# Lecture IV:
# Convolutional Neural Network (CNN)

# Three Steps for Deep Learning

| Step 1: Neural Network | Step 2: goodness of function | Step 3: pick the best function |
|---|---|---|

Deep Learning is so simple ……

Now If you want to find a function

If you have lots of function input/output (?) as training data
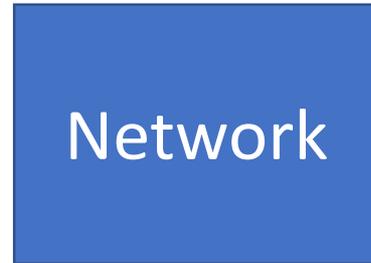
⟶ You can use deep learning

# For example, you can do …….

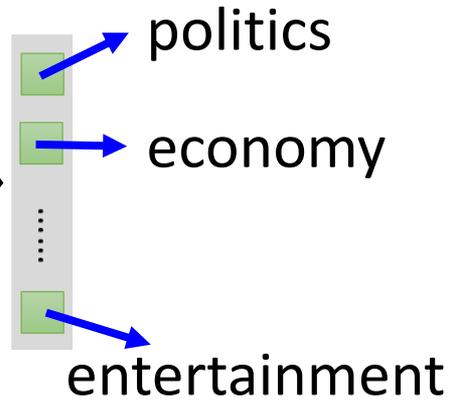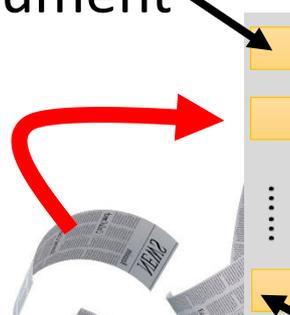**_Spam filtering_**

"Talk" in e-mail

"free" in e-mail

Network

1/0
(Yes/No)

1 (Yes)

0 (No)

(http://spam-filter-review.toptenreviews.com/)

# For example, you can do …….



"stock" in document

Machine

politics

economy

entertainment

"president" in document

http://top-breaking-news.com/

entertainment    politics    economy

# For example, you can do .......

- Image Recognition



"monkey"

"cat"

"dog"

"monkey"

"cat"

"dog"

Network

Convolutional Neural Network (CNN) is usually considered.

# Why CNN for Image?

- When processing image, the first layer of fully connected network would be very large



100

100

100 x 100 x 3

$3 \times 10^7$

1000

Softmax

Can the fully connected network be simplified by considering the properties of image processing?

# Why CNN for Image

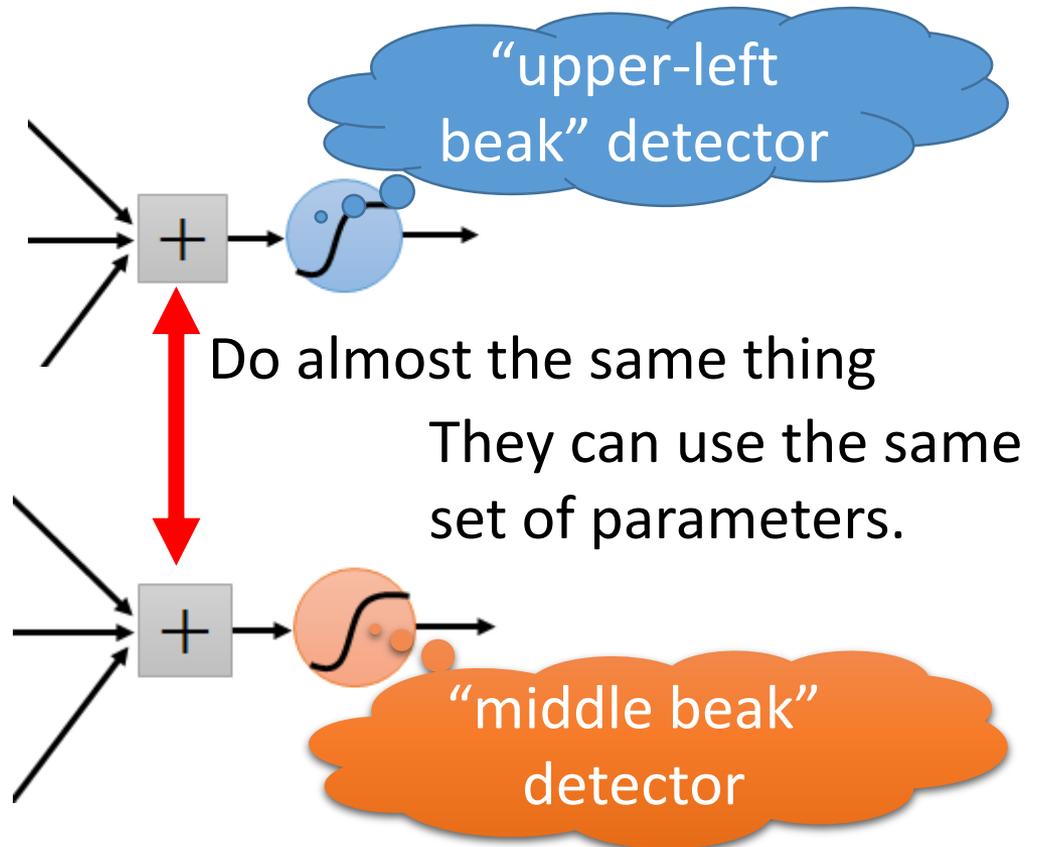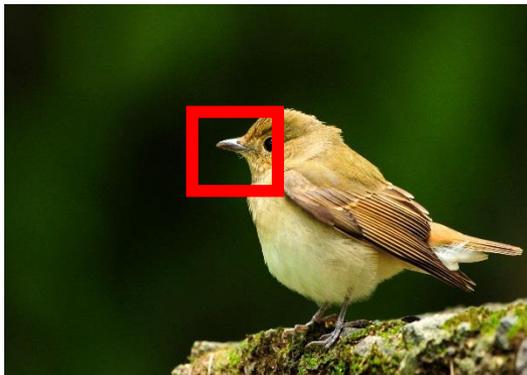- Some patterns are much smaller than the whole image

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

Connecting to small region with less parameters



"beak" detector

# Why CNN for Image

- The same patterns appear in different regions.



"upper-left beak" detector

Do almost the same thing

They can use the same set of parameters.

"middle beak" detector

# Why CNN for Image

- Subsampling the pixels will not change the object

bird



subsampling

bird
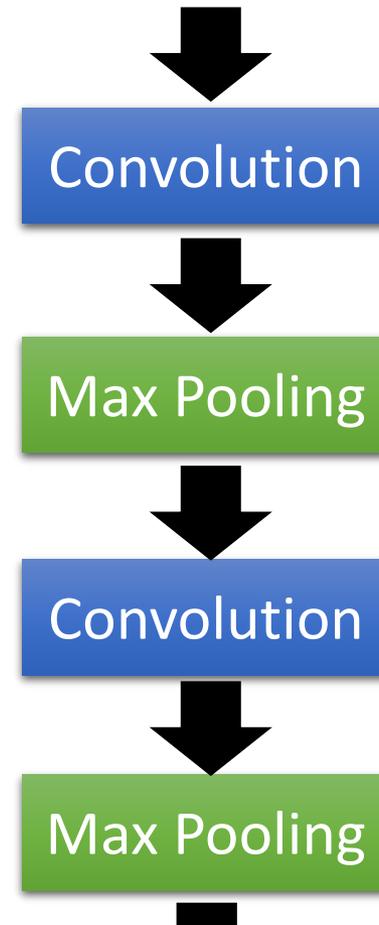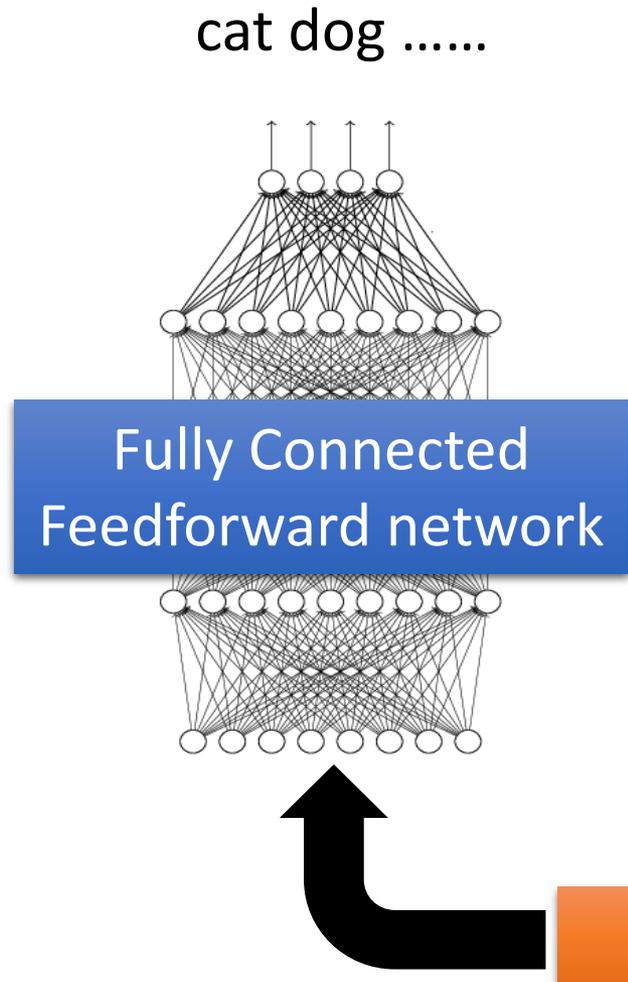
We can subsample the pixels to make image smaller

 Less parameters for the network to process the image

# 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
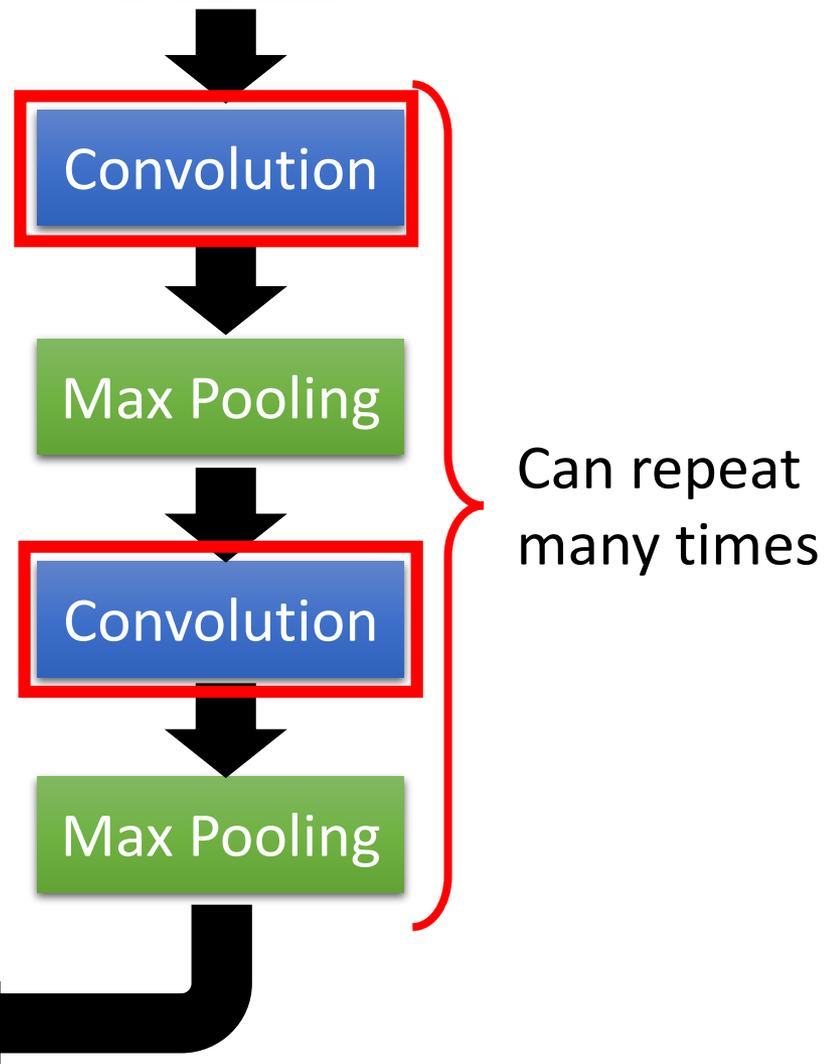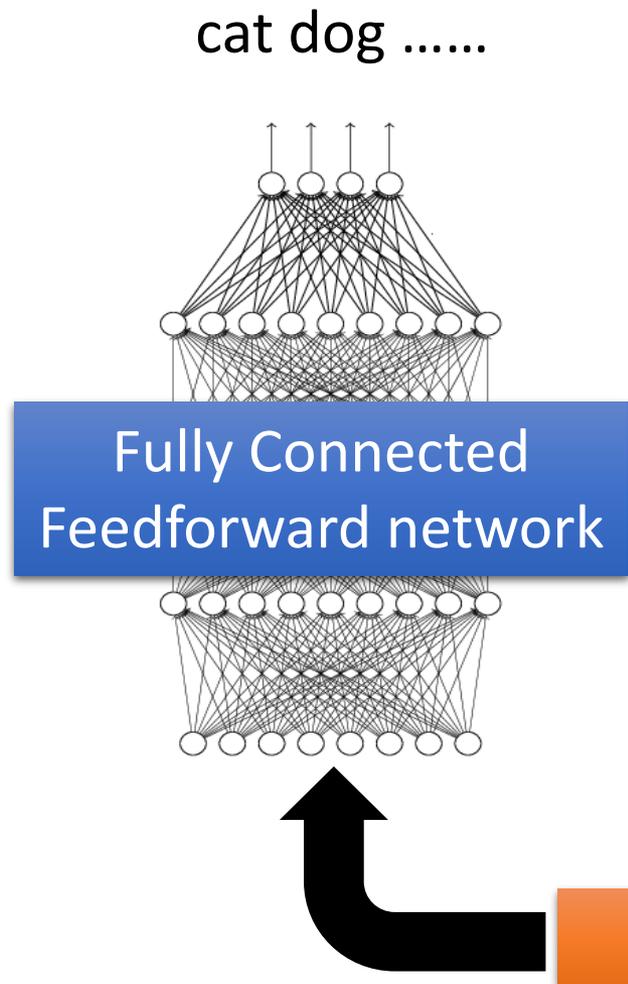
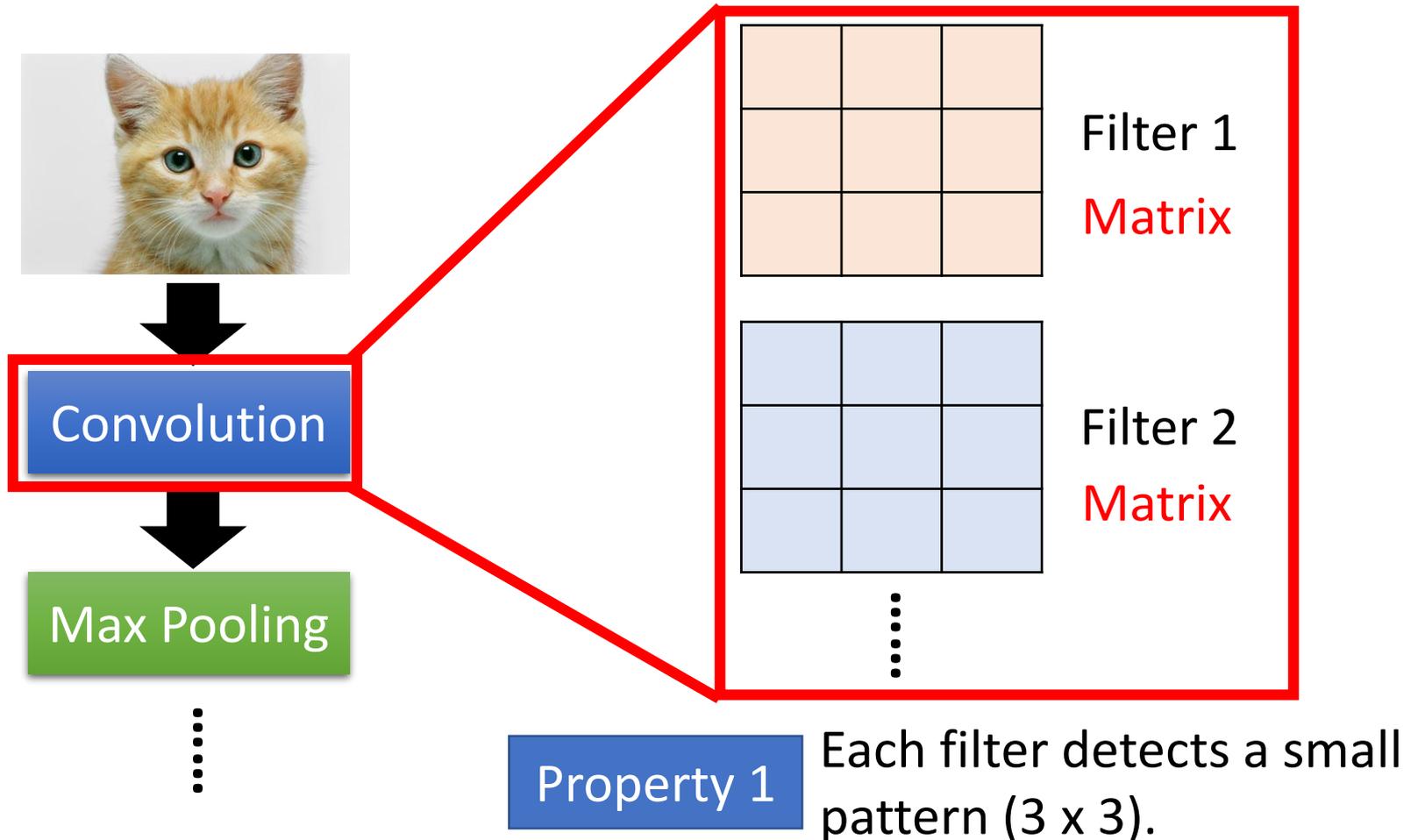Max Pooling

Convolution

Max Pooling

Can repeat many times

Flatten

# The whole CNN



cat dog ……

Fully Connected Feedforward network

Flatten

Convolution

Max Pooling

Convolution

Max Pooling

Can repeat many times

# CNN – Convolution

**The values in the matrices are learned from training data.**



Convolution

Max Pooling

Filter 1

Matrix

Filter 2

Matrix

Property 1

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

# CNN – Convolution

**The values in the matrices are learned from training data.**

| 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

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

Filter 1

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

stride=1

6 x 6 image

| 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 |

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

Property 2

# CNN – Convolution

Filter 2

| -1 | 1 | -1 |
|----|---|----|
| -1 | 1 | -1 |
| -1 | 1 | -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

Do the same process for every filter

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

Feature Map

4 x 4 image

# CNN – Colorful image

Colorful image



| | | | |
|---|---|---|---|
| 1 | -1 | -1 |
| -1 | 1 | -1 |
| -1 | -1 | 1 |

Filter 1

| | | | |
|---|---|---|---|
| -1 | 1 | -1 |
| -1 | 1 | -1 |
| -1 | 1 | -1 |

Filter 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 |

# Convolution v.s. Fully Connected



image

convolution

Fully-connected

Filter 1

6 x 6 image

**Less parameters!**

1: 1
2: 0
3: 0
4: 0
⋮
7: 0
8: 1
9: 0
10: 0
⋮
13: 0
14: 0
15: 1
16: 1
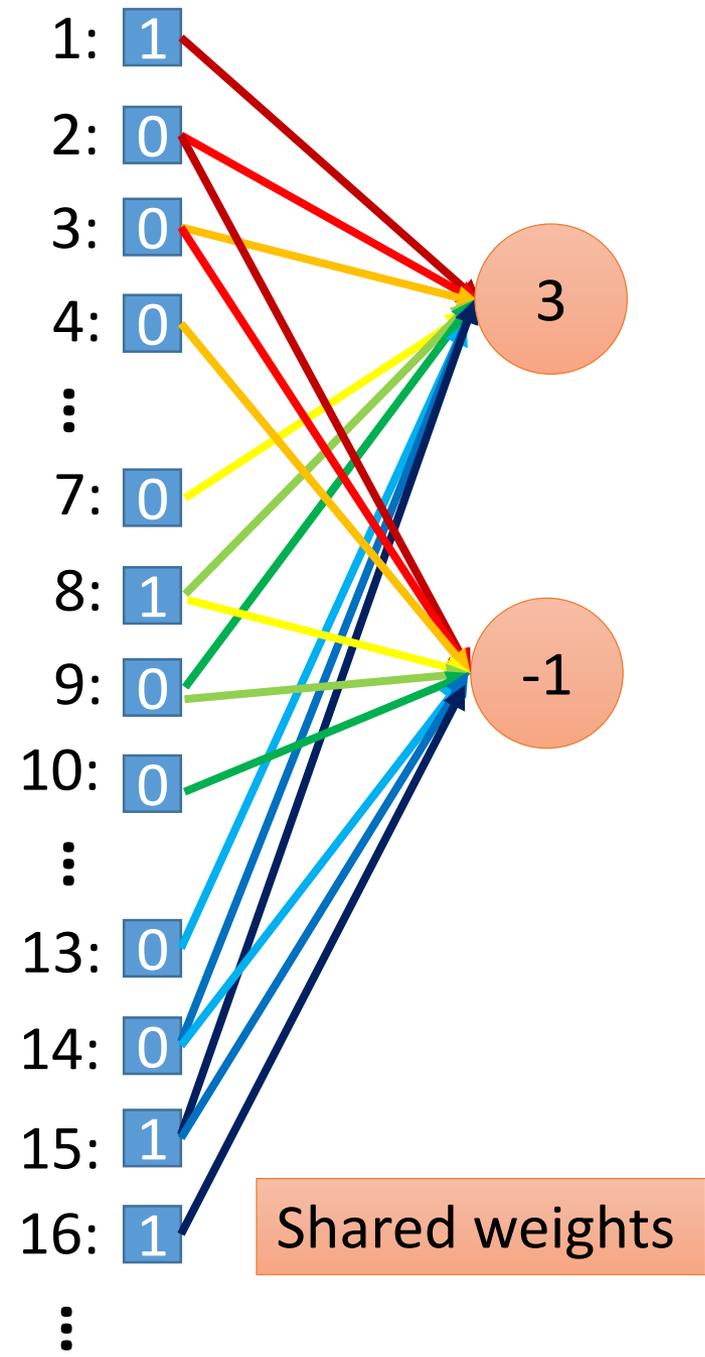⋮

3

Only connect to 9 input, not fully connected

Filter 1

6 x 6 image
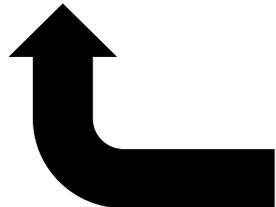
Less parameters!

Even less parameters!

1: 1
2: 0
3: 0
4: 0
⋮
7: 0
8: 1
9: 0
10: 0
⋮
13: 0
14: 0
15: 1
16: 1
⋮

3

-1

Shared weights

# The whole CNN



cat dog ......

Fully Connected
Feedforward network

Convolution

Max Pooling

Convolution

Max Pooling

Can repeat
many times

Flatten

# 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 |
| 3 | -2 |

| | |
|---|---|
| 0 | 1 |
| -2 | -1 |

| | |
|---|---|
| -1 | -1 |
| -1 | -1 |

| | |
|---|---|
| -1 | -1 |
| -2 | 1 |

| | |
|---|---|
| -1 | -1 |
| -1 | 0 |

| | |
|---|---|
| -2 | 1 |
| -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

| -1 | 1 |
|----|---|
| 0  | 3 |

2 x 2 image
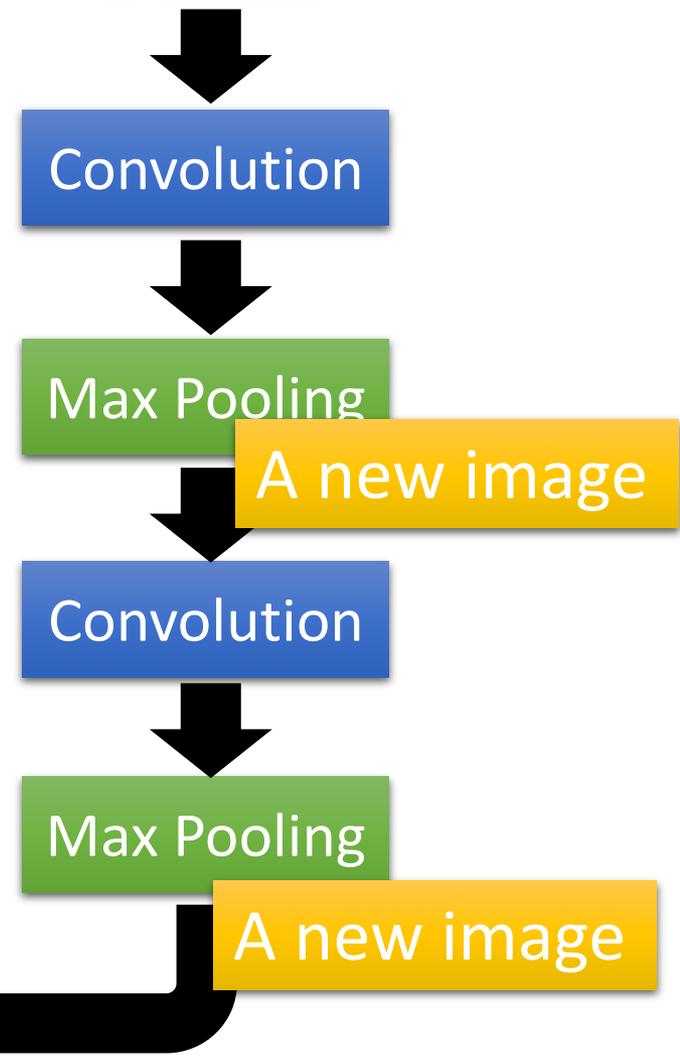
Each filter is a channel

# The whole CNN



-1 1
0 3

A new image

Smaller than the original image

The number of the channel is the number of filters

Convolution

Max Pooling

Convolution

Max Pooling

Can repeat many times
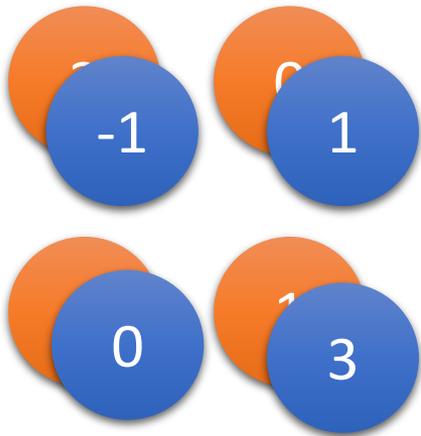
# The whole CNN

# Flatten

Flatten



Fully Connected Feedforward network
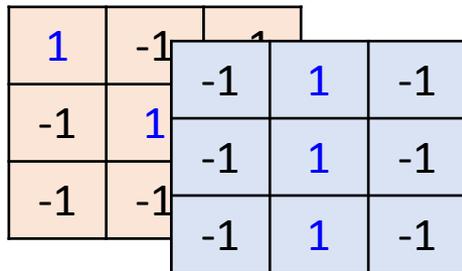
# *CNN in Keras*

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

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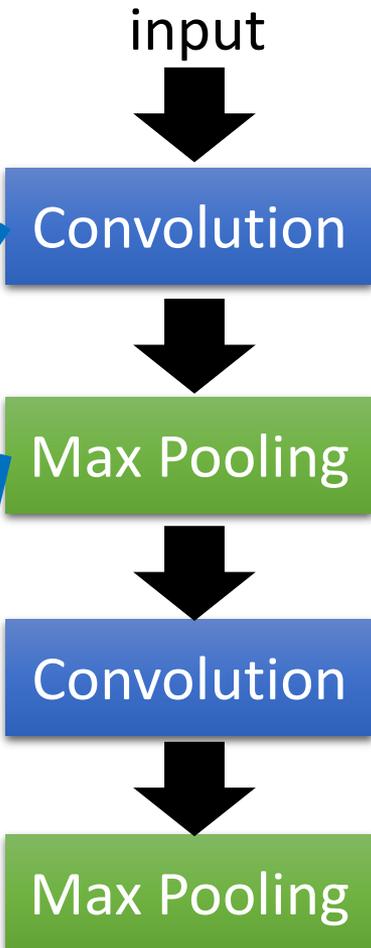


There are 25 3x3 filters.

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)**

input

1 x 28 x 28

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

Convolution

How many parameters for each filter?

**9**    25 x 26 x 26

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

Max Pooling

25 x 13 x 13

```
model2.add(Convolution2D(50,3,3))
```

Convolution

How many parameters for each filter?

**225**    50 x 11 x 11

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

Max Pooling

50 x 5 x 5

# CNN in Keras

input

1 x 28 x 28

Convolution

25 x 26 x 26

Max Pooling

25 x 13 x 13

Convolution

50 x 11 x 11

Max Pooling

50 x 5 x 5

output

Fully Connected Feedforward network

```
model2.add(Dense(output_dim=100))
model2.add(Activation('relu'))
model2.add(Dense(output_dim=10))
model2.add(Activation('softmax'))
```

1250
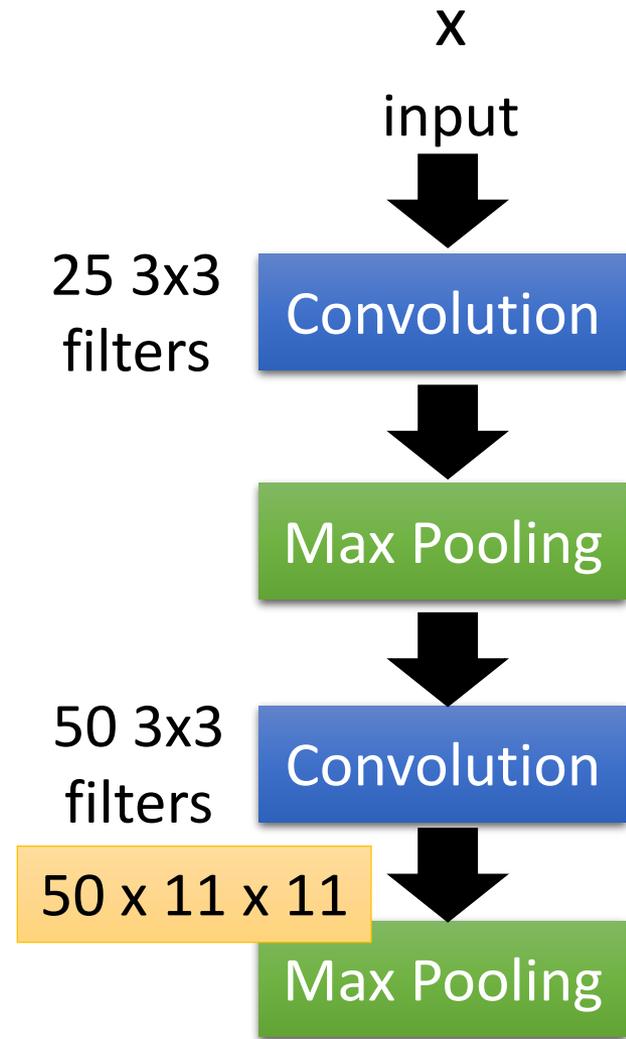
Flatten

```
model2.add(Flatten())
```

# Live Demo

# *What does CNN learn?*

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

Degree of the activation of the k-th filter:

$$a^k = \sum_{i=1}^{11} \sum_{j=1}^{11} a_{ij}^k$$

$x^* = arg \max_x a^k$ (gradient ascent)

11



11

$a_{ij}^k$

| | | | |
|---|---|---|---|
| 3 | -1 | …… | -1 |
| -3 | 1 | …… | -3 |
| 3 | -2 | …… | -1 |

x
input

25 3x3 filters → Convolution

Max Pooling

50 3x3 filters → Convolution
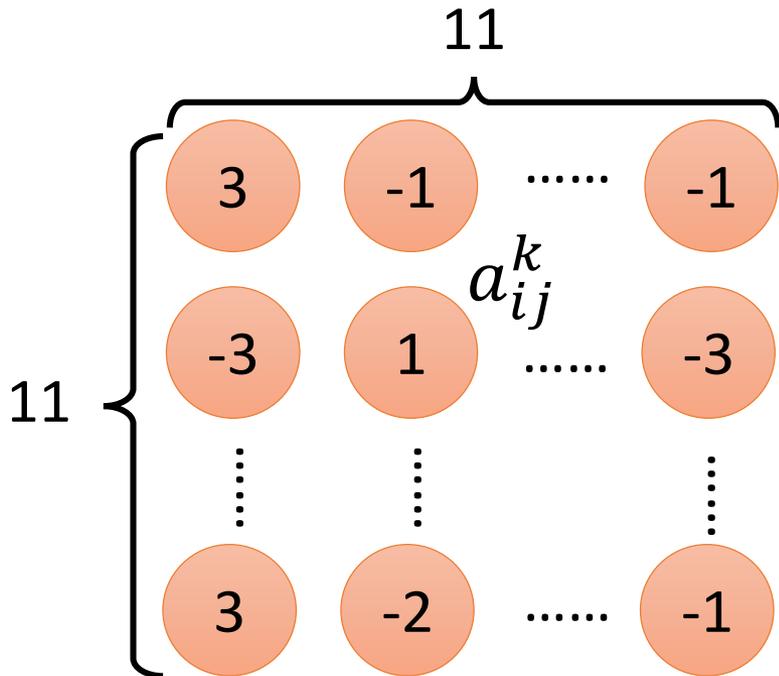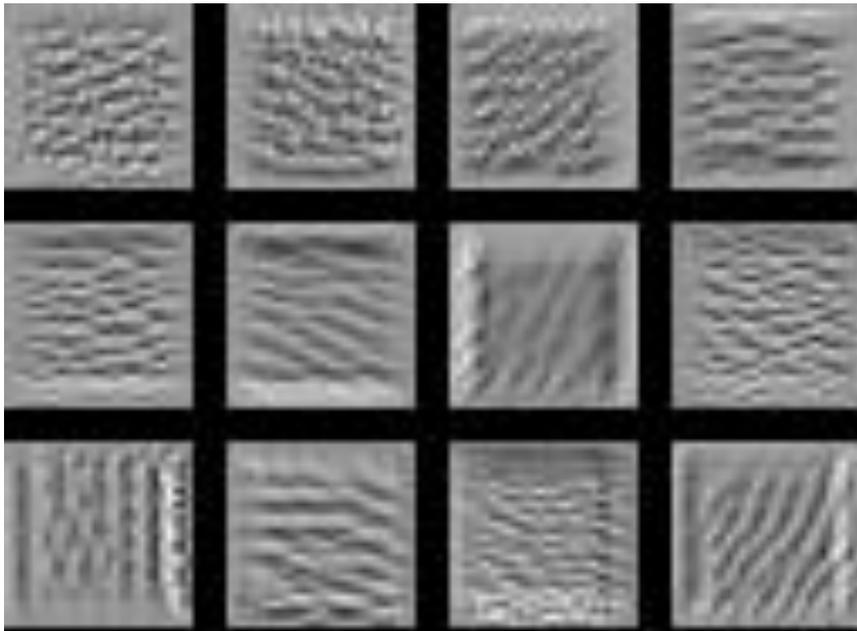
50 x 11 x 11

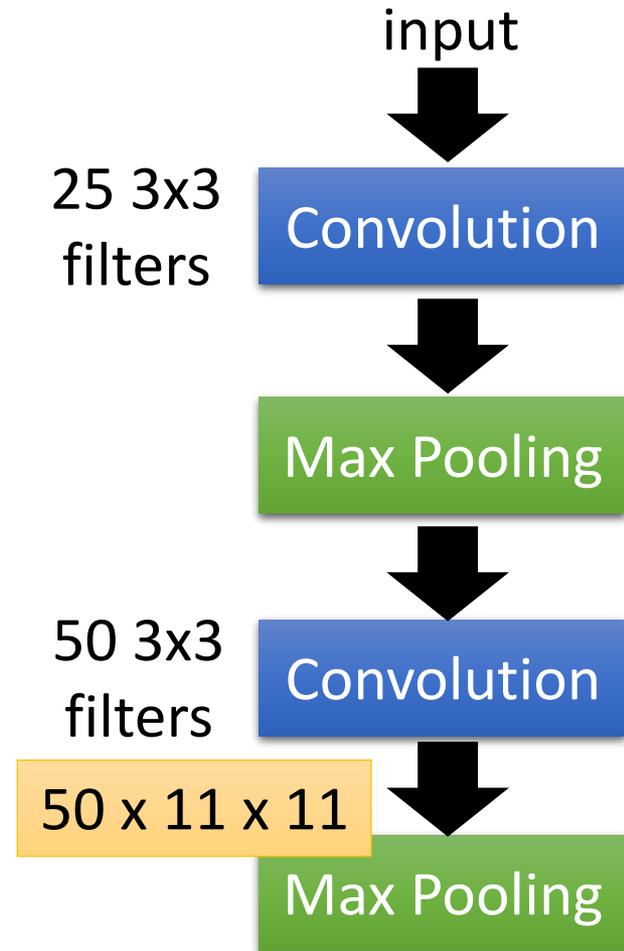Max Pooling

# What does CNN learn?

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

Degree of the activation of the k-th filter:

$$a^k = \sum_{i=1}^{11} \sum_{j=1}^{11} a_{ij}^k$$

$$x^* = arg \max_x a^k \quad \text{(gradient ascent)}$$

input



Each small figure corresponds to a filter.

25 3x3 filters

**Convolution**

**Max Pooling**
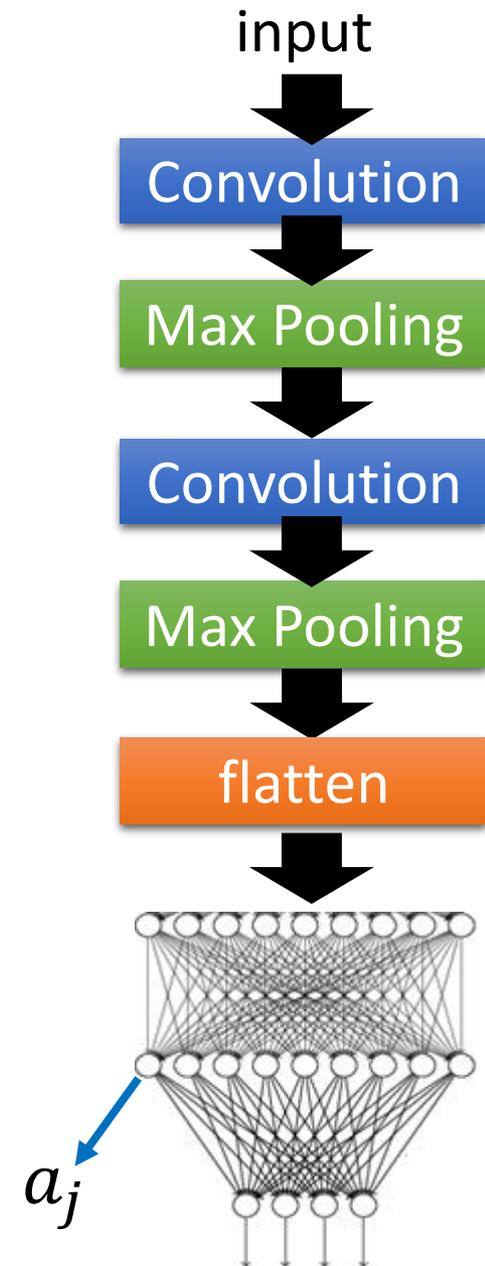
50 3x3 filters

**Convolution**

50 x 11 x 11

**Max Pooling**

# *What does CNN learn?*

Find an image maximizing the output of neuron:

$$x^* = arg \max_x a^j$$



Each figure corresponds to a neuron

input

Convolution

Max Pooling

Convolution

Max Pooling

flatten

$a_j$

# *What does CNN learn?*

$$x^* = \arg \max_x y^i$$

Can we see digits?



input

Convolution

Max Pooling

Convolution

Max Pooling

flatten

$y_i$

Deep Neural Networks are Easily Fooled
https://www.youtube.com/watch?v=M2IebCN9Ht4

Evolving AI Lab

# What does CNN learn?

$$x^* = arg \max_x y^i$$

$$x^* = arg \max_x \left( y^i + \boxed{\sum_{i,j} |x_{ij}|} \right)$$

Over all
pixel values

# Deep Dream

Modify image → CNN

$$\begin{bmatrix} 3.9 \\ -1.5 \\ 2.3 \\ \vdots \end{bmatrix}$$

- Given a photo, machine adds what it sees ……

CNN exaggerates what it sees

http://deepdreamgenerator.com/

# Deep Dream

- Given a photo, machine adds what it sees ……



http://deepdreamgenerator.com/

# Deep Style

- Given a photo, make its style like famous paintings



https://dreamscopeapp.com/

# Deep Style

- Given a photo, make its style like famous paintings



https://dreamscopeapp.com/

# *Deep Style*



A Neural Algorithm of Artistic Style

https://arxiv.org/abs/1508.06576

content

style

CNN

?

# Application: Playing Go



19 x 19 matrix (image)

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.

# *Training*

## Collecting records of many previous plays



## Machine mimics human player

# Why CNN for Go playing?

- 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 Go playing?

- 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 convolve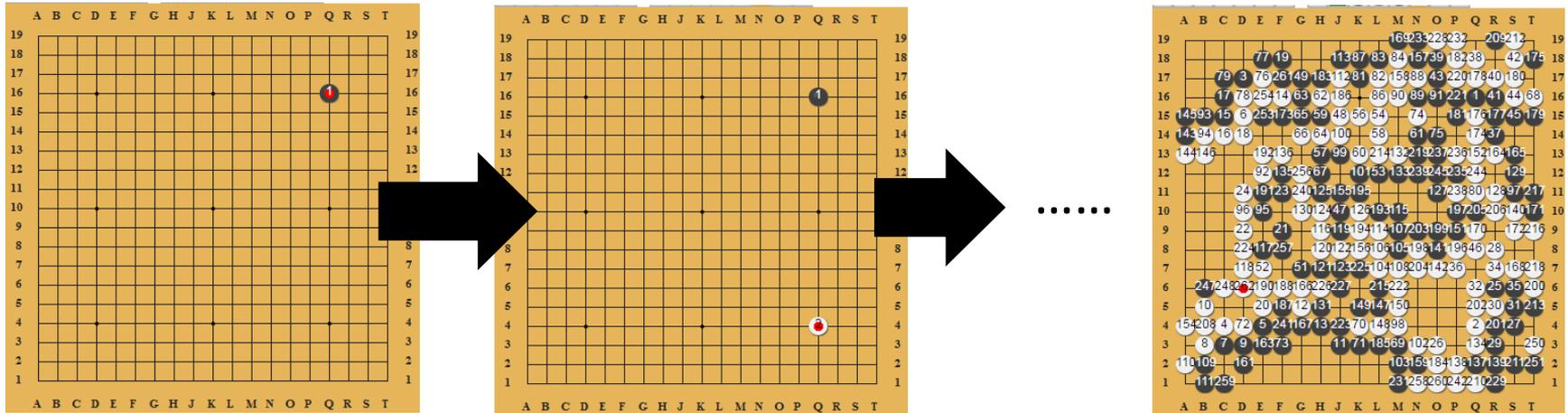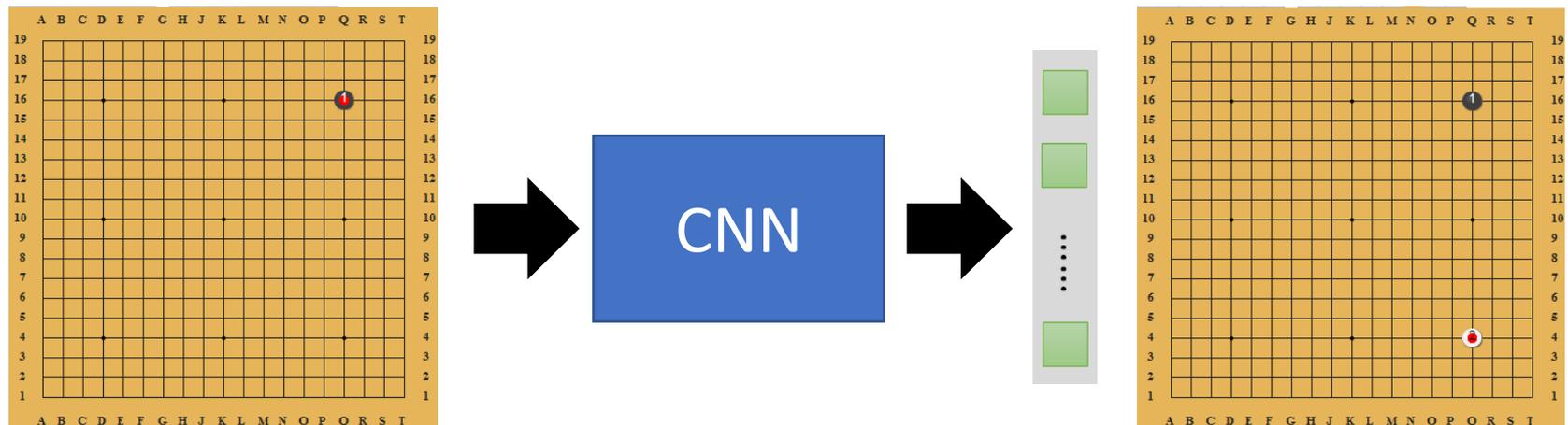s $k$ filters of kernel size $3 \times 3$ with stride 1, again followed by a rectifier nonlinearity. The final layer convolves 1 filter of kernel size $1 \times 1$ with stride 1, with a different bias for each position, and applies a softmax func-
tion. The **Alpha Go does not use Max Pooling ......** Extended Data Table 3 additionally show the results of training with $k = 128$, 256 and 384 filters.

# More Application: Speech



CNN

The filters move in the frequency direction.

Frequency

Image

Time

**Spectrogram**

# More Application: Text



Source of image:
http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.703.6858&rep=rep1&type=pdf

# Lecture V:
# Recurrent Neural Network (RNN)

Neural Network with Memory

# Example Application

- Slot Filling

I would like to arrive Taipei on November 2nd.

ticket booking system

Slot
Destination: Taipei

time of arrival: November 2nd

# 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.

Each dimension corresponds to a word in the lexicon

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

apple = [ 1   0   0   0   0]

bag   = [ 0   1   0   0   0]

cat   = [ 0   0   1   0   0]

dog   = [ 0   0   0   1   0]

elephant   = [ 0   0   0   0   1]

# Beyond 1-of-N encoding

## *Dimension for "Other"*

| | |
|---|---|
| apple | 0 |
| bag | 0 |
| cat | 0 |
| dog | 0 |
| elephant | 0 |
| ⋮ | |
| "other" | 1 |

w = "Gandalf"    w = "Sauron"

## *Word hashing*

| | |
|---|---|
| a-a-a | 0 |
| a-a-b | 0 |
| ⋮ | ⋮ |
| a-p-p | 1 |
| ⋮ | ⋮ |
| p-l-e | 1 |
| ⋮ | ⋮ |
| p-p-l | 1 |
| ⋮ | ⋮ |

**26 X 26 X 26**

w = "apple"

# Example Application

Solving slot filling by Feedforward network?

Input: a word

(Each word is represented as a vector)
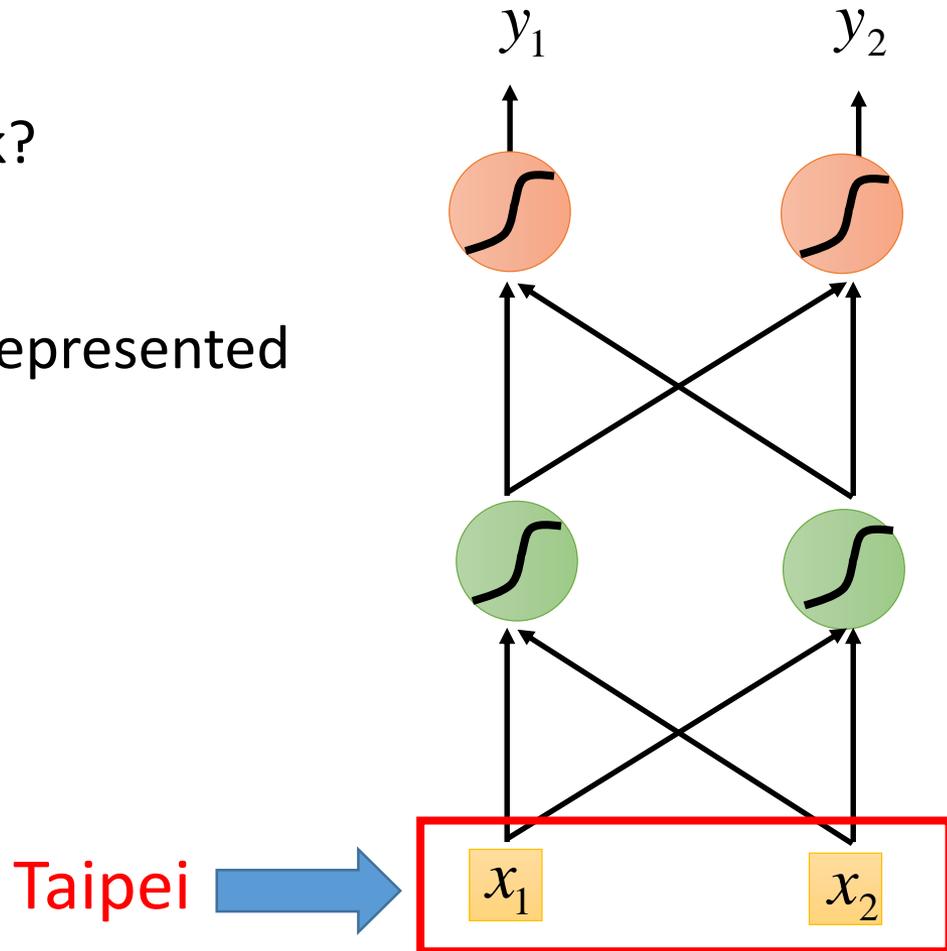
Output:

Probability distribution that the input word belonging to the slots

dest

time of departure

$y_1$

$y_2$

Taipei

$x_1$

$x_2$

# Example Application

# Three Steps for Deep Learning

| Step 1: Recurrent Neural Network | → | Step 2: goodness of function | → | Step 3: pick the best function |
|---|---|---|---|---|

Deep Learning is so simple ……

# Recurrent Neural Network (RNN)

The output of hidden layer are stored in the memory.

store

Memory can be considered as another input.

$y_1$ $y_2$

$a_1$ $a_2$

$x_1$ $x_2$

# Example

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix} \ldots \ldots$

output sequence: $\begin{bmatrix} 4 \\ 4 \end{bmatrix}$



given Initial values

All the weights are "1", no bias

All activation functions are linear

# Example

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ $\begin{bmatrix} 2 \\ 2 \end{bmatrix}$ ......

output sequence: $\begin{bmatrix} 4 \\ 4 \end{bmatrix}$ $\begin{bmatrix} 12 \\ 12 \end{bmatrix}$



All the weights are "1", no bias

All activation functions are linear

# Example

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix}$ ......

output sequence: $\begin{bmatrix} 4 \\ 4 \end{bmatrix} \begin{bmatrix} 12 \\ 12 \end{bmatrix} \begin{bmatrix} 32 \\ 32 \end{bmatrix}$

**Changing the sequence order will change the output.**

$y_1$ 32

$y_2$ 32

store

6      6      16      16

All the weights are "1", no bias

2      2

All activation functions are linear

# RNN

The same network is used again and again.

# RNN



Different

Prob of "leave" in each slot

Prob of "Taipei" in each slot

Prob of "arrive" in each slot

Prob of "Taipei" in each slot

store

store

leave

Taipei

arrive

Taipei

The values stored in the memory is different.

# Of course it can be deep ...

# Bidirectional RNN

# Long Short-term Memory (LSTM)

$$a = h(c')f(z_o)$$

$z_o$ → Output Gate

$f(z_o)$

multiply

$h(c')$

Forget Gate

$c$  $f(z_f)$

$c'$

$cf(z_f)$

$z_f$

$z_i$ → Input Gate

$f(z_i)$  $g(z)f(z_i)$

multiply

$g(z)$

Block

$z$

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

*Multiple-layer LSTM*

I will not implement this!

This is quite standard now ...

https://img.komicolle.org/2015-09-20/src/14426967627131.gif

# Three Steps for Deep Learning

Step 1: define a set of function → Step 2: goodness of function → Step 3: pick the best function

Deep Learning is so simple ......

# *Learning Target*

# Three Steps for Deep Learning

| Step 1: define a set of function | → | Step 2: goodness of function | → | Step 3: pick the best function |
|---|---|---|---|---|

Deep Learning is so simple ……

# Learning



Backpropagation through time (BPTT)

copy

$a_1$

$a_2$

$w$

$w \leftarrow w - \eta \partial L / \partial w$

$y_1$

$y_2$

$x_1$

$x_2$

RNN Learning is difficult in practice.

# Unfortunately ……

- RNN-based network is not always easy to learn

Real experiments on Language modeling

# The error surface is rough.



The error surface is either very flat or very steep.

Clipping

Total Loss

0.35
0.30
0.25
0.20
0.15
0.10
0.05

$w_2$

4.6
4.8
5.0
5.2
5.4

−2.8  −2.6  −2.4  −2.2  −2.0

$w_1$

[Razvan Pascanu, ICML'13]

# Why?

| | |
|---|---|
| $w = 0.01 \Rightarrow y^{1000} \approx 0$ | small $\partial L / \partial w$ $\Rightarrow$ Large Learning rate? |
| $w = 0.99 \Rightarrow y^{1000} \approx 0$ | |
| $w = 1 \Rightarrow y^{1000} = 1$ | Large $\partial L / \partial w$ $\Rightarrow$ Small Learning rate? |
| $w = 1.01 \Rightarrow y^{1000} \approx 20000$ | |

$= w^{999}$

**_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 opened.)

Gated Recurrent Unit (GRU): simpler than LSTM



Output Gate

Forget Gate

Cell

add

Input Gate

Block

[Cho, EMNLP'14]

# Helpful Techniques

Clockwise RNN

Structurally Constrained
Recurrent Network (SCRN)



[Jan Koutnik, JMLR'14]

[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 "Taipei" in each slot

Probability of "on" in each slot



$y^1$  $y^2$  $y^3$

$a^3$

$a^2$

Input and output are both sequences with the same length

RNN can do more than that!

$x^1$  $x^2$  $x^3$

arrive  Taipei  on  November  2nd

# Many to one

Keras Example:
https://github.com/fchollet/keras/blob/master/examples/imdb_lstm.py

- Input is a **_vector sequence_**, but output is only **_one vector_**

**_Sentiment Analysis_**

Positive

Negative

Very Positive

Very Positive
Positive
Neutral
Negative
Very Negative

I think this …… I like it

# Many to Many

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

machine

learning

Containing all information about input sequence

# 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→機器學習)

# 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→機器學習)



Add a symbol "**<END>**"

[Ilya Sutskever, NIPS'14][Dzmitry Bahdanau, arXiv'15]

# Many to Many: Title Generation

[Alexander M Rush, EMNLP 15][Chopra, NAACL 16][Lopyrev, arXiv 2015][Shen, arXiv 2016][Yu & Lee,SLT 2016]

Training Data

Title 1

Title 2

Title 3

Title

Input: a document (word sequence), output: its title (shorter word sequence)

https://arxiv.org/pdf/1512.01712v1.pdf

# Many to Many:
# Video Caption Generation



Video

A girl is running.

A group of people is knocked by a tree.

A group of people is walking in the forest.

# Many to Many:
# Video Caption Generation

- Can machine describe what it see from video?

# One to Many: Image Caption Generation

- Input an image, but output a sequence of words

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



A vector for whole image

CNN

Input image

a    woman    is    <END>

***Caption Generation***

# One to Many:
# Image Caption Generation

- Can machine describe what it see from image?

Project of MTK

# Attention-based Model



What you learned in these lectures

Breakfast today

What is deep learning?

summer vacation 10 years ago

Answer ← Organize

http://henrylo1605.blogspot.tw/2015/05/blog-post_56.html

# Attention-based Model



Input → DNN/RNN → output

Reading Head Controller

Reading Head

Machine's Memory

# Attention-based Model v2



Input → DNN/RNN → output

Reading Head Controller

Writing Head Controller

Writing Head

Reading Head

Machine's Memory

Neural Turing Machine

# Reading Comprehension



Query → DNN/RNN → answer

Reading Head Controller

Semantic Analysis

Each sentence becomes a vector.

# Reading Comprehension

- End-To-End Memory Networks. S. Sukhbaatar, A. Szlam, J. Weston, R. Fergus. NIPS, 2015.

The position of reading head:

| Story (16: basic induction) | Support | Hop 1 | Hop 2 | Hop 3 |
|---|---|---|---|---|
| Brian is a frog. | yes | 0.00 | 0.98 | 0.00 |
| Lily is gray. | | 0.07 | 0.00 | 0.00 |
| Brian is yellow. | yes | 0.07 | 0.00 | 1.00 |
| Julius is green. | | 0.06 | 0.00 | 0.00 |
| Greg is a frog. | yes | 0.76 | 0.02 | 0.00 |
| **What color is Greg?  Answer: yellow** | **Prediction: yellow** | | | |

Keras has example:

https://github.com/fchollet/keras/blob/master/examples/babi_memnn.py

# Visual Question Answering

# Visual Question Answering

# Visual Question Answering

- Huijuan Xu, Kate Saenko. Ask, Attend and Answer: Exploring Question-Guided Spatial Attention for Visual Question Answering. arXiv Pre-Print, 2015



**Is there a red square on the bottom of the cat?**
GT: yes                          Prediction: yes

# Speech Question Answering

- **TOEFL Listening Comprehension Test by Machine**
- Example:

Audio Story: (The original story is 5 min long.)

Question: " What is a possible origin of Venus' clouds? "

Choices:

(A) gases released as a result of volcanic activity

(B) chemical reactions caused by high surface temperatures

(C) bursts of radio energy from the plane's surface

(D) strong winds that blow dust into the atmosphere

# Model Architecture



Everything is learned from training examples

Answer

Attention

Select the choice most similar to the answer

Question Semantics

Attention

...... It be quite possible that this be due to volcanic eruption because volcanic eruption often emit gas. If that be the case volcanism could very well be the root cause of Venus 's thick cloud cover. And also we have observe burst of radio energy from the planet 's surface. These burst be similar to what we see when volcano erupt on earth ......

Semantic Analysis

Question: "what is a possible origin of Venus' clouds?"

Speech Recognition

Semantic Analysis

Audio Story:

# Model Architecture - Attention Mechanism

## Understand the question

**Vs : consider both Question and Story with attention weight α**

Concatenate the output of last hidden layer in bi-directional GRU

$$\alpha = \frac{V_Q \cdot S_t}{|V_Q| \cdot |S_t|}$$

(similarity score)

$$\Sigma \quad V_s = \sum_{t=1}^{8} \propto_t * S_t$$

**Module for Vector Representation**

$\alpha_1$  $\alpha_2$  $\alpha_3$  $\alpha_4$  $\alpha_5$  $\alpha_6$  $\alpha_7$  $\alpha_8$

$V_Q$

Concatenate the output of hidden layer at each time step

A bi-directional GRU

$y_f(1)$  $y_f(2)$  $\cdots$  $y_f(T)$

$y_b(1)$ ← $y_b(2)$  $\cdots$ ← $y_b(T)$

$W_1$  $W_2$  $\cdots$  $W_T$

**Question**

$V_Q$ : vector representation for question

$S_1$  $S_2$  $S_3$  $S_4$  $S_5$  $S_6$  $S_7$  $S_8$

$w_1$  $w_2$  $w_3$  $w_4$  $w_5$  $w_6$  $w_7$  $w_8$

Sentence 1      Sentence 2

**Story (through ASR)**

# Model Architecture

**Attention Mechanism...** **Compare similarity between choices and $V_{QH}$**

Take the choice with the highest score as answer

Process Question by VecRep module

Get Vs through attention module

To keep question info, add $V_Q$ and $V_s$

A hop means the machine considers question and story jointly once

Do n... for co... story again

**Vs : final vector considering both Question and Story**
word-level: consider every word information

$V_s = \sum_{t=1}^{8} \alpha_t * S_t$

(similarity score)

**Module for Vector Representation**

**VecRep**

$V_Q$

$y_b(1)$  $y_f(T)$

$y_f(1)$  $y_f(2)$  ...  $y_f(T)$

$y_b(1)$  $y_b(2)$  $y_b(T)$

$w_1$  $w_2$  ...  $w_T$

**Question**

$V_Q$ : vector representation for question

word-level Attention

$\alpha_1$ $\alpha_2$ $\alpha_3$ $\alpha_4$ $\alpha_5$ $\alpha_6$ $\alpha_7$ $\alpha_8$

**0.6**  **0.2**  **0.1**  **0.1**

$V_A$  **Att**  $V_B$  $V_C$  $V_D$

VecRep  VecRep  VecRep  VecRep

**Choice A**  **Choice B**  **Choice C**  **Choice D**

$V_{Q_n}$

**hop 1**  **hop 2**  **hop n**

VecRep → $V_{Q_0}$ → (+) → $V_{Q_1}$ → (+) → $V_{Q_2}$  ......  (+)

**Question**

$V_{s_1}$  $V_{s_2}$  ......  $V_{s_n}$

Att  Att  Att
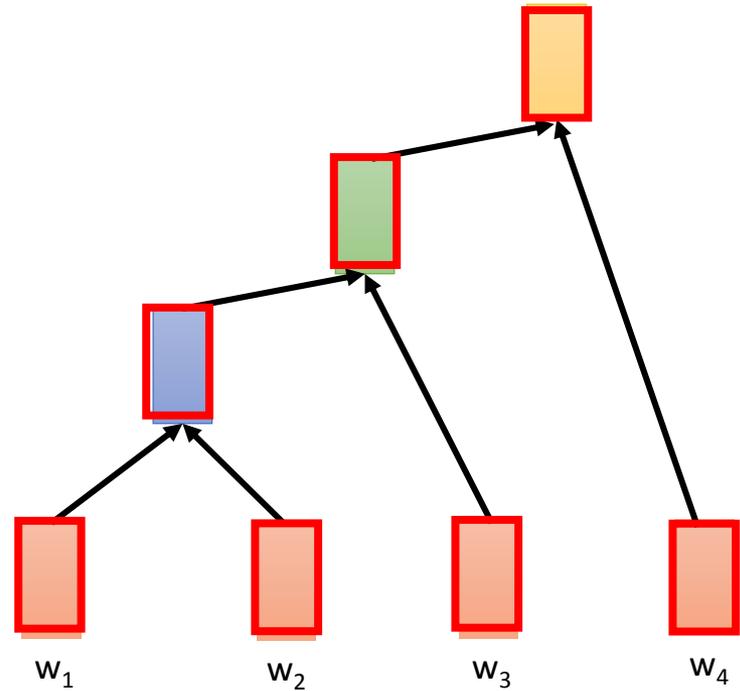
**Story (through ASR)**

$S_4$  $S_5$  $w_4$  $w_5$

**VecRep**

# Sentence Representation

Bi-directional RNN

Tree-structured Neural Network



Attention on all phrases

# Simple Baselines

(2) select the ***shortest*** choice as answer

(4) the choice with semantic most similar to others

random

(1)  (2)  (3)  (4)  (5)
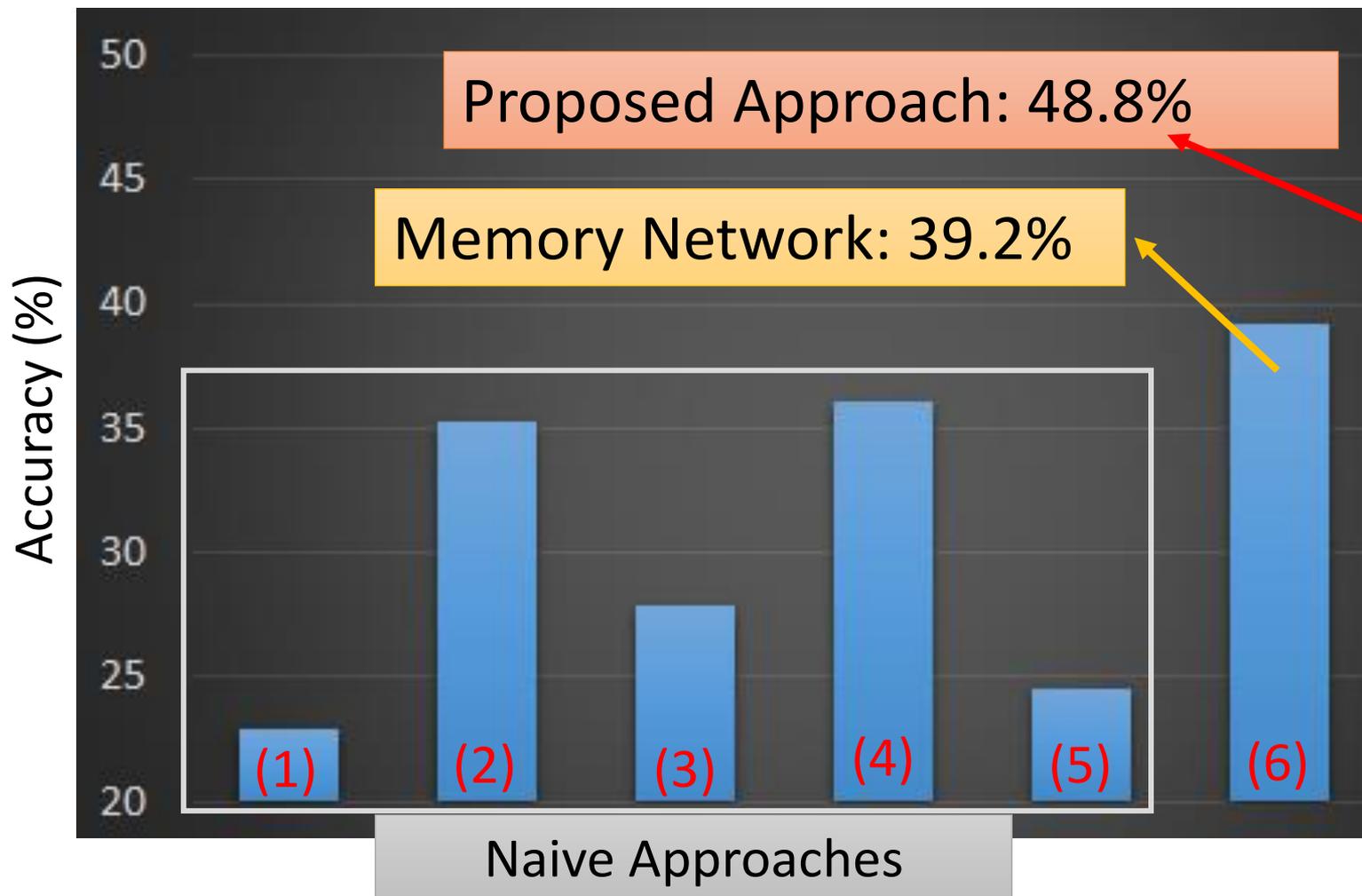
Accuracy (%)

Naive Approaches

# Memory Network

# Proposed Approach

[Tseng & Lee, Interspeech 16]
[Fang & Hsu & Lee, SLT 16]

# Analysis

**Type 1**: Comprehension
**Type 2**: Pragmatic understanding
**Type 3**: Connecting information, making inference and drawing conclusions