech 16]

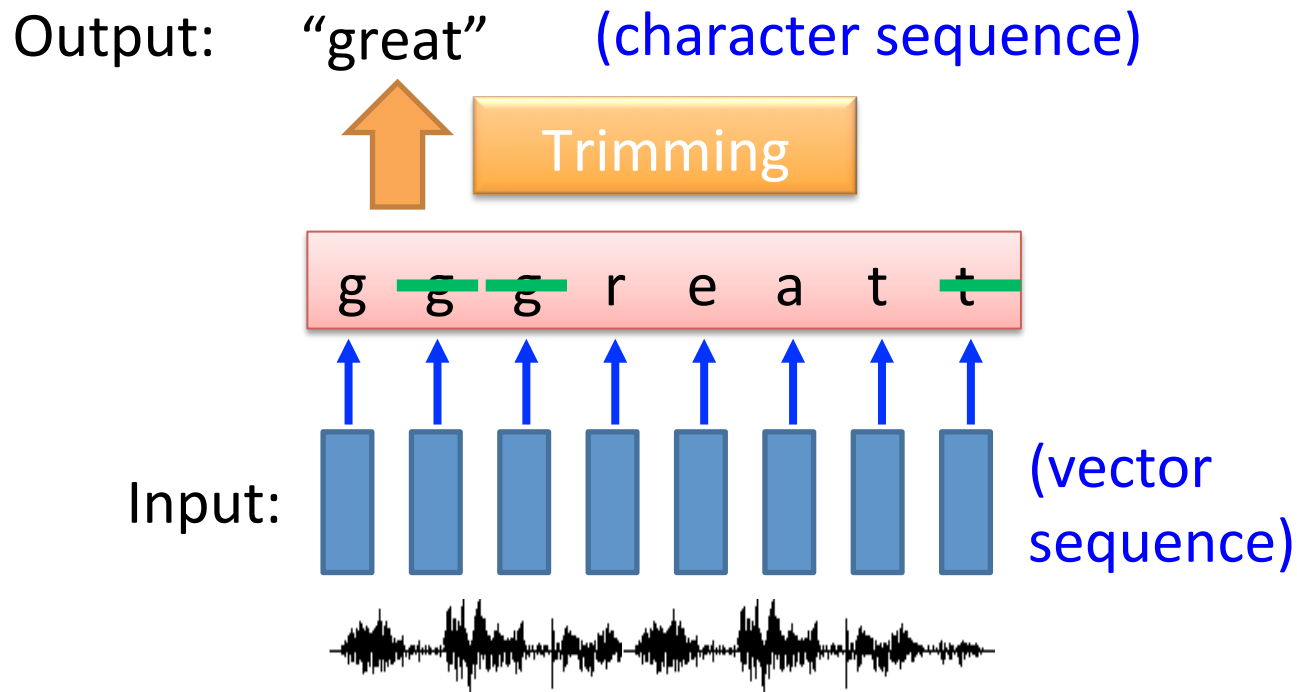
- Input is a vector sequence, but output is only one vector

## Key Term Extraction



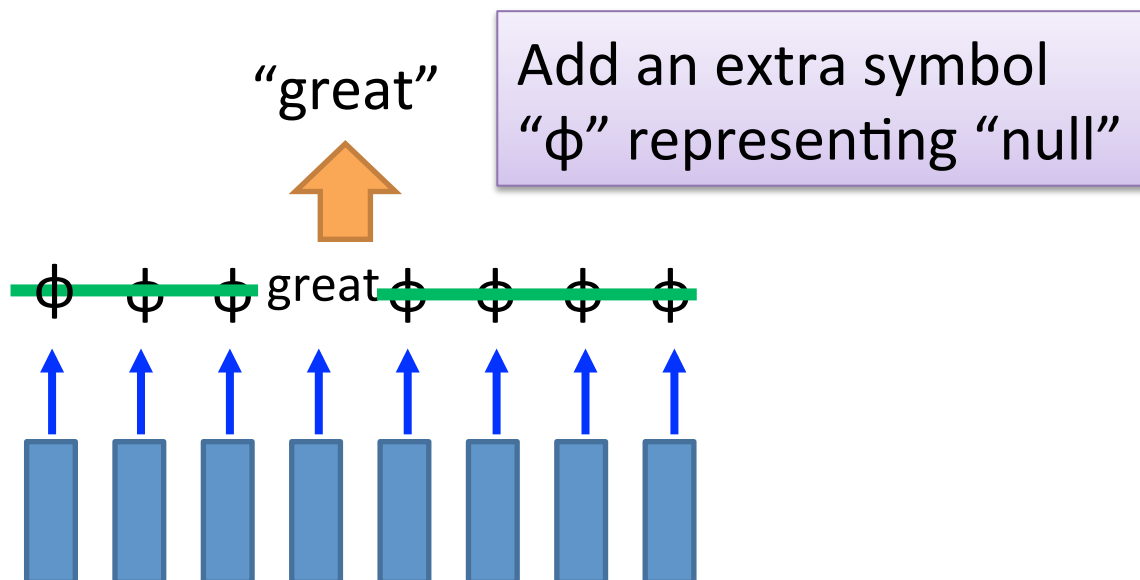
# Many to Many (Output is shorter)

- Both input and output are both sequences, **but the output is shorter.**
  - E.g. **Speech Recognition**



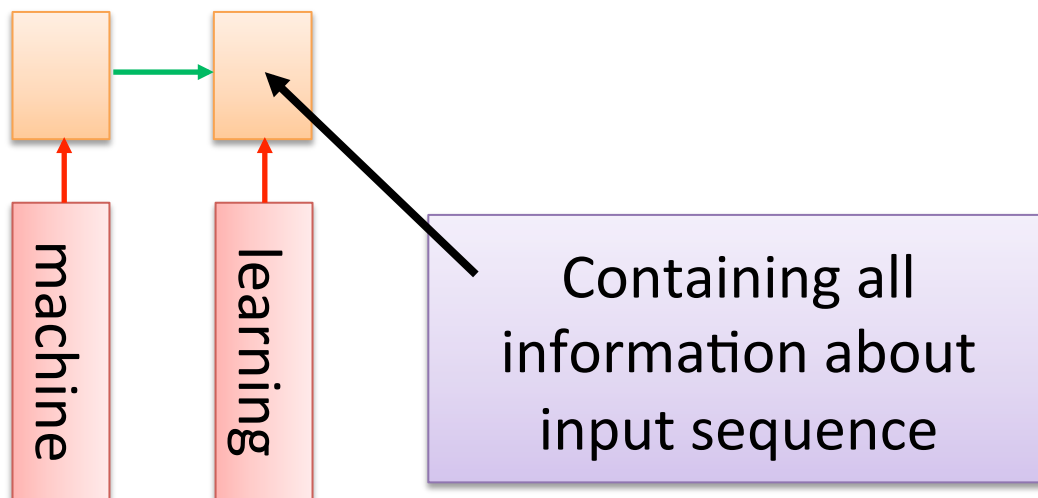
# Many to Many (Output is shorter)

- Both input and output are both sequences, **but the output is shorter.**
- Connectionist Temporal Classification (CTC) [Alex Graves, ICML'06]  
[Alex Graves, ICML'14][Haşim Sak, Interspeech'15][Jie Li, Interspeech'15][Andrew Senior, ASRU'15]



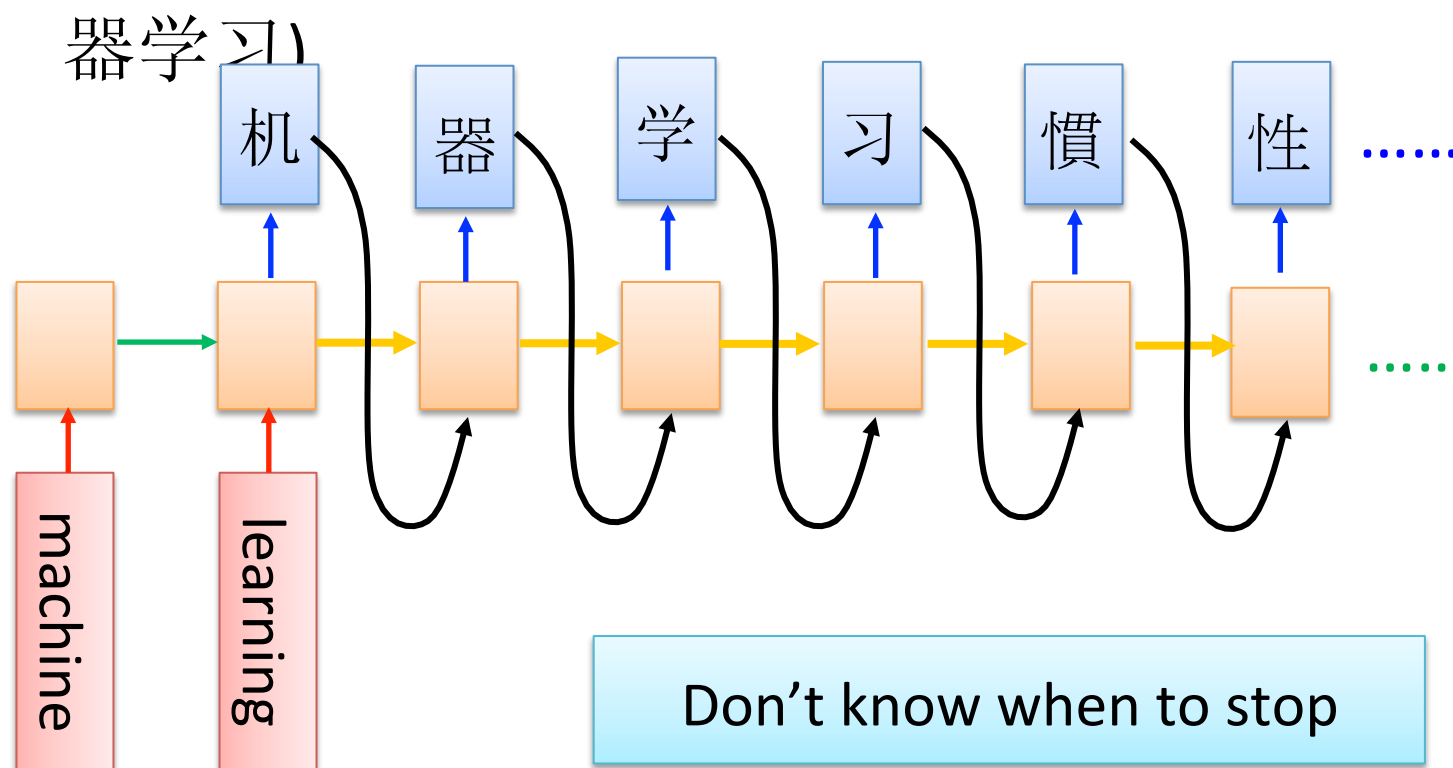
# Many to Many (No Limitation)

- Both input and output are both sequences *with different lengths*. → *Sequence to sequence learning*
  - E.g. *Machine Translation* (machine learning → 机器学习)



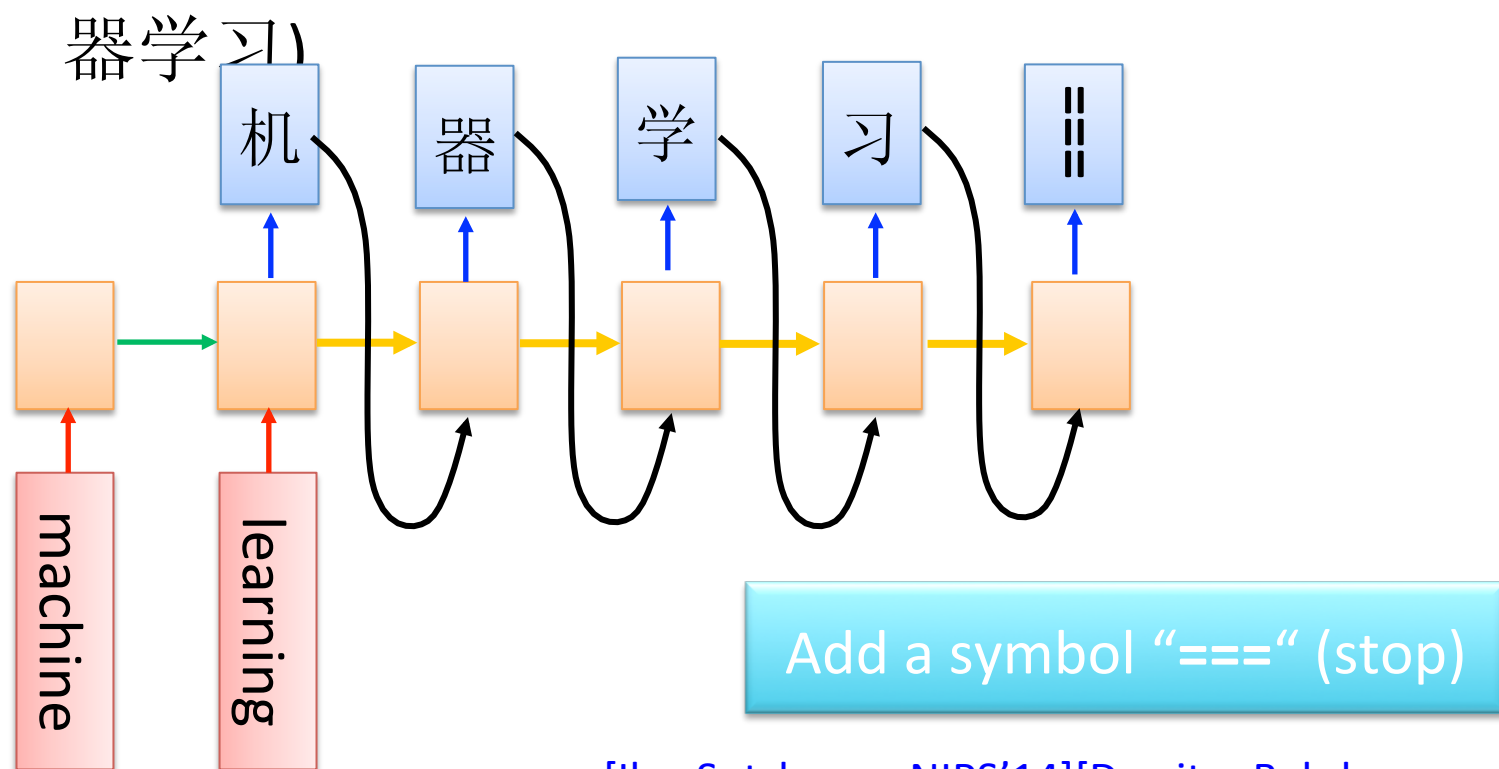
# Many to Many (No Limitation)

- Both input and output are both sequences *with different lengths*. → *Sequence to sequence learning*
  - E.g. *Machine Translation* (machine learning → 机器学习)



# Many to Many (No Limitation)

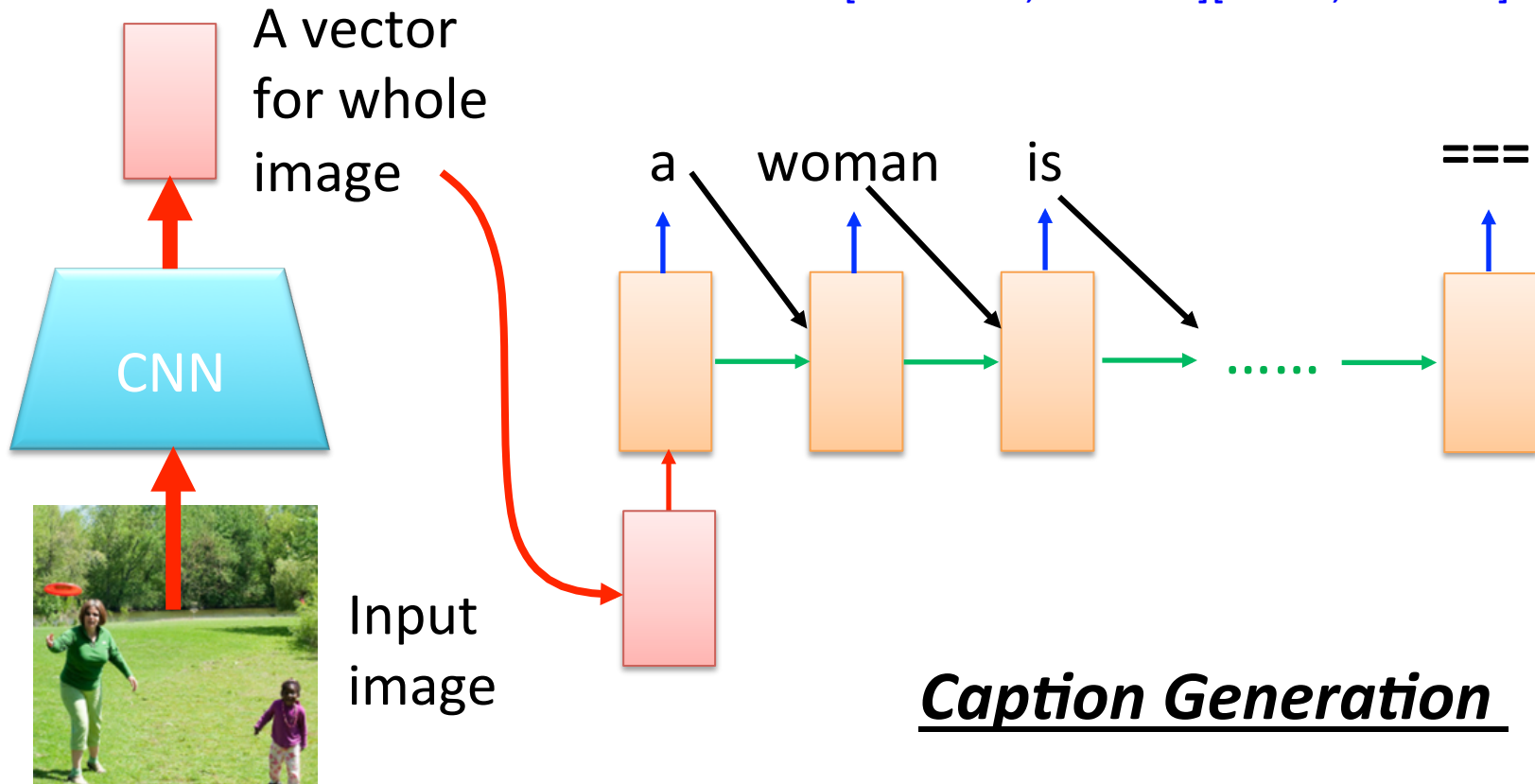
- Both input and output are both sequences *with different lengths*. → *Sequence to sequence learning*
  - E.g. *Machine Translation* (machine learning → 机器学习)



# Image Caption Generation

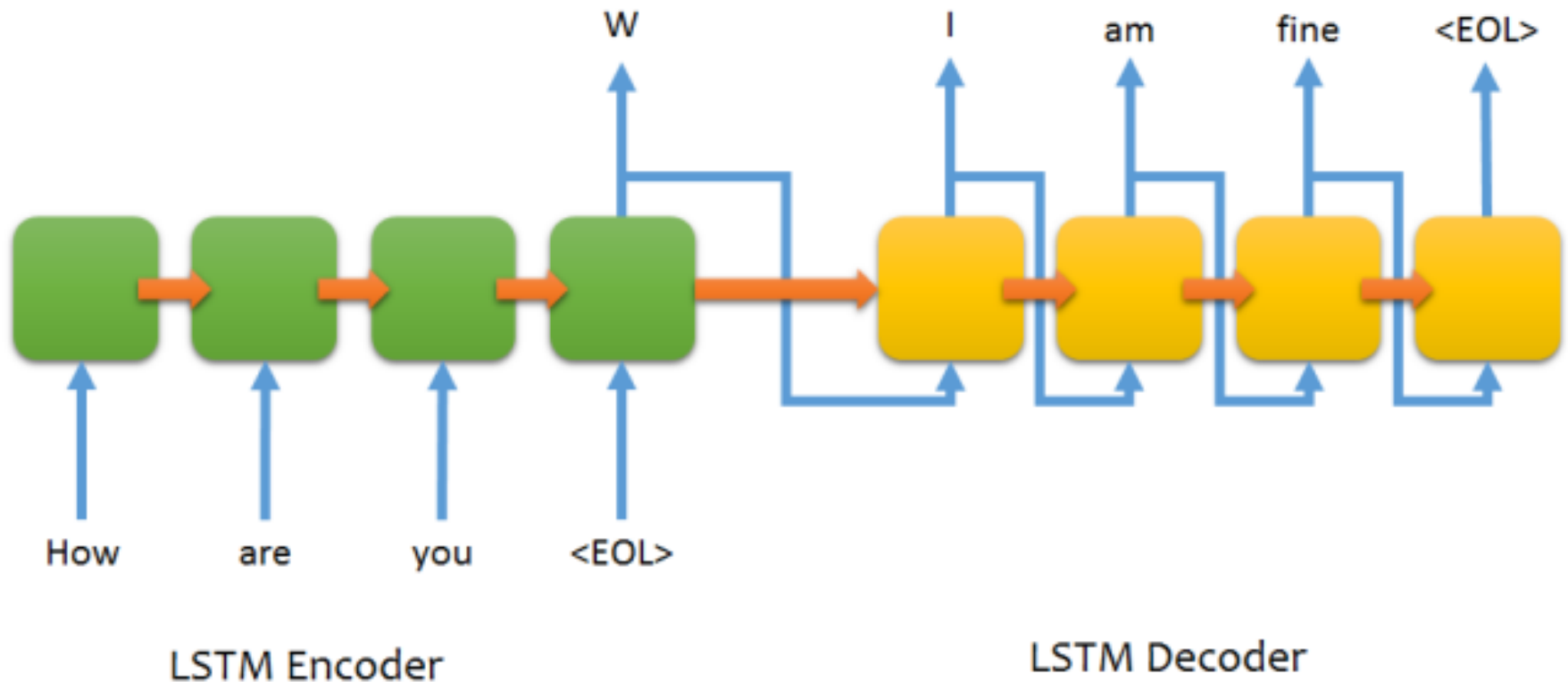
- Input an image, but output a sequence of words

[Kelvin Xu, arXiv'15][Li Yao, ICCV'15]





# Chat-bot



movie (~40,000 sentences), presidential debate...