



# Introduction to Deep Learning

# Acknowledgments



This presentation is heavily based on:

- <http://cs.nyu.edu/~fergus/pmwiki/pmwiki.php>
- <http://deeplearning.net/reading-list/tutorials/>
- <http://deeplearning.net/tutorial/lenet.html>
- [http://ufldl.stanford.edu/wiki/index.php/UFLDL\\_Tutorial](http://ufldl.stanford.edu/wiki/index.php/UFLDL_Tutorial)

... and many other

# Types of Learning

**Supervised:** Learning with a **labeled training** set

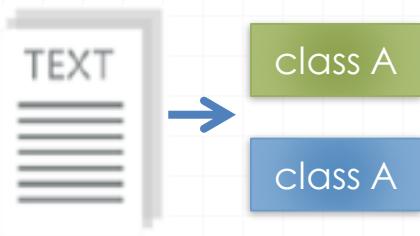
Example: email **classification** with already labeled emails

**Unsupervised:** Discover **patterns** in **unlabeled** data

Example: **cluster** similar documents based on text

**Reinforcement learning:** learn to **act** based on **feedback/reward**

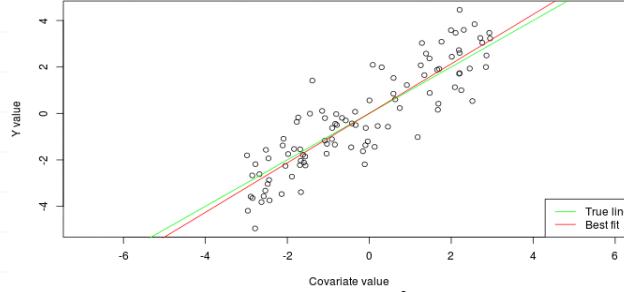
Example: learn to play Go, reward: **win or lose**



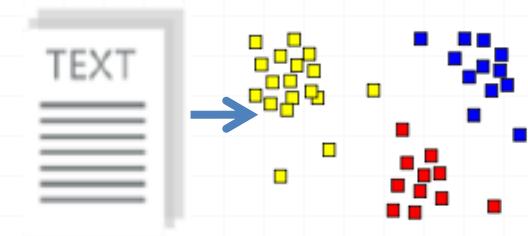
Classification

Anomaly Detection  
Sequence labeling

...



Regression



Clustering

# Buzz...

HOME ▾ MENU ▾ CONNECT THE LATEST POPULAR MOST SHARED

MIT Technology Review

## 10 BREAKTHROUGH TECHNOLOGIES 2013

Introduction The 10 Technologies Past Years

**Deep Learning**

With massive amounts of computational power, machines can now recognize objects and translate speech in real time. Artificial intelligence is finally getting smart.

**Temporary Social Media**

Messages that quickly self-destruct could enhance the privacy of online communications and make people freer to be spontaneous.

**Prenatal DNA Sequencing**

Reading the DNA of fetuses will be the next frontier of the genomic revolution. But do you really want to know about the genetic problems or musical aptitude of your unborn child?

**Additive Manufacturing**

Skeptical about 3-D printing? GE, the world's largest manufacturer, is on the verge of using the technology to make jet parts.

**Baxter: The Blue-Collar Robot**

Rodney Brooks's newest creation is easy to interact with, but the complex innovations behind the robot show just how hard it is to get along with people.

**Memory Implants**

**Smart Watches**

**Ultra-Efficient Solar Power**

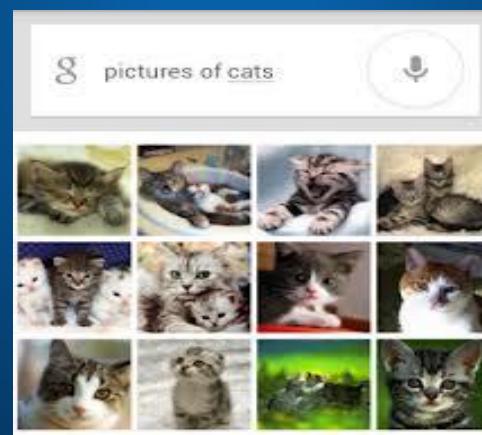
**Big Data from Cheap Phones**

**Supergrids**

Deep Learning is highlighted with a red box.

MIT Technology Review, April 23<sup>rd</sup>, 2013

# Deep Learning – from Research to Technology



Deep Learning - breakthrough in visual and speech recognition

# Classical Computer Vision Pipeline



# Classical Computer Vision Pipeline.

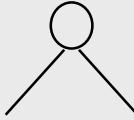
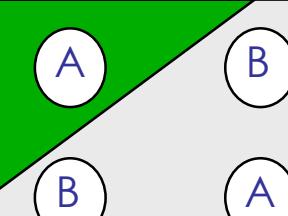
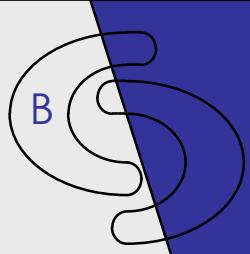
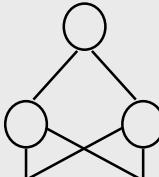
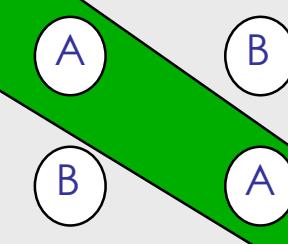
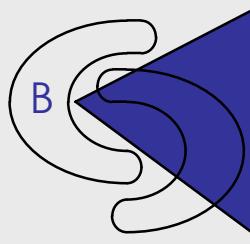
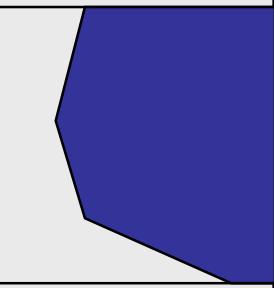
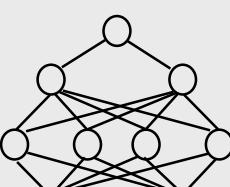
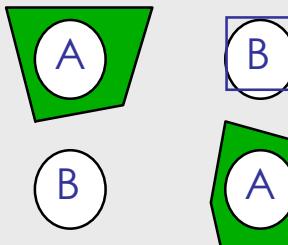
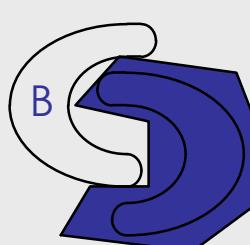
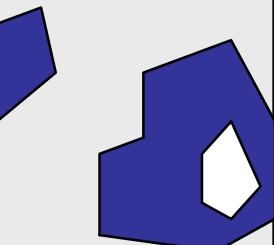
CV experts

1. Select / develop features: SURF, HoG, SIFT, RIFT, ...
2. Add on top of this Machine Learning for multi-class recognition and train classifier



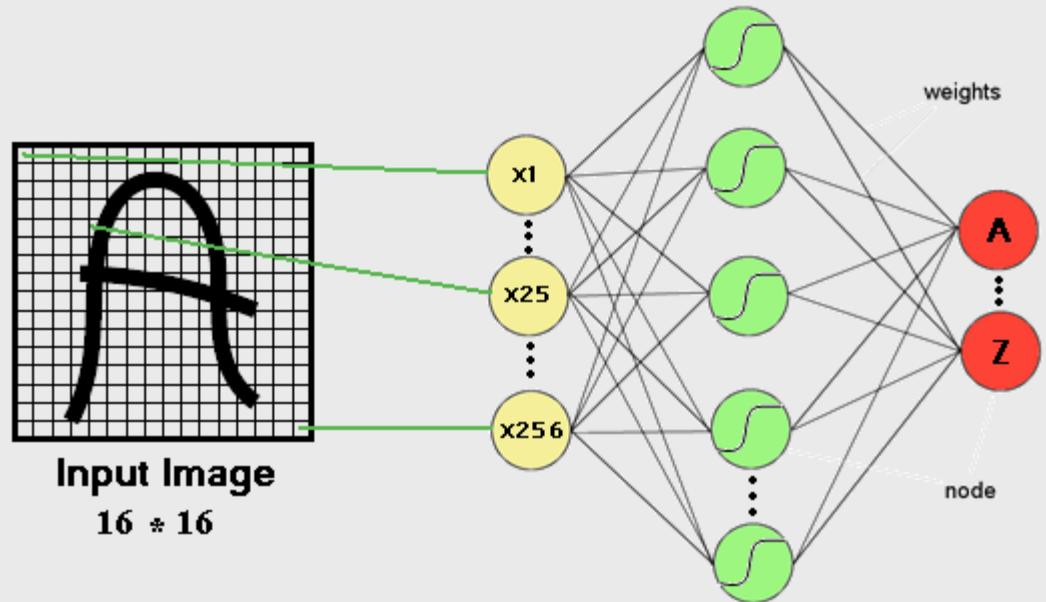
Classical CV feature definition is domain-specific and time-consuming

# Behavior of multilayer neural networks

Structure	Types of Decision Regions	Exclusive-OR Problem	Classes with Meshed regions	Most General Region Shapes
Single-Layer 	Half Plane Bounded By Hyperplane			
Two-Layer 	Convex Open Or Closed Regions			
Three-Layer 	Arbitrary (Complexity Limited by No. of Nodes)			

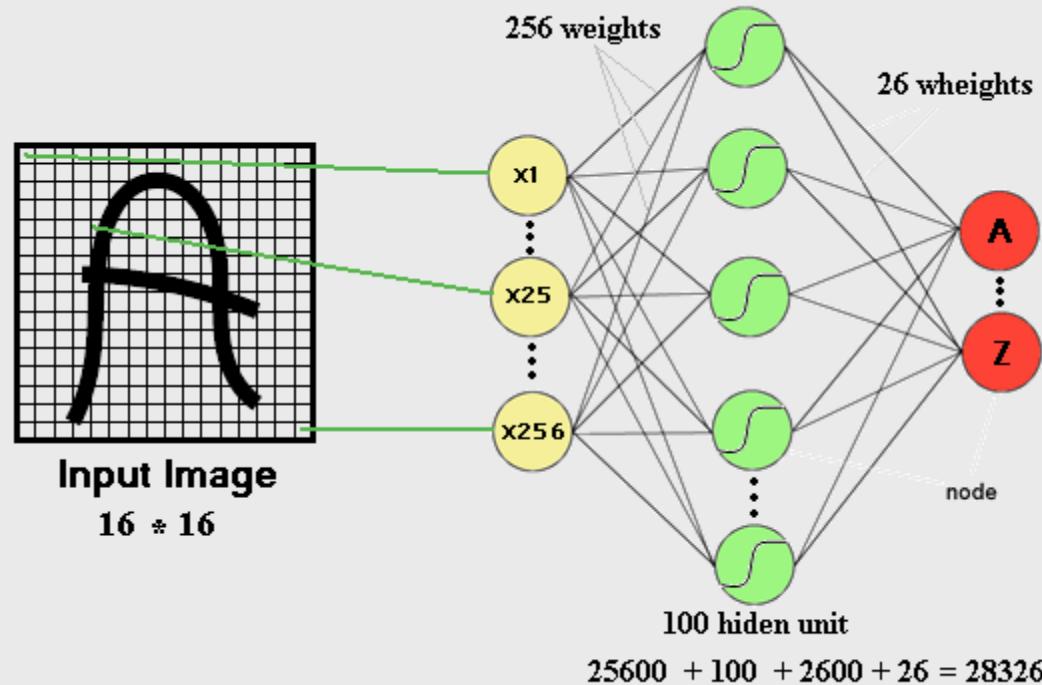
# Multi-layer perceptron and image processing

- ⦿ One or more hidden layers
- ⦿ Sigmoid activations functions



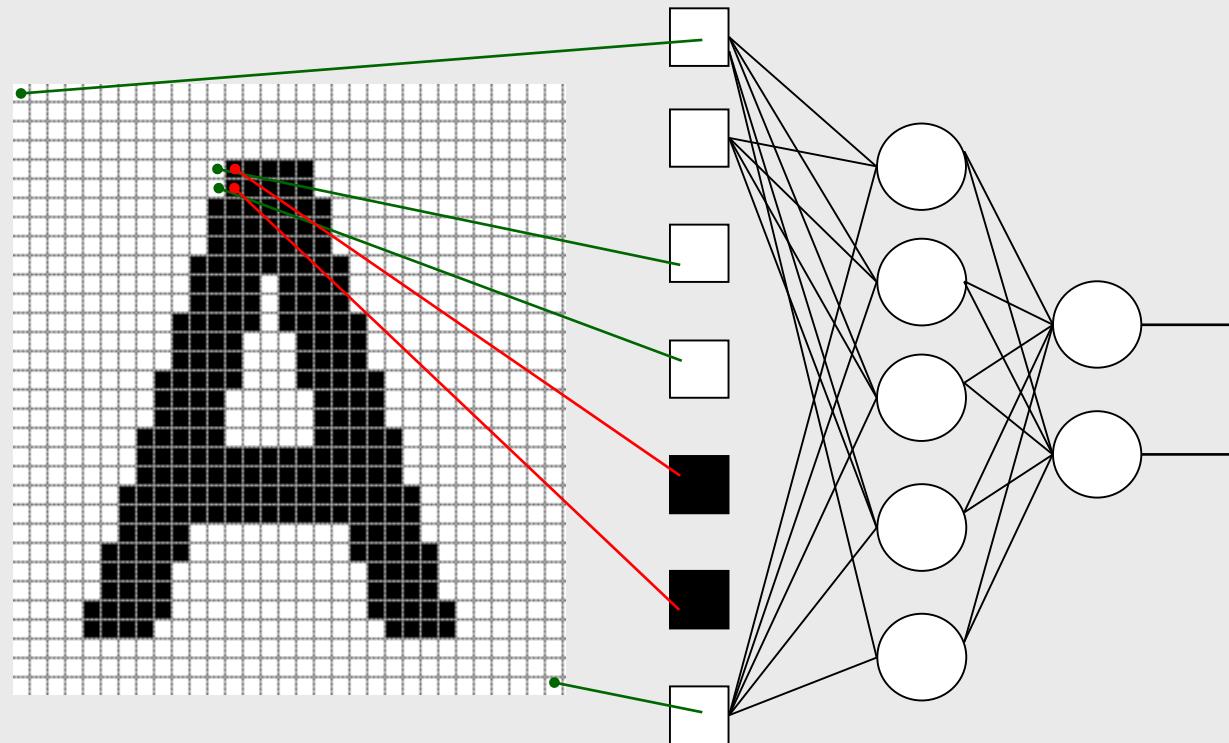
# Drawbacks of previous neural networks

- the number of trainable parameters becomes extremely large



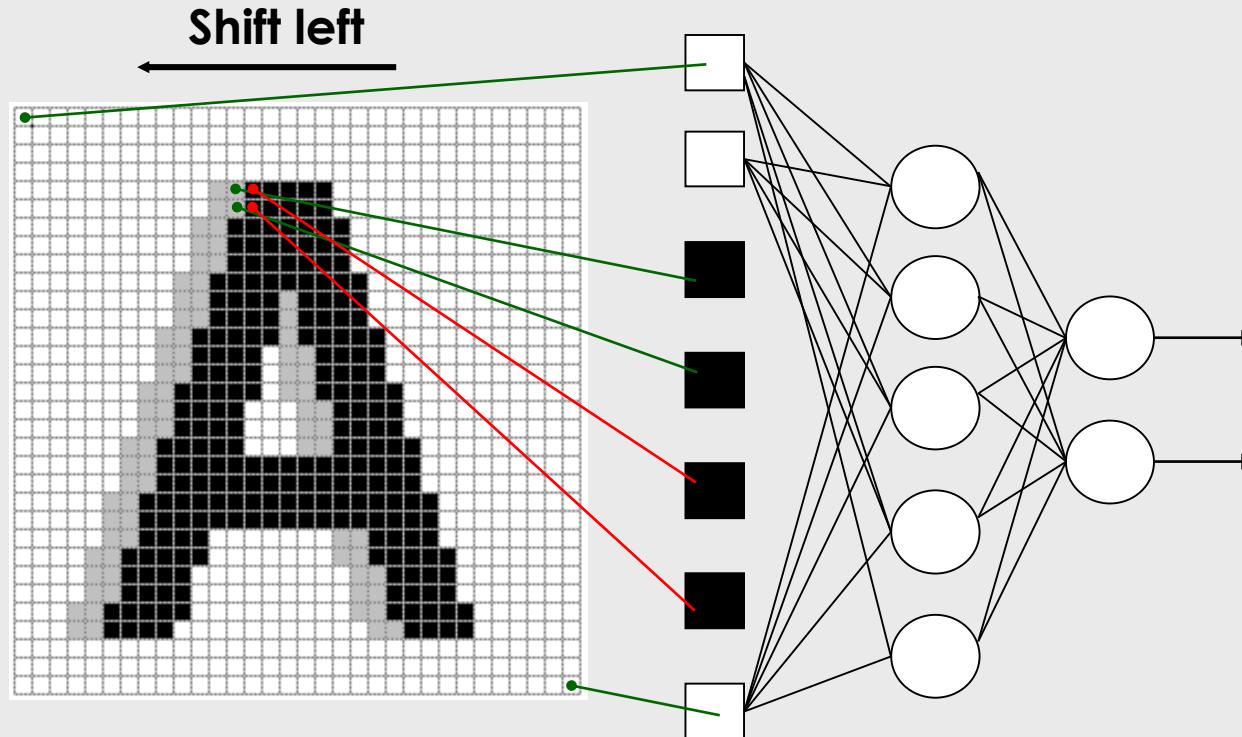
# Drawbacks of previous neural networks

- ⌚ Little or no invariance to shifting, scaling, and other forms of distortion

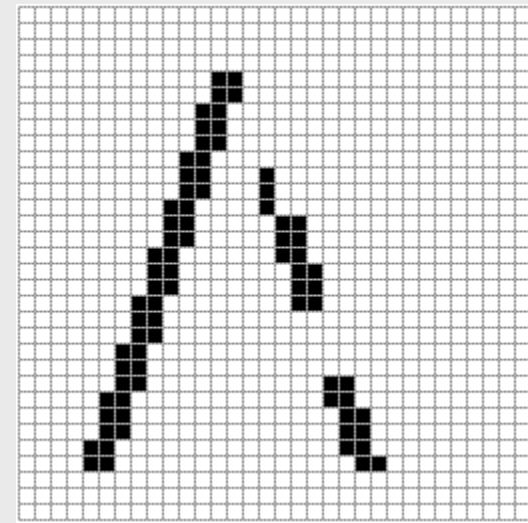
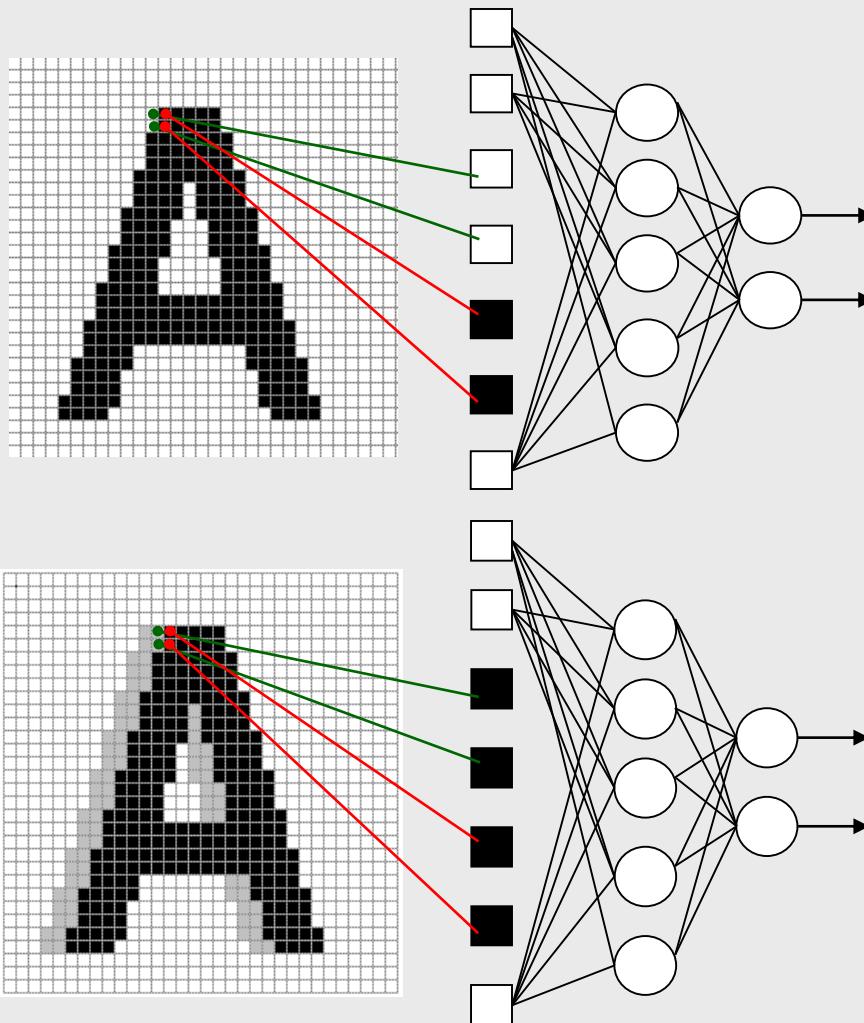


# Drawbacks of previous neural networks

- ⌚ Little or no invariance to shifting, scaling, and other forms of distortion



# Drawbacks of previous neural networks



**154 input change  
from 2 shift left**

**t**

**77 : black to white**

**77 : white to black**

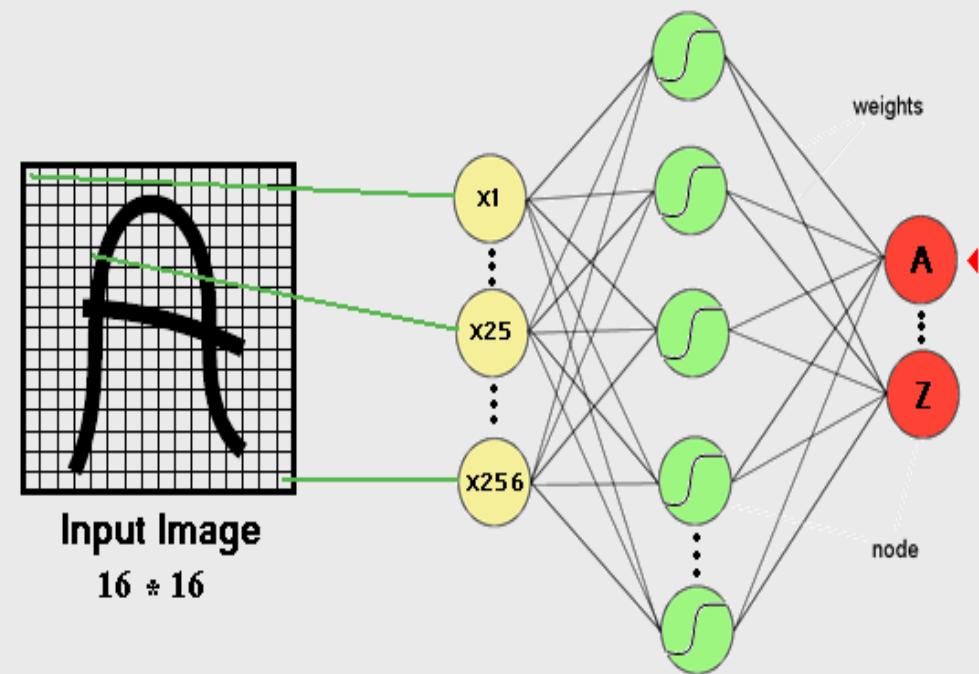
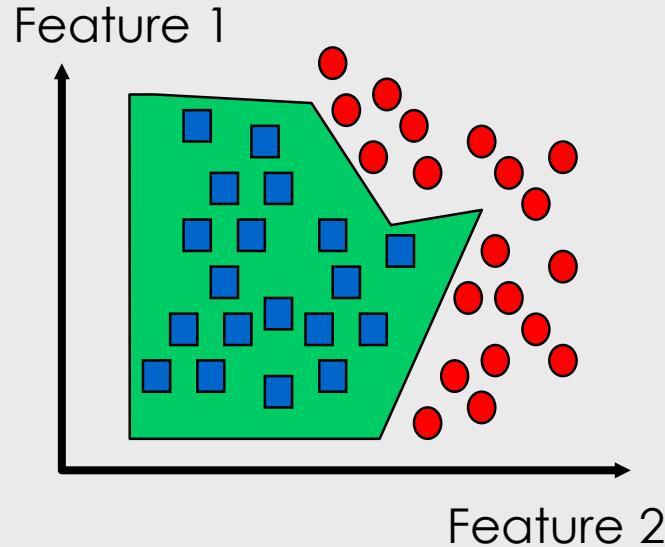
# Drawbacks of previous neural networks

- ⌚ scaling, and other forms of distortion



# Drawbacks of previous neural networks

- the topology of the input data is completely ignored
- work with raw data.



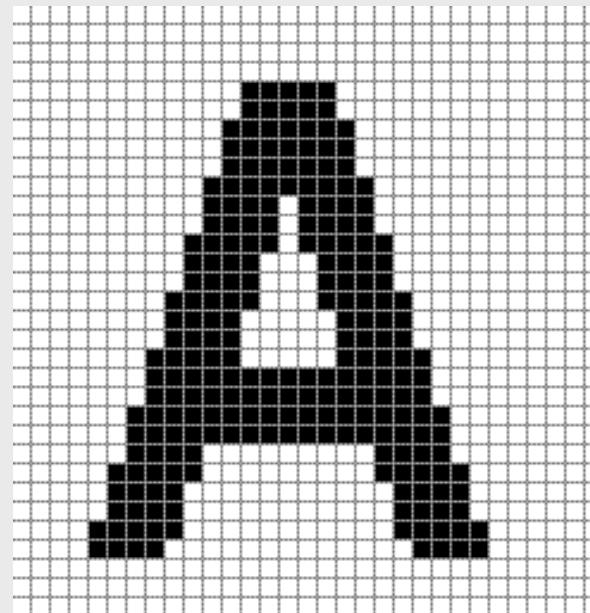
# Drawbacks of previous neural networks

**Black and white patterns:**

$$2^{32 \times 32} = 2^{1024}$$

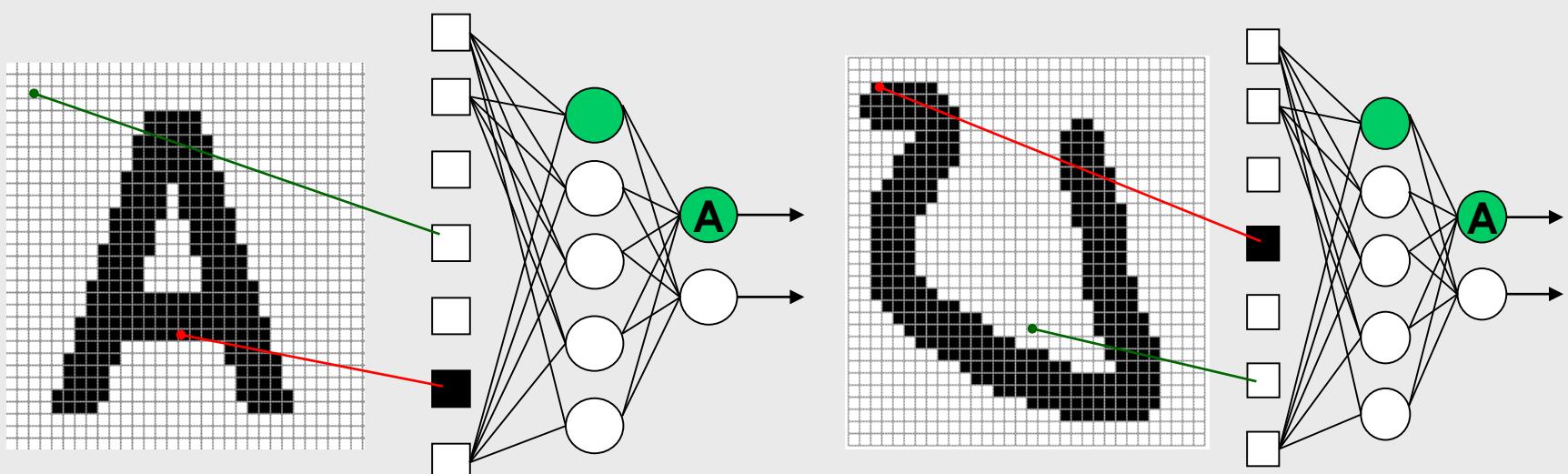
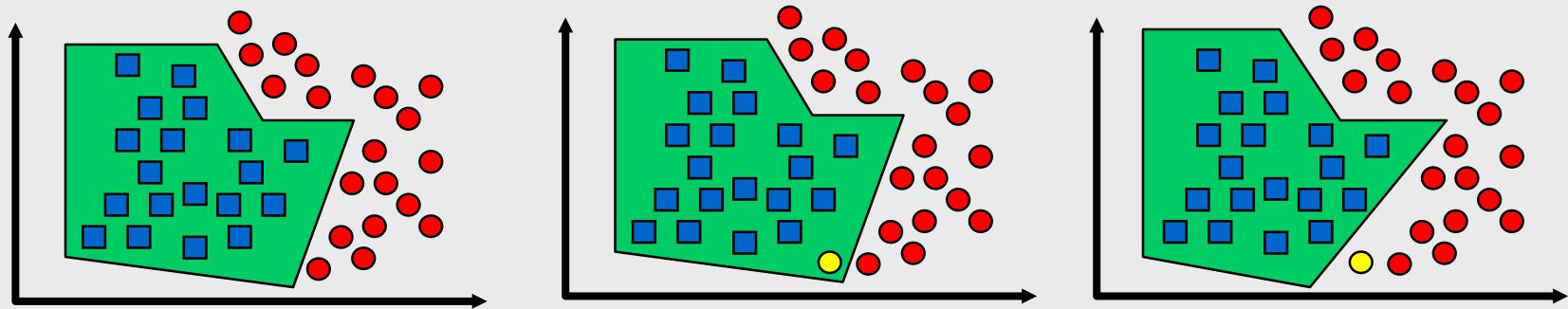
**Gray scale patterns :**

$$256^{32 \times 32} = 256^{1024}$$



**32 \* 32 input image**

# Drawbacks of previous neural networks



# Improvement

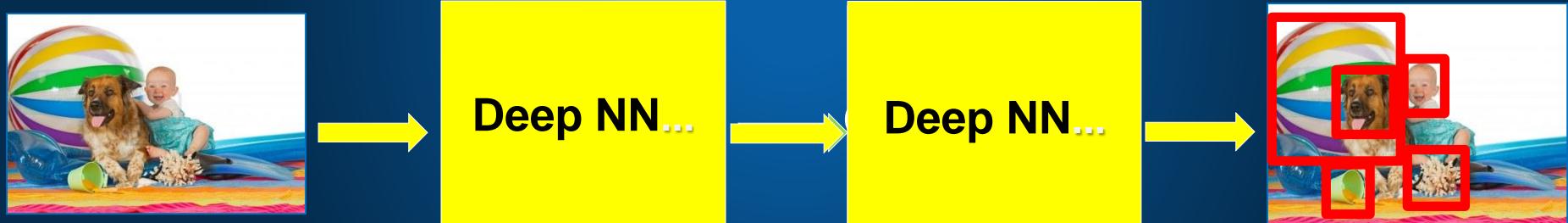
- ⌚ Fully connected network of sufficient size can produce outputs that are invariant with respect to such variations.
  - ⌚ Training time
  - ⌚ Network size
  - ⌚ Free parameters

# Deep Learning –based Vision Pipeline.

Deep Learning:

- Build features automatically based on training data
- Combine feature extraction and classification

DL experts: define NN topology and train NN



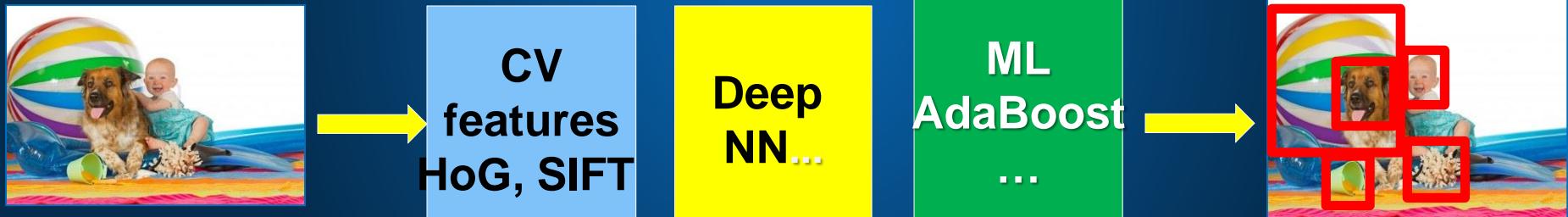
Deep Learning promise:  
train good feature automatically,  
same method for different domain

# Computer Vision +Deep Learning + Machine Learning

We want to combine Deep Learning + CV + ML

- Combine pre-defined features with learned features;
- Use best ML methods for multi-class recognition

CV+DL+ML experts needed to build the best-in-class

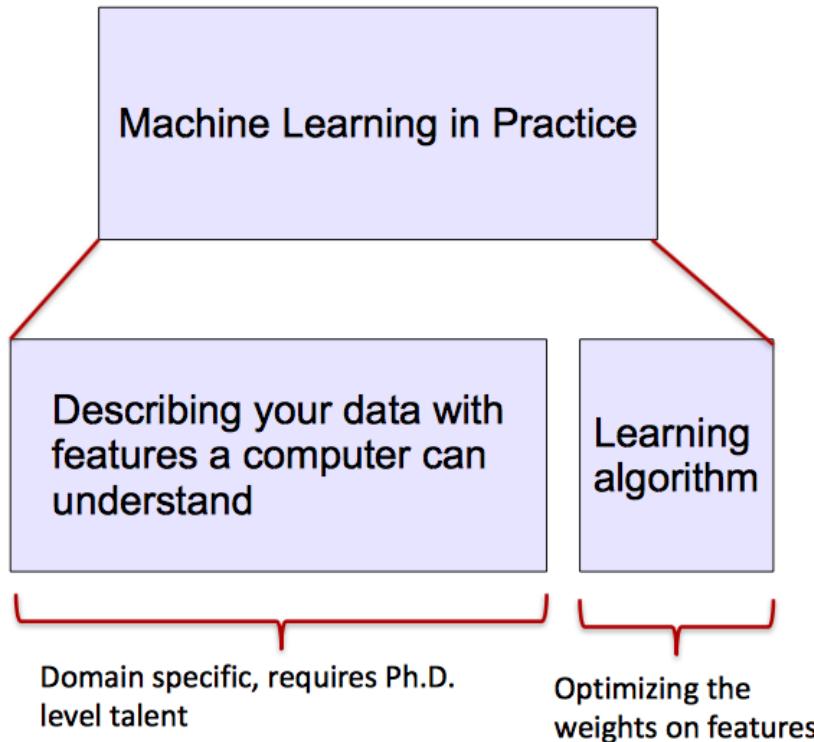


Combine best of Computer Vision  
Deep Learning and Machine Learning

# ML vs. Deep Learning

Most machine learning methods work well because of **human-designed representations** and **input features**

ML becomes just **optimizing weights** to best make a final prediction



Feature	NER
Current Word	✓
Previous Word	✓
Next Word	✓
Current Word Character n-gram	all
Current POS Tag	✓
Surrounding POS Tag Sequence	✓
Current Word Shape	✓
Surrounding Word Shape Sequence	✓
Presence of Word in Left Window	size 4
Presence of Word in Right Window	size 4

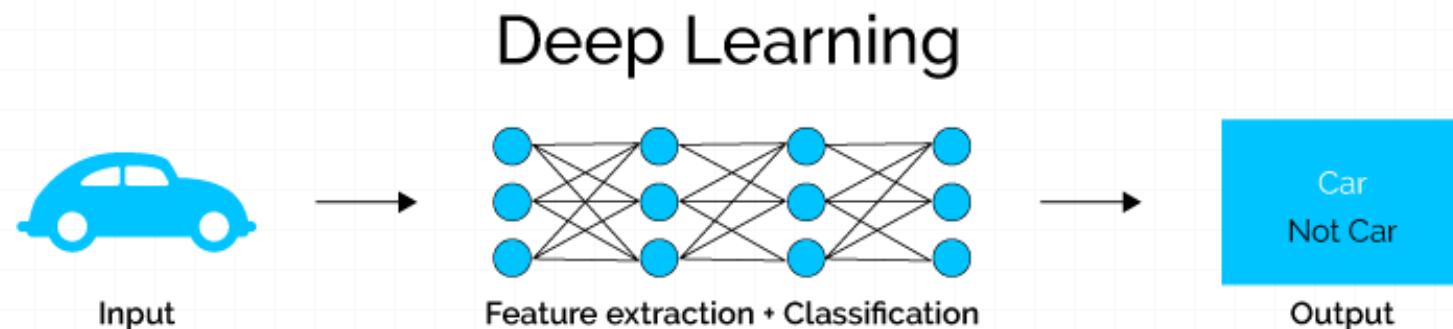
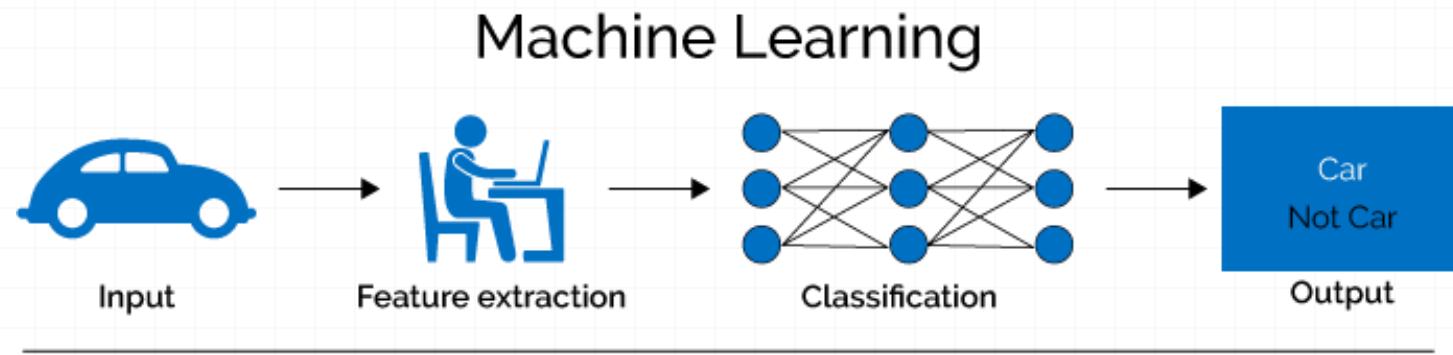
# What is Deep Learning (DL) ?

A machine learning subfield of learning **representations** of data.

Exceptional effective at **learning patterns**.

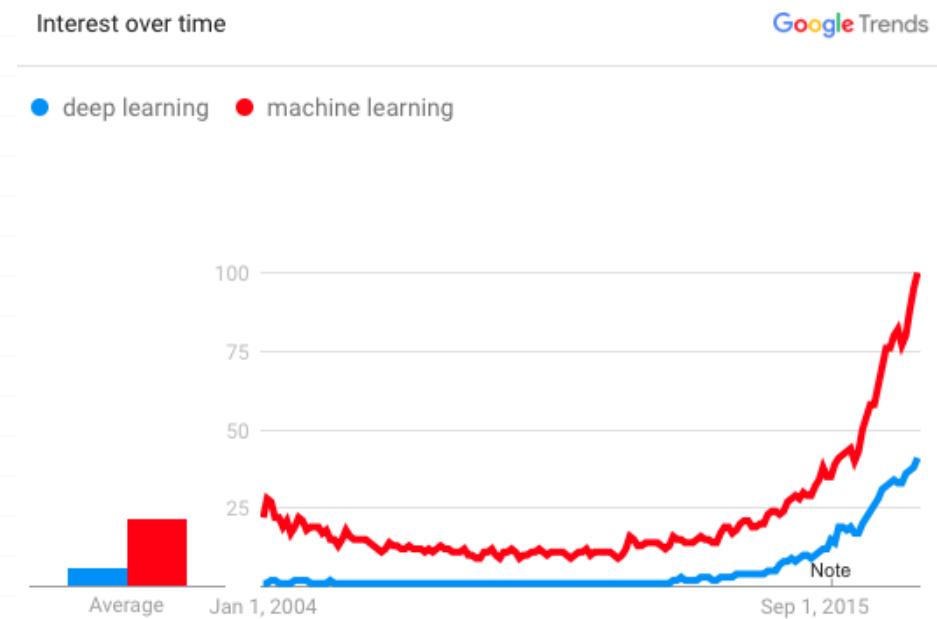
Deep learning algorithms attempt to learn (multiple levels of) representation by using a **hierarchy of multiple layers**

If you provide the system **tons of information**, it begins to understand it and respond in useful ways.



# Why is DL useful?

- Manually designed features are often **over-specified, incomplete** and take a **long time to design** and validate
- Learned Features are **easy to adapt, fast** to learn
- Deep learning provides a very **flexible**, (almost?) **universal**, learnable framework for representing world, visual and linguistic information.
- Can learn both unsupervised and supervised
- Effective **end-to-end** joint system learning
- Utilize large amounts of training data



In ~2010 DL started outperforming other ML techniques first in speech and vision, then NLP

# Consider learning an image:

- Some patterns are much smaller than the whole image

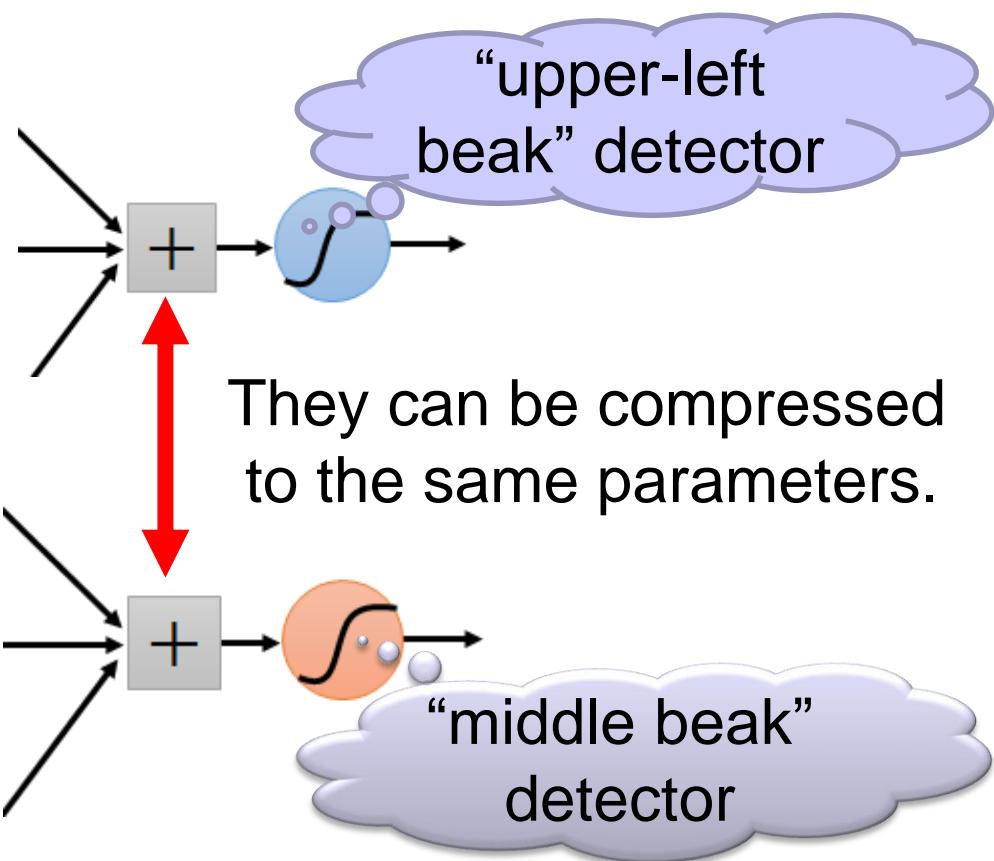
Can represent a small region with fewer parameters



Same pattern appears in different places:

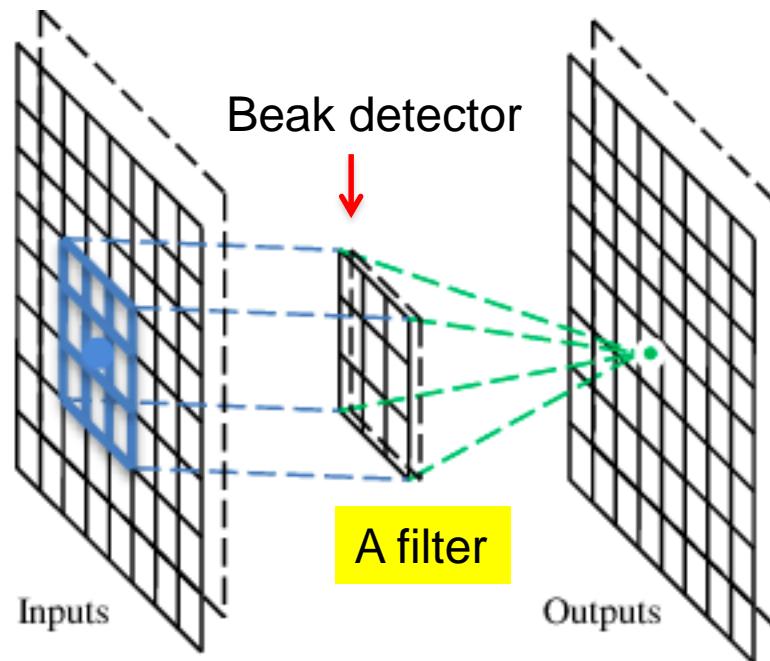
They can be compressed!

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



# 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

⋮ ⋮

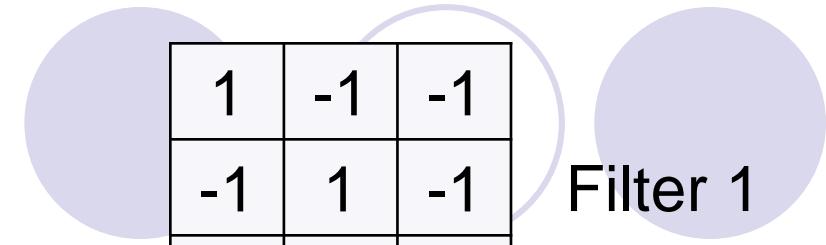
Each filter detects a small pattern (3 x 3).

# 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

Dot  
product



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

3

-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

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

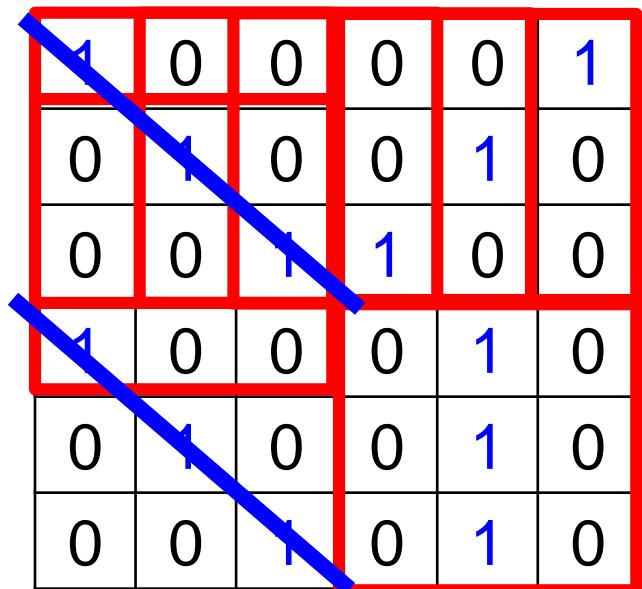
Filter 1

3

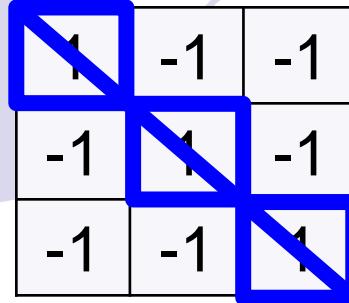
-3

# Convolution

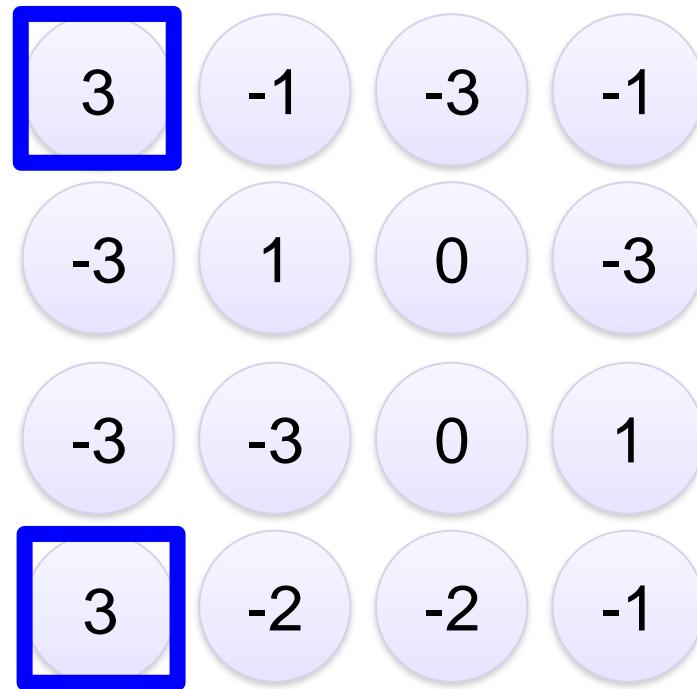
stride=1



6 x 6 image



Filter 1



# Convolution

stride=1

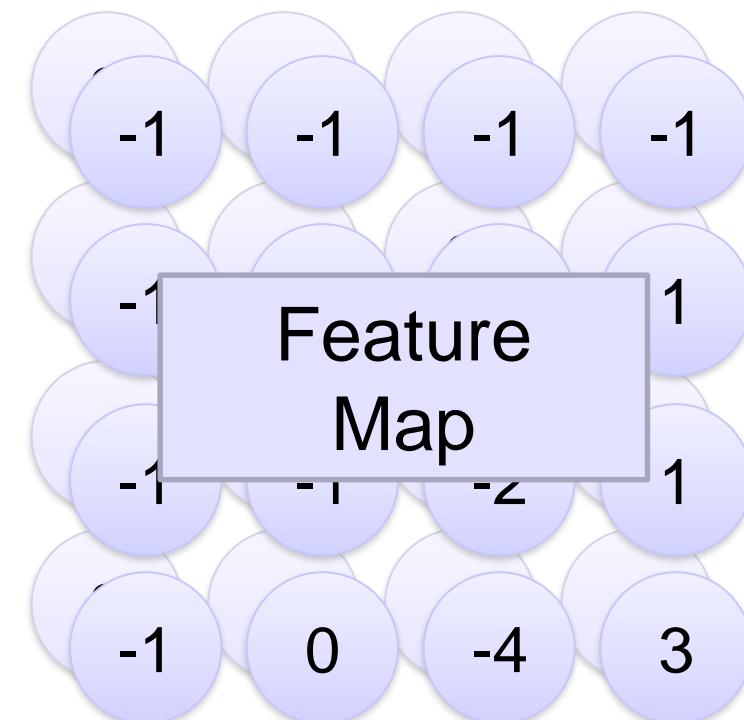
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

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

Filter 2

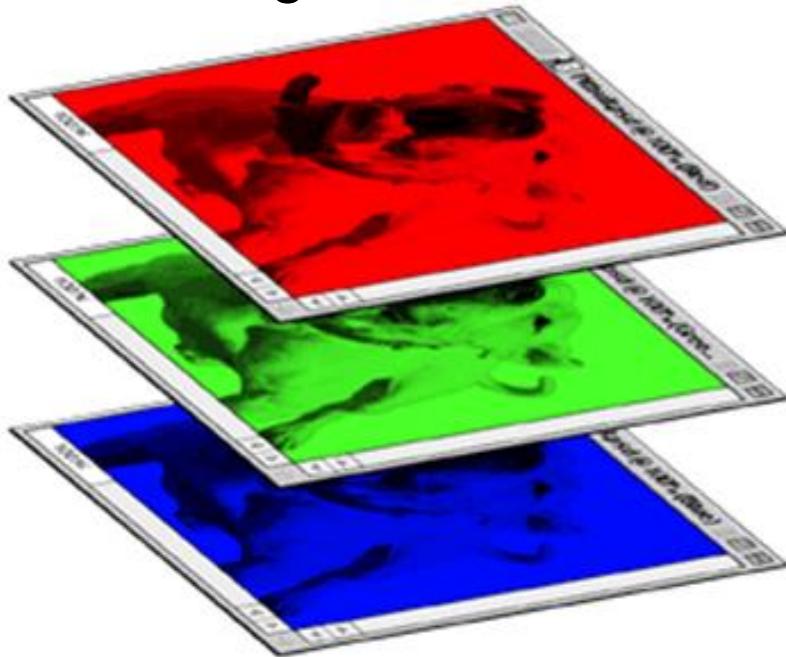
Repeat this for each filter



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

# Color image: RGB 3 channels

Color 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

# Deep Learning Taxonomy

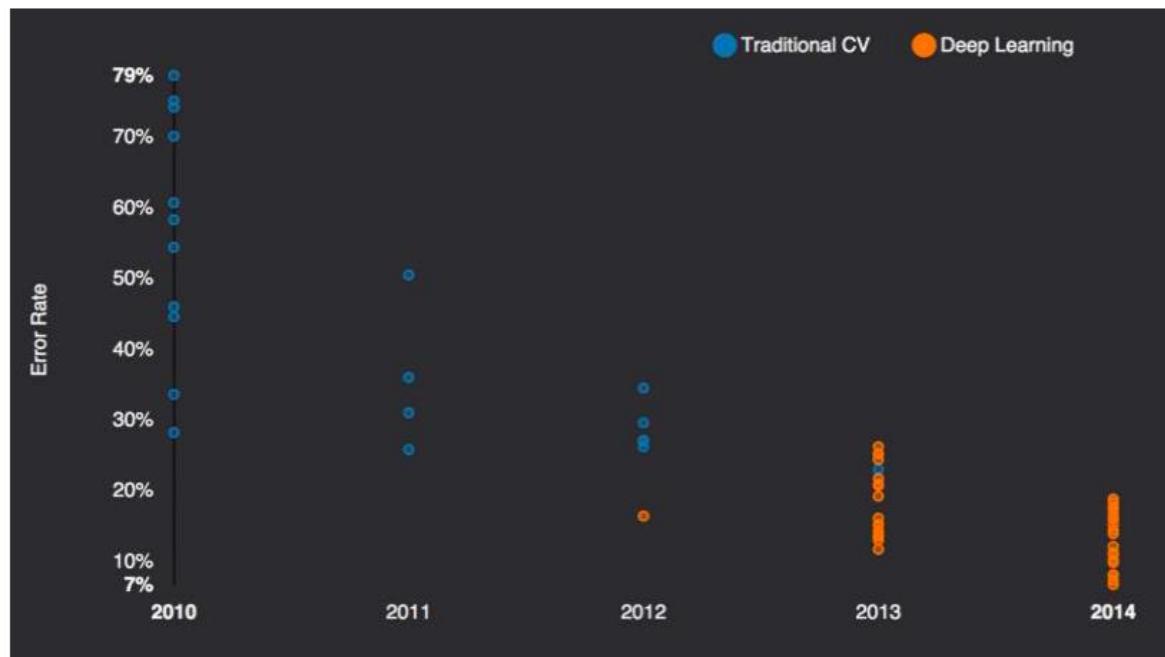
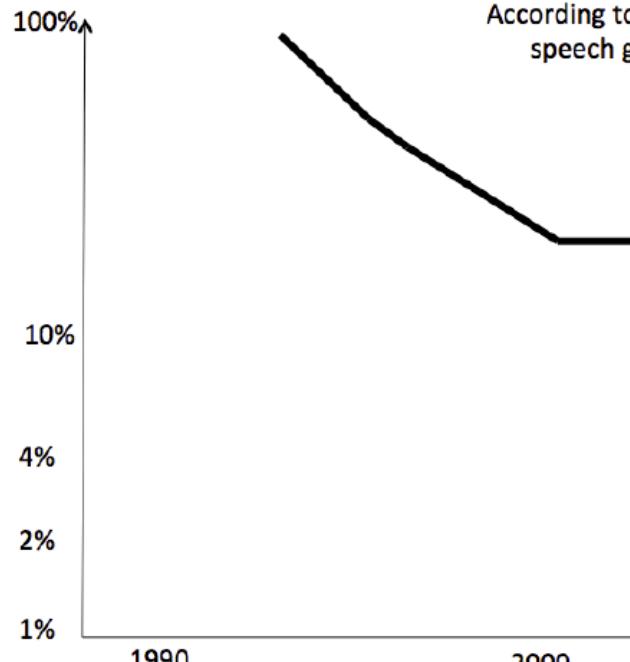
## Supervised:

- Convolutional NN ( LeCun)
- Recurrent Neural nets (Schmidhuber )

## Unsupervised

- Deep Belief Nets / Stacked RBMs (Hinton)
- Stacked denoising autoencoders (Bengio)
- Sparse AutoEncoders ( LeCun, A. Ng, )

# State of the art in ...



Several big improvements in recent years in NLP

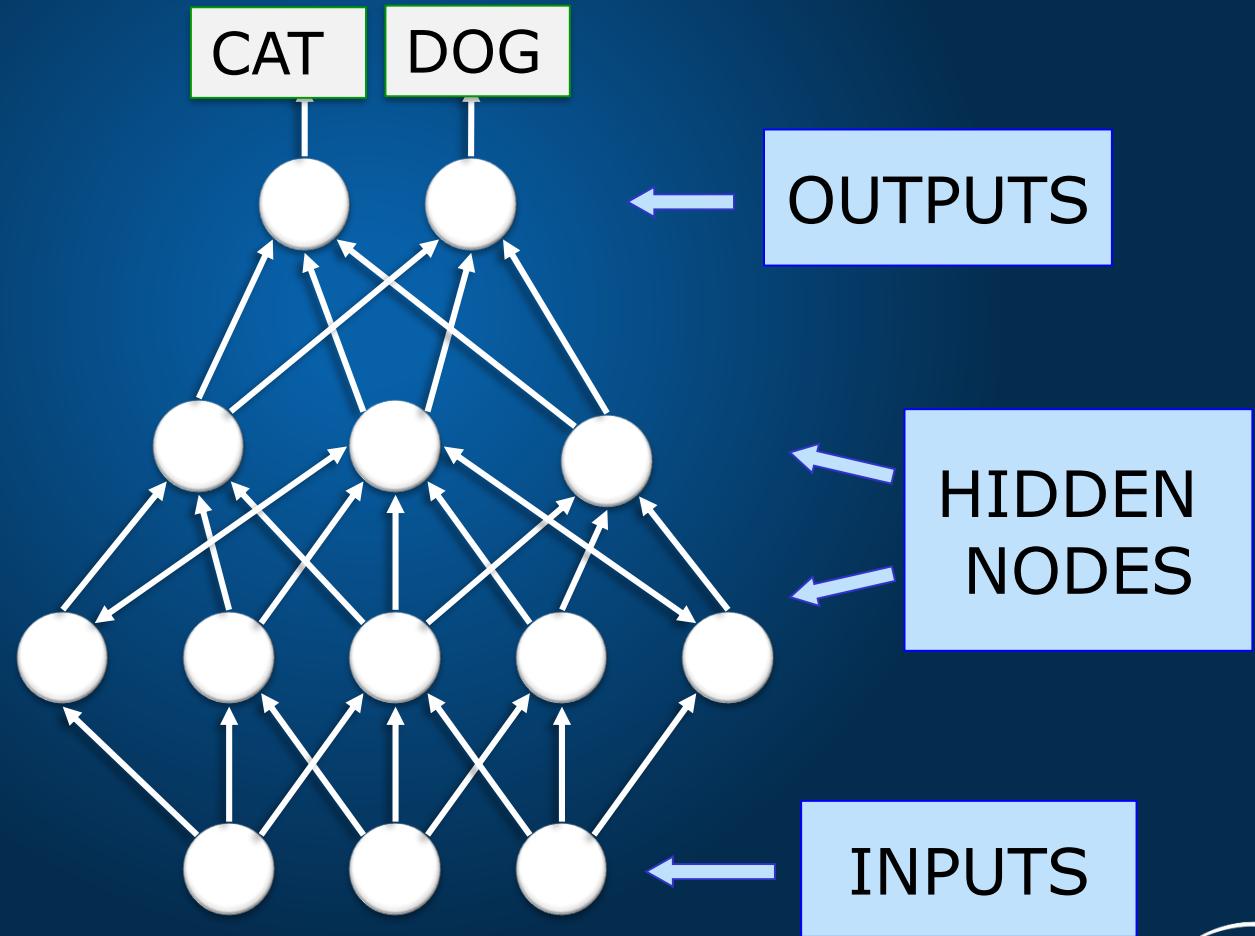
- ✓ Machine Translation
- ✓ Sentiment Analysis
- ✓ Dialogue Agents
- ✓ Question Answering
- ✓ Text Classification ...

Leverage different levels of representation

- words & characters
- syntax & semantics

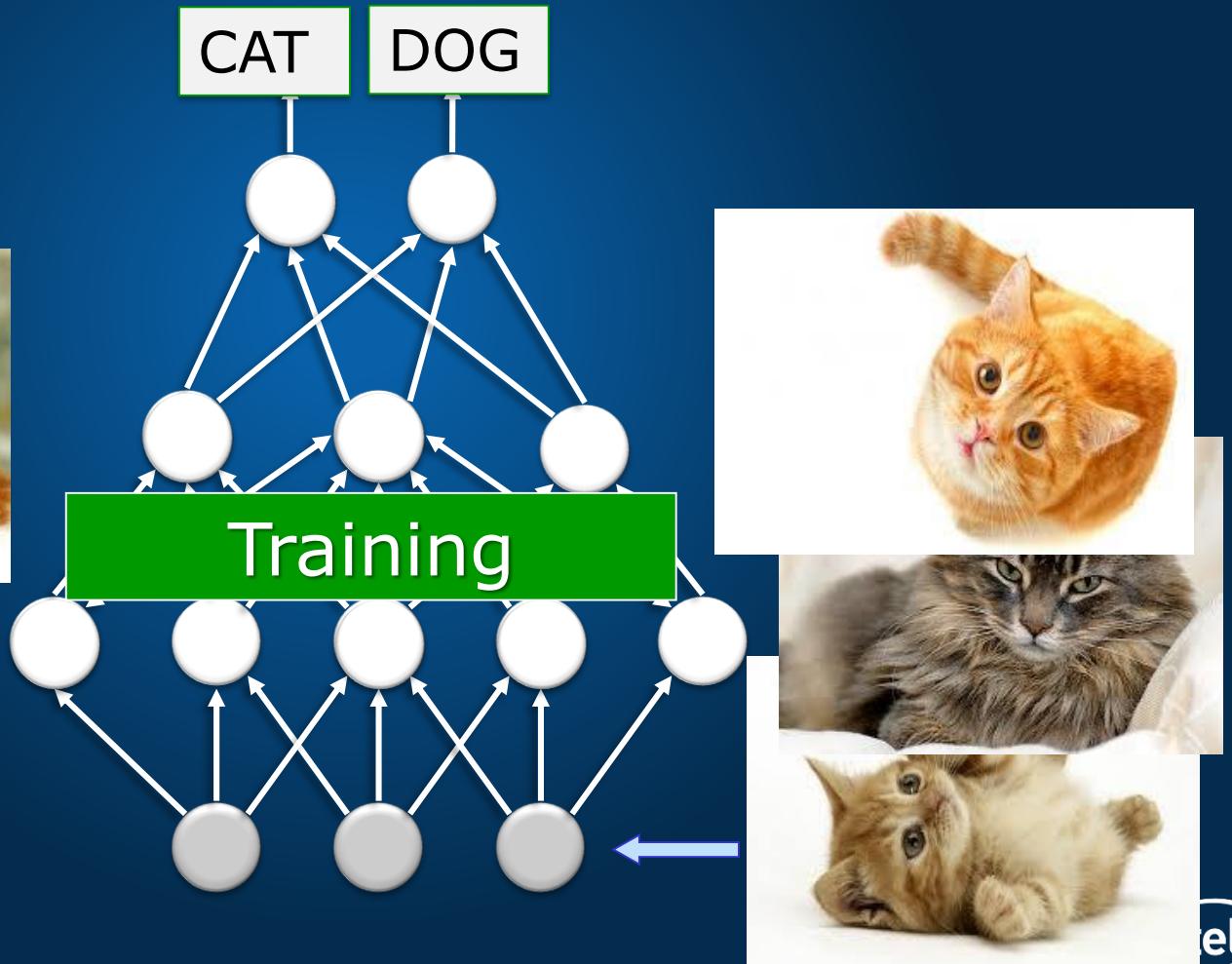
# Deep Learning Basics

Deep Learning – is a set of machine learning algorithms based on multi-layer networks



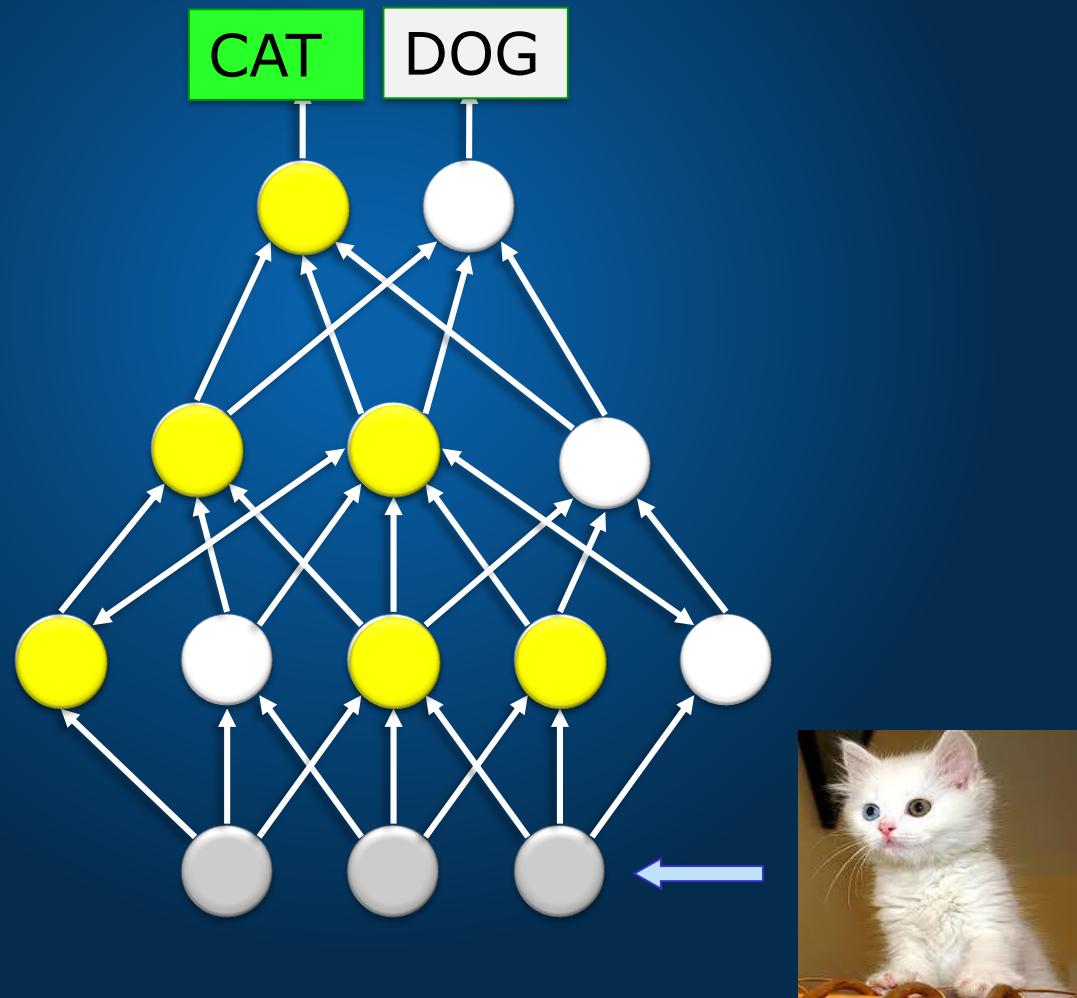
# Deep Learning Basics

Deep Learning – is a set of machine learning algorithms based on multi-layer networks



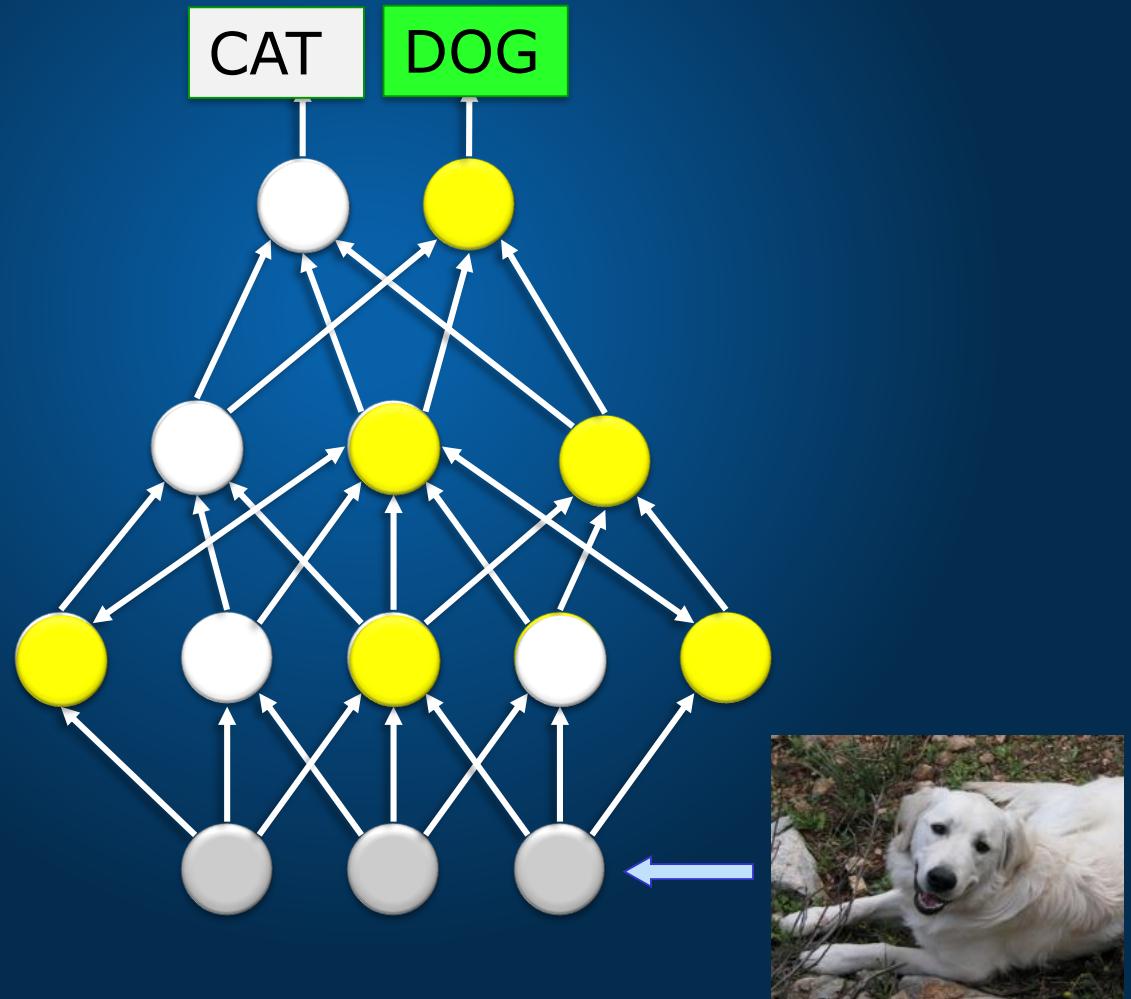
# Deep Learning Basics

Deep Learning – is a set of machine learning algorithms based on multi-layer networks



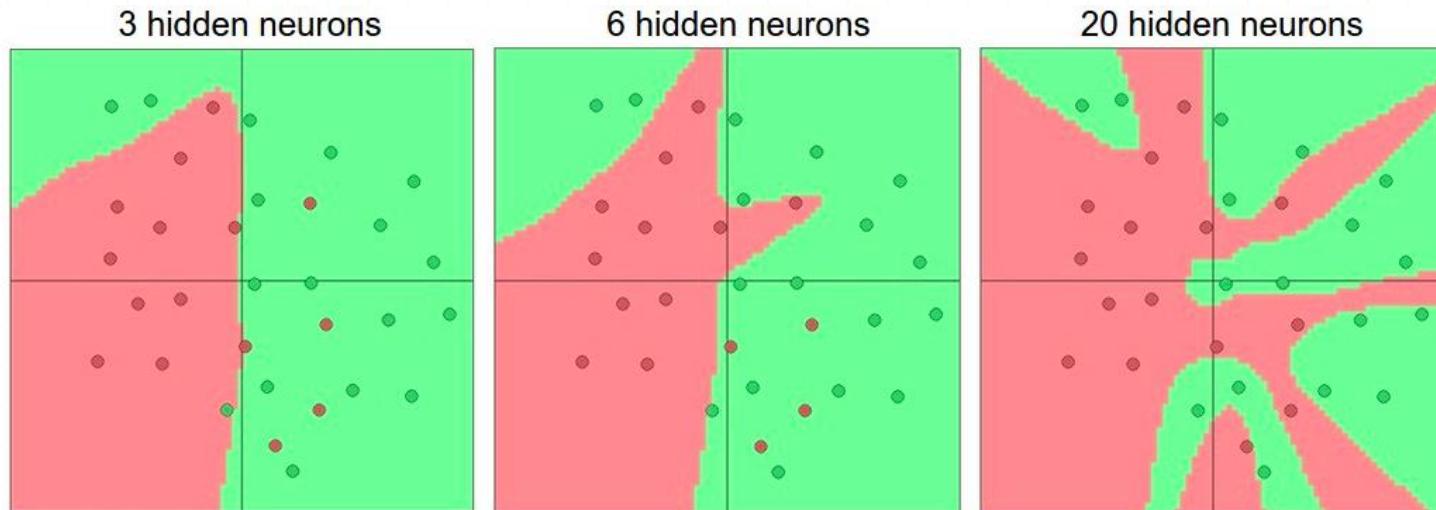
# Deep Learning Basics

Deep Learning – is a set of machine learning algorithms based on multi-layer networks



# Activation functions

Non-linearities needed to learn complex (non-linear) representations of data, otherwise the NN would be just a linear function  $W_1 W_2 x = Wx$



[http://cs231n.github.io/assets/nn1/layer\\_sizes.jpeg](http://cs231n.github.io/assets/nn1/layer_sizes.jpeg)

More layers and neurons can approximate more complex functions

Full list: [https://en.wikipedia.org/wiki/Activation\\_function](https://en.wikipedia.org/wiki/Activation_function)

# Convolutional Neural Networks (CNNs)

Main CNN idea for text:

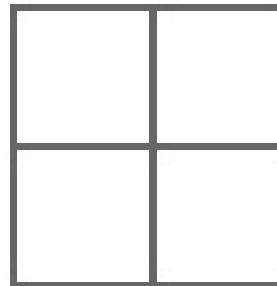
Compute vectors for n-grams and **group them afterwards**

Feature Map

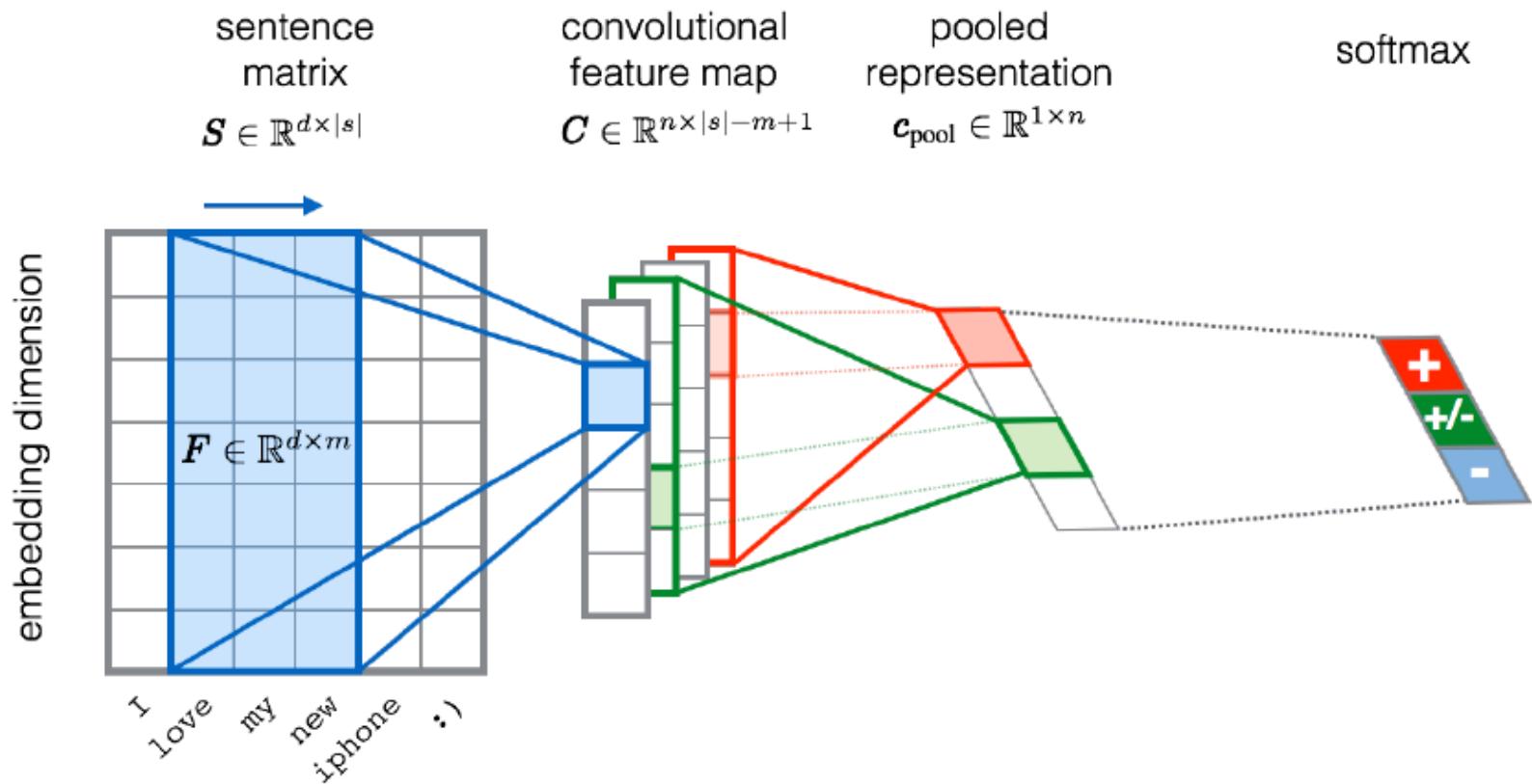
6	4	8	5
5	4	5	8
3	6	7	7
7	9	7	2

max pool  
2x2 filters  
and stride 2

Max-Pooling

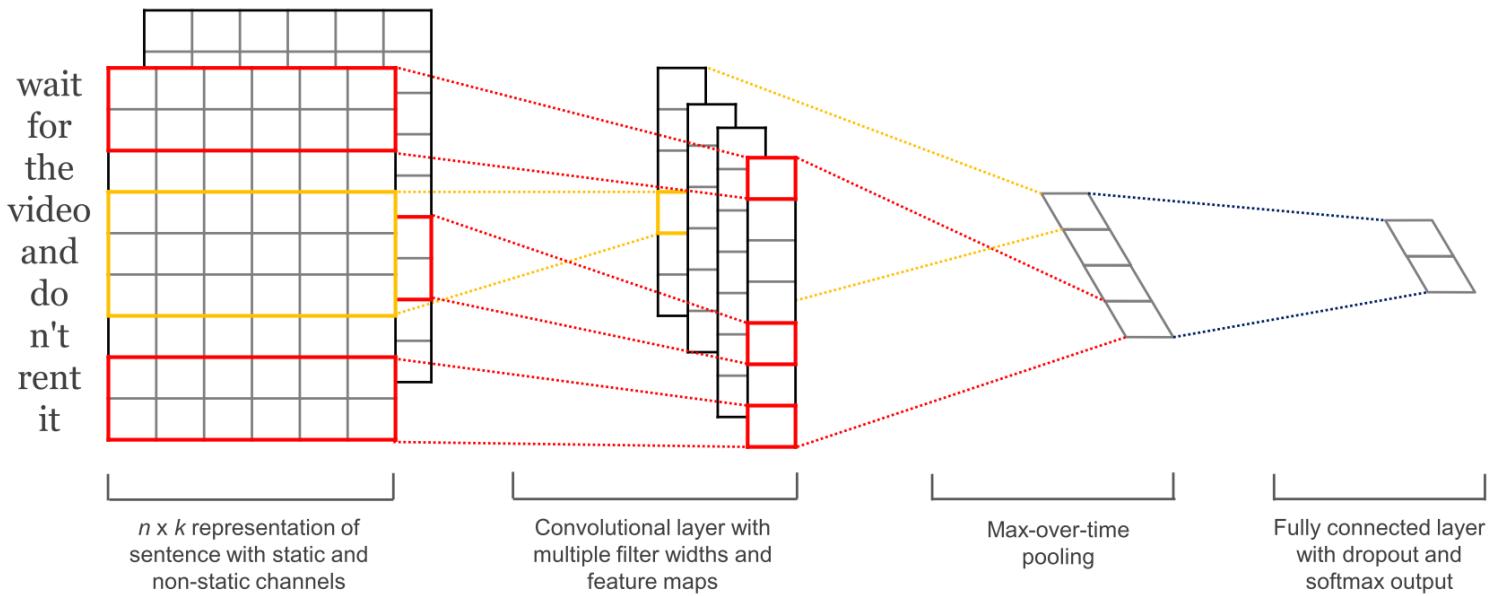


# CNN for text classification



Severyn, Aliaksei, and Alessandro Moschitti. "UNITN: Training Deep Convolutional Neural Network for Twitter Sentiment Classification." SemEval@ NAACL-HLT. 2015.

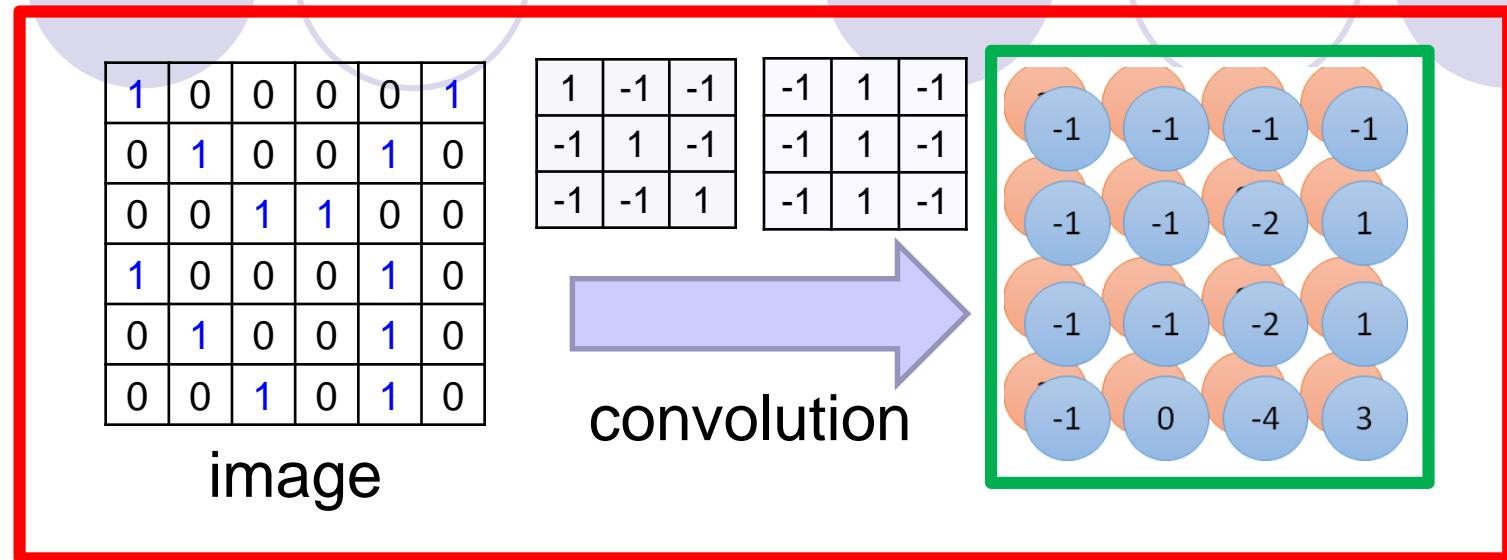
# CNN with multiple filters



Kim, Y. "Convolutional Neural Networks for Sentence Classification", EMNLP (2014)

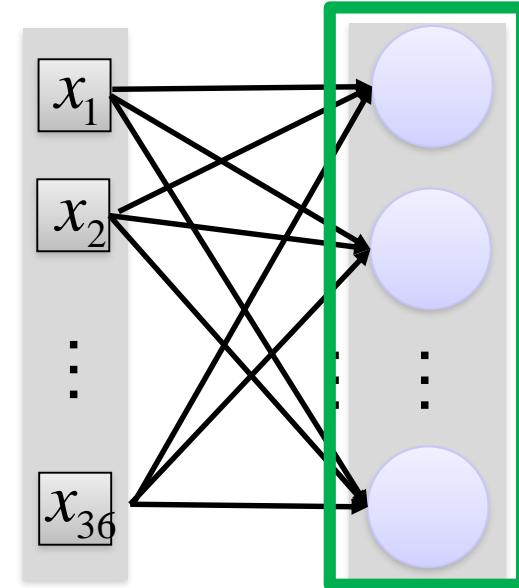
sliding over 3, 4 or 5 words at a time

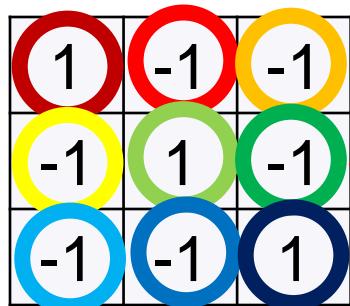
# Convolution v.s. Fully Connected



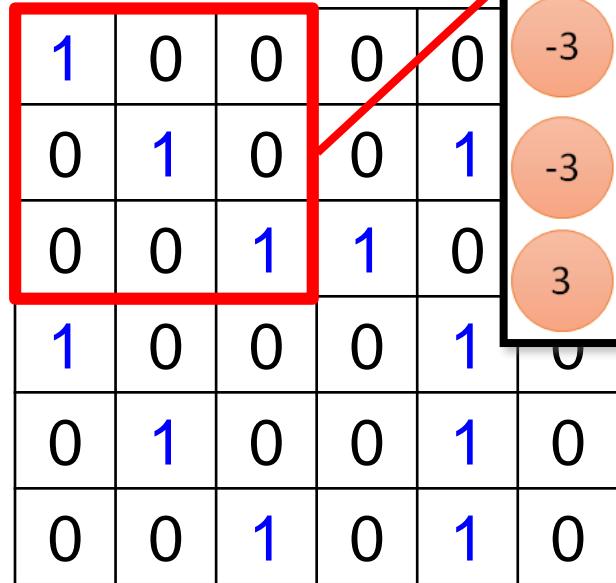
Fully-  
connected

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

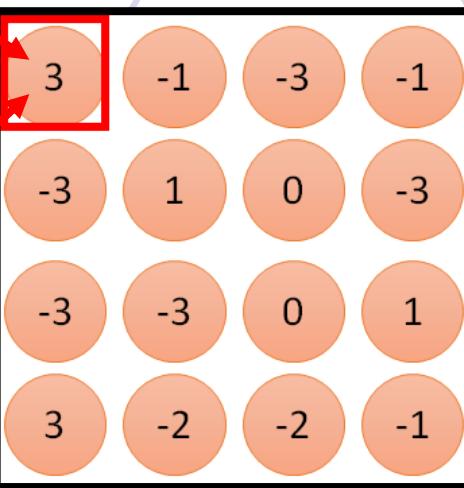




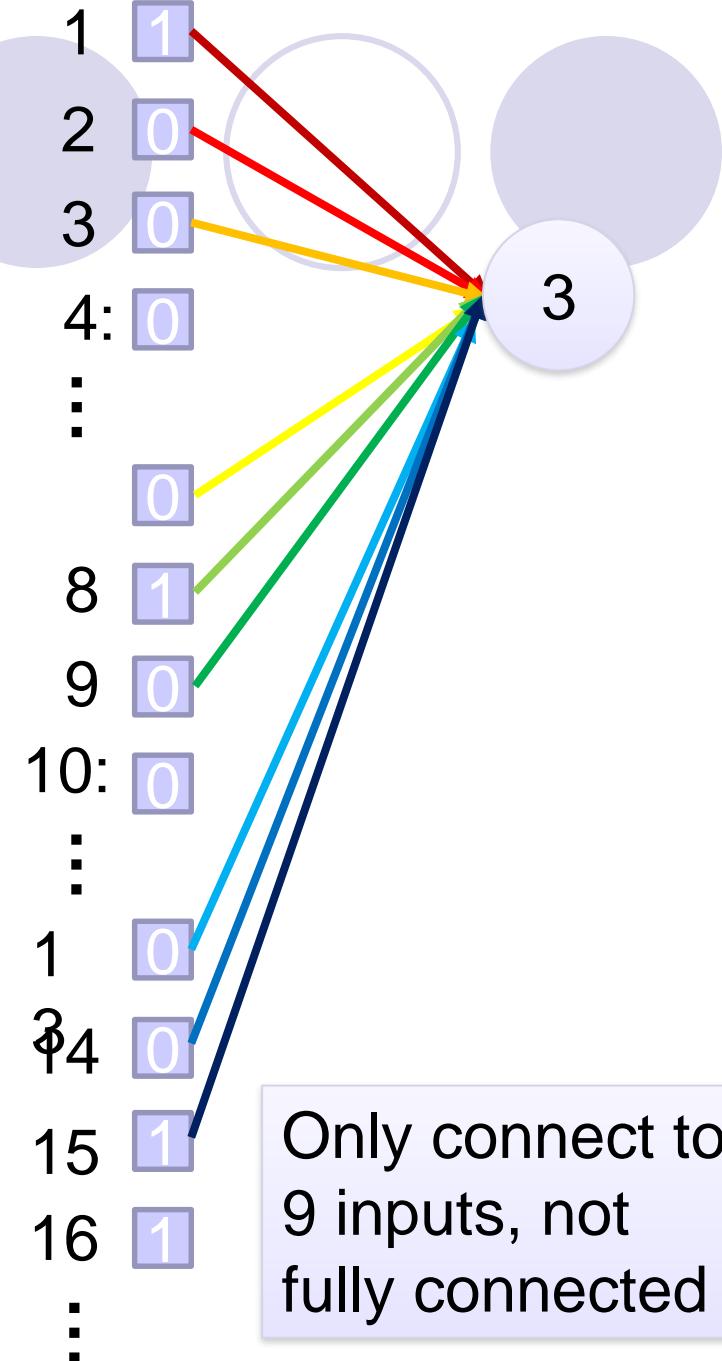
Filter 1

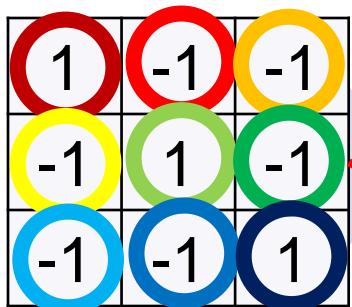


$6 \times 6$  image

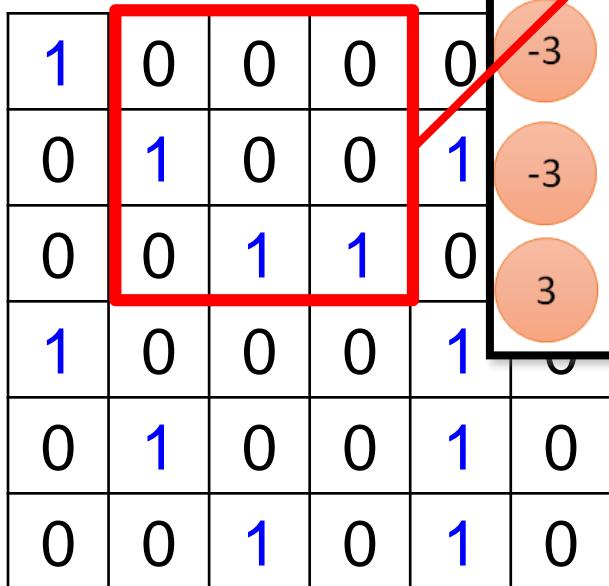


fewer parameters!





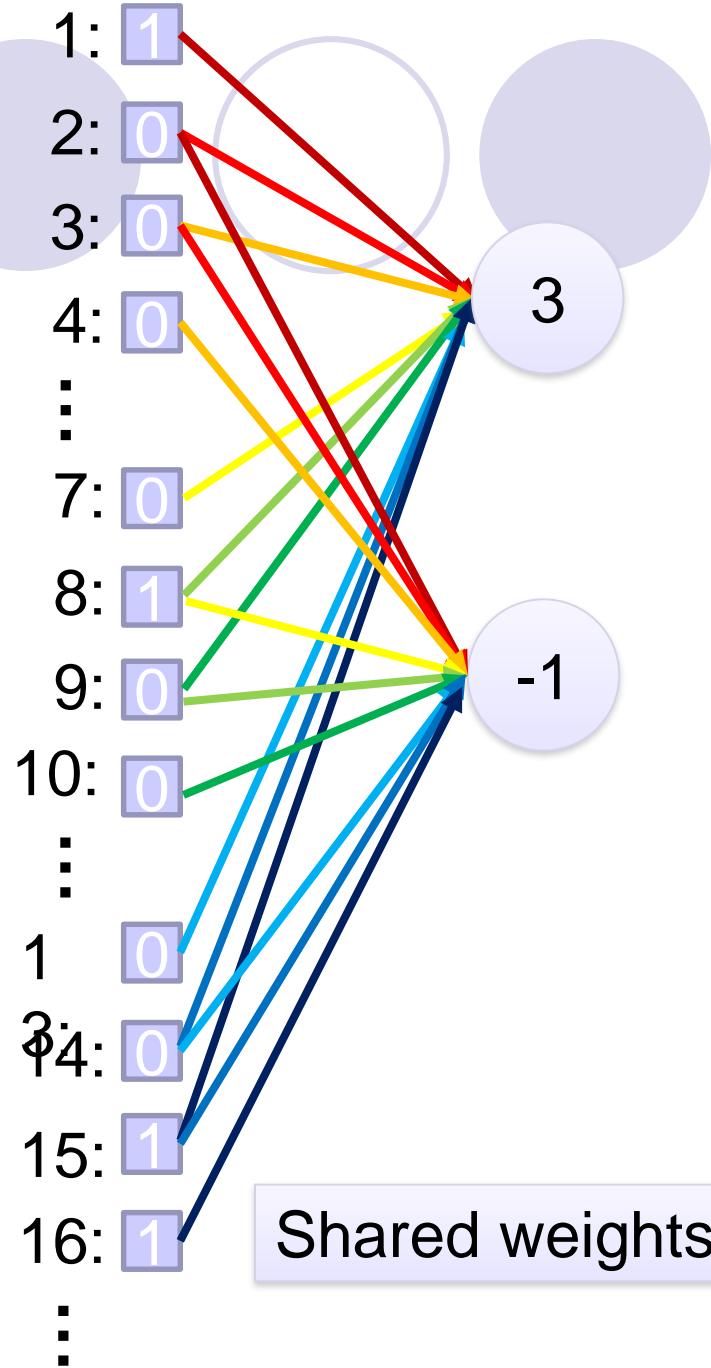
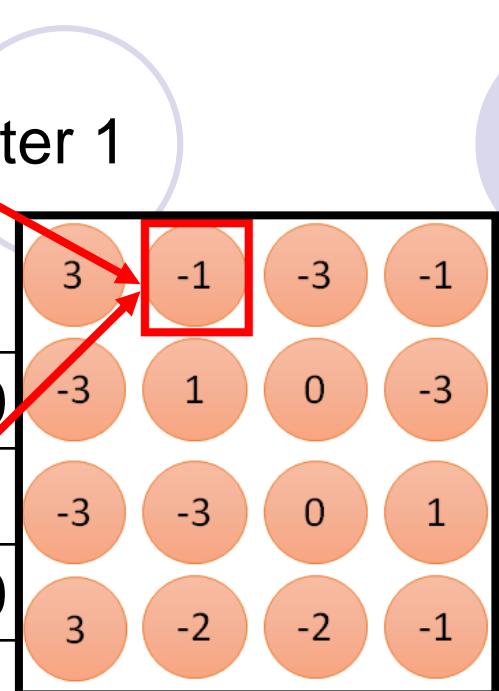
Filter 1



$6 \times 6$  image

Fewer parameters

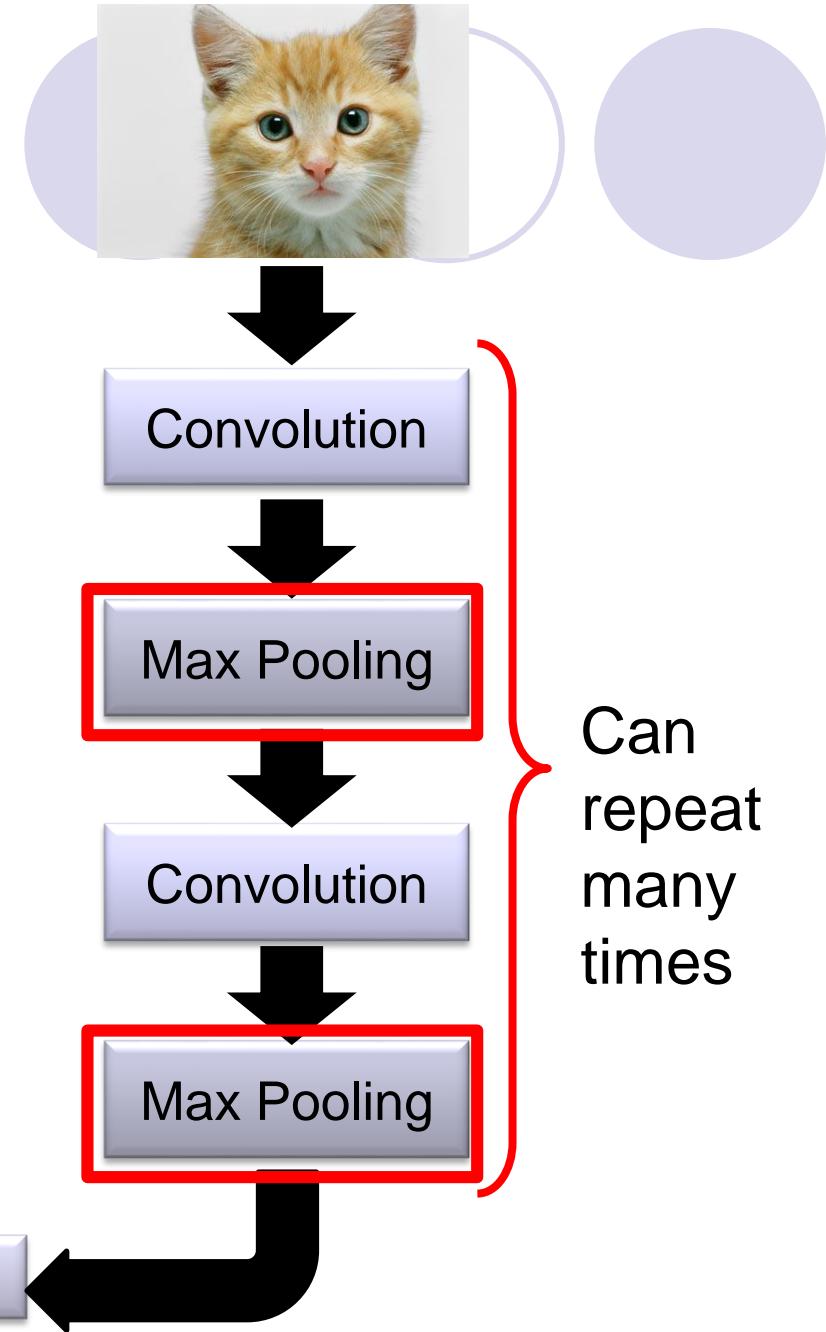
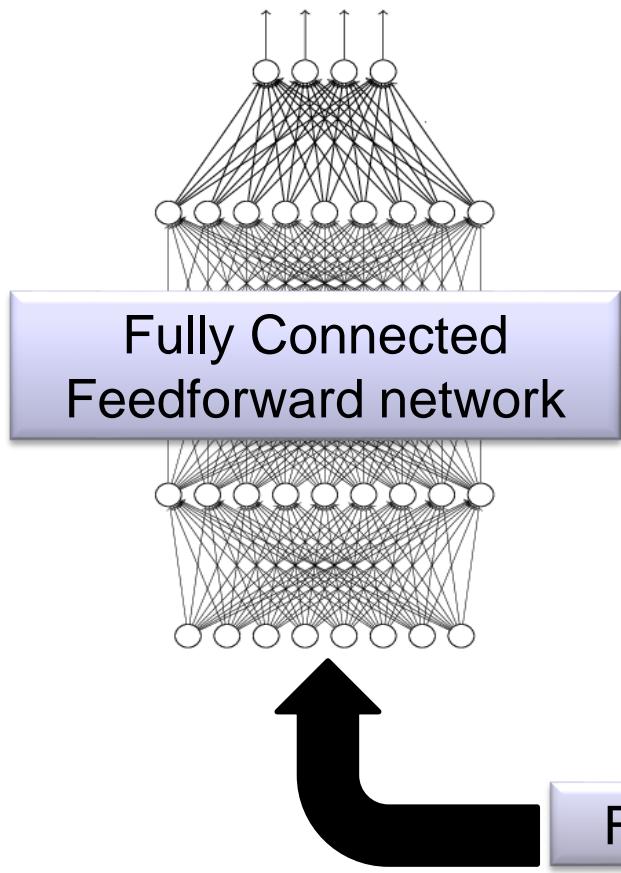
Even fewer parameters



Shared weights

# The whole CNN

cat dog .....



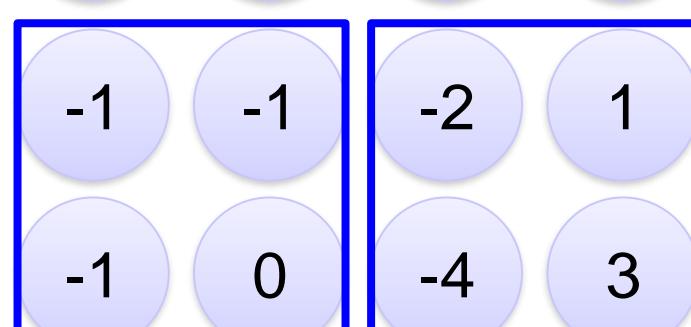
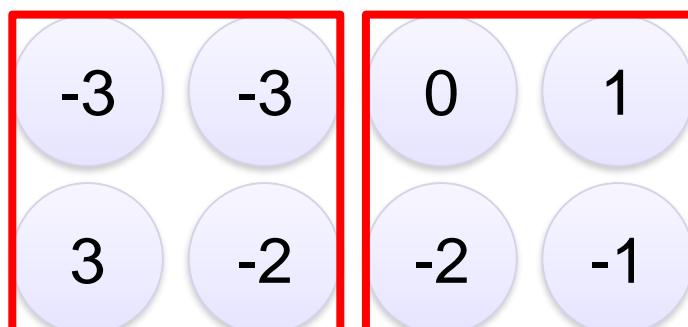
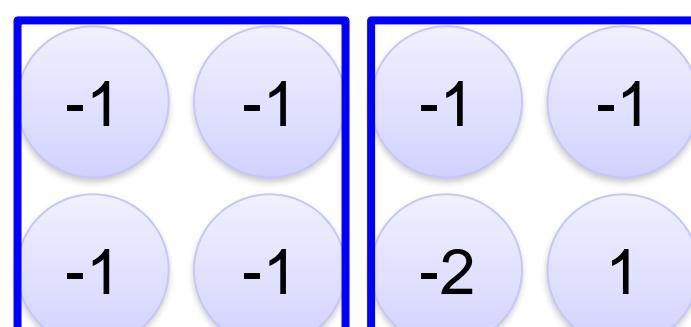
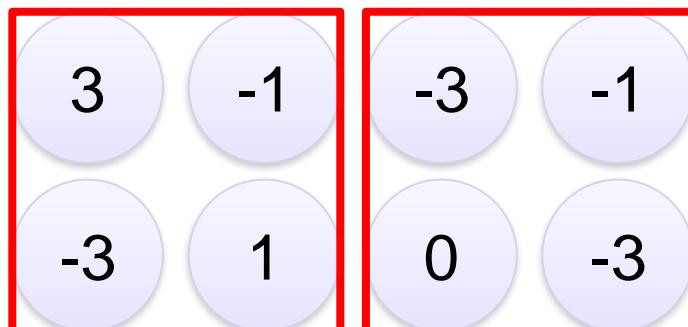
# Max Pooling

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

Filter 1

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

Filter 2



# Why Pooling

- Subsampling pixels will not change the object  
bird



Subsampling



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

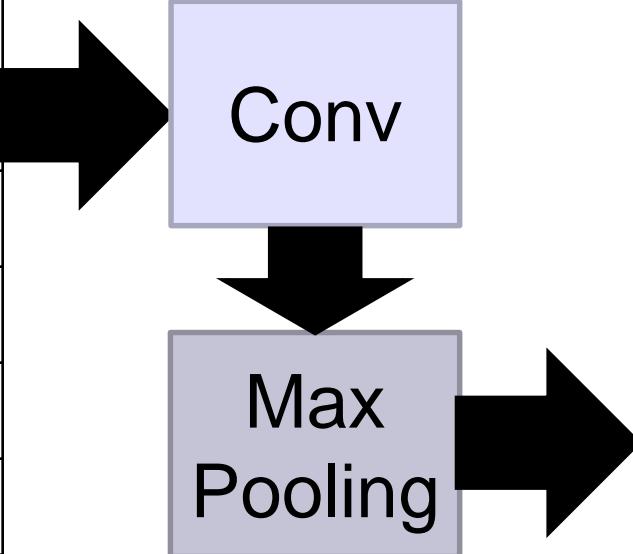
A CNN compresses a fully connected network in two ways:

- Reducing number of connections
- Shared weights on the edges
- Max pooling further reduces the complexity

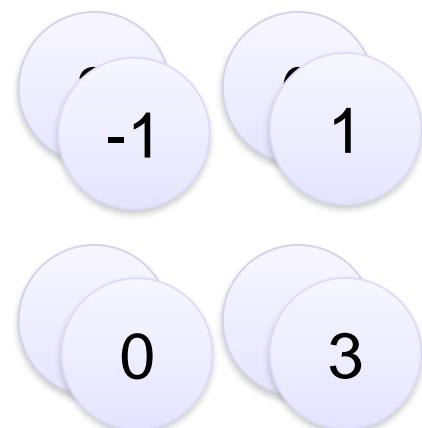
# Max Pooling

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



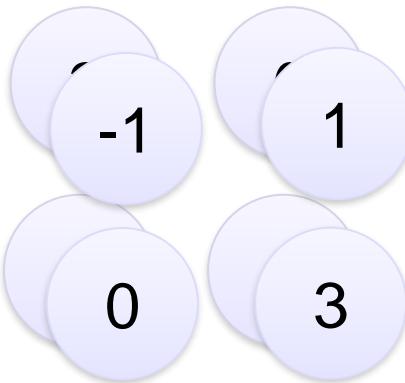
New image  
but smaller



2 x 2 image

Each filter  
is a channel

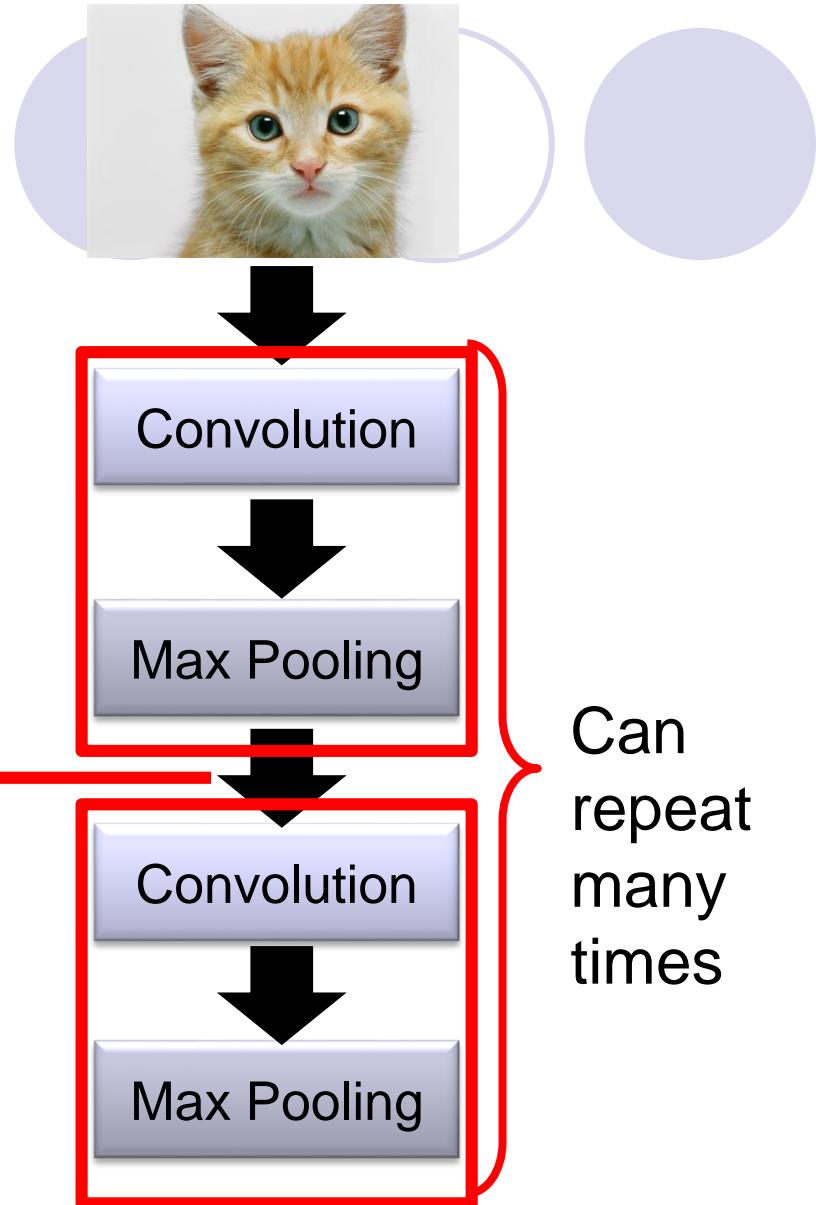
# The whole CNN



A new image

Smaller than the original image

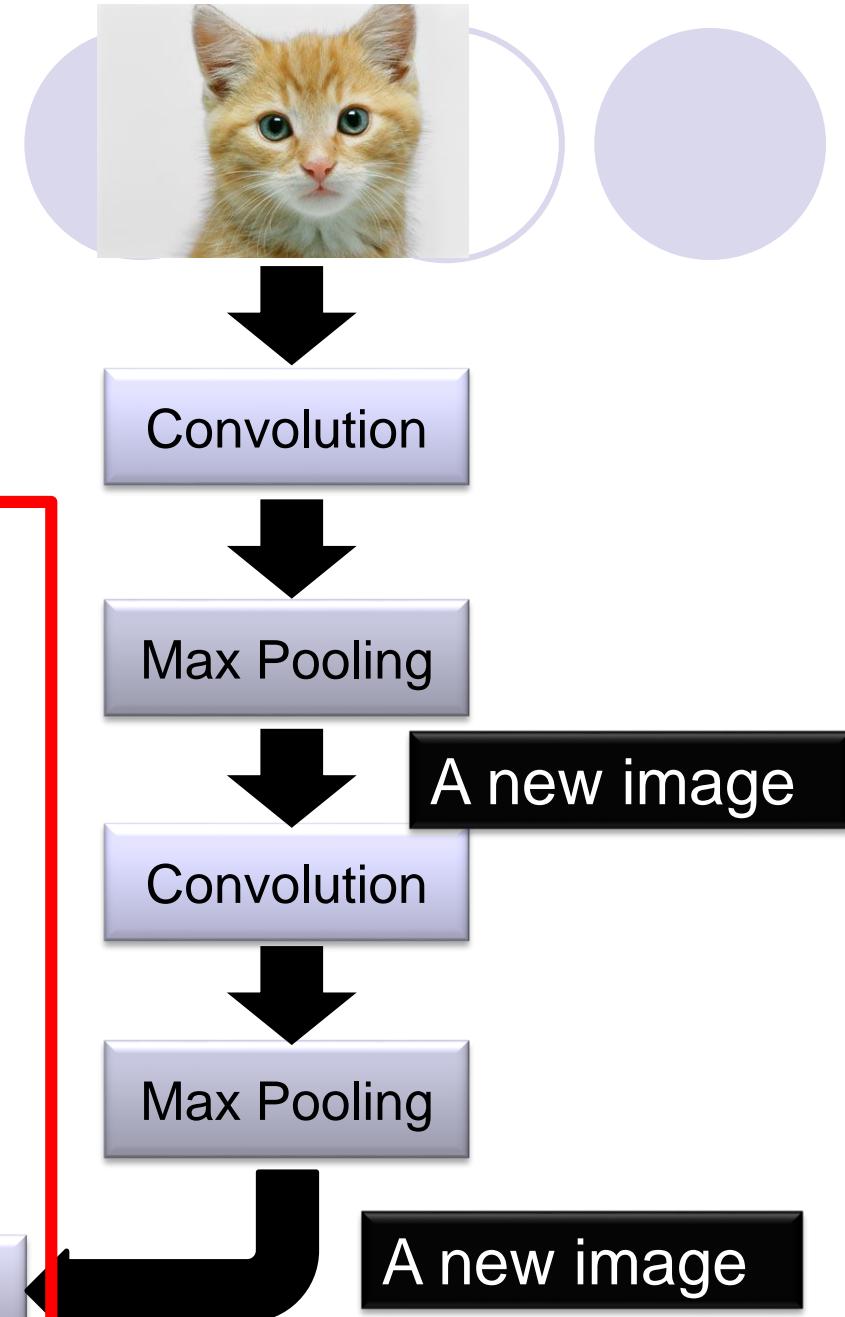
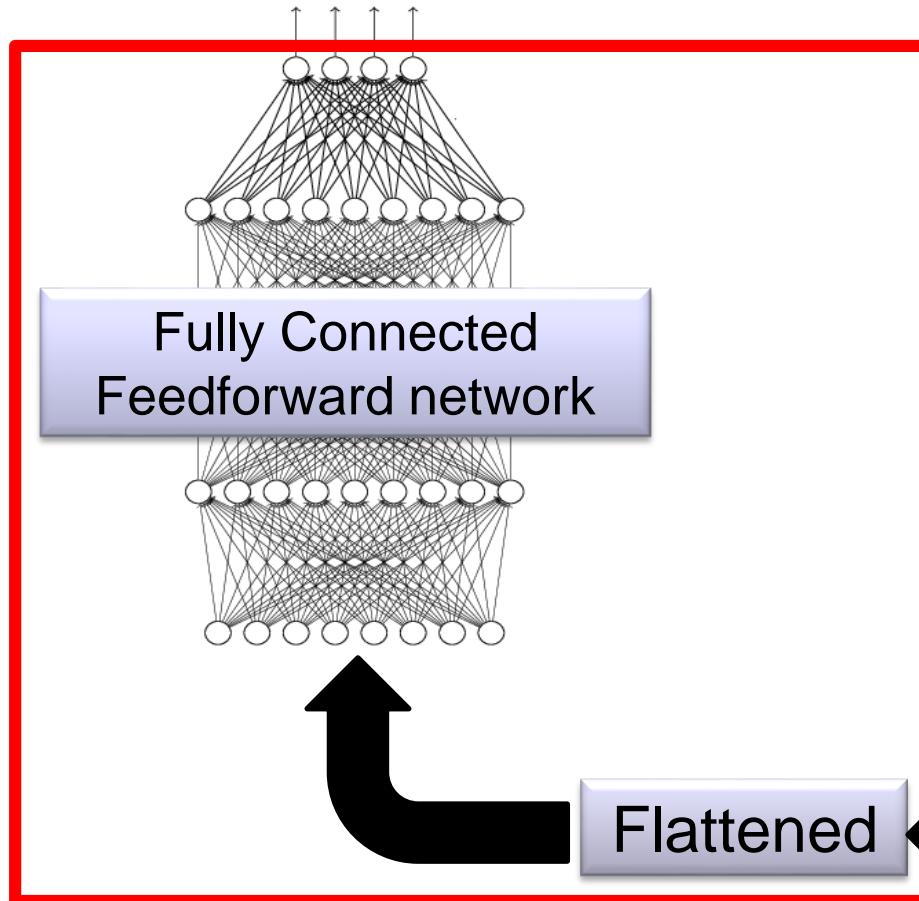
The number of channels  
is the number of filters



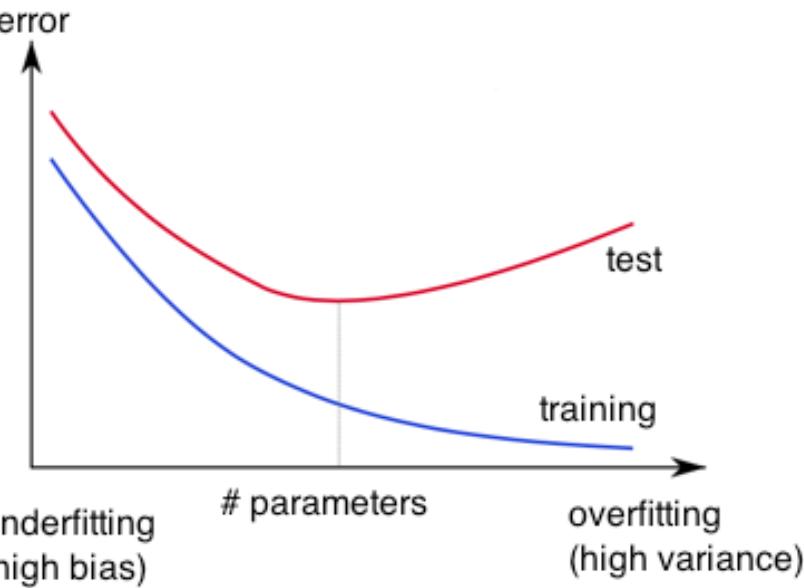
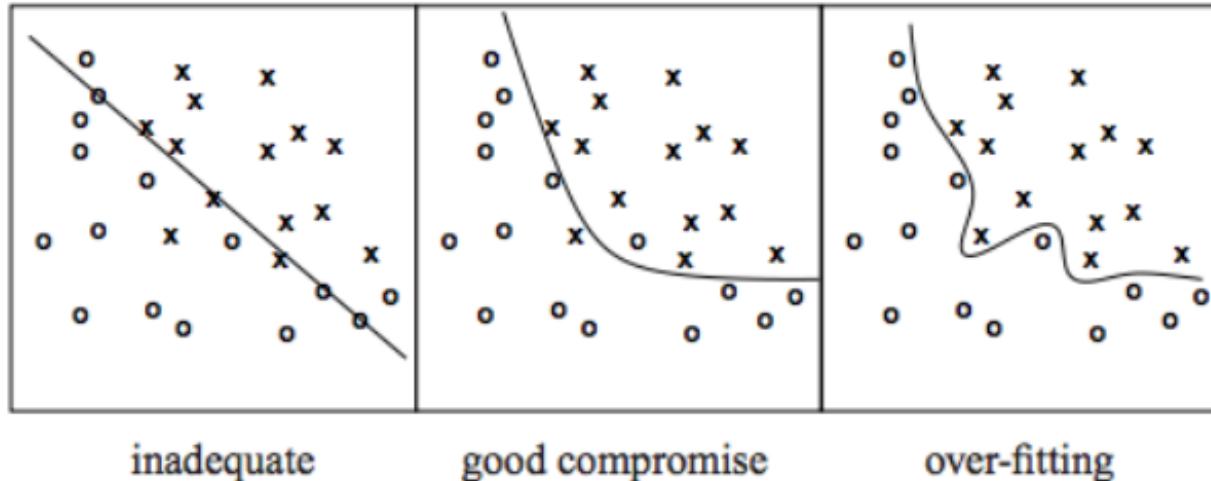
Can  
repeat  
many  
times

# The whole CNN

cat dog .....



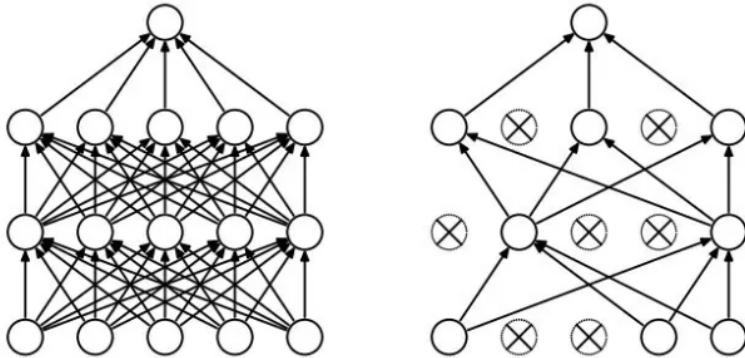
# Overfitting



<http://wiki.bethancrane.com/overfitting-of-data>

Learned hypothesis may **fit** the training data very well, even outliers (**noise**) but fail to **generalize** to new examples (test data)

# Regularization



## Dropout

- Randomly drop units (along with their connections) during training
- Each unit retained with fixed probability  $p$ , independent of other units
- **Hyper-parameter**  $p$  to be chosen (tuned)

Srivastava, Nitish, et al. ["Dropout: a simple way to prevent neural networks from overfitting."](#) Journal of machine learning research (2014)

## L2 = weight decay

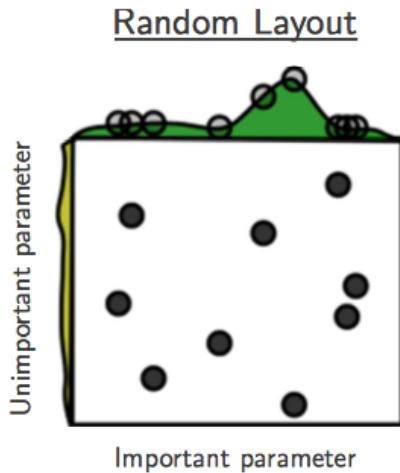
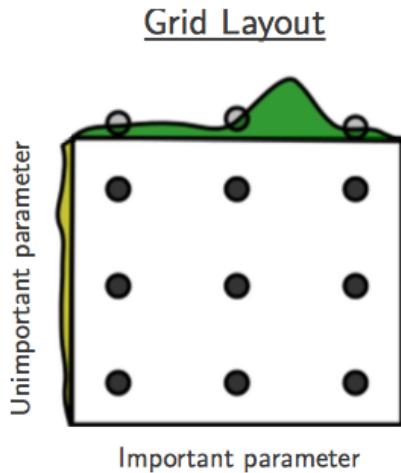
- Regularization term that penalizes big weights, added to the objective
- Weight decay value determines how dominant regularization is during gradient computation
- Big weight decay coefficient → big penalty for big weights

$$J_{reg}(\theta) = J(\theta) + \lambda \sum_k \theta_k^2$$

## Early-stopping

- Use validation error to decide when to stop training
- Stop when monitored quantity has not improved after  $n$  subsequent epochs
- $n$  is called patience

# Tuning hyper-parameters



$$g(x) \approx g(x) + h(y)$$

$g(x)$  shown in green  
 $h(y)$  is shown in yellow

Bergstra, James, and Yoshua Bengio. "[Random search for hyper-parameter optimization.](#)" *Journal of Machine Learning Research*, Feb (2012)

"Grid and random search of 9 trials for optimizing function  $g(x) \approx g(x) + h(y)$   
With grid search, nine trials only test  $g(x)$  in three distinct places.  
With random search, all nine trials explore distinct values of  $g$ . "

Both try configurations randomly and **blindly**  
Next trial is independent to all the trials done before

**Bayesian optimization for hyper-parameter tuning:**

Library available!

Make smarter choice for the next trial, minimize the number of trials

1. Collect the performance at several configurations
2. Make inference and decide what configuration to try next

# Loss functions and output

## Classification

Training examples

$R^n \times \{\text{class\_1, ..., class\_n}\}$   
(one-hot encoding)

Output Layer

Soft-max  
[map  $R^n$  to a probability distribution]

$$P(y = j | \mathbf{x}) = \frac{e^{\mathbf{x}^\top \mathbf{w}_j}}{\sum_{k=1}^K e^{\mathbf{x}^\top \mathbf{w}_k}}$$

Cost (loss) function

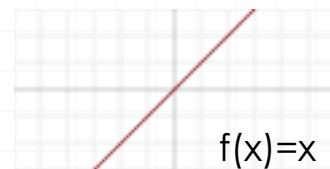
$$J(\theta) = -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K \left[ y_k^{(i)} \log \hat{y}_k^{(i)} + (1 - y_k^{(i)}) \log (1 - \hat{y}_k^{(i)}) \right]$$

List of loss functions

## Regression

$R^n \times R^m$

Linear (Identity)  
or Sigmoid



Mean Squared Error

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2$$

Mean Absolute Error

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n |y^{(i)} - \hat{y}^{(i)}|$$

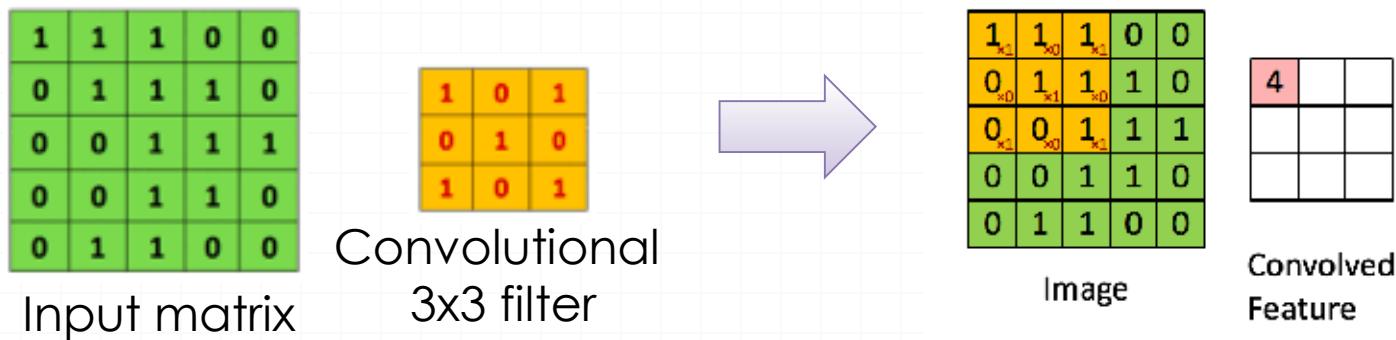
# Convolutional Neural Networks (CNNs)

Main CNN idea for text:

**Compute vectors for n-grams** and group them afterwards

Example: “this takes too long” compute vectors for:

This takes, takes too, too long, this takes too, takes too long, this takes too long



[http://deeplearning.stanford.edu/wiki/index.php/Feature\\_extraction\\_using\\_convolution](http://deeplearning.stanford.edu/wiki/index.php/Feature_extraction_using_convolution)



# D i s a d v a n t a g e s

- ▪ ▪ ▪ ▪ ▪

# Disadvantages

- ⌚ From a **memory** and **capacity** standpoint the **CNN** is not much bigger than a **regular two layer network**.
- ⌚ At runtime the **convolution** operations are computationally expensive and take up about **67%** of the time.
- ⌚ CNN's are about **3X** slower than their fully connected equivalents (size-wise).

# Disadvantages

## ⦿ Convolution operation

- ⦿ 4 nested loops ( 2 loops on input image & 2 loops on kernel)

## ⦿ Small kernel size

- ⦿ make the inner loops very inefficient as they frequently JMP.

## ⦿ Cash unfriendly memory access

- ⦿ Back-propagation require both row-wise and column-wise access to the input and kernel image.
- ⦿ 2-D Images represented in a row-wise-serialized order.
- ⦿ Column-wise access to data can result in a high rate of cash misses in memory subsystem.



# A p p l i c a t i o n



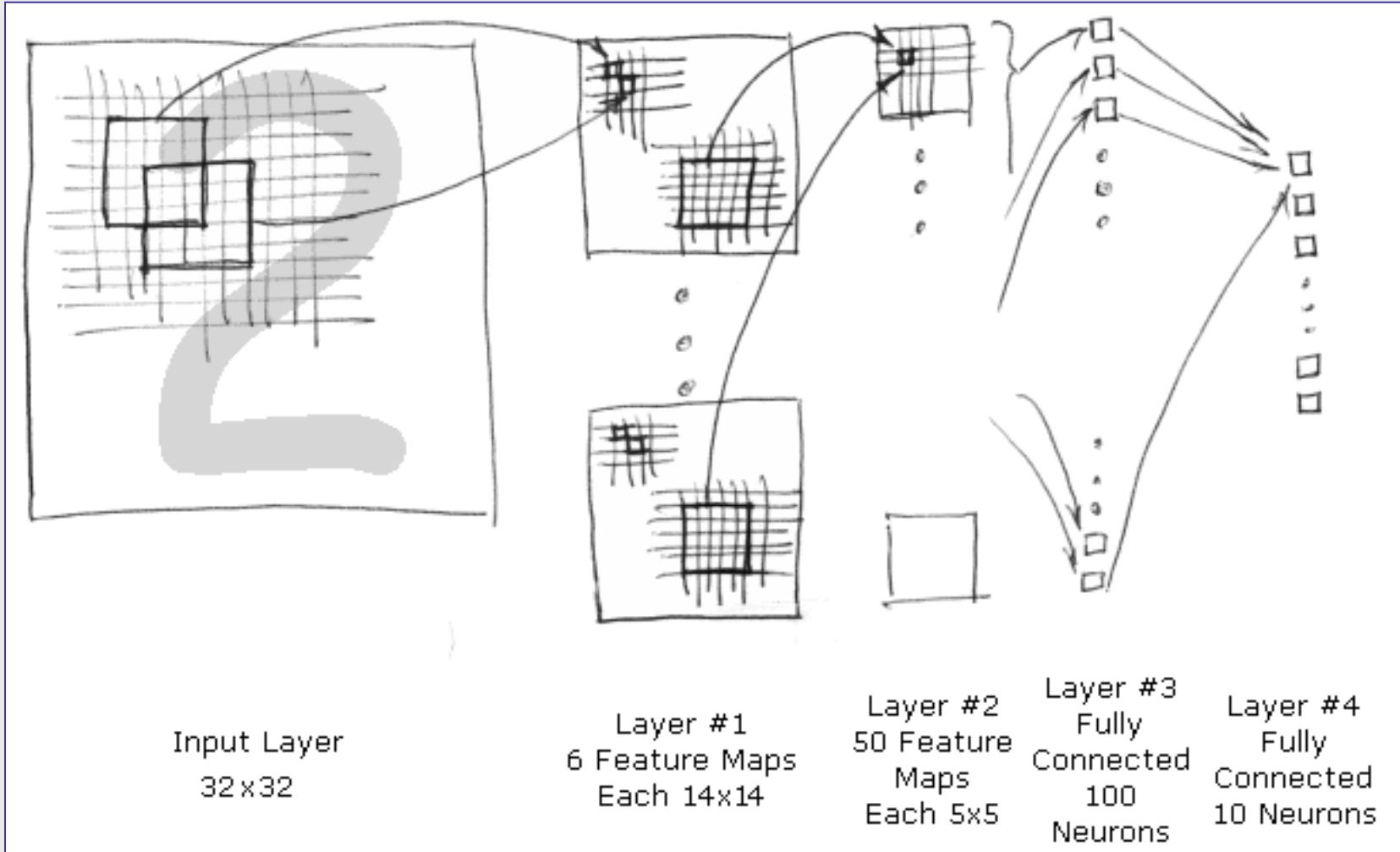
# Application



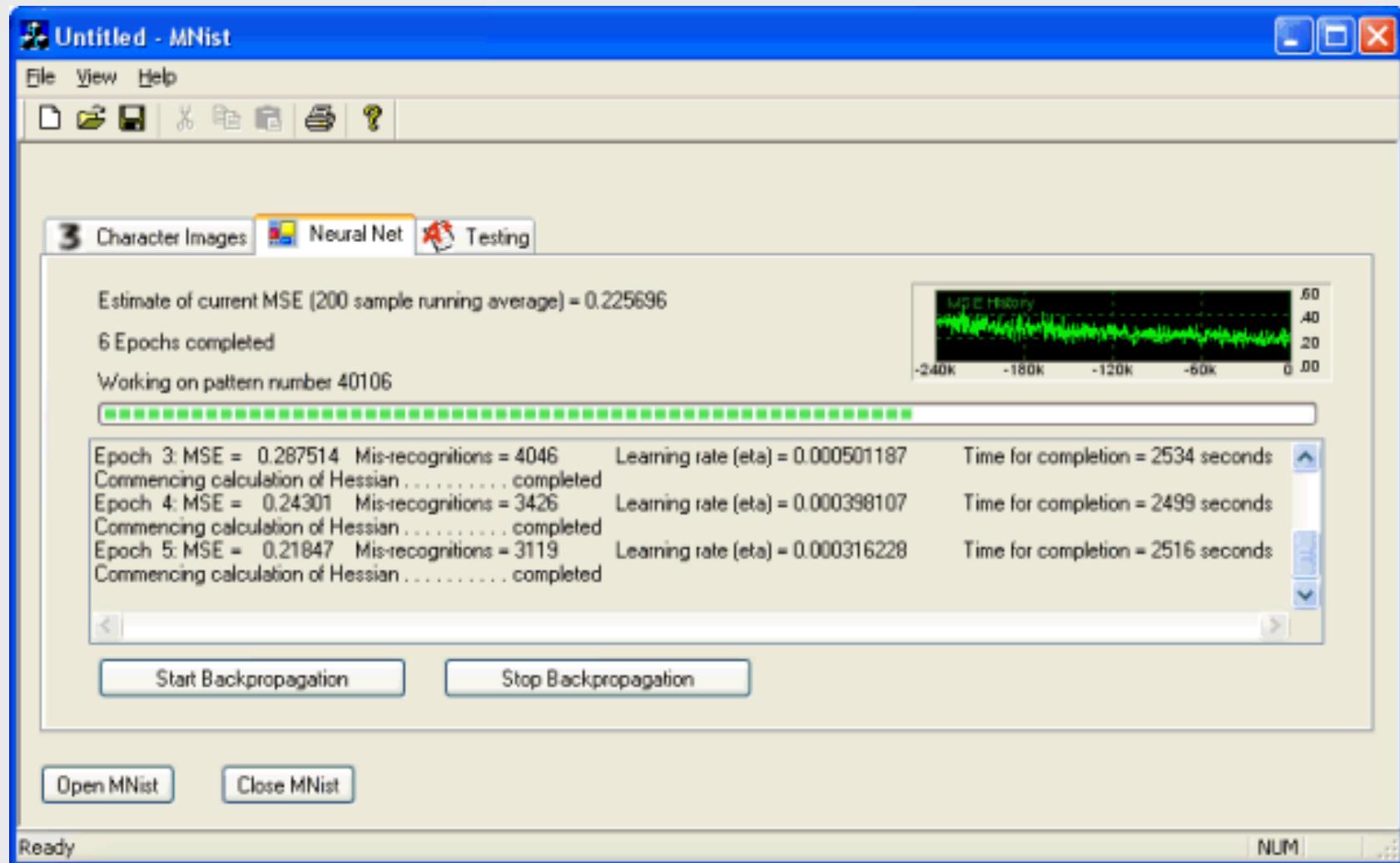
Mike O'Neill

- ➊ A Convolutional neural network achieves 99.26% accuracy on a modified NIST database of hand-written digits.
- ➋ MNIST database : Consist of 60,000 hand written digits uniformly distributed over 0-9.

# Application



# Application



# Application

				
3818	6597			
0 => 6	0 => 7			
				
2018	2182	5457		
1 => 7	1 => 3	1 => 8		
				
4176	8059	8094	9664	
2 => 7	2 => 1	2 => 8	2 => 7	
				
1681	2280	4740		
3 => 7	3 => 5	3 => 5		
				
247	2130	8520	8527	9792
4 => 6	4 => 9	4 => 9	4 => 9	4 => 9

# Application

											
340	674	1299	1737	2035	2040	2597	3558	4360	5937	9729	9770
5 => 3 5 => 3 5 => 3 5 => 3 5 => 6 5 => 3 5 => 0 5 => 3 5 => 3 5 => 6 5 => 0											
											
2135	2654	3365	3422	3762	4699	4838	6558	8287	9627	9679	9698
6 => 1 6 => 1 6 => 1 6 => 0 6 => 8 6 => 1 6 => 5 6 => 3 6 => 8 6 => 5 6 => 5 6 => 2											
											
282	1226	3225	3808	9009	9015	9024					
7 => 3 7 => 2 7 => 9 7 => 2 7 => 2 7 => 2 7 => 2											
											
184	582	947	1033	1068	1319	1782	1878	4497	4879	4956	6555
8 => 3 8 => 2 8 => 9 8 => 1 8 => 4 8 => 0 8 => 9 8 => 3 8 => 7 8 => 6 8 => 4 8 => 9 8 => 5											
											
1247	1709	1901	2582	2939	3503	3850	3869	4369	4761	6571	6632
9 => 5 9 => 5 9 => 4 9 => 7 9 => 5 9 => 1 9 => 4 9 => 4 9 => 4 9 => 8 9 => 7 9 => 8 9 => 8											

# References

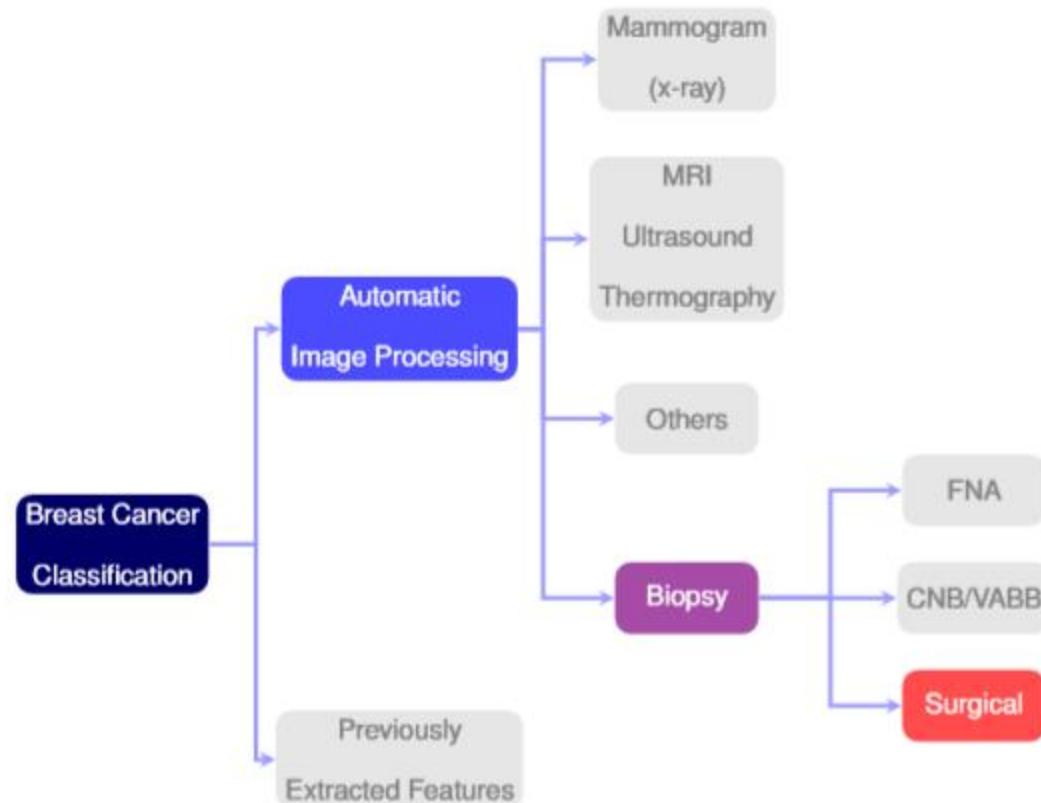
- [1].Y. LeCun and Y. Bengio.“**Convolutional networks for images, speech, and time-series.**” In M. A. Arbib, editor, *The Handbook of Brain Theory and Neural Networks*. MIT Press, 1995.
- [2].Fabien Lauer, ChingY. Suen, Gérard Bloch,”**A trainable feature extractor for handwritten digit recognition**”,Elsevier, october 2006.
- [3].Patrice Y. Simard, Dave Steinkraus, John Platt, "**Best Practices for Convolutional Neural Networks Applied to Visual Document Analysis,**" International Conference on Document Analysis and Recognition (ICDAR), IEEE Computer Society, Los Alamitos, pp. 958-962, 2003.

# Questions





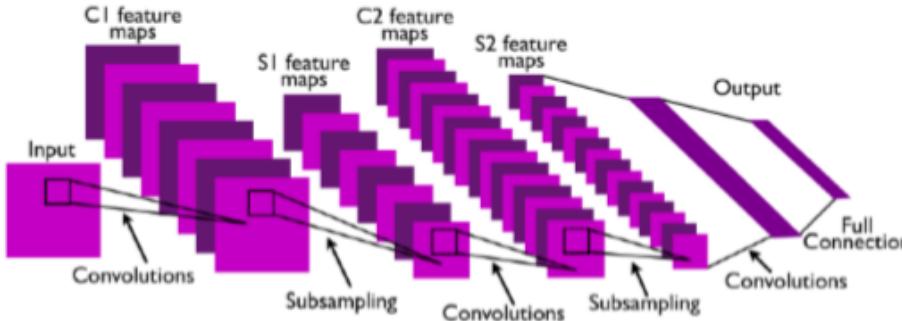
Pairs of blurred license plates and their  
respective reconstructions computed by CNN.



**Figure 3.1:** Categories of related works on breast cancer classification.

**Source:** The author (2015).

There are three main types of layers used to build CNN architectures: *convolutional layer*, *pooling layer*, and *fully-connected layer*. Normally, a full CNN architecture is obtained by stacking several of these layers. An example of typical CNN architecture with two feature stages is shown in Figure 2.8.

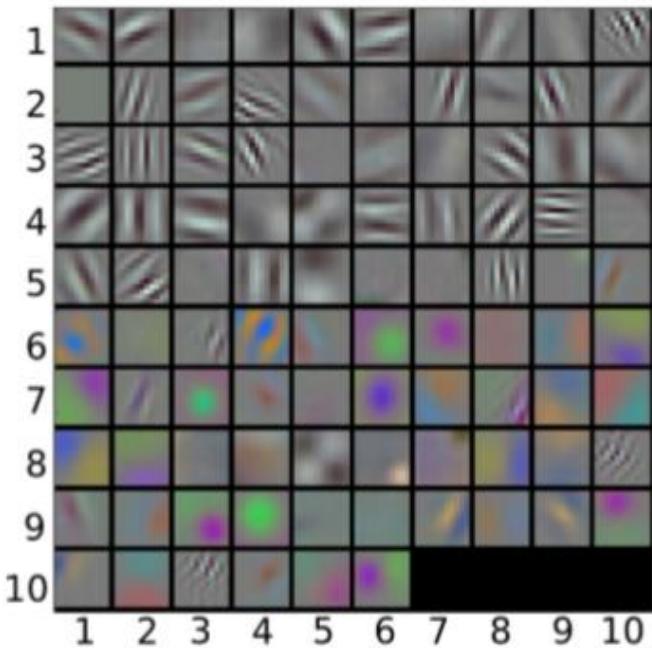


**Figure 2.8:** An example of a typical CNN architecture with two feature stages.

**Source:** Extracted from [168].

In a CNN, the key computation is the convolution of a feature detector with an input signal. Convolutional layer computes the output of neurons connected to local regions in the input, each one computing a dot product between their weights and the region they are connected to in the input volume. The set of weights which is convolved with the input is called **filter** or **kernel**. Every filter is spatially small (width and height), but extends through the full depth of the input volume. For inputs such as images, typical filters are small areas (e.g.,  $3 \times 3$ ,  $5 \times 5$ , or  $8 \times 8$ ) and each neuron is only connected to this area in the previous layer. The weights are shared across neurons, leading the filters to learn frequent patterns that occur in any part of the image. The distance between the applications of filters are called **stride**. If the stride hyperparameter is smaller than the filter size, the convolution is applied in overlapping windows.

Convolution with a collection of filters, like the learned filters (also named feature maps or activation maps) in Figure 2.9, improves the representation: at the first layer of a CNN the features go from individual pixels to simple primitives, like horizontal and vertical lines, circles, and patches of color. We can see that the model learns the filters for horizontal and vertical edges, and also learns the filters that resemble Gabor filters (edge detectors) [36, 98]. In contrast to conventional single-channel image processing filters, these CNN filters are computed across all of the input channels. Due to its translation-invariant property, convolutional filters yield a high response wherever a feature is detected.



**Figure 2.9:** The first layer of learned convolutional filters in CaffeNet. This Caffe [141] reference model is based on AlexNet [158]. These filters are tuned to edges of different orientations, frequency, phase and colors. The filter outputs expand the dimensionality of the visual representation from the three color channels of the image to these 96 primitives. Deeper layers further enrich the representation.

**Source:** The author (2015).

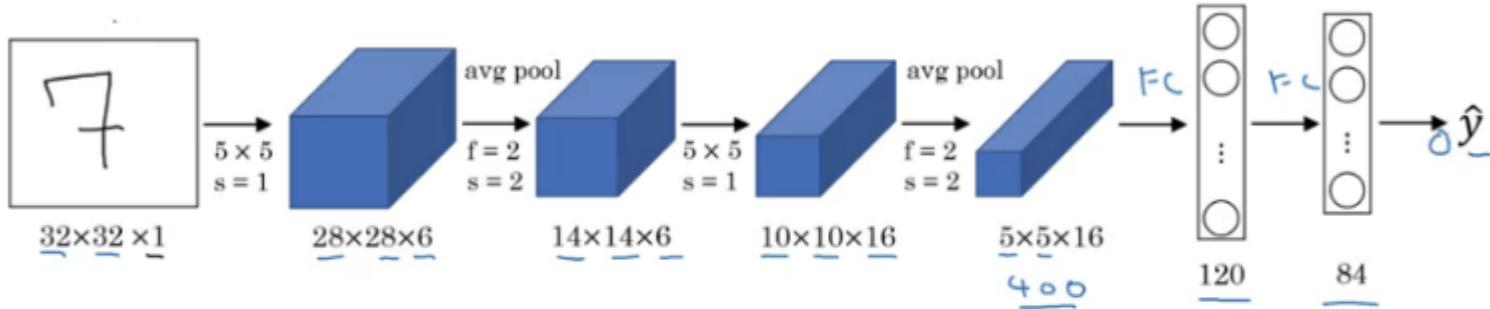
When different features can be learned in different spatial locations the weights sharing scheme of the convolutional layer is relaxed and this layer is then referred to as *locally-connected layer*.

The insertion of a pooling (subsampling) layer between two successive convolutional layers is common. The main objective of this practice is to progressively reduce the spatial size of the representation. Thus, reducing the number of parameters and computations required by the network helps the overfitting control. The pooling layer downsamples the volume spatially in each depth slice of the input volume independently. Thus, the pool operator resizes the input along the width and the height, discarding the activations. In practice, the *max pooling* function, which applies a window function to the input patch, and computes the maximum in that neighborhood, has shown better results [248]. Furthermore, the pooling units can perform other functions like *L2-norm pooling* or *average pooling*.

In a fully-connected layer, neurons have full connections to all activations in the previous layer and their activations can be computed using a matrix multiplication followed by a bias offset. This type of layer is standard in a regular NN. The last fully-connected layer holds the net output, such as probability distributions over classes [90, 158].

## LeNet-5

Let's start with LeNet-5:

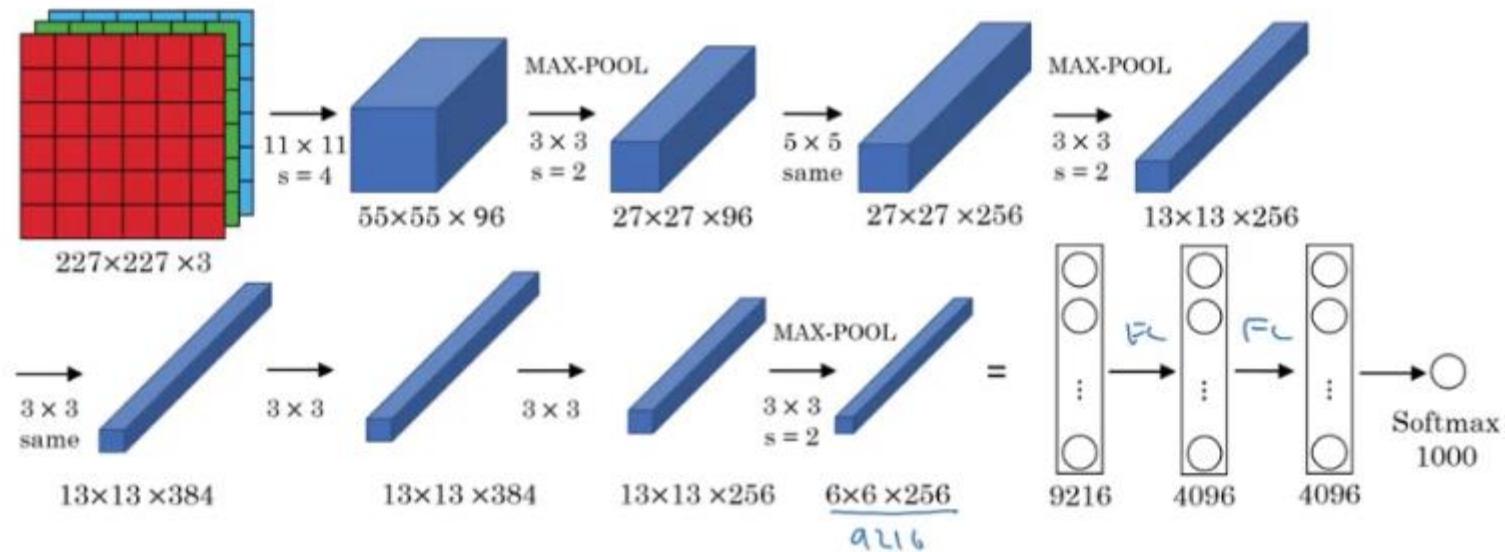


It takes a grayscale image as input. Once we pass it through a combination of convolution and pooling layers, the output will be passed through fully connected layers and classified into corresponding classes. The total number of parameters in LeNet-5 are:

- **Parameters:** 60k
- **Layers flow:** Conv -> Pool -> Conv -> Pool -> FC -> FC -> Output
- **Activation functions:** Sigmoid/tanh and ReLu

## AlexNet

An illustrated summary of AlexNet is given below:



This network is similar to LeNet-5 with just more convolution and pooling layers:

- **Parameters:** 60 million
- **Activation function:** ReLu

