# Convolutional Neural Networks

Jiahui Chen
Department of Mathematical Sciences
University of Arkansas
Reference: Ming Li's notes

# Introduction

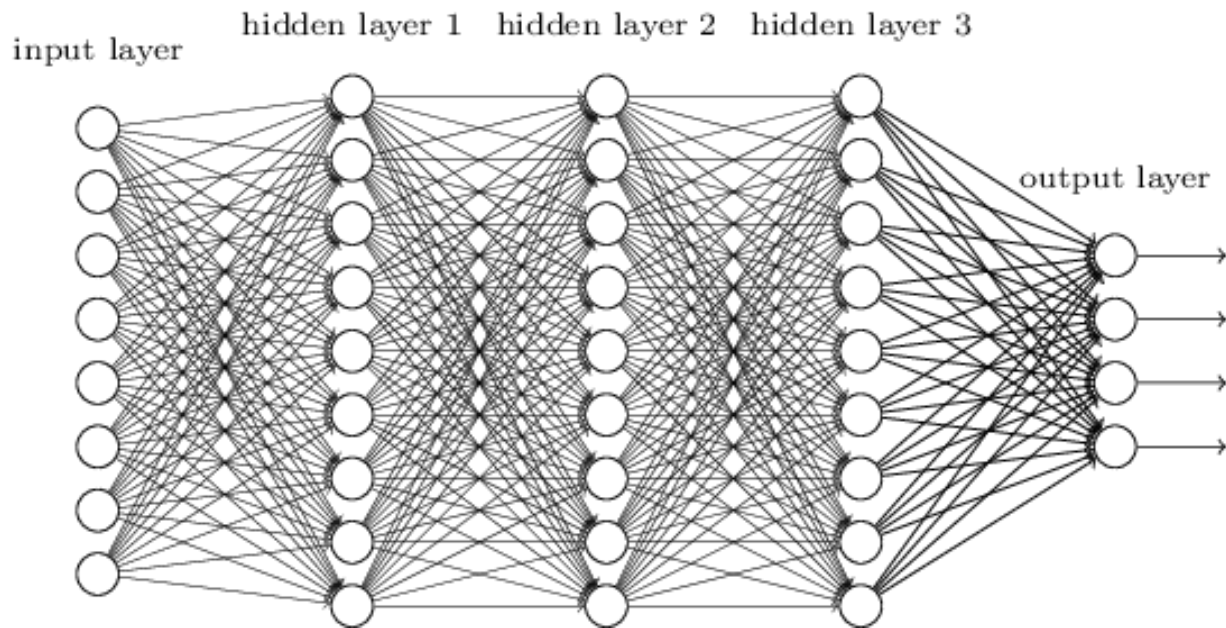- Fukushima (1980) – Neo-Cognitron; LeCun (1998) – Convolutional Neural Networks (CNN, or ConvNet);…

**Motivation:**

Images typically have $1000^2$ pixels, which give rise to $1000^2$ data points or features, leading to an intractably high dimension of the weight space. However, not all of them are essential due to the spatial patterns. Many of them have shared properties. Therefore, the weight space dimension can be dramatically reduced if an appropriate pre-processing of the image data is carried out to analyze or extract spatial correlations or patterns in images.

# Motivation

- We know it is good to learn a small model.
- From this fully connected model, do we really need all the edges?
- Can some of these be shared?

input layer    hidden layer 1    hidden layer 2    hidden layer 3
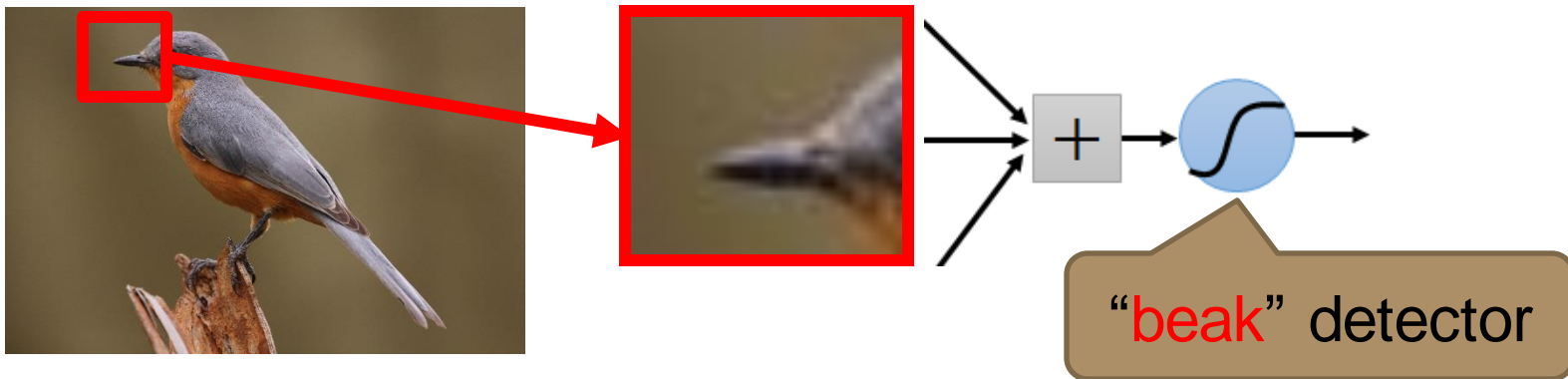
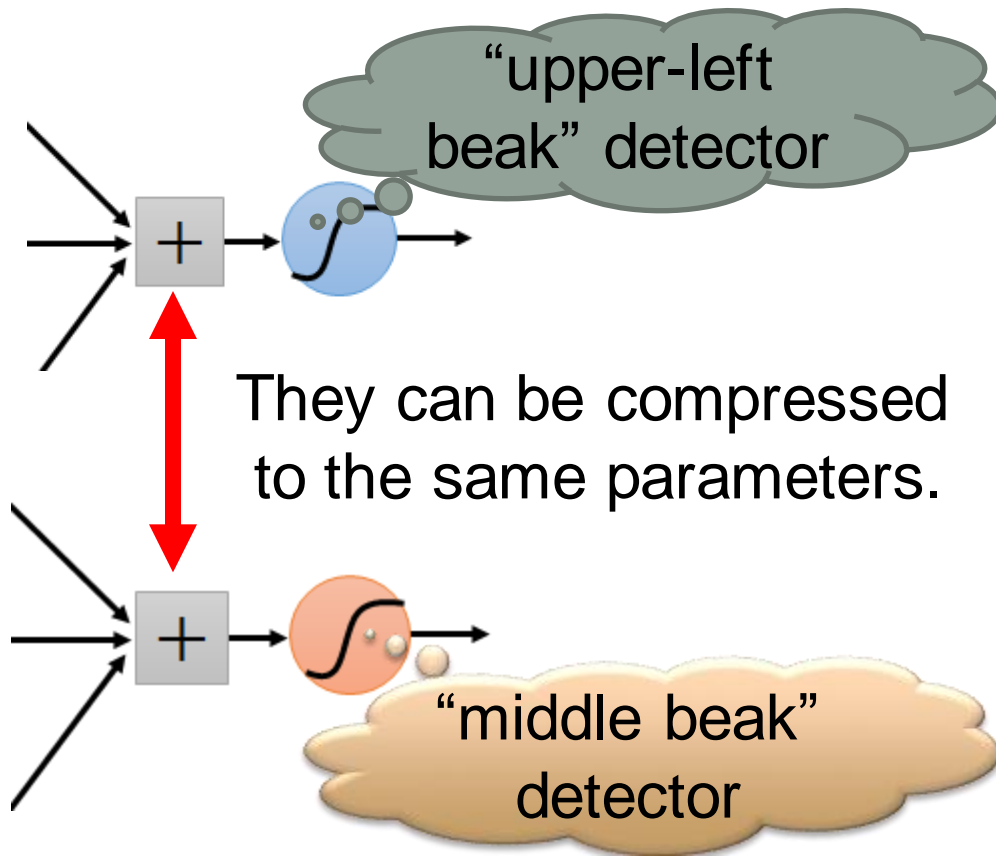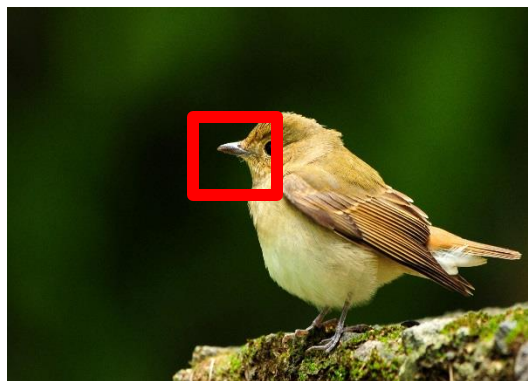output layer

# Motivation

Consider learning an image:

Some patterns are much smaller than the whole image

Can represent a small region with fewer parameters
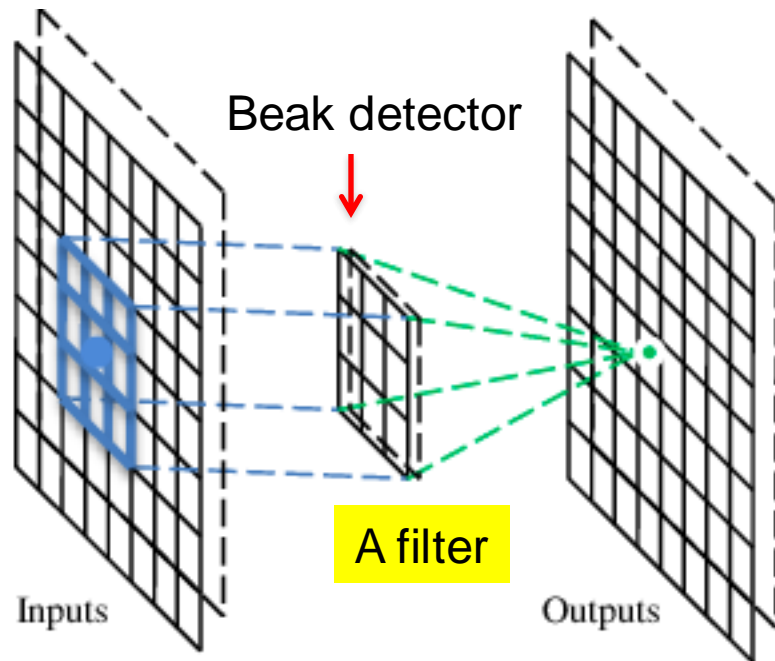


"beak" detector

Same pattern appears in different places:
They can be compressed!
What about training a lot of such "small" detectors and each detector must "move around".



"upper-left beak" detector

They can be compressed to the same parameters.

"middle beak" detector

# 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.

Beak detector

A filter

Inputs

Outputs

# Convolution

These are the network parameters to be learned.

| 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

| -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 |

Dot product

3    -1

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

Filter 1

6 x 6 image

# 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

3    -3

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

Filter 1

stride=1

Filter 1

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

Filter 2

stride=1

Repeat this for each filter

Feature Map

6 x 6 image

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

Filter 1

Filter 2

Color 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 |

image

convolution

Fully-connected

Filter 1

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

6 x 6 image

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

fewer parameters!

1  **1**
2  **0**
3  **0**
4  **0**
⋮
8  **1**
9  **0**
10: **0**
⋮
1  **0**
3  **0**
15  **1**
16  **1**
⋮

3

Only connect to 9 inputs, not fully connected

# The whole 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 |

# Why Pooling

- Subsampling pixels will not change the object
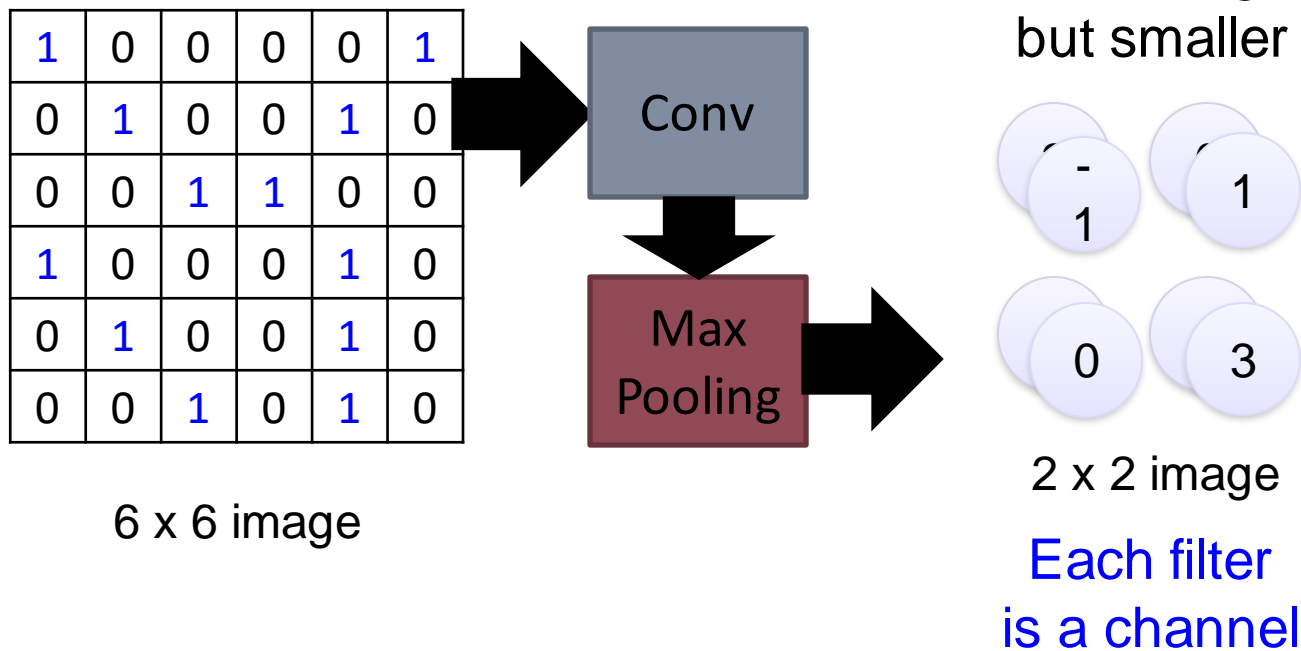
bird



bird

Subsampling

We can subsample the pixels to make image smaller
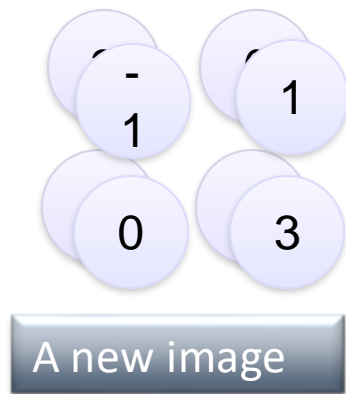
fewer parameters to characterize the image

# A CNN compresses a fully connected network in two ways:

- Reducing number of connections
- Shared weights on the edges
- 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

-1   1

0   3

2 x 2 image

Each filter is a channel

# The whole CNN



Convolution

Max Pooling

A new image

Smaller than the original image

The number of channels is the number of filters

Convolution

Max Pooling

Can repeat many times

# The whole CNN



cat dog ……

**Fully Connected Feedforward network**

Convolution

Max Pooling

A new image

Convolution

Max Pooling

Flattened

A new image

# Flattening



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=(28,28,1)) )
```

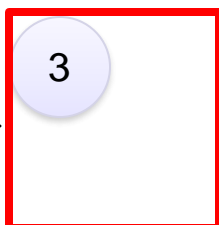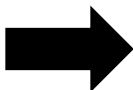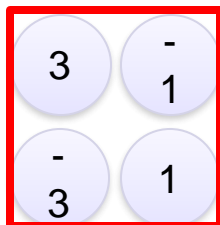| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

... 

... 

There are
**25 3x3**
filters.

Input_shape = ( 28 , 28 , 1)

28 x 28 pixels        1: black/white, 3: RGB
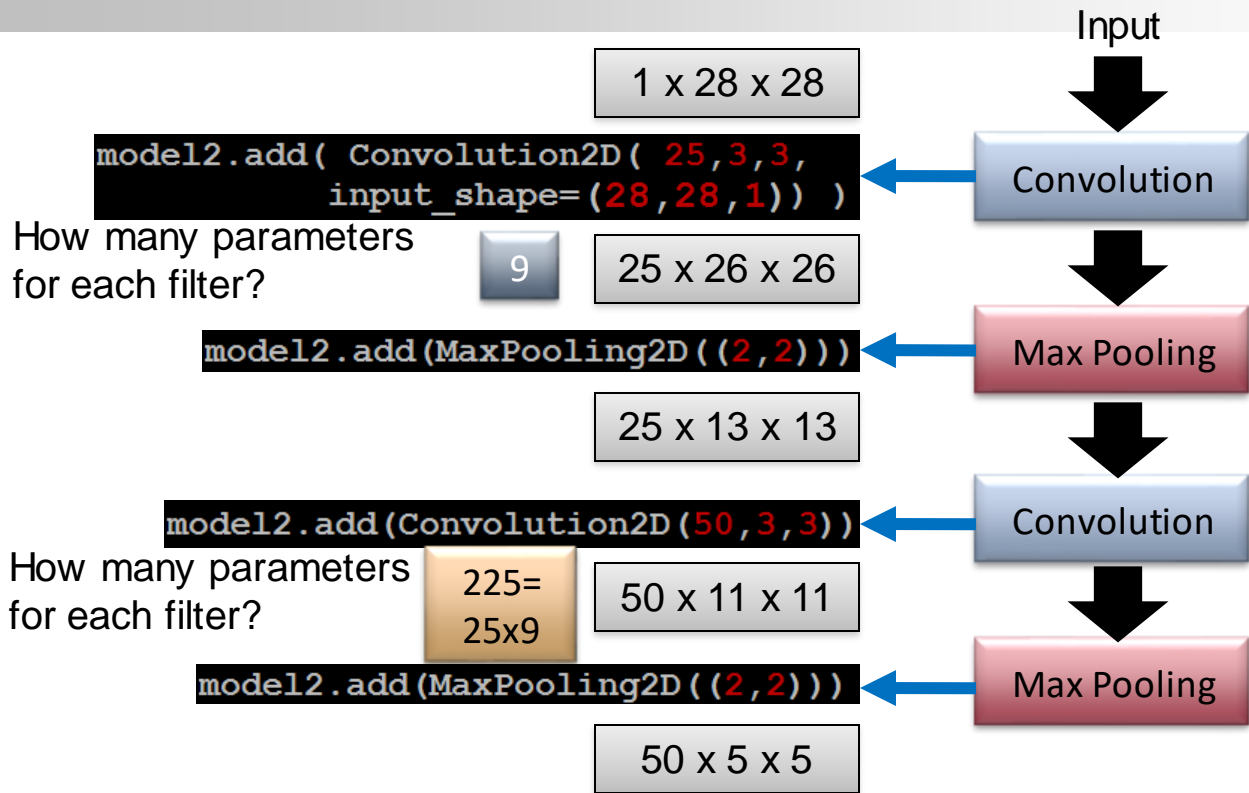
```
model2.add(MaxPooling2D((2,2)))
```

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

→ 3

input

Convolution

Max Pooling

Convolution

Max Pooling

# CNN in Keras

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

Input

1 x 28 x 28

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

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=
25x9

50 x 11 x 11

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

Max Pooling

50 x 5 x 5

# AlphaGo



19 x 19 matrix

Black: 1

white: -1

none: 0

Neural Network

Next move
(19 x 19 positions)

Fully-connected feedforward network can be used
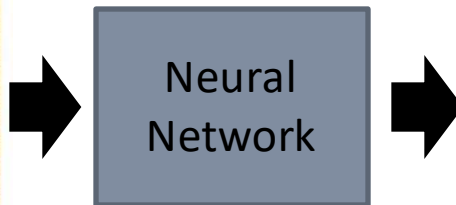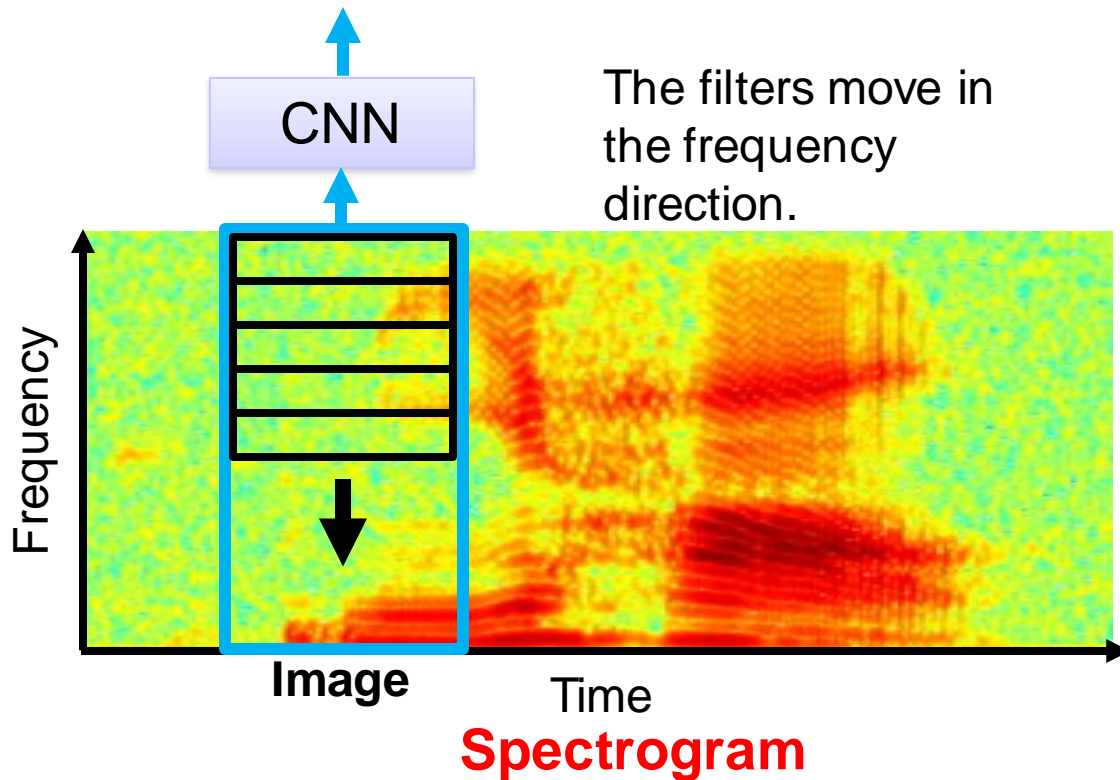
But CNN performs much better

# AlphaGo's policy network
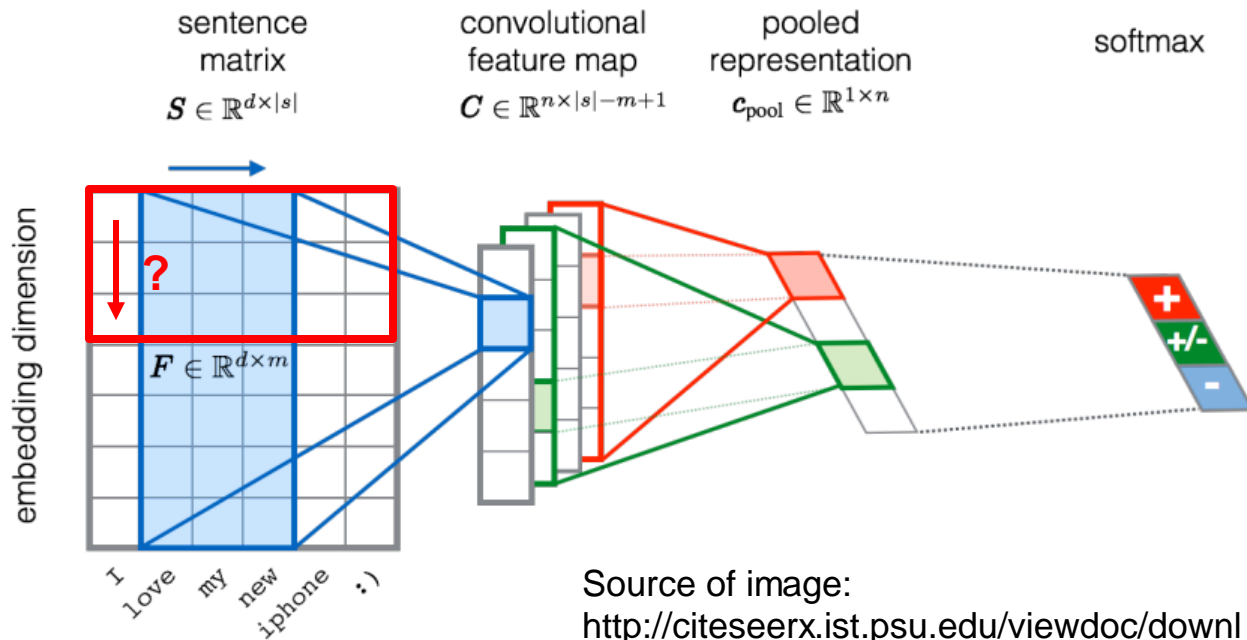
The following is quotation from their Nature article:

> Note: AlphaGo does not use Max Pooling.

**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 match version of AlphaGo used $k = 192$ filters; Fig. 2b and Extended Data Table 3 additionally show the results of training with $k = 128, 256$ and 384 filters.

# CNN in speech recognition



CNN

The filters move in the frequency direction.

Frequency

**Image**

Time

**Spectrogram**

# CNN in text classification



sentence matrix $S \in \mathbb{R}^{d \times |s|}$

convolutional feature map $C \in \mathbb{R}^{n \times |s|-m+1}$

pooled representation $c_{\text{pool}} \in \mathbb{R}^{1 \times n}$

softmax

embedding dimension

$F \in \mathbb{R}^{d \times m}$

I love my new iphone :)

+
+/-
-

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