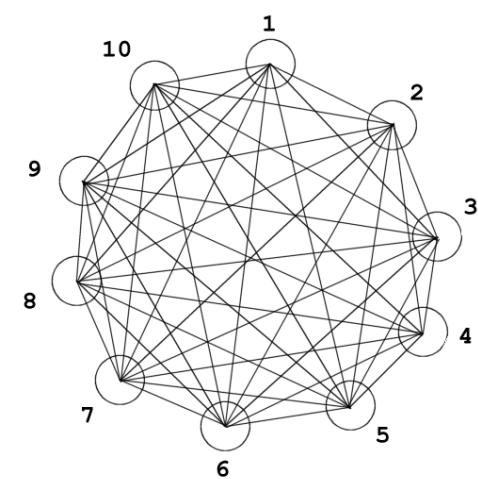
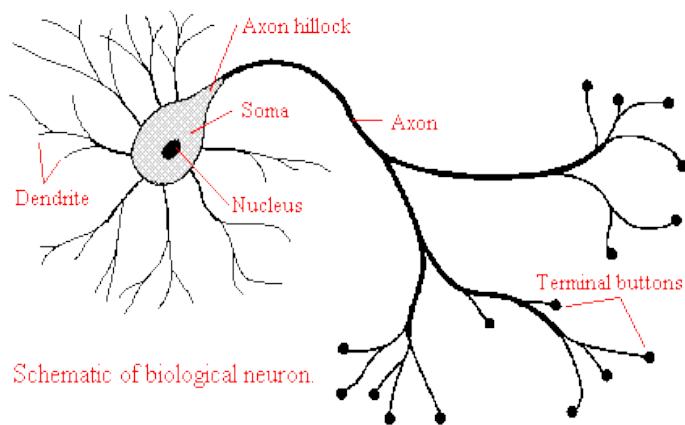


Artificial Neural Network

Anna Helena Reali Costa

PCS



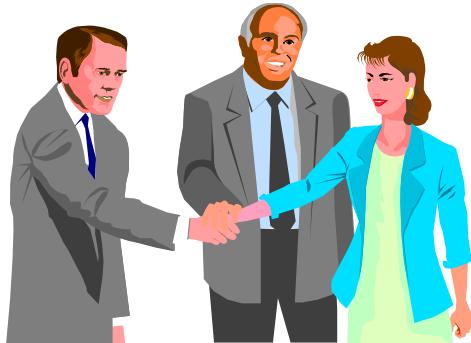
ANN

- A different style of computation: **parallel distributed processing**
- A **universal computational architecture**: the same structure carries out many different functions
- It can learn new knowledge, therefore it is **adaptive**

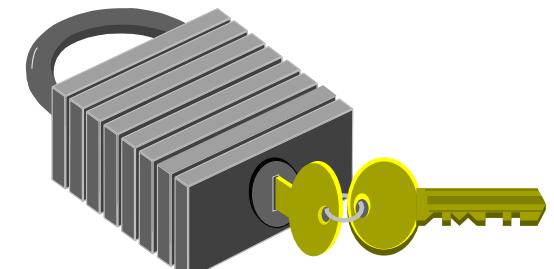
History

- McCulloch and Pitts introduced the artificial neuron in 1943.
 - Simplified model of a biological neuron
- Fell out of favor in the late 1960's
 - Perceptron limitations (Minsky and Papert)
- Resurgence in the mid 1980's
 - Nonlinear Neuron Functions
 - Backpropagation training (Werbos)
- Currently resounding success with Deep NN

Applications

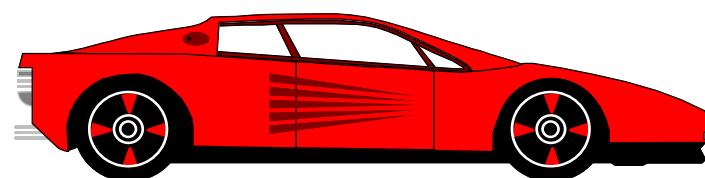


Face / Speech
Recognition



User Authentication.

Identification of military targets:
B-52, Boeing 747, Space Shuttle



Autonomous Vehicle Navigation



Oil exploration:
Lithology, etc.

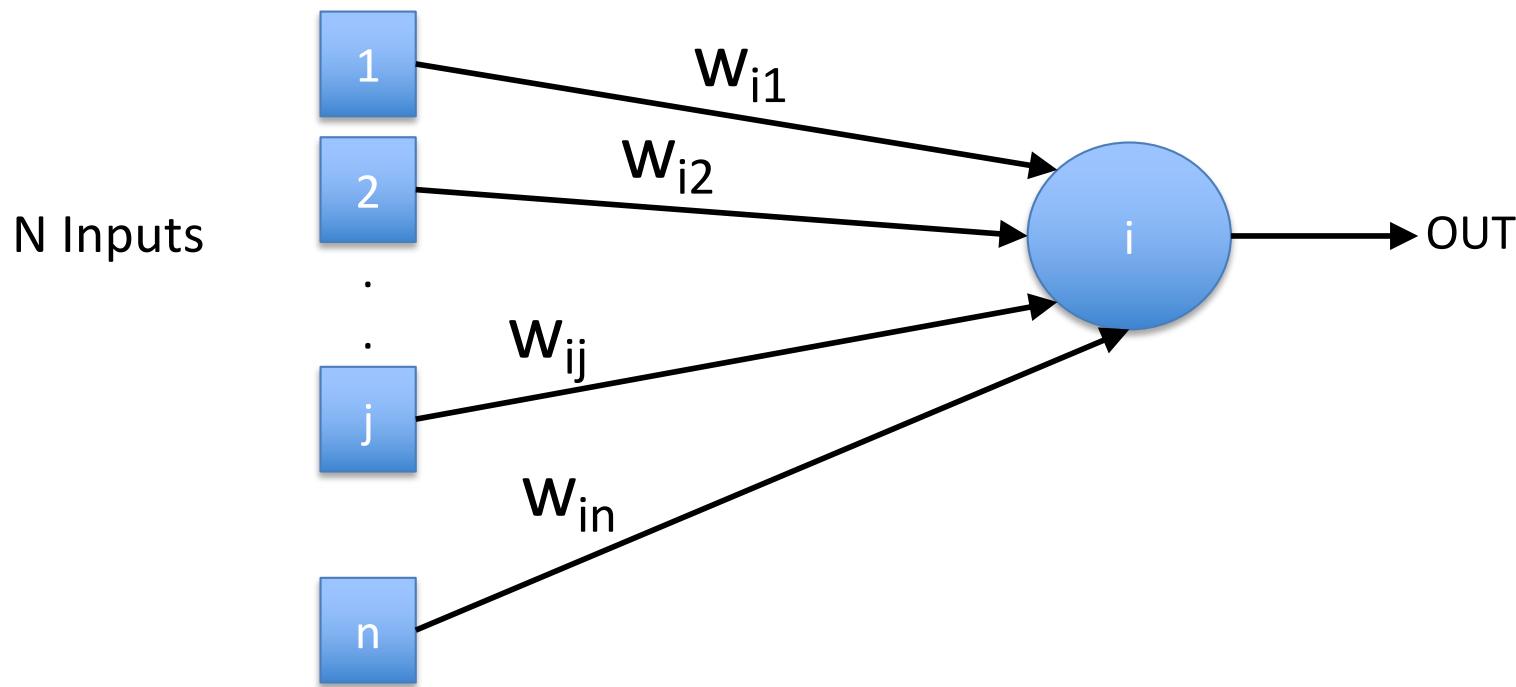


Prediction on
Financial Market

Design of an ANN

- The design of an ANN involves determining the following elements:
 - Neurons and **activation function**.
 - Connections and arrangement of neurons: **topology** (network architecture).
 - Synaptic **weights**: values (in the case of learned weights) or a **training algorithm** to be used and its parameters.
 - **Recall**: procedure to be used for the network to calculate the outputs for given new inputs
- Unfortunately, there is no single "recipe" ...

Perceptron (Frank Rosenblatt, 1958)



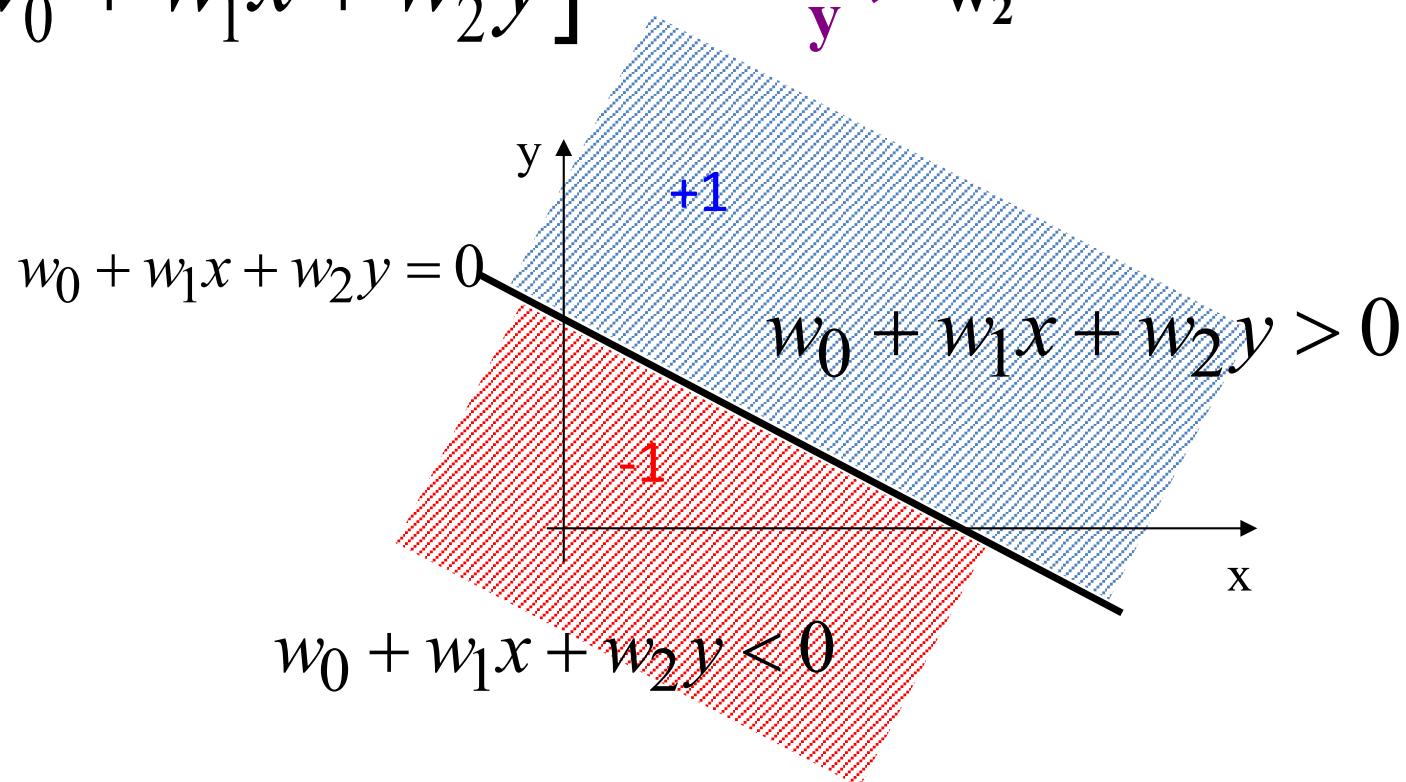
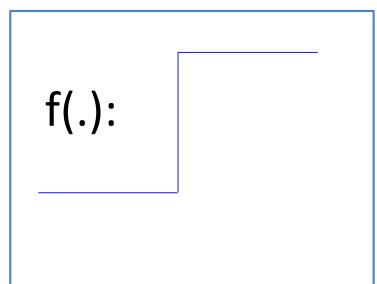
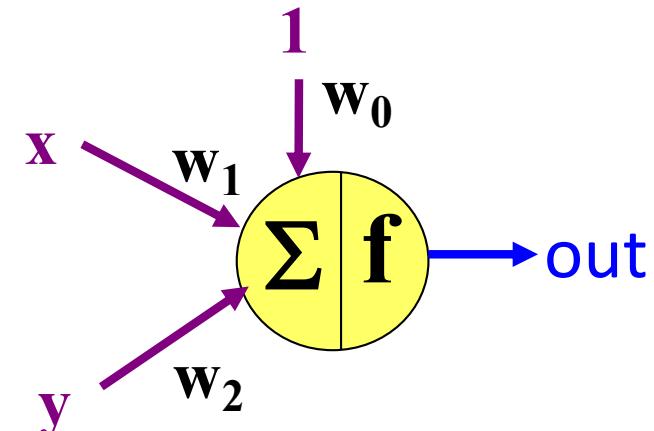
$$OUT = f(net) = \begin{cases} 1 : net > 0 \\ -1 : net \leq 0 \end{cases} \quad \text{with } net = \sum_j w_{ij} In_j$$

↑
Activation function

Perceptron: example

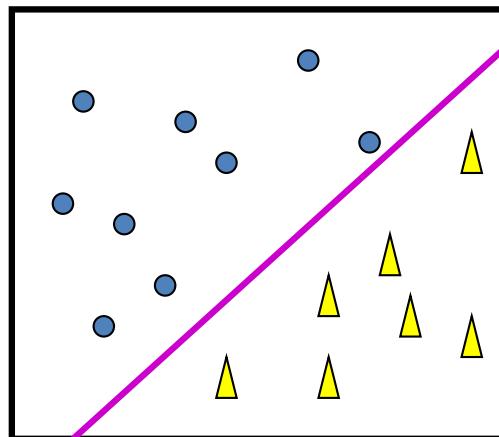
- Two inputs: x, y
- Two weights: w_1, w_2
- Bias input: w_0

$$out = f[w_0 + w_1x + w_2y]$$



Representational power of Perceptrons

Perceptron can represent any **linearly separable** function.

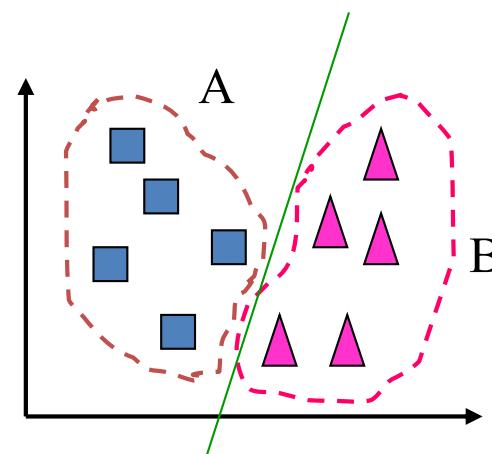


It represents a hyperplane decision surface in the n-dimensional space of instances.

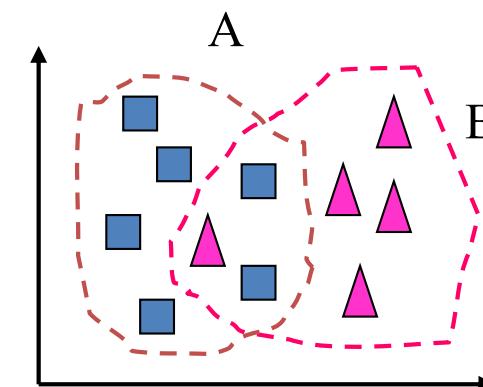
1 neuron: 2 classes +1 and -1

2 neurons: 4 classes +1+1, +1-1, -1+1, -1-1

N neurons: 2^N classes



← OK
Not OK →



Exercise

Table 1

- Design a perceptron (Figure 1) which calculates the logical implication function $y = x_1 \rightarrow x_2$, described in Table 1.

x_1	x_2	y
0	0	1
0	1	1
1	0	0
1	1	1

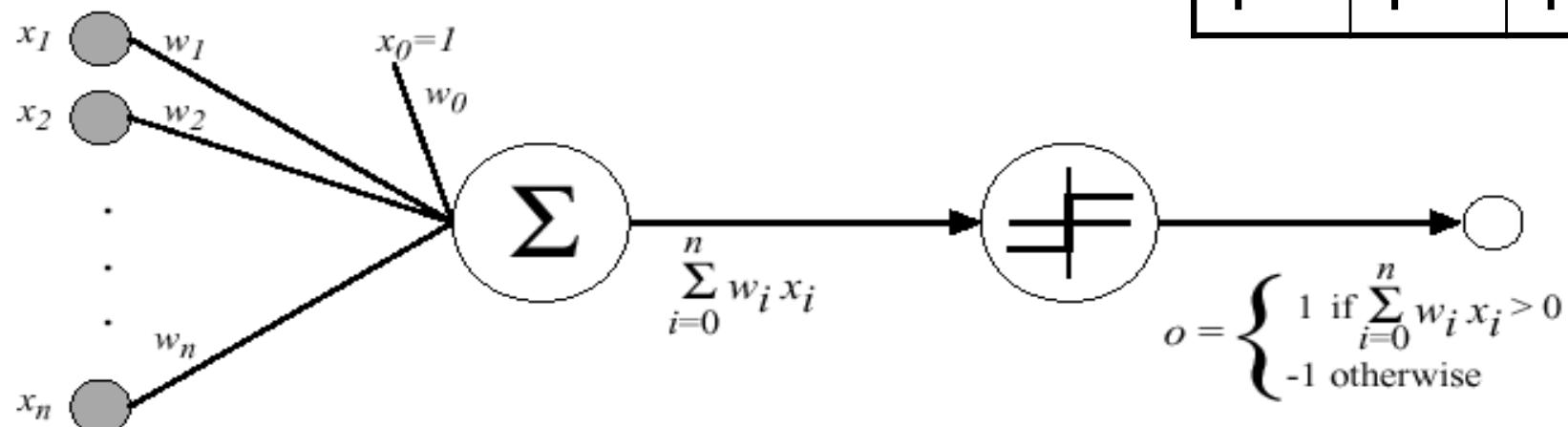


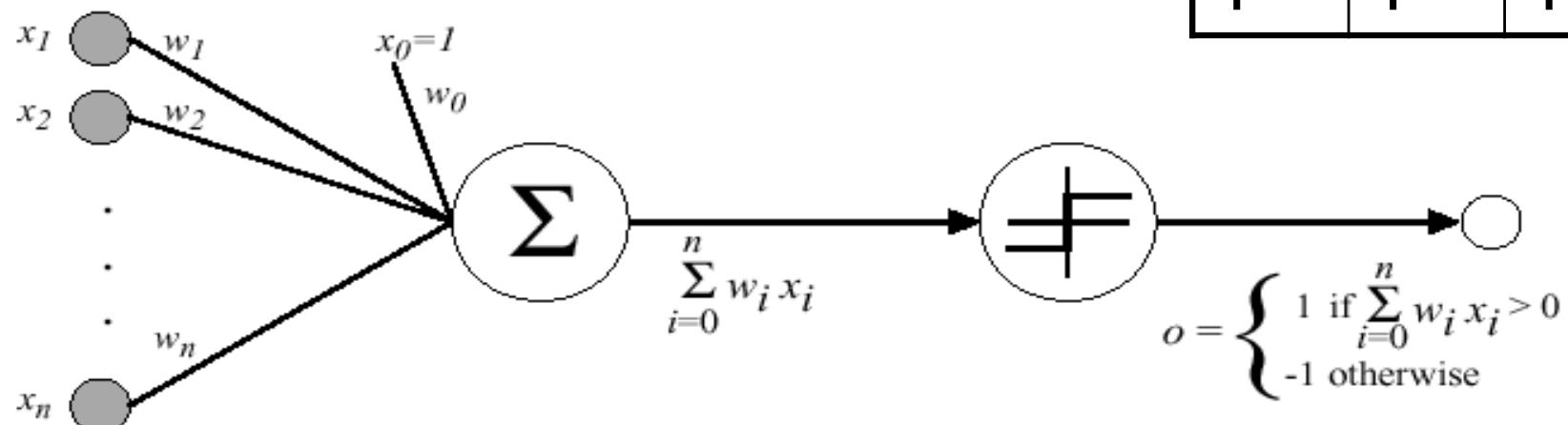
Figure 1 – Perceptron.

Exercise

Table 1

- Design a perceptron (Figure 1) which calculates the logical implication function $y = x_1 \rightarrow x_2$, described in Table 1.

x_1	x_2	y
0	0	1
0	1	1
1	0	0
1	1	1



**Uma resposta: $w_0 = 1, w_1 = -2, w_2 = 2$
(0 representa NL0, 1 representa NL1)**

The Perceptron training rule

Delta Rule

- Let t be the target output for the current training example, o be the output generated by the perceptron, and η be a positive constant called the **learning rate**:

Delta Rule:

$$w_i \leftarrow w_i + \Delta w_i$$

and

$$\Delta w_i = \eta(t - o)x_i$$

The Perceptron training algorithm

- Initialize weight matrix \mathbf{w}
- For each training sample i , with $(\mathbf{x}_i, s(\mathbf{x}_i)) = (\mathbf{x}_i, \mathbf{t}_i)$:
 - Use current \mathbf{w} to calculate $s'(\mathbf{x}_i) = \mathbf{o}_i$
 - If $|\mathbf{t}_i - \mathbf{o}_i| > \varepsilon$ then update:
$$\mathbf{w} \leftarrow \mathbf{w} + \eta \cdot (\mathbf{t}_i - \mathbf{o}_i) \cdot \mathbf{x}_i$$
- Stop when $\varepsilon \geq |\mathbf{t}_i - \mathbf{o}_i|$ for all samples $(\mathbf{x}_i, \mathbf{t}_i)$

If $|t_i - o_i| = |\delta| \leq \varepsilon \rightarrow \mathbf{w}$ does not change

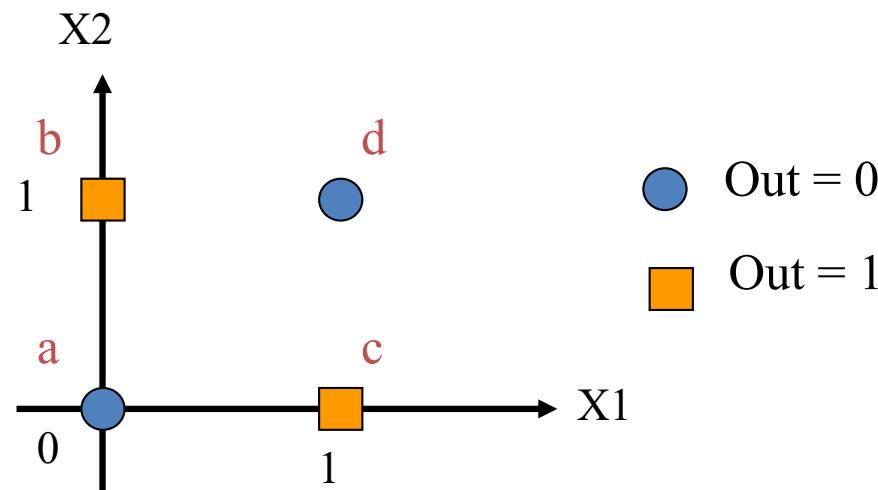
If $\delta > 0$ and $|\delta| > \varepsilon \rightarrow \mathbf{w}$ increases (because o_i is very small)

If $\delta < 0$ and $|\delta| > \varepsilon \rightarrow \mathbf{w}$ decreases (because o_i is very large)

Limitations of Perceptron with a single layer

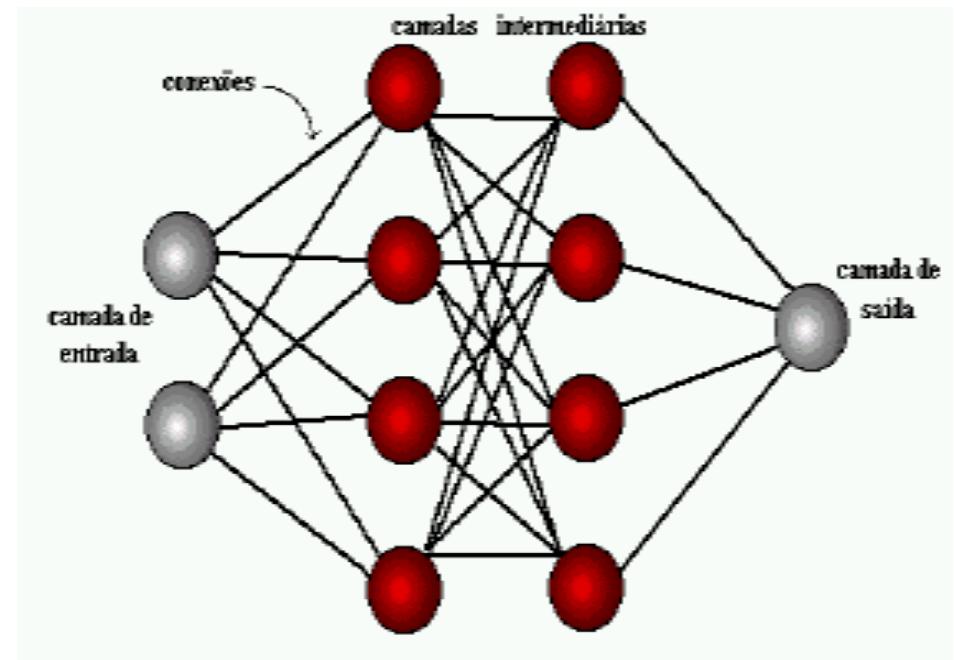
- Unfortunately, various functions of interest are not linearly separable.
- For example, the Perceptron can not represent XOR (exclusive OR).

	X1	X2	XOR
a	0	0	0
b	0	1	1
c	1	0	1
d	1	1	0

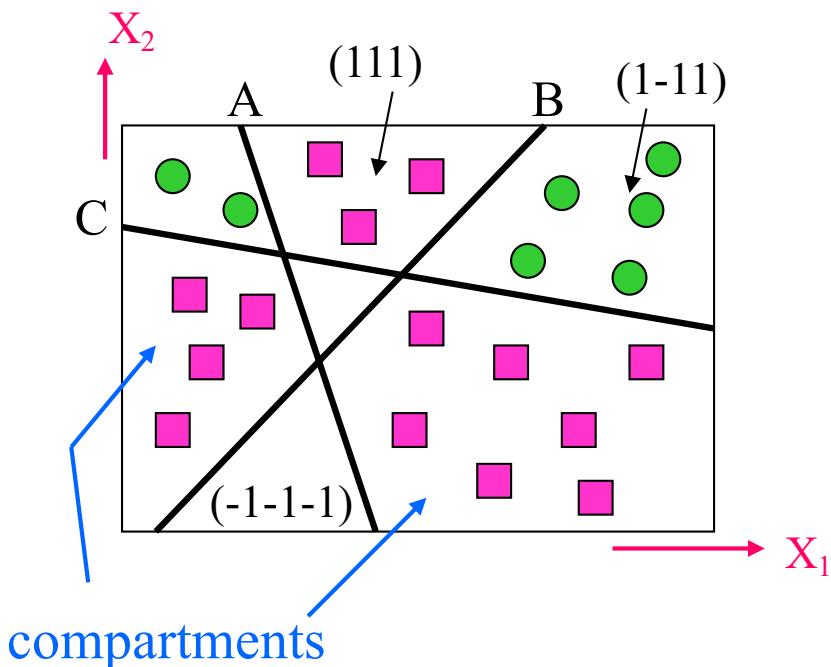


MLP: Multilayer Perceptrons

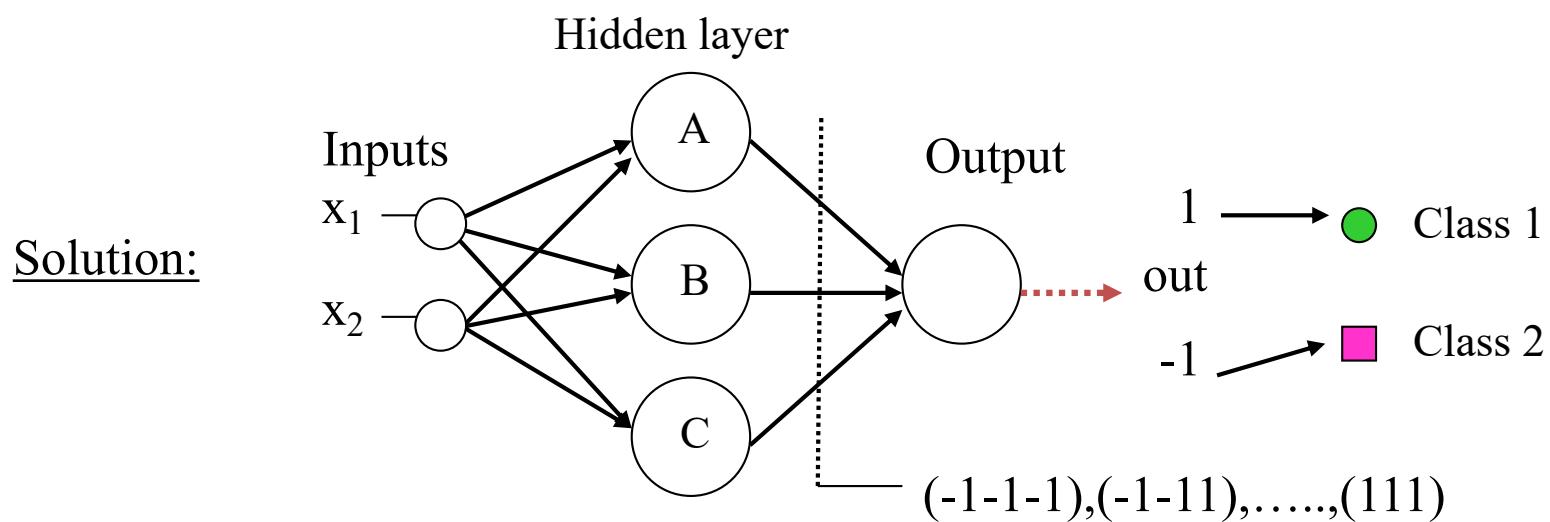
- Sets of perceptrons arranged in several layers.
- At least one hidden layer.
- An intermediate layer is sufficient to approximate any continuous function.
- Two intermediate layers are sufficient to approximate any mathematical function.



MLP: Multi-layer Perceptrons



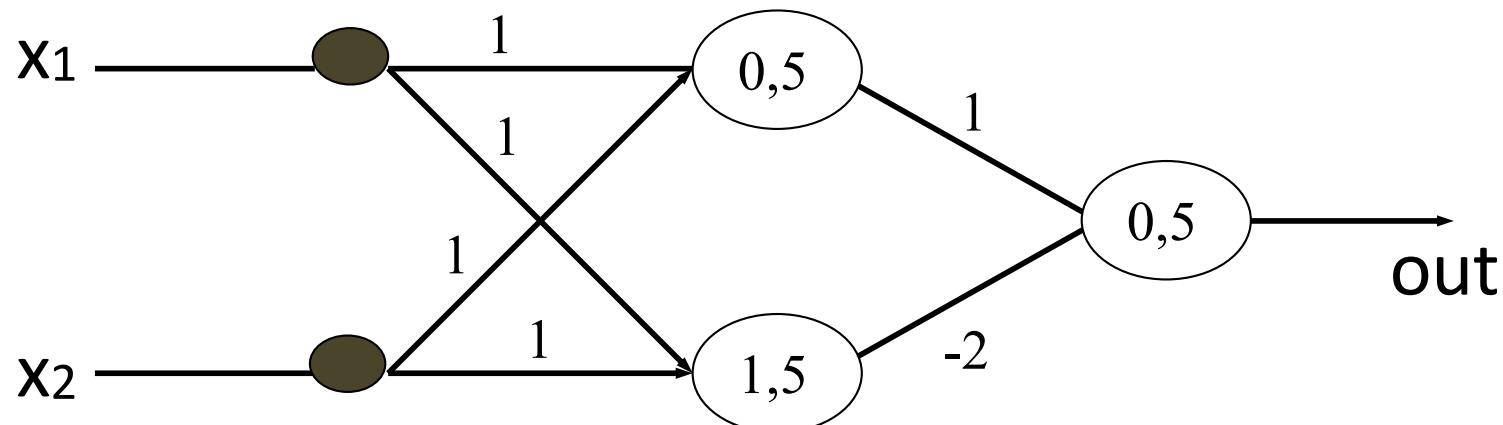
Classification requires 3 straight lines that create 7 compartments and two decision regions: one for ■ and one for ●



Example of an MLP

- $x_1 = x_2$ = binary inputs.
- $w_1 = w_2 = w_3 = w_4 = w_5 = 1$ e $w_6 = -2$
- $f_1(x \cdot w) = 1$ if the activation level ≥ 0.5 ,
0 otherwise.
- $f_2(x \cdot w) = 1$ if the activation level ≥ 1.5 ,
0 otherwise.

x_1	x_2	Out
0	0	
0	1	
1	0	
1	1	

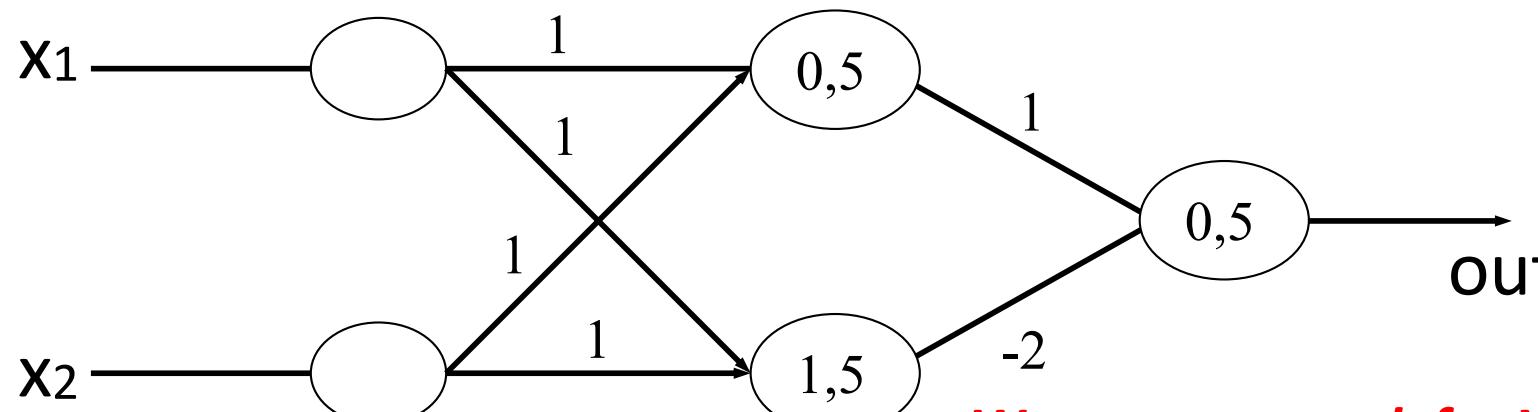


Make the recall of this network, filling the table.

Example of an MLP

- $x_1 = x_2$ = binary inputs.
- $w_1 = w_2 = w_3 = w_4 = w_5 = 1$ e $w_6 = -2$
- $f_1(x \cdot w) = 1$ if the activation level ≥ 0.5 ,
0 otherwise.
- $f_2(x \cdot w) = 1$ if the activation level ≥ 1.5 ,
0 otherwise.

x_1	x_2	Out
0	0	0
0	1	1
1	0	1
1	1	0

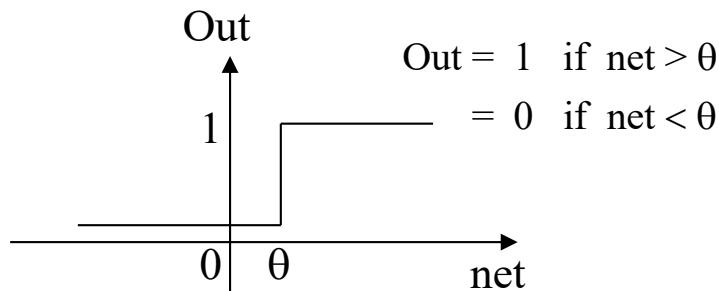


We got a network for XOR !!!

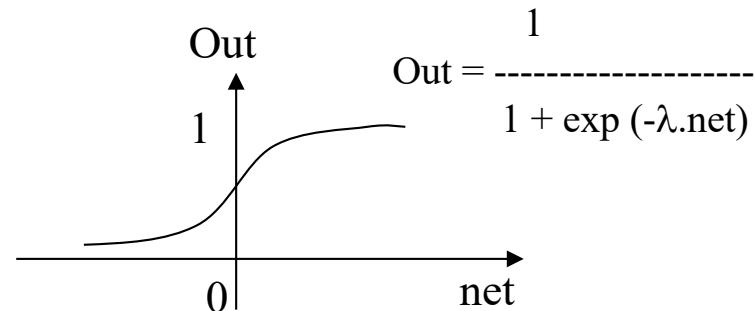
We now need an algorithm to train the weights ...

Types of Activation Function

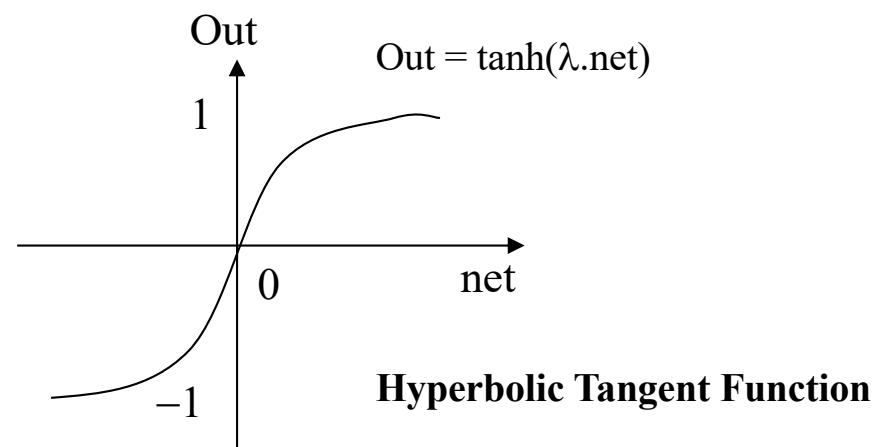
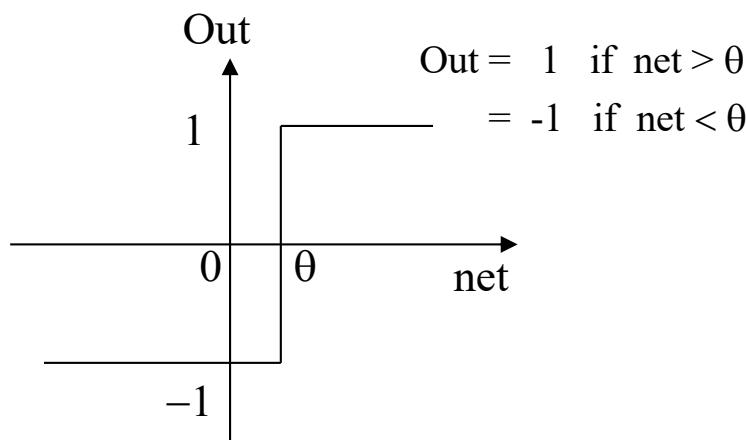
- **Linear-Threshold functions:**



- **Sigmoid functions:**



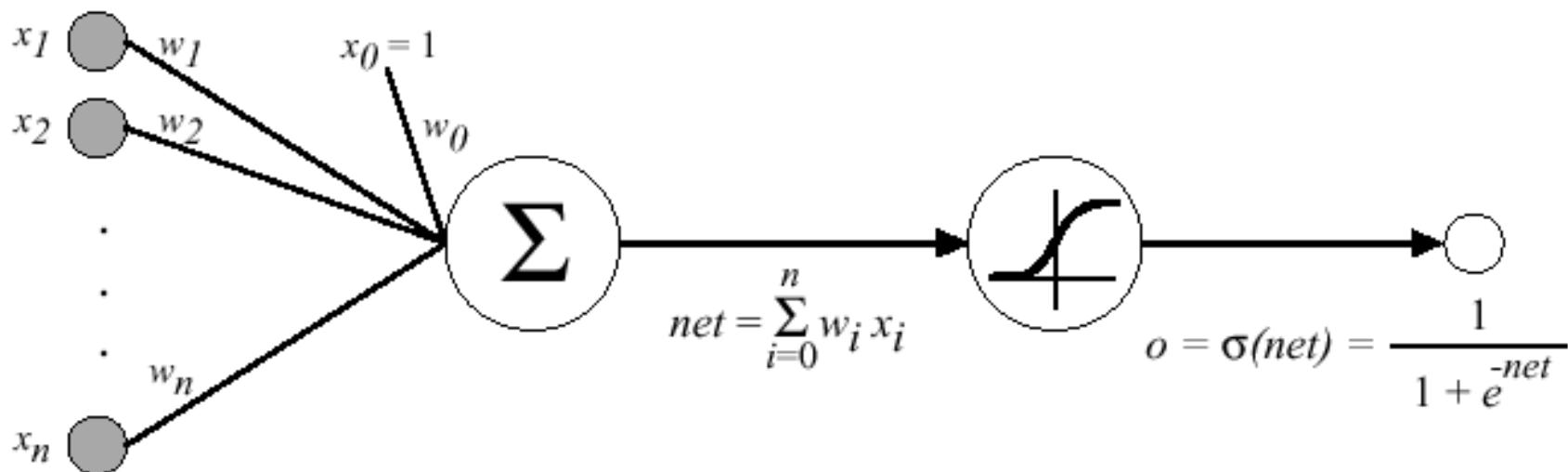
Logistics function



Hyperbolic Tangent Function

MLP / Sigmoid

$$\sigma(\text{net}) = \frac{1}{1 + e^{-\text{net}}} \quad \text{or} \quad \sigma(\text{net}) = \tanh(\text{net})$$



$$\sigma(x) = \frac{1}{1 + e^{-x}} \Rightarrow \sigma'(x) = \sigma(x)(1 - \sigma(x))$$

$$\sigma(x) = \tanh(x) \Rightarrow \sigma'(x) = \frac{(1 - \sigma^2(x))}{2}$$

derivative function

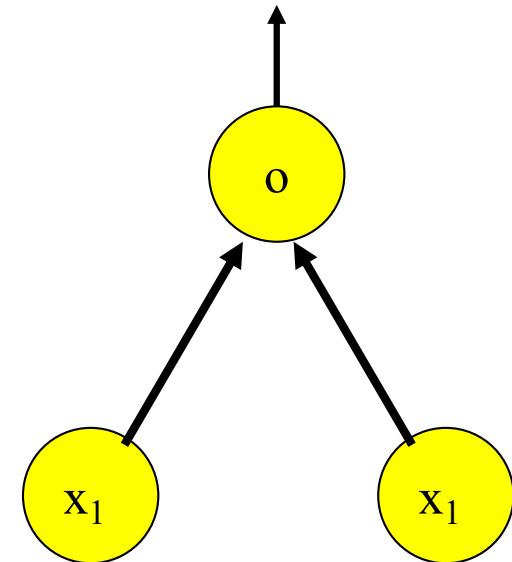
**Interesting
property ...**

MLP training

- **Key idea:** use of **gradient descent** to search the hypothesis space of possible weight vectors to find the weights that best fit the training examples.
- Learn w_i 's that minimize **squared error**:

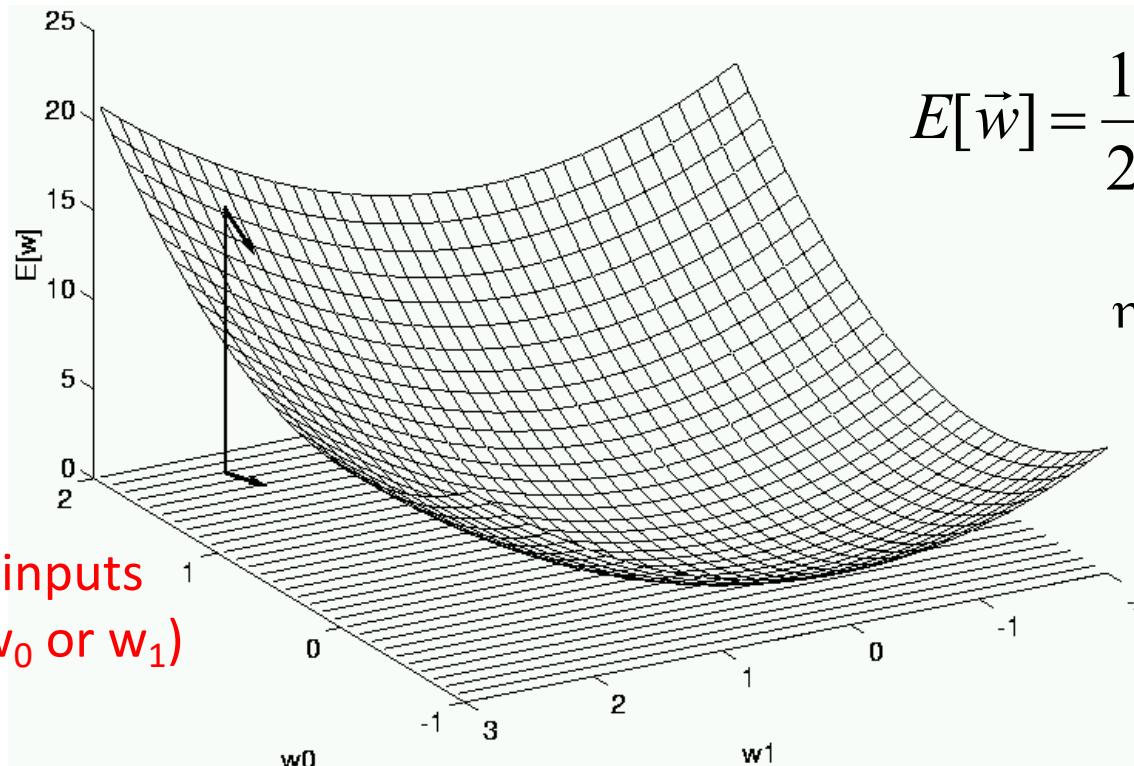
$$E[\vec{w}] = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

$D = \text{training data}$



Gradient Descent

1 neuron with two inputs
(and two weights w_0 or w_1)



$$E[\vec{w}] = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

η positive

Gradient:

$$\nabla E[\vec{w}] = \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

Rule:

$$\Delta \vec{w} = -\eta \nabla E[\vec{w}] \quad \Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

Because it is desired to move the weight vector in the direction that the error E decreases

Gradient Descent (one layer)

$$\begin{aligned}\frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_d (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_d \frac{\partial}{\partial w_i} (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_d 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d) \\ &= \sum_d (t_d - o_d) \frac{\partial}{\partial w_i} (t_d - \sigma(\vec{w} \cdot \vec{x}_d)) \\ &= \sum_d (t_d - o_d) \left(-\frac{\partial}{\partial w_i} \sigma(\vec{w} \cdot \vec{x}_d) \right) \quad \textcolor{red}{\sigma \text{ is the logistics function}} \\ &= -\sum_d (t_d - o_d) \sigma(\vec{w} \cdot \vec{x}_d) (1 - \sigma(\vec{w} \cdot \vec{x}_d)) x_{i,d}\end{aligned}$$

Gradient Descent: multiple Outputs

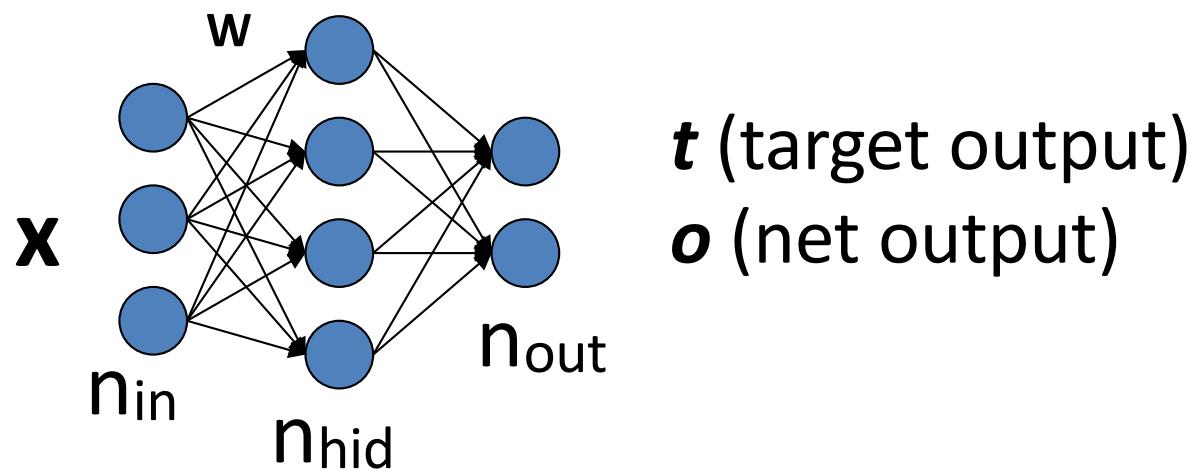
The error should now be redefined to k outputs:

$$E[\vec{w}] = \frac{1}{2} \sum_{d \in D} \sum_{k \in \text{outputs}} (t_{kd} - o_{kd})^2$$



The *Backpropagation* algorithm

- A feed-forward network with one hidden layer with n_{hid} sigmoidal units, n_{out} outputs (sigmoidal units), n_{in} inputs, several training samples $\langle x, t \rangle$



Input parameters of the backpropagation algorithm:
 $\{\langle x_1, t_1 \rangle, \langle x_2, t_2 \rangle, \dots\}, \eta, n_{\text{hid}}, n_{\text{in}}, n_{\text{out}}$

The *Backpropagation* algorithm

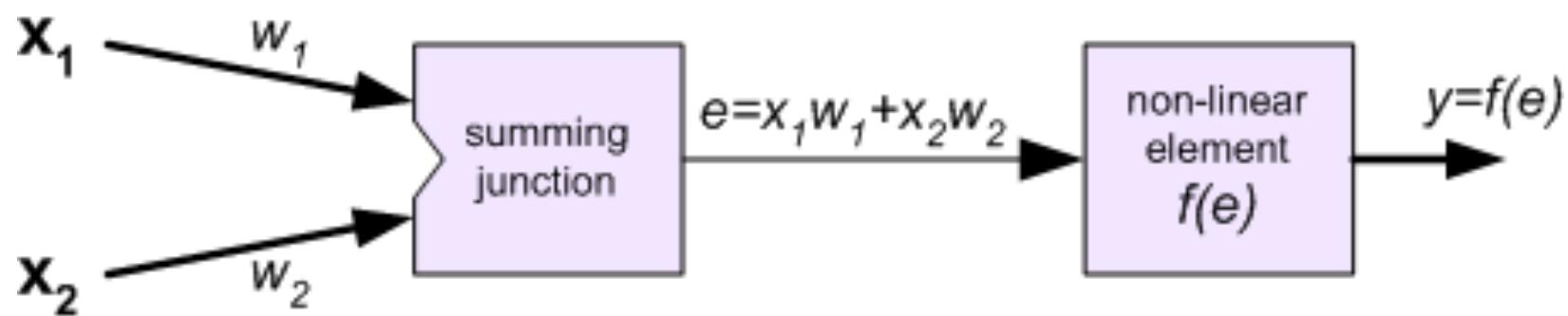
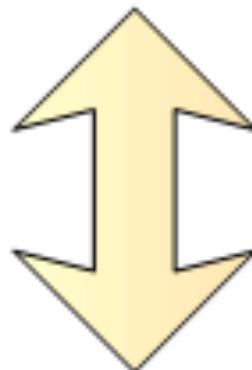
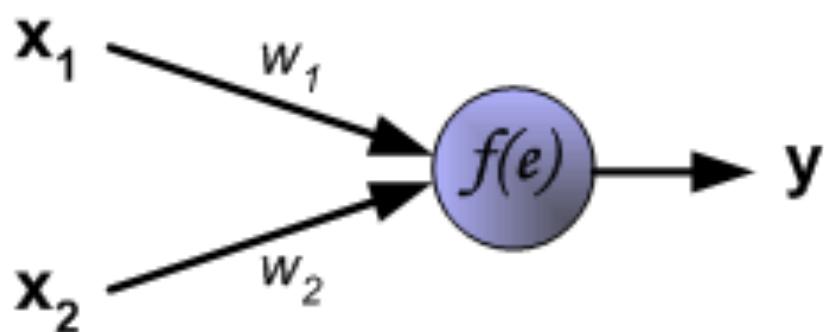
- Initialize w to small random numbers
- Until the termination condition is met, do:
 - For each $\langle x, t \rangle$ in training-examples, do:
 - //Propagate the input forward:
 1. Propagate the input forward and compute each o_k .
 - //Propagate the error backward:
 2. For each output unit k, calculate the error δ_k :
 - $$\delta_k \leftarrow o_k (1 - o_k) (t_k - o_k)$$
 - 3. For each hidden unit h, calculate the error δ_h :

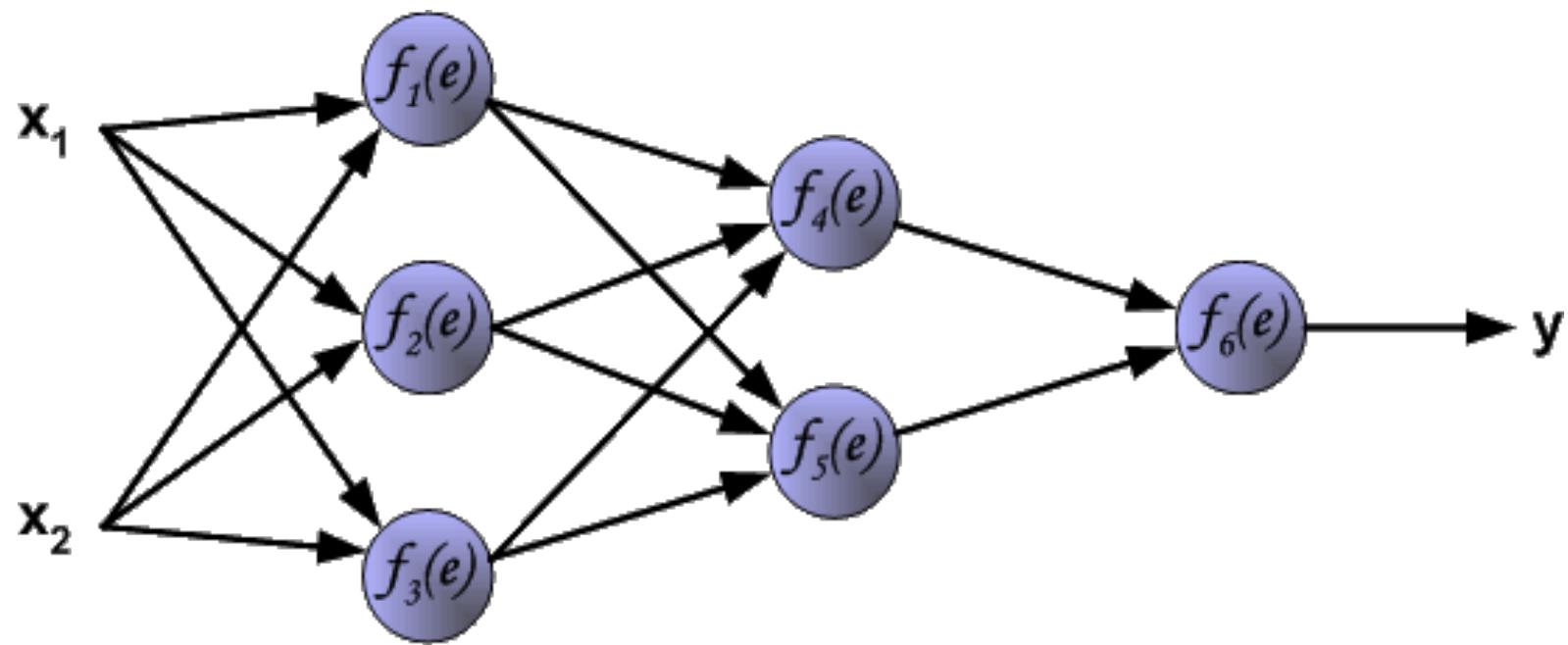
$$\delta_h \leftarrow o_h (1 - o_h) \sum_{k \in \text{saidas}} W_{hk} \cdot \delta_k$$

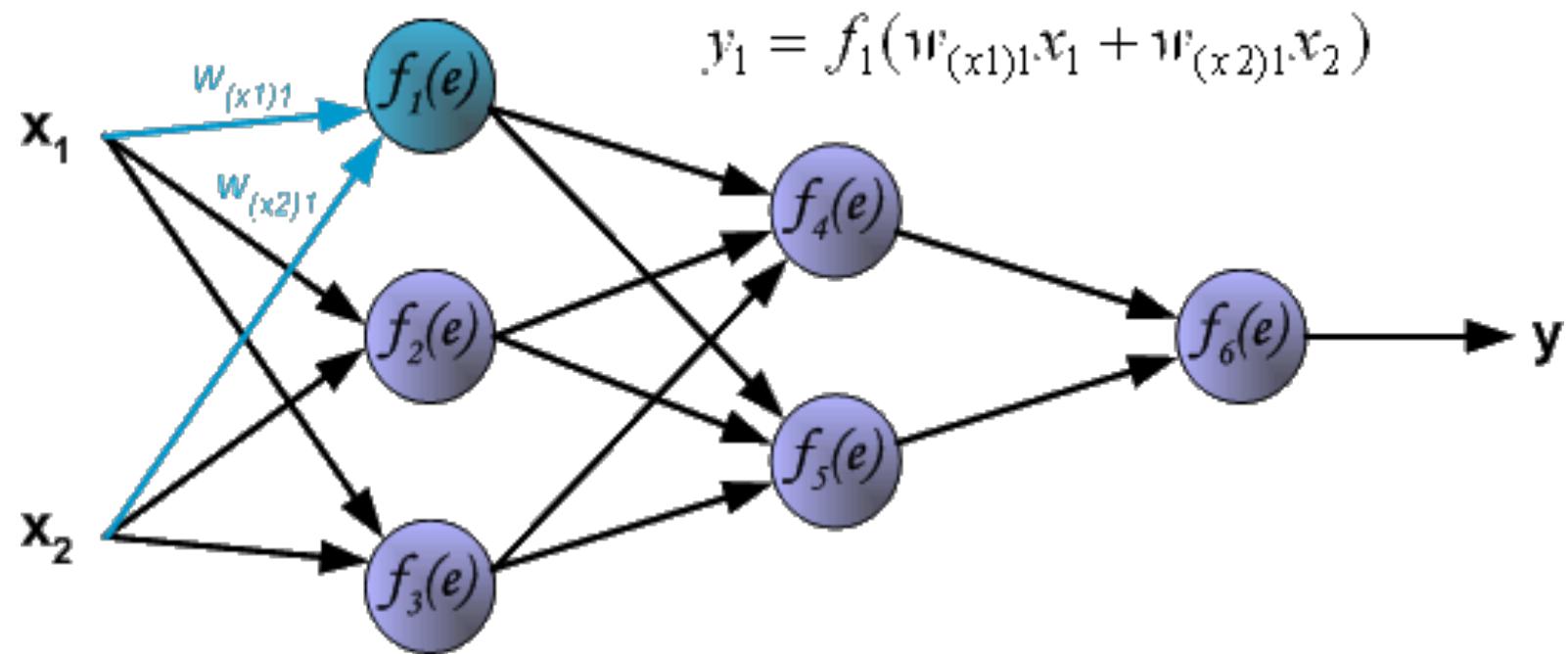
4. Update each network weight w_{ij} :

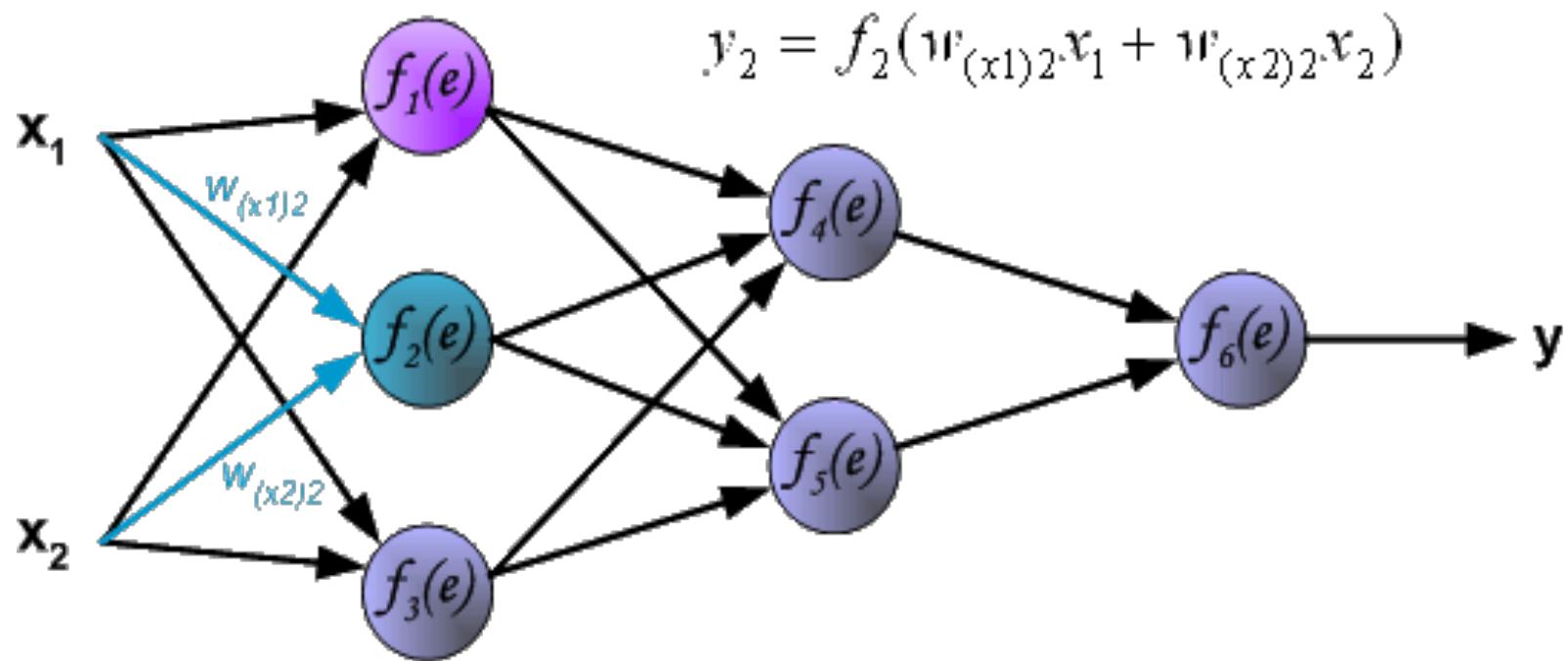
$$w_{ij} \leftarrow w_{ij} + \eta \delta_j x_{ij}$$

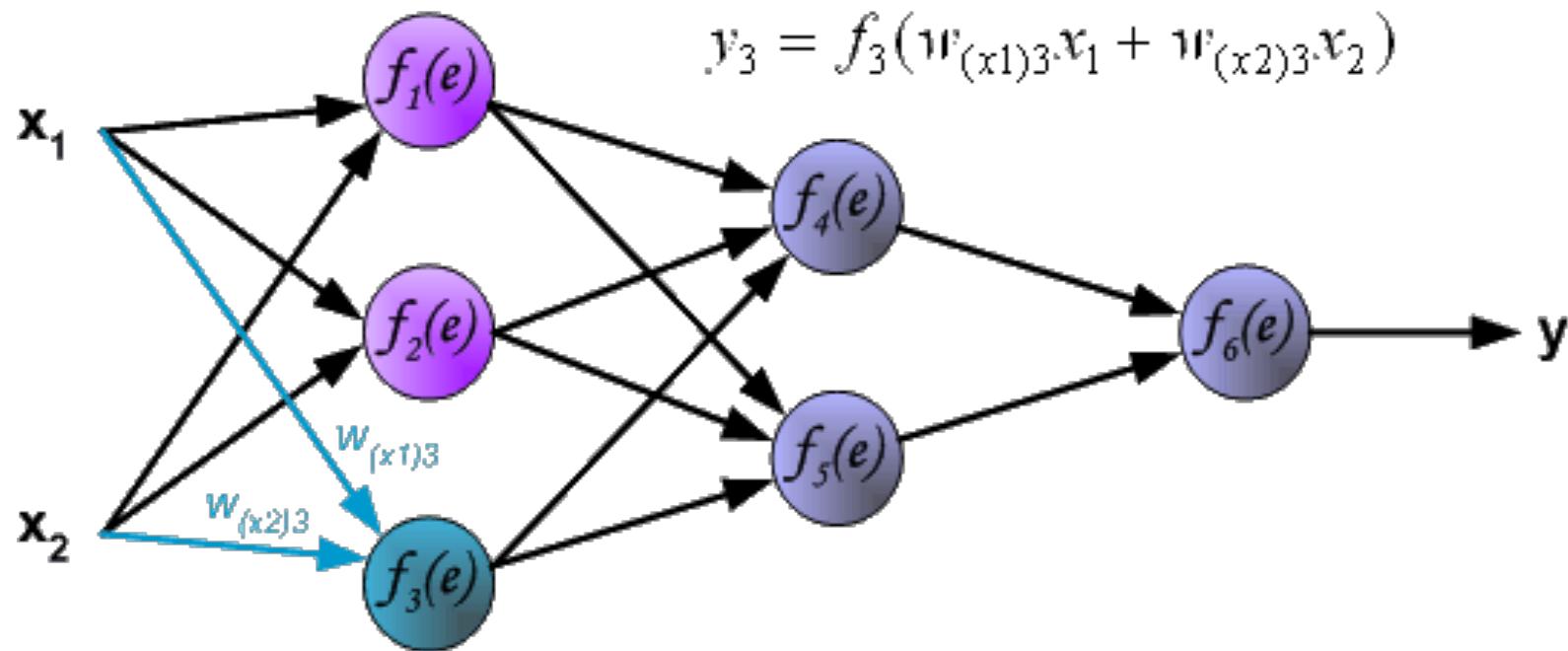
Input from i to j

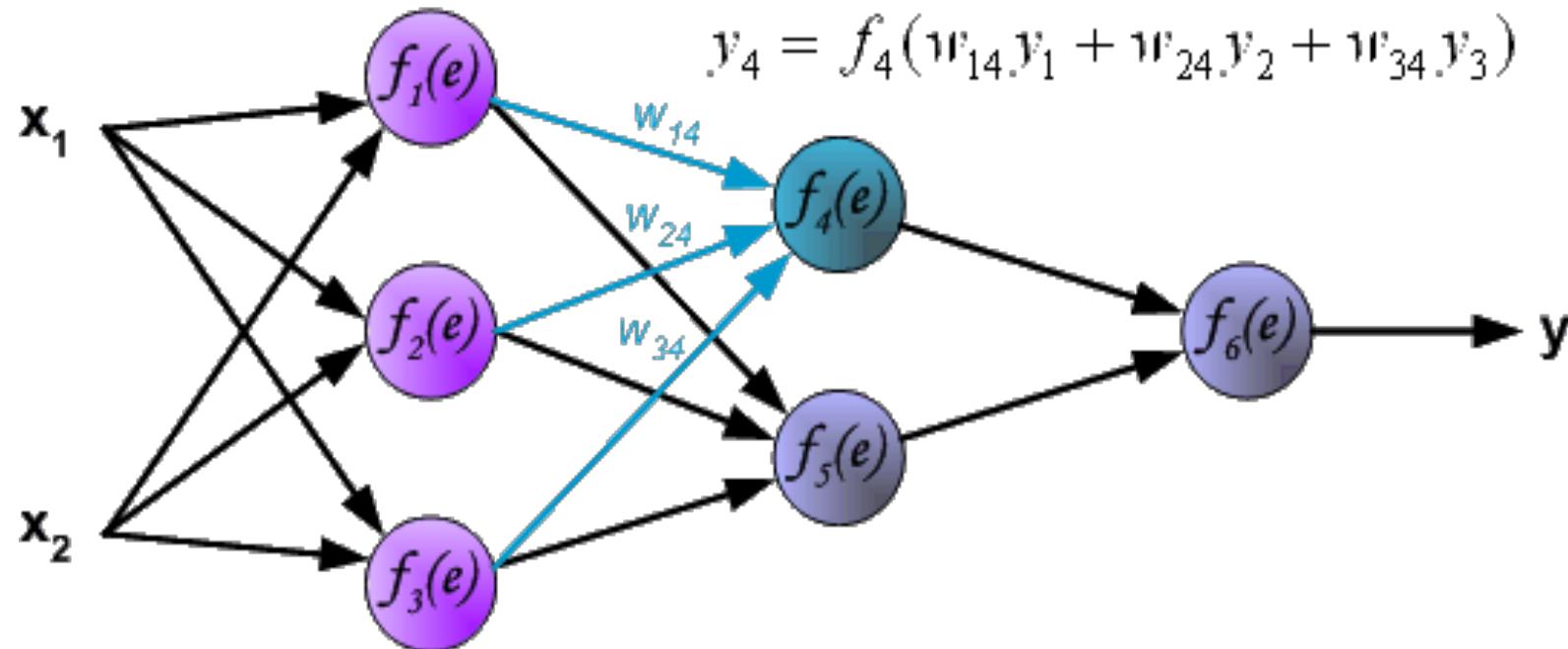




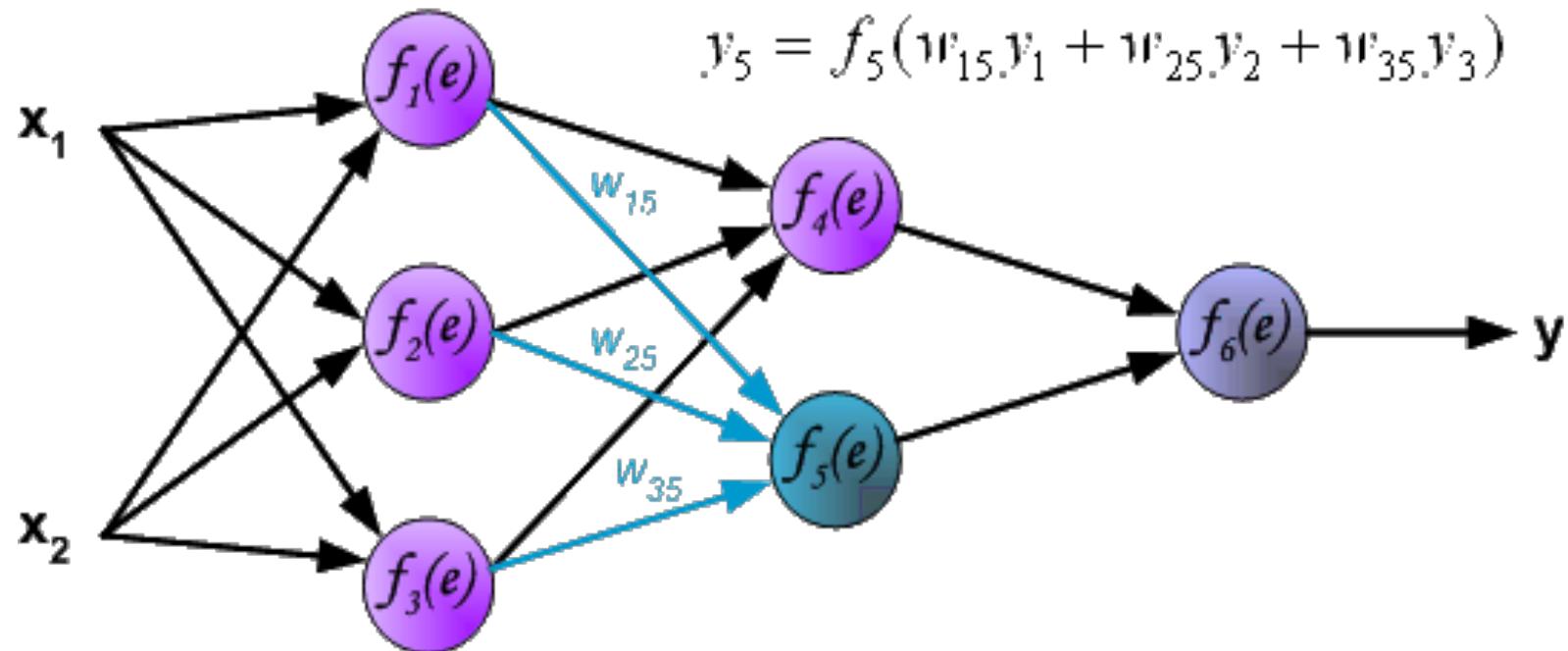


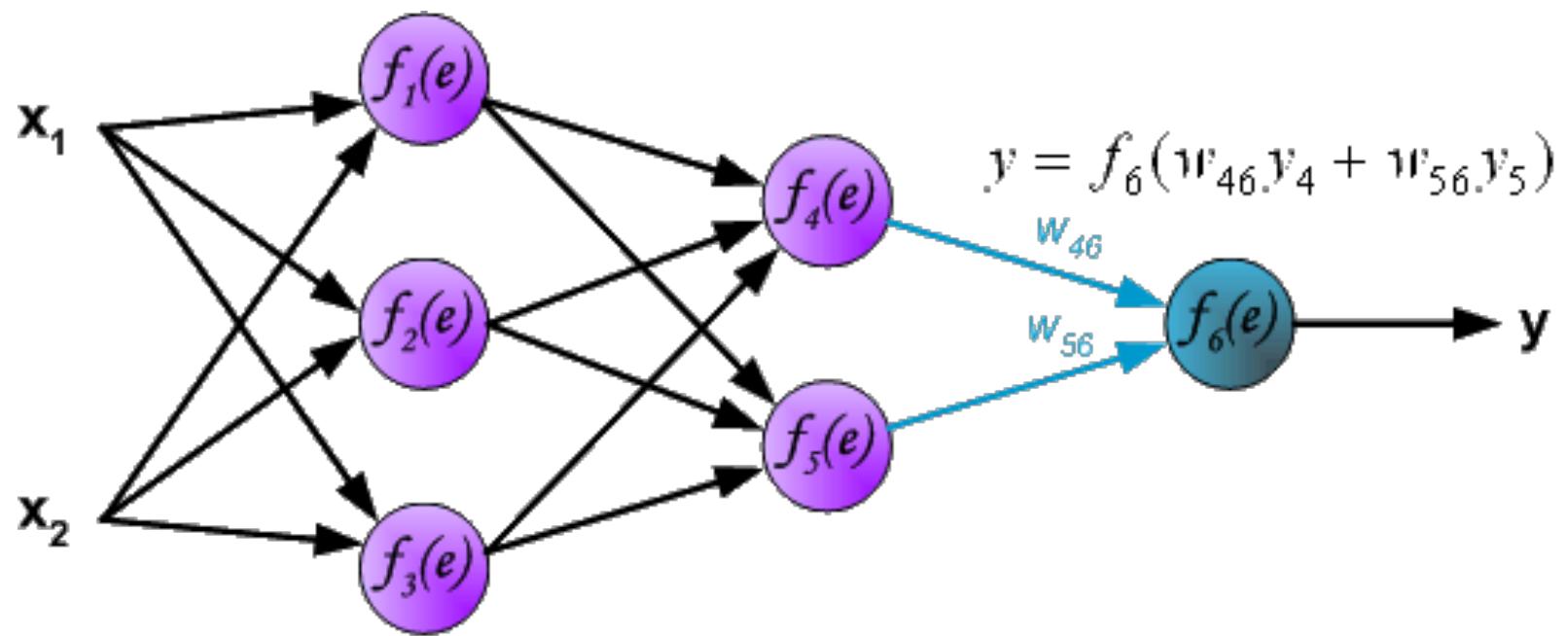


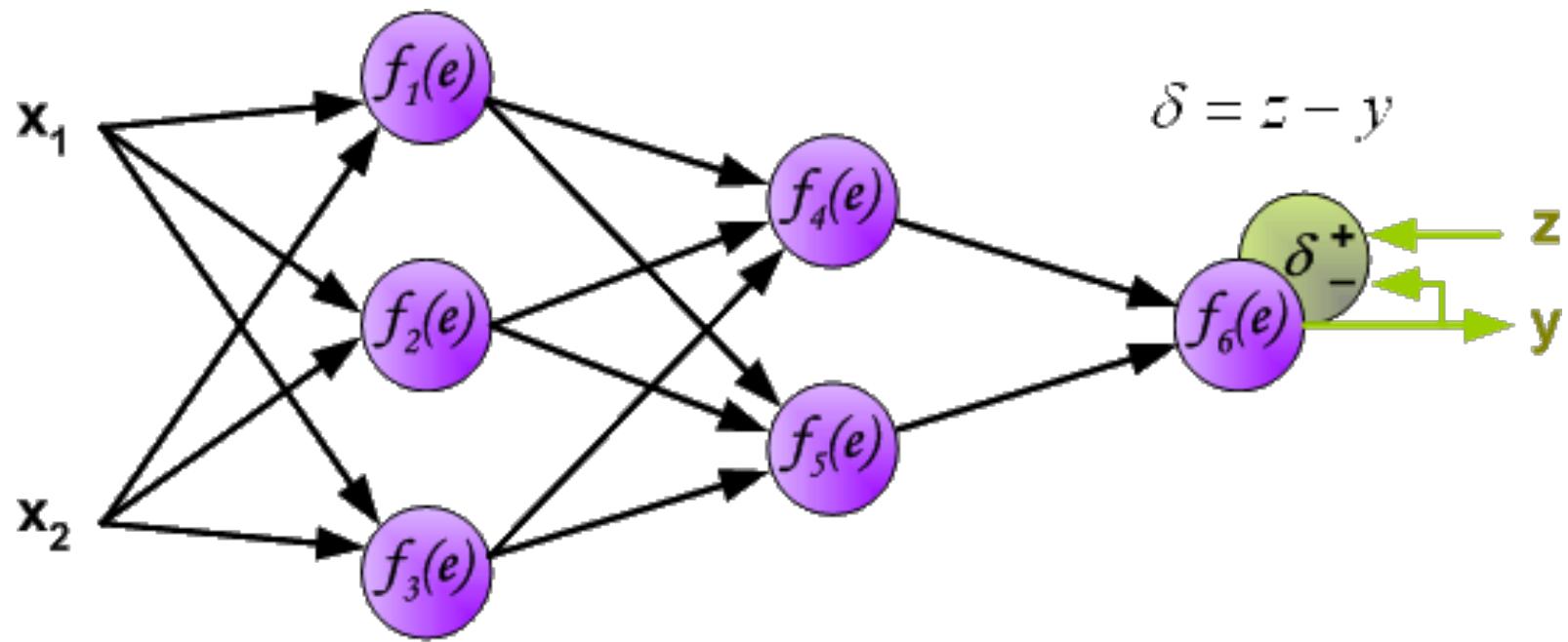


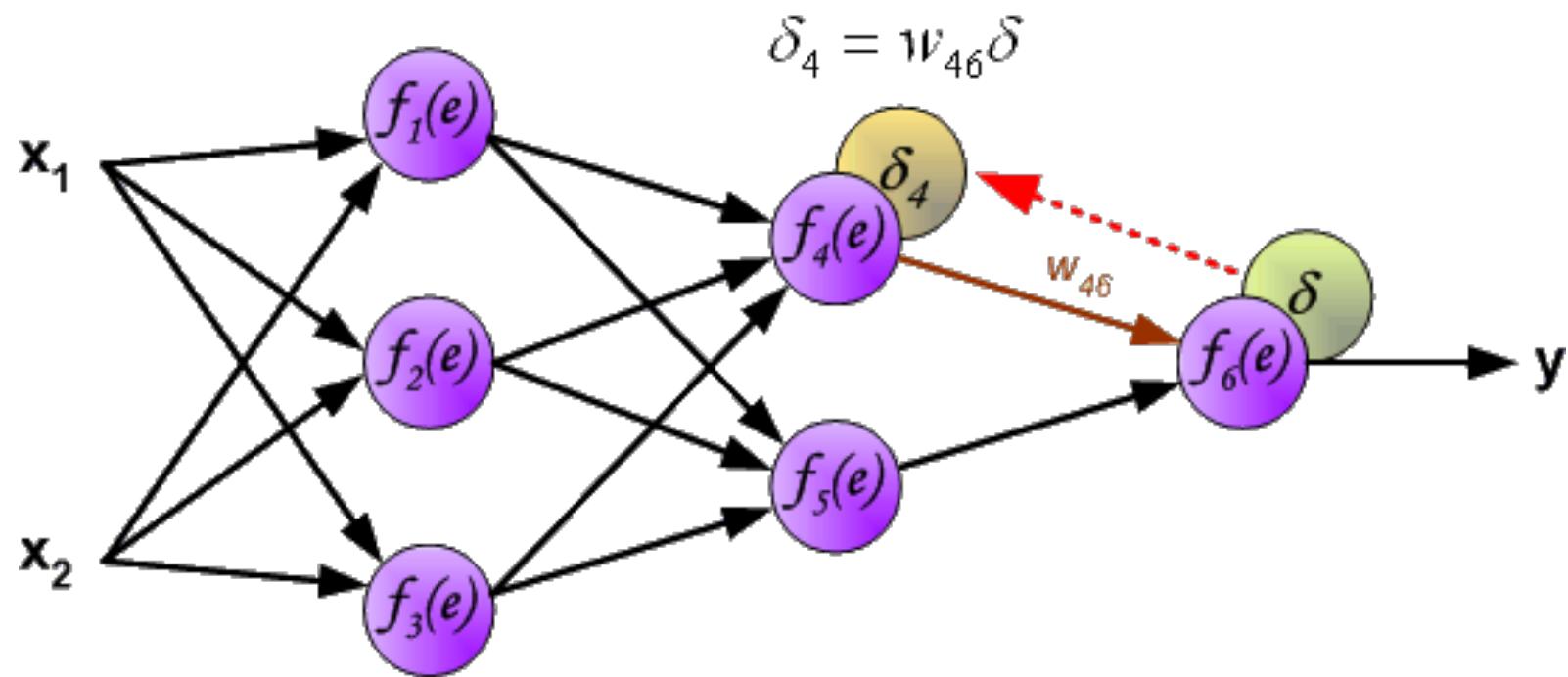


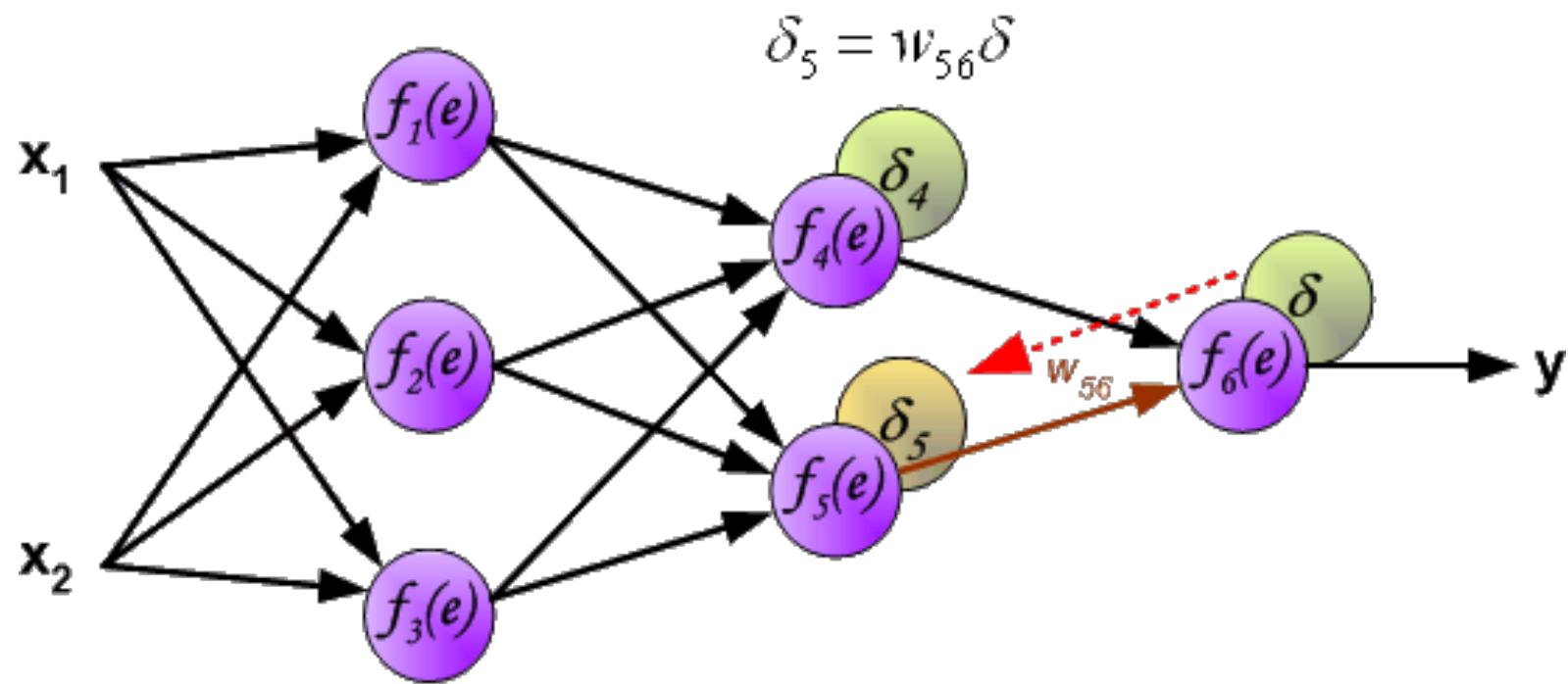
$$y_4 = f_4(w_{14}y_1 + w_{24}y_2 + w_{34}y_3)$$

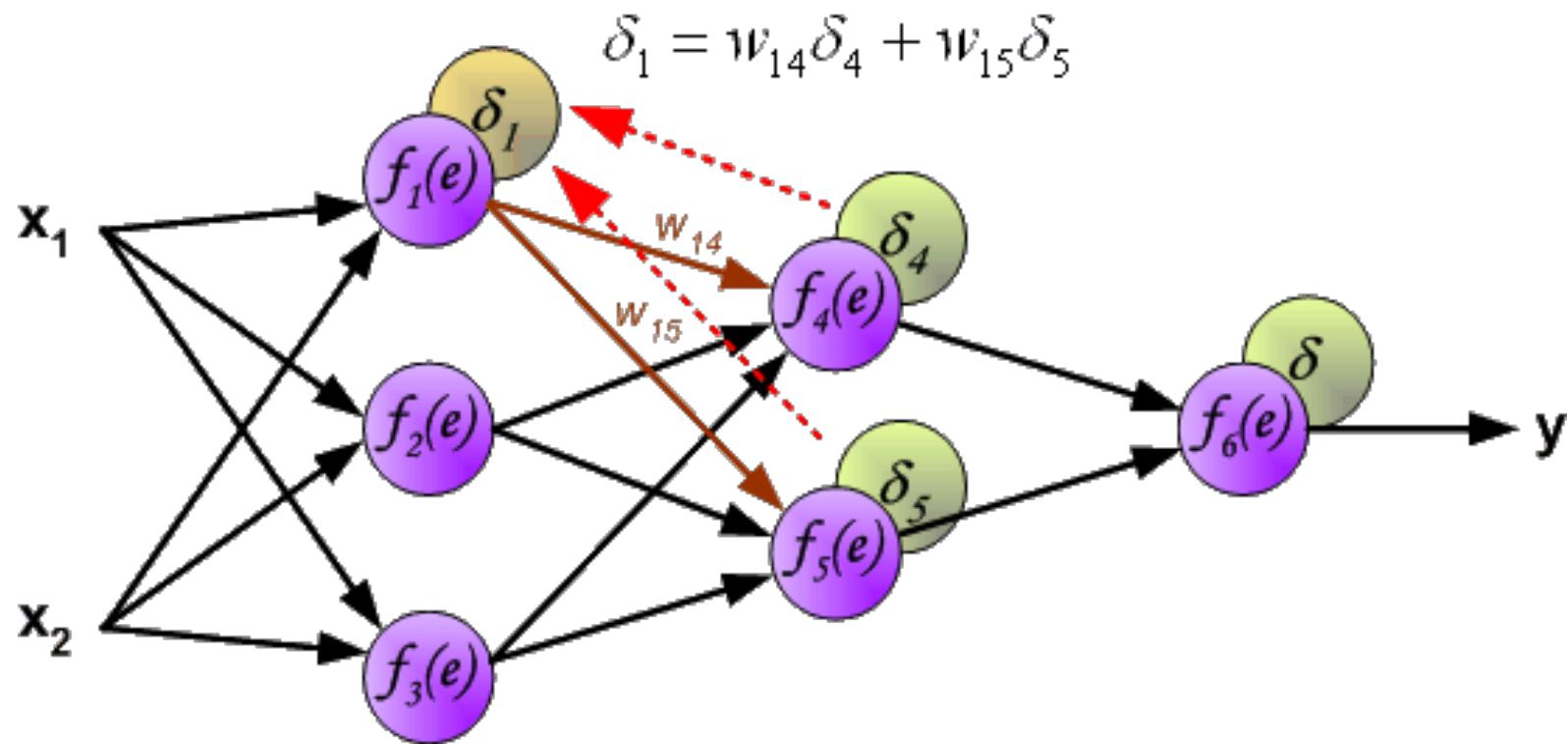


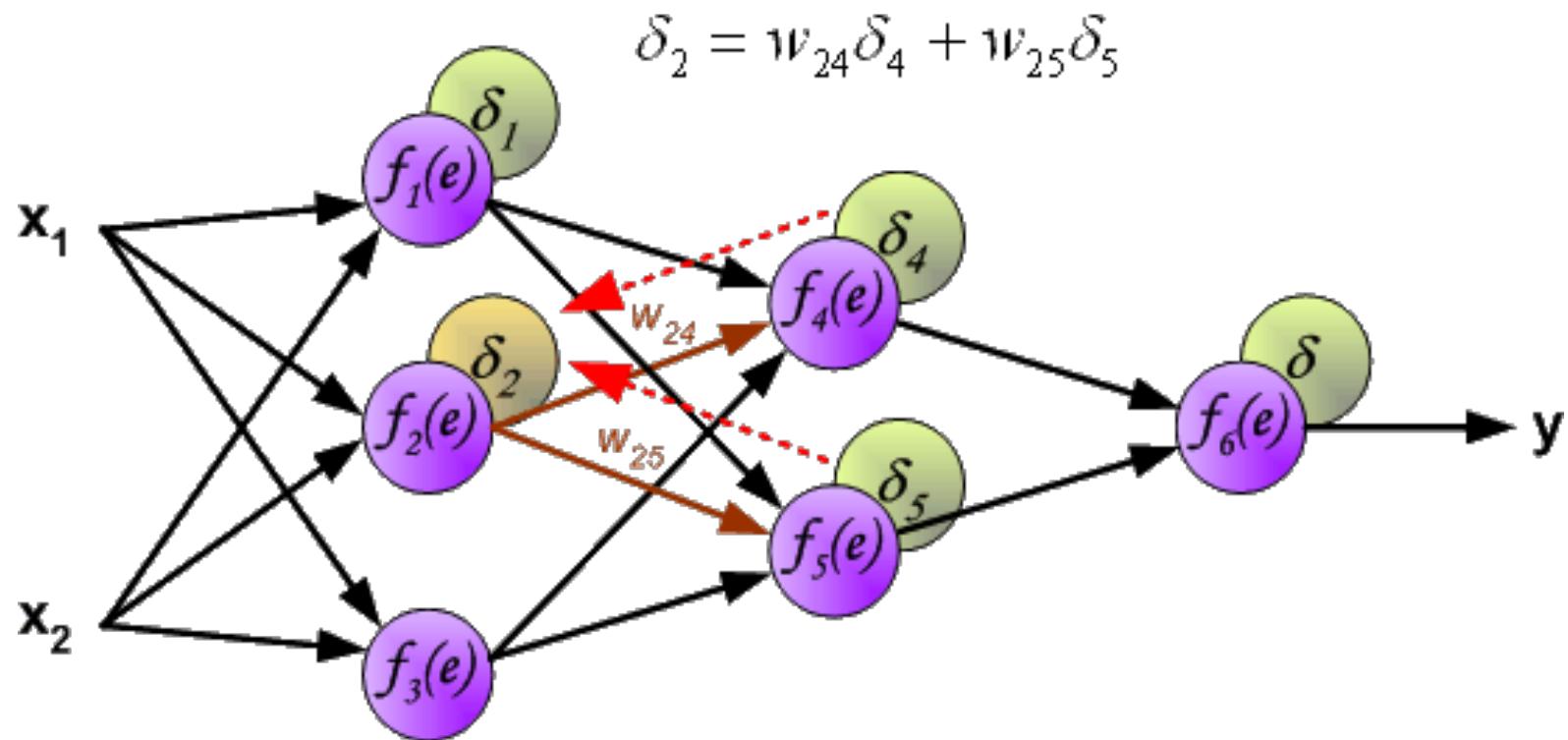


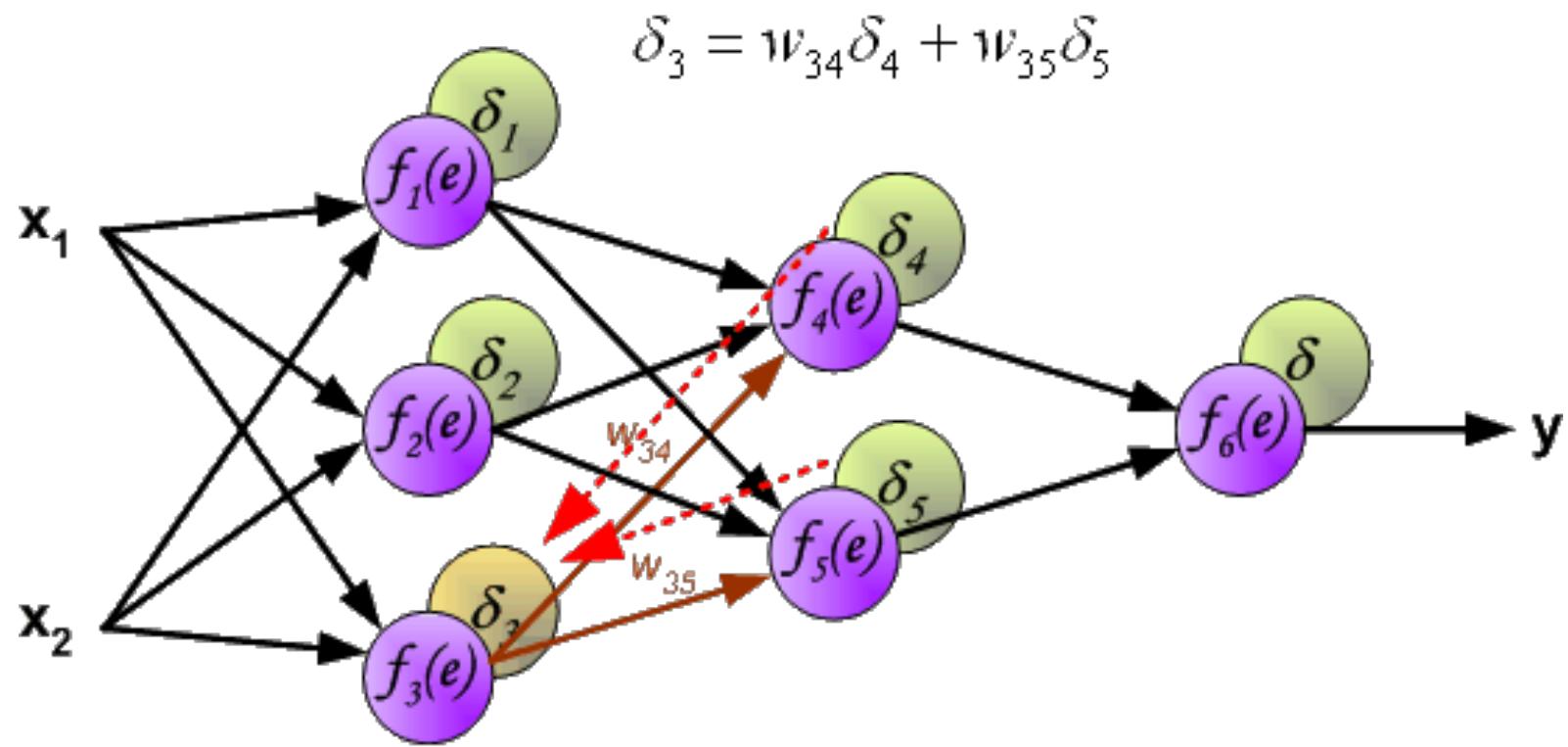


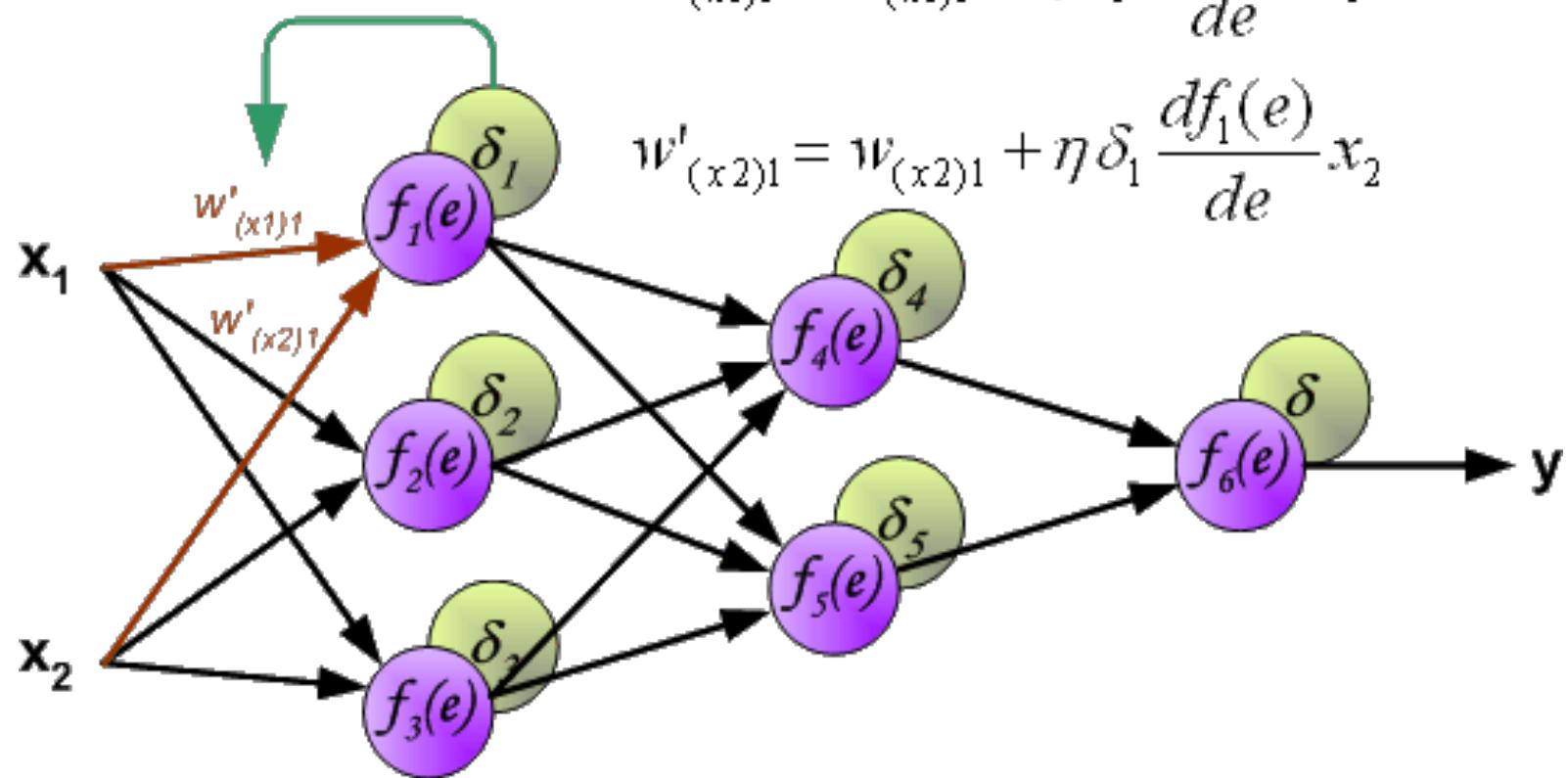










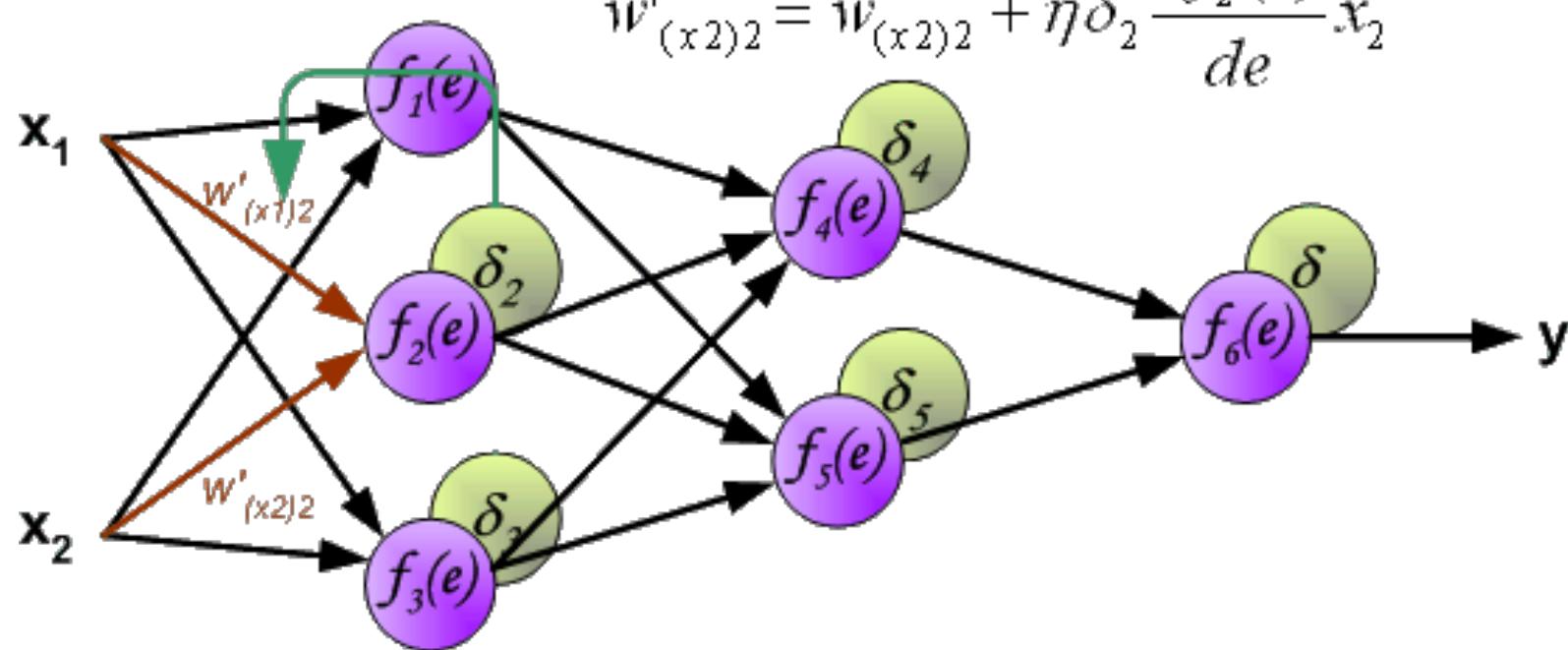


$$w'_{(x1)1} = w_{(x1)1} + \eta \delta_1 \frac{df_1(e)}{de} x_1$$

$$w'_{(x2)1} = w_{(x2)1} + \eta \delta_1 \frac{df_1(e)}{de} x_2$$

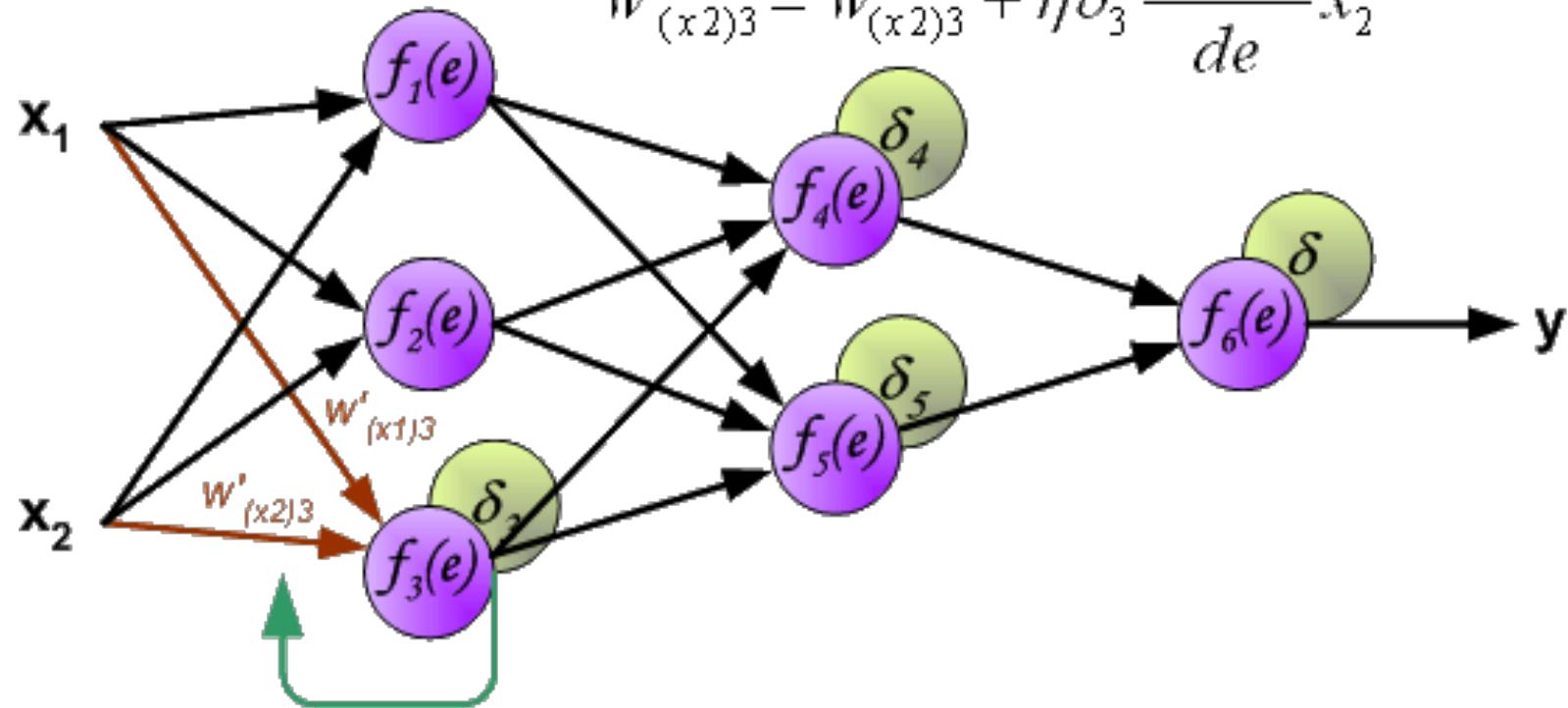
$$w'_{(x1)2} = w_{(x1)2} + \eta \delta_2 \frac{df_2(e)}{de} x_1$$

$$w'_{(x2)2} = w_{(x2)2} + \eta \delta_2 \frac{df_2(e)}{de} x_2$$



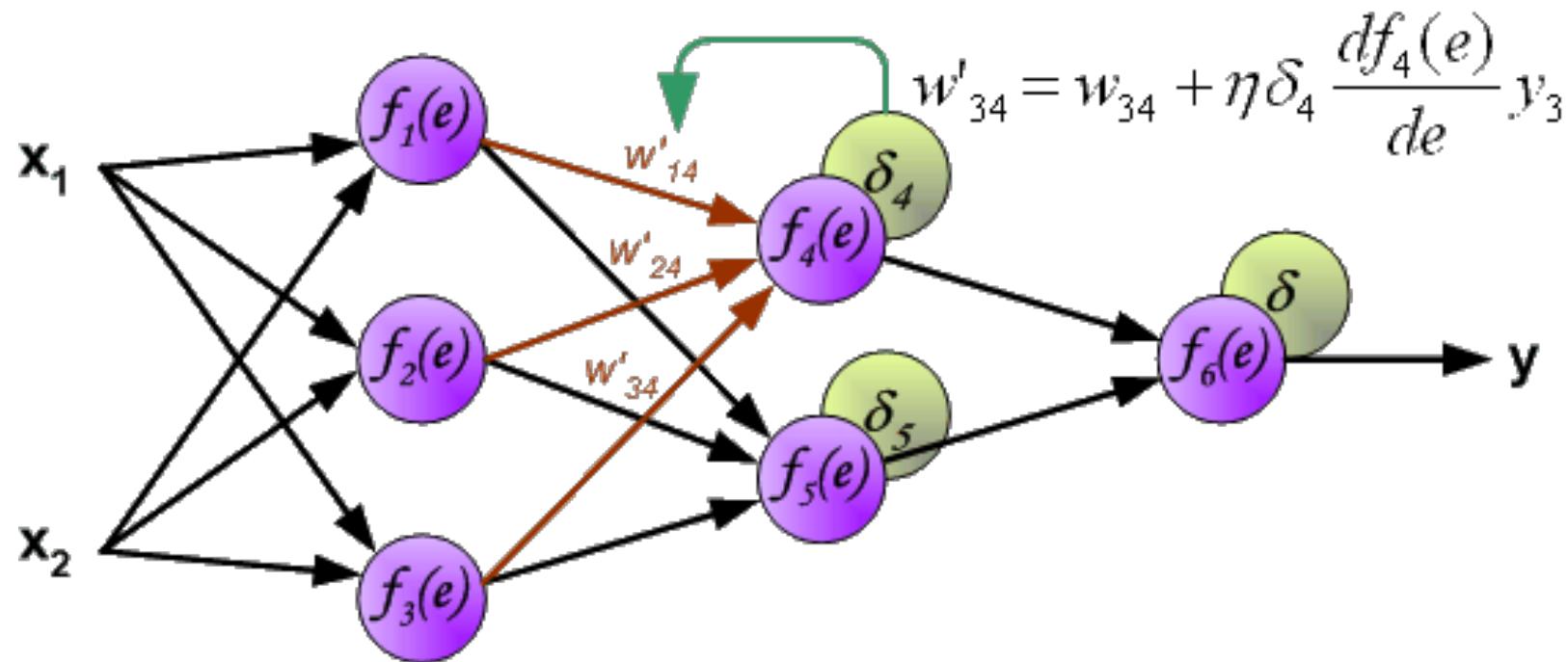
$$w'_{(x1)3} = w_{(x1)3} + \eta \delta_3 \frac{df_3(e)}{de} x_1$$

$$w'_{(x2)3} = w_{(x2)3} + \eta \delta_3 \frac{df_3(e)}{de} x_2$$



$$w'_{14} = w_{14} + \eta \delta_4 \frac{df_4(e)}{de} y_1$$

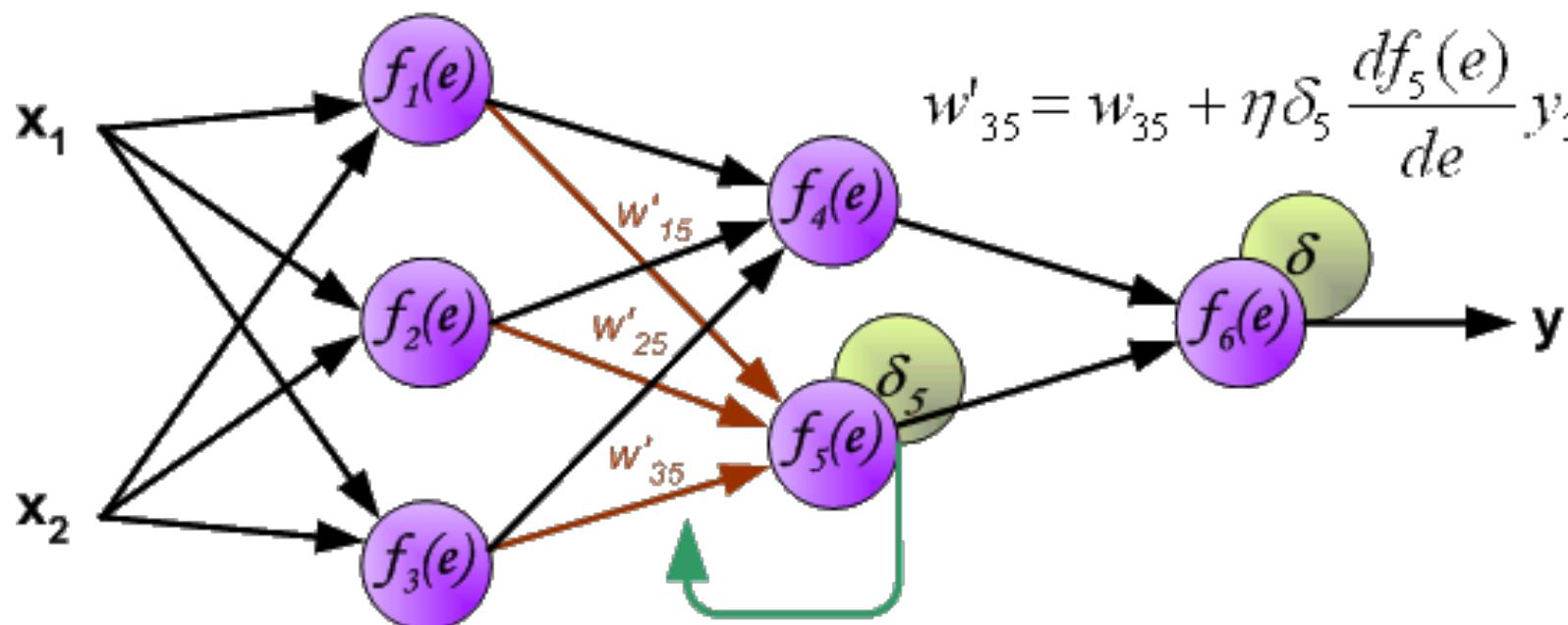
$$w'_{24} = w_{24} + \eta \delta_4 \frac{df_4(e)}{de} y_2$$



$$w'_{15} = w_{15} + \eta \delta_5 \frac{df_5(e)}{de} y_1$$

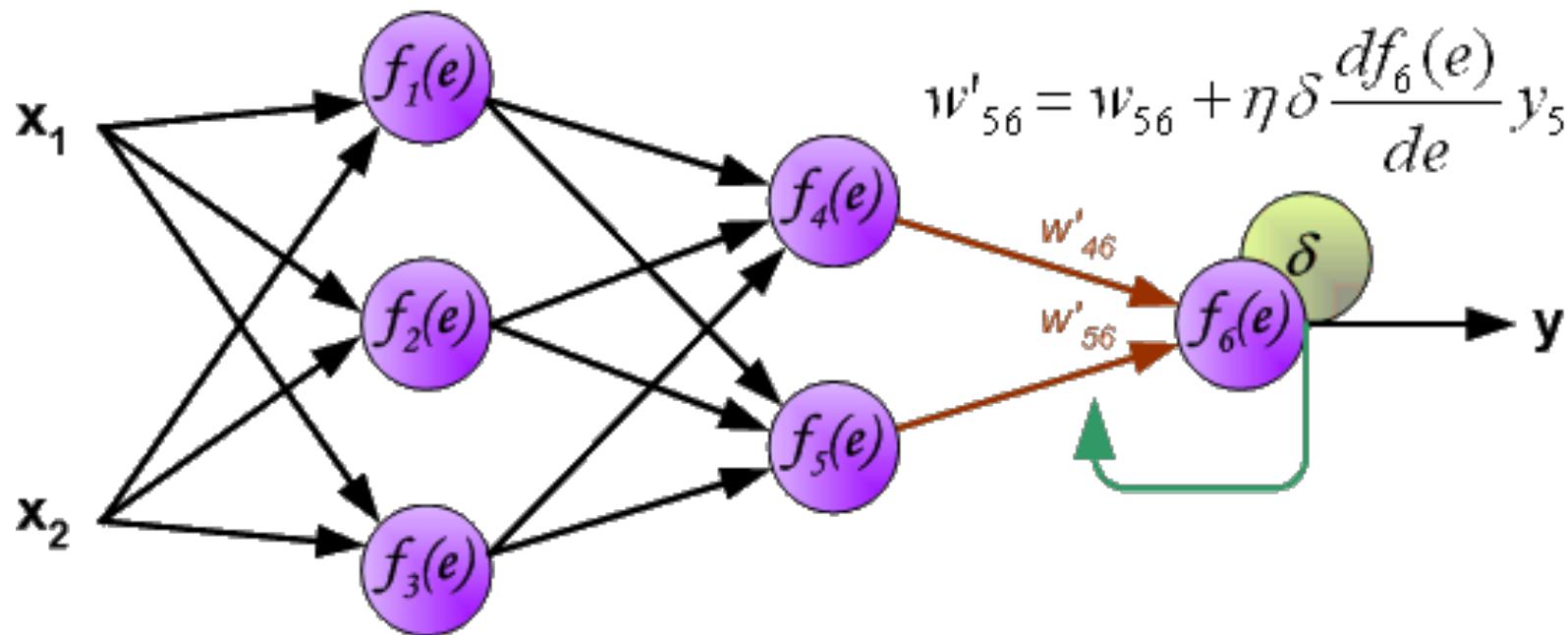
$$w'_{25} = w_{25} + \eta \delta_5 \frac{df_5(e)}{de} y_2$$

$$w'_{35} = w_{35} + \eta \delta_5 \frac{df_5(e)}{de} y_3$$

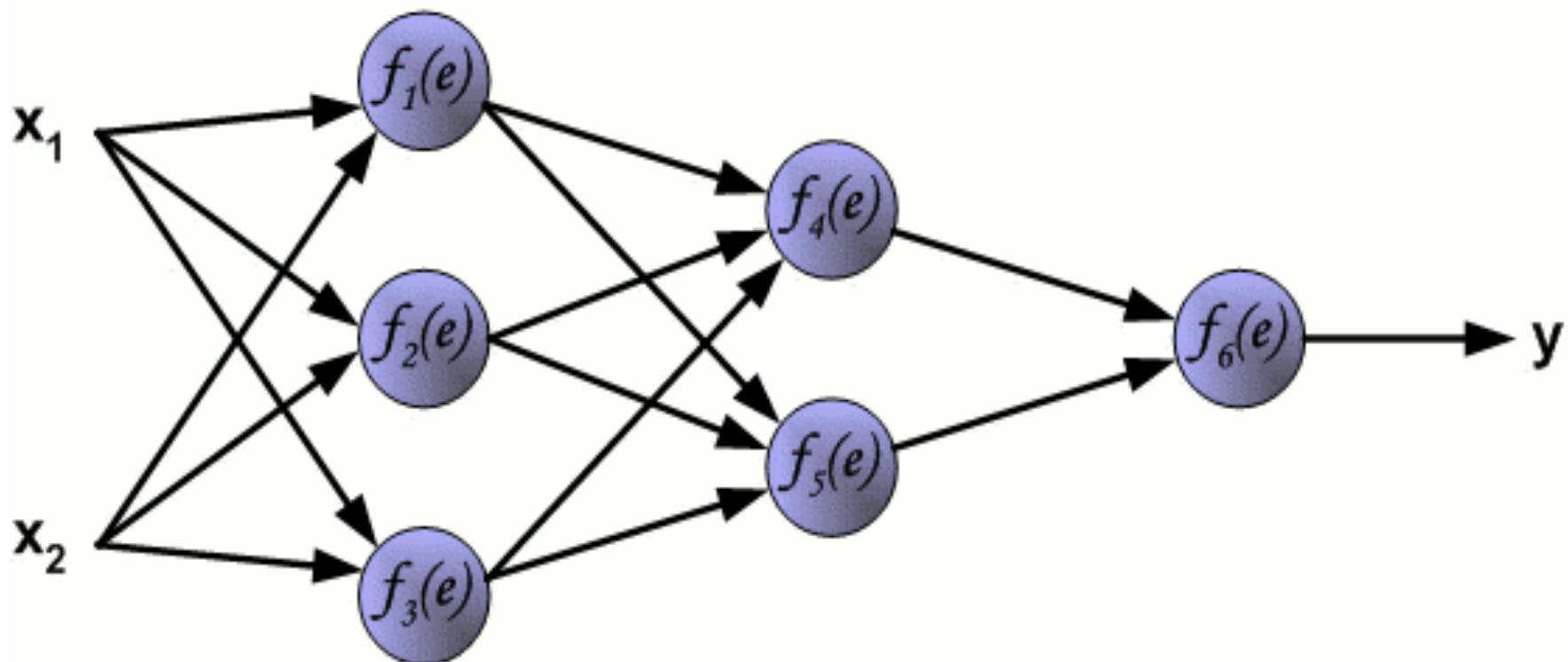


$$w'_{46} = w_{46} + \eta \delta \frac{df_6(e)}{de} y_4$$

$$w'_{56} = w_{56} + \eta \delta \frac{df_6(e)}{de} y_5$$



Backpropagation: animation

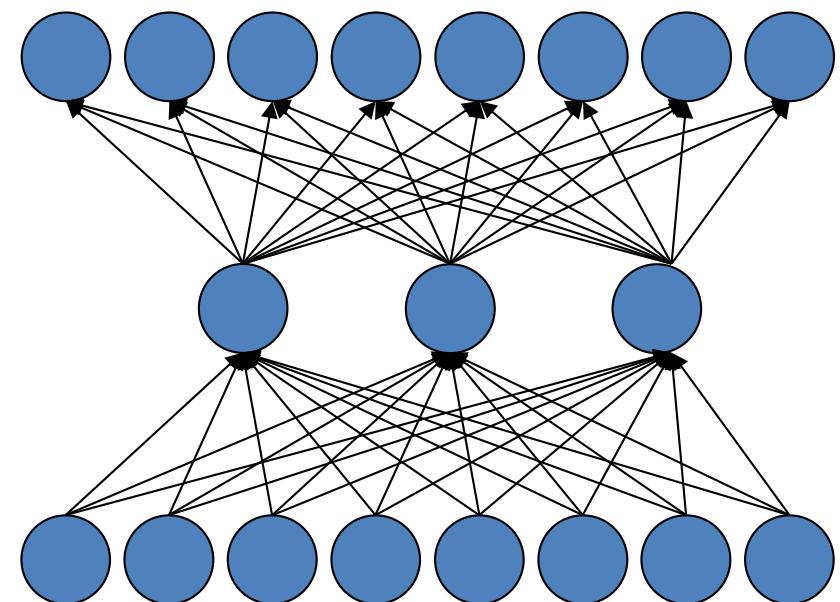


http://www.trapexit.org/images/b/ba/Animate_ANN.gif

Example:

initial weights: (-0.1 a +0.1); $\eta = 0.3$

Entrada		Saída
10000000	→	10000000
01000000	→	01000000
00100000	→	00100000
00010000	→	00010000
00001000	→	00001000
00000100	→	00000100
00000010	→	00000010
00000001	→	00000001

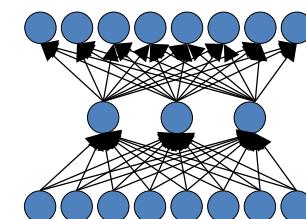


Learning the representation (inner layer)

Entrada		Saída hs		Saída
10000000	→	.89 .04 .08	→	10000000
01000000	→	.15 .99 .99	→	01000000
00100000	→	.01 .97 .27	→	00100000
00010000	→	.99 .97 .71	→	00010000
00001000	→	.03 .05 .02	→	00001000
00000100	→	.01 .11 .88	→	00000100
00000010	→	.80 .01 .98	→	00000010
00000001	→	.60 .94 .01	→	00000001

1 0 0
0 1 1

Intermediate representation "discovers" binary code !!!



Application × ANN type

- There are many other types of ANN:
 - For **classification** or **prediction** tasks, usually feed-forward networks (such as MLP) are used.
 - For **clustering** tasks, the types of network used are: Simple Competitive Networks, Adaptive Resonance Theory (ART) networks, Kohonen Self-Organizing Maps (SOM).
 - In **association** tasks, an ANN can be trained to "remember" a number of patterns; the type of ANN usually used for this task is Hopfield network.

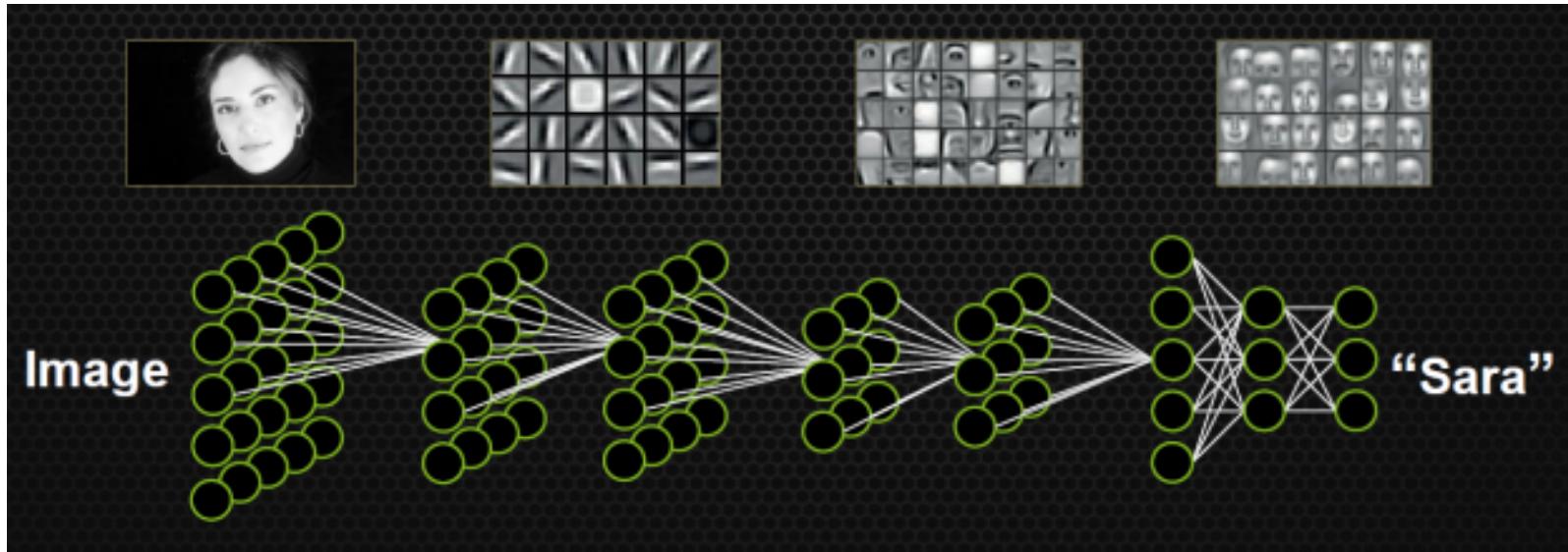
Recent advances and future applications of ANNs

- Integration of fuzzy logic into neural networks
- Pulsed neural networks
- Hardware specialized for neural networks
- Deep Learning:
 - enable much deeper (and larger) networks (5 to 10 hidden layers): have the ability of building up a complex hierarchy of concepts.
 - one can show that there are functions which a k -layer network can represent compactly (with a number of hidden units that is **polynomial** in the number of inputs), that a $(k - 1)$ -layer network cannot represent unless it has an **exponentially** large number of hidden units.

References

- Russel, S.; Norvig, P. Artificial Intelligence: a modern approach. 2nd.edition. Prentice Hall, 2003. Cap. 20.5.
- Mitchell, T.M. Machine Learning. WCB/McGraw-Hill, 1997. Cap.4.
- *Simon Haykin. Neural Networks: A Comprehensive Foundation.*
- *Livro em Português: Braga, Ludermir e Carvalho. Redes Neurais Artificiais. LTC.*

There are very interesting stuff on the web.



Deep Learning

Anna Helena Reali Costa

“more intuitive” applications



Face recognition



“Baseball player is throwing ball in game.”

Generating natural language descriptions



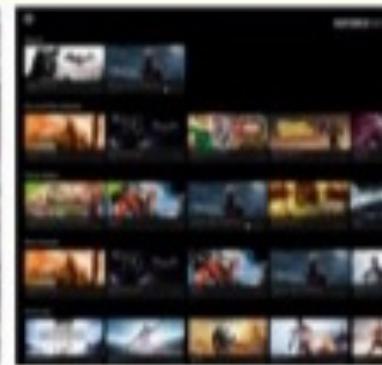
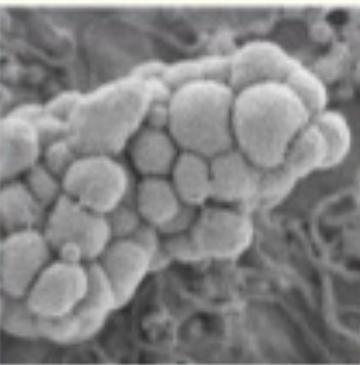
Dog or Mop?



Chiuaua or Muffin?

Deep Learning Applications

DEEP LEARNING EVERYWHERE



INTERNET & CLOUD

- Image Classification
- Speech Recognition
- Language Translation
- Language Processing
- Sentiment Analysis
- Recommendation

MEDICINE & BIOLOGY

- Cancer Cell Detection
- Diabetic Grading
- Drug Discovery

MEDIA & ENTERTAINMENT

- Video Captioning
- Video Search
- Real Time Translation

SECURITY & DEFENSE

- Face Detection
- Video Surveillance
- Satellite Imagery

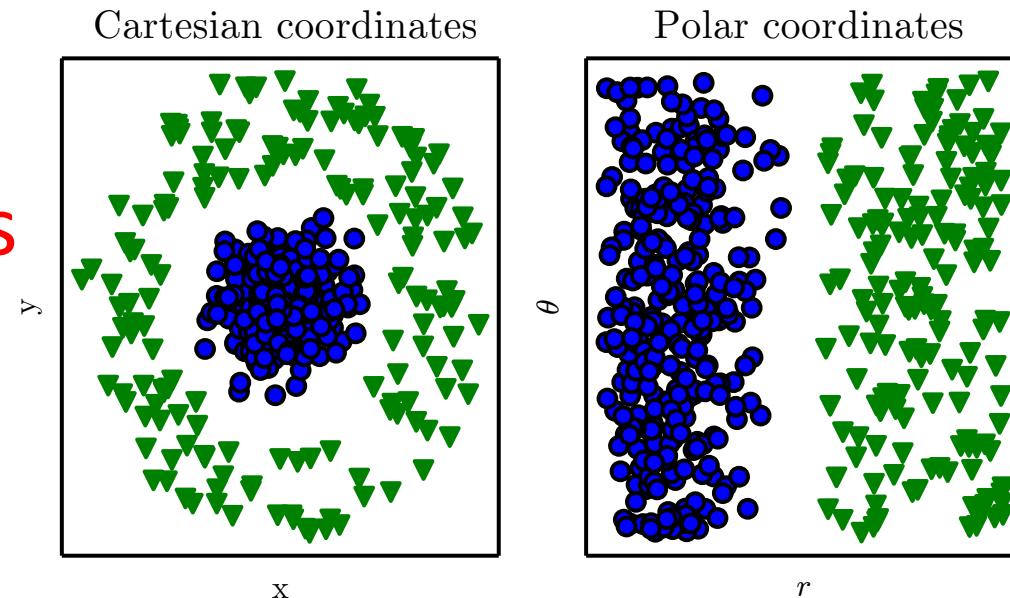
AUTONOMOUS MACHINES

- Pedestrian Detection
- Lane Tracking
- Recognize Traffic Sign

Machine Learning depends on Representation

- The performance of simple machine learning algorithms depends heavily on the **representation** of the data they are given (*features!*).

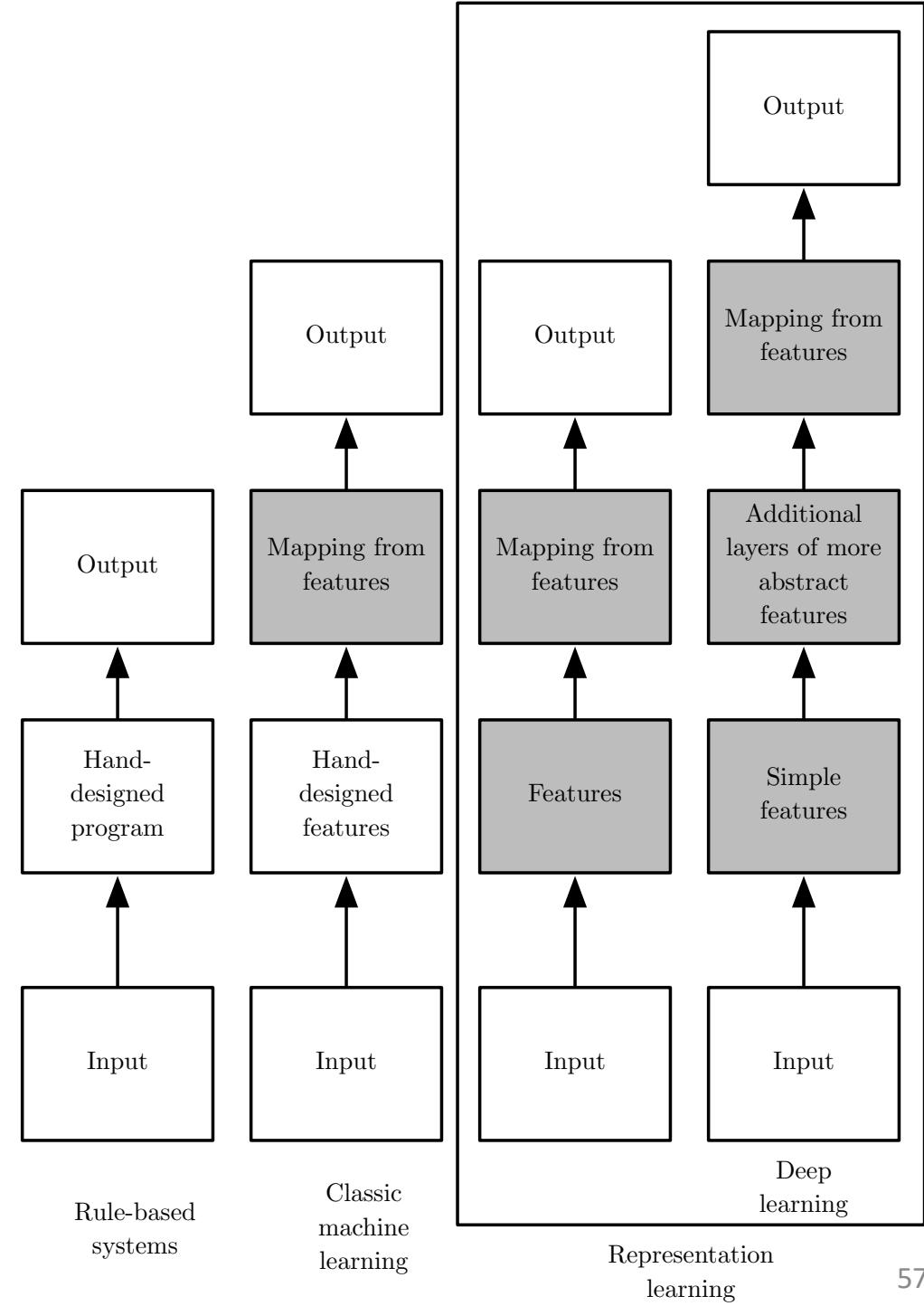
Representations Matter!!



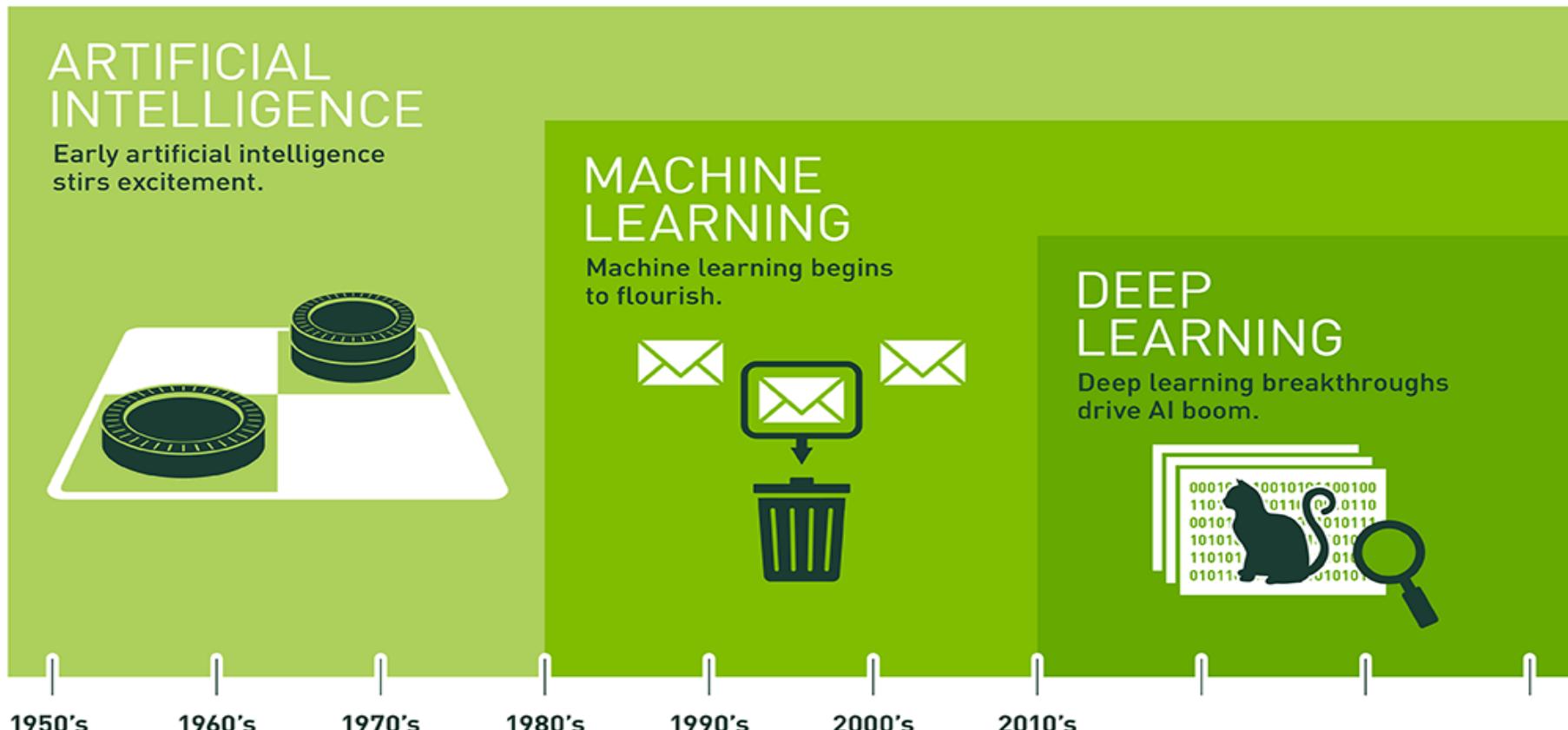
Representation Learning

- Solution: to use machine learning to **discover** not only the mapping from representation to output but also the representation itself.
 - This approach is known as **representation learning**.
- **Deep learning** allows the computer to build complex representations out of simpler representations.

Learning Multiple Components



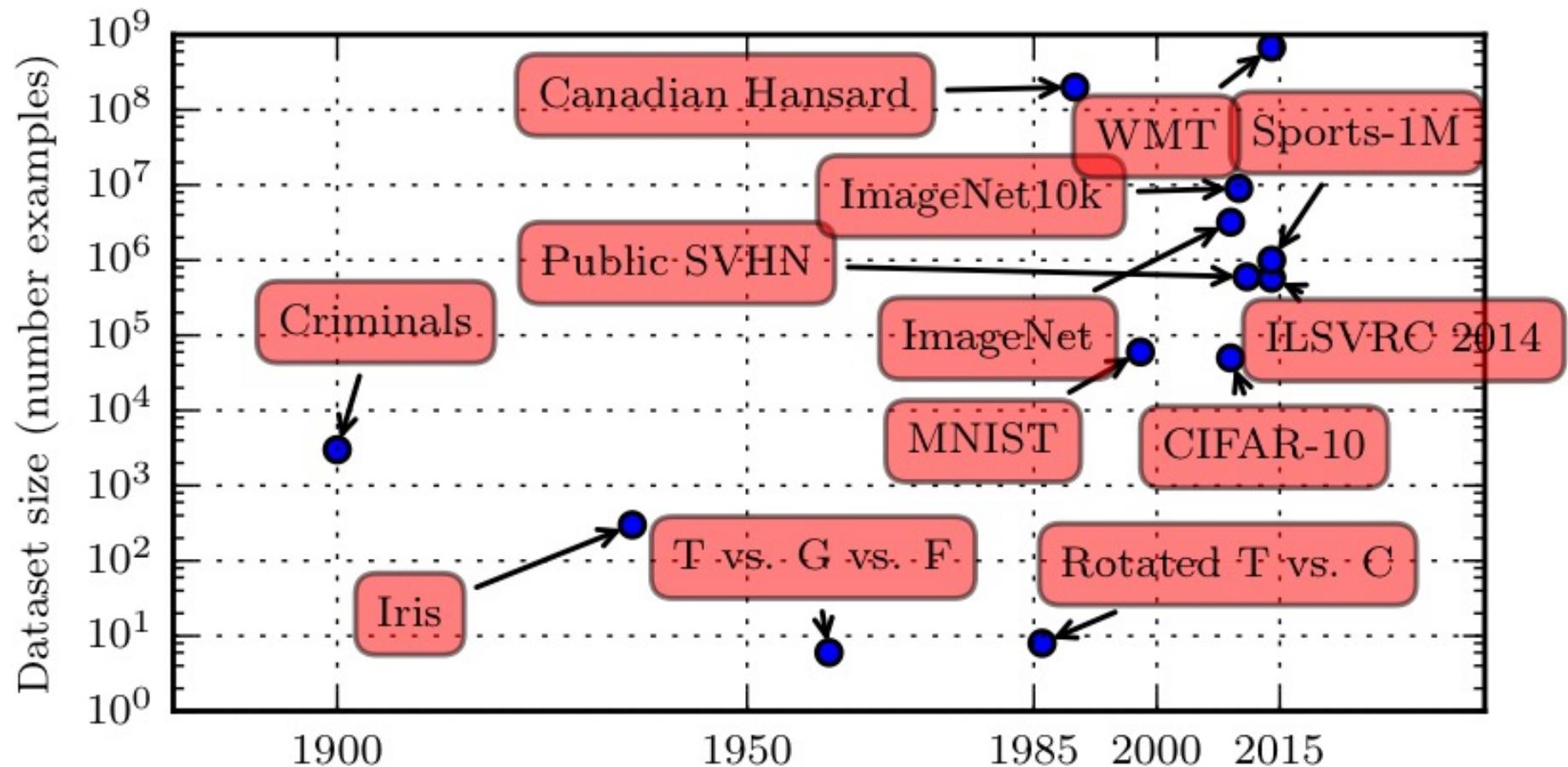
Deep Learning



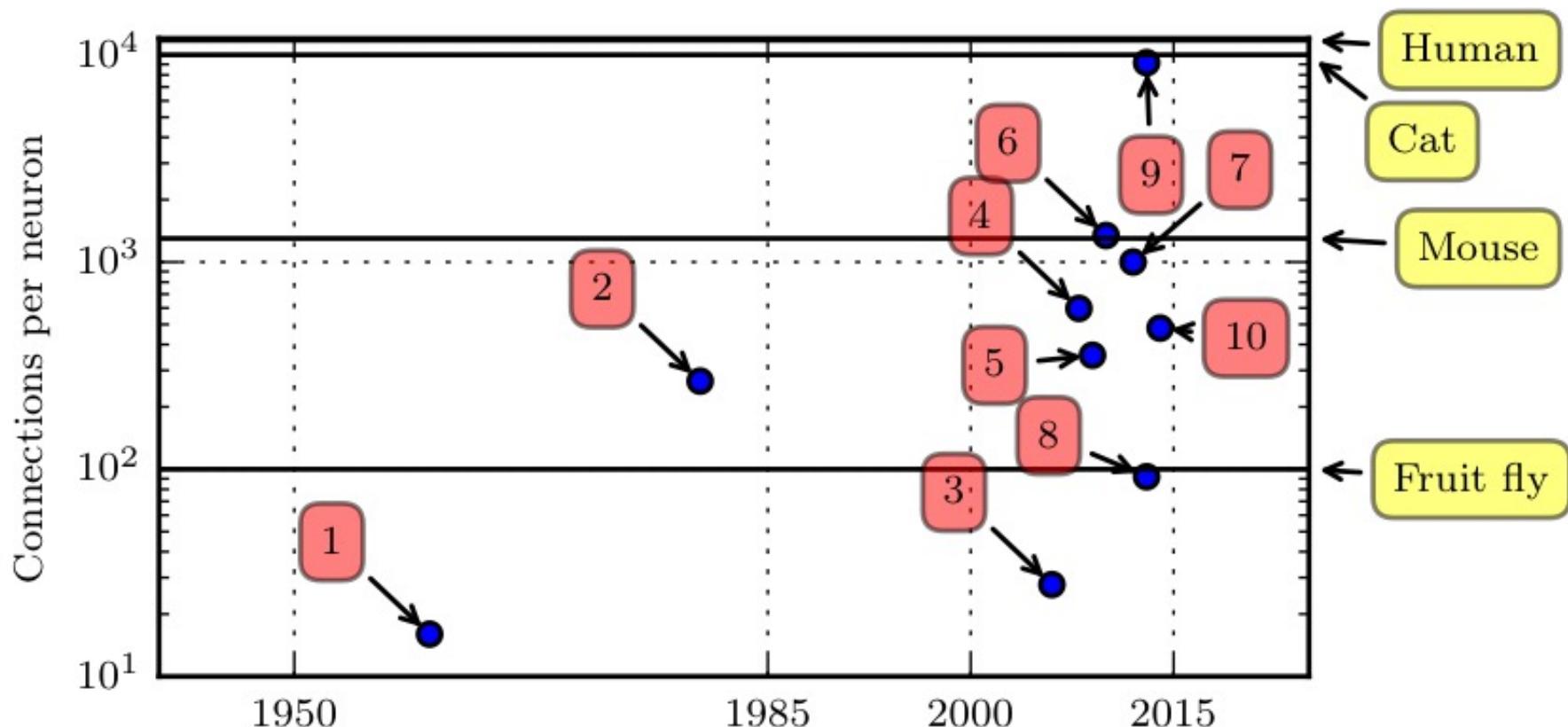
Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.

<https://blogs.nvidia.com/blog/2016/07/29/whats-difference-artificial-intelligence-machine-learning-deep-learning-ai/>

Historical Trends: Growing Datasets



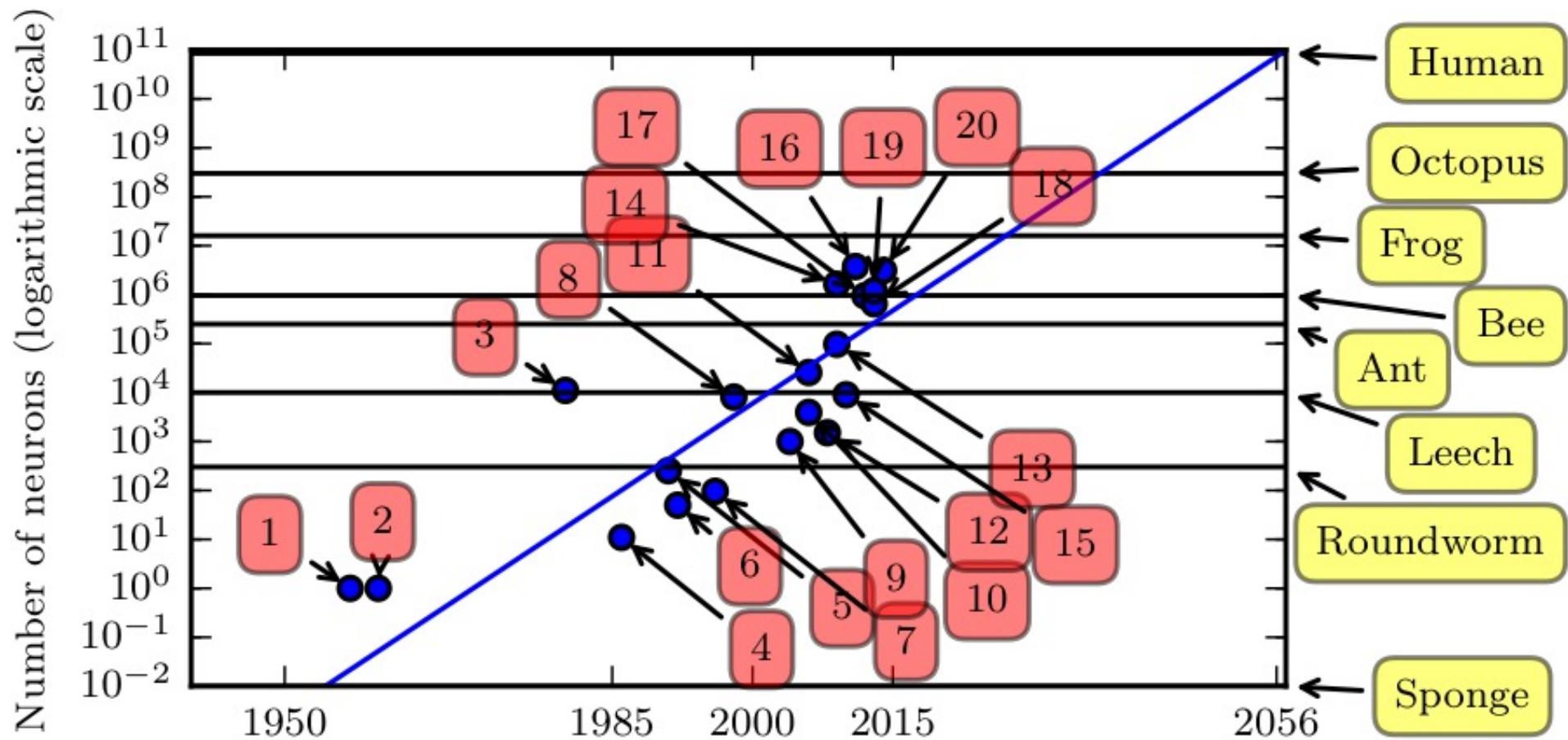
Historical Trends: Growing Connections per Neuron



9: COTS HPC unsupervised CNN (Coates *et al.*, 2013)

10: GoogLeNet (Szegedy *et al.*, 2014)

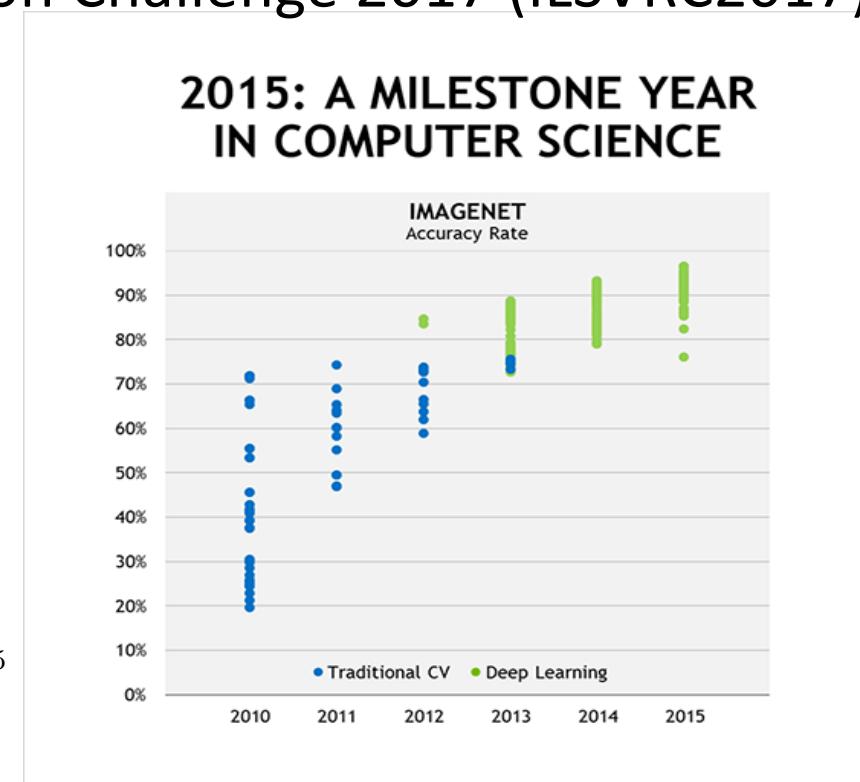
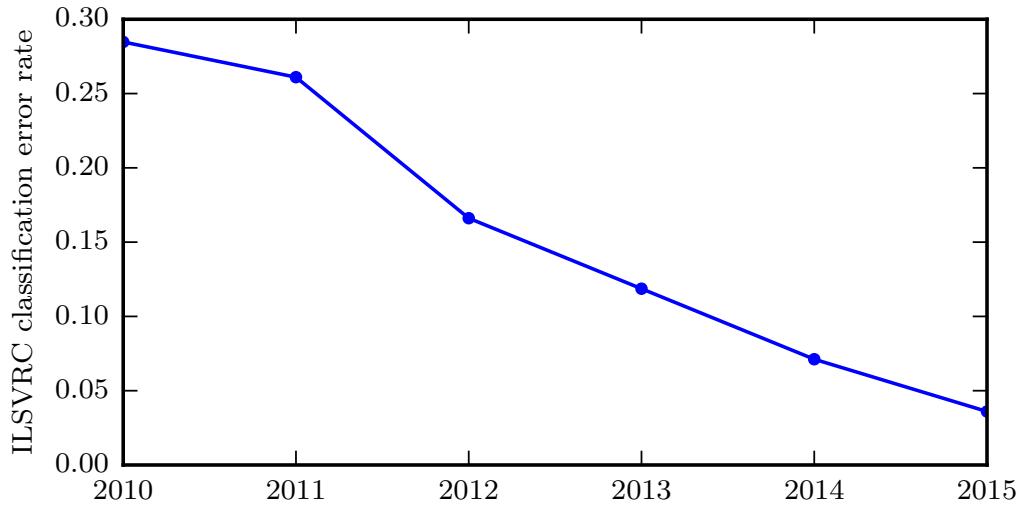
Historical Trends: Growing Number of Neurons



9: Echo state network (Jaeger and Haas, 2004)
20: GoogLeNet (Szegedy *et al.*, 2014)

Historical Trends: Increasing Accuracy

- DL has solved increasingly complicated applications with increasing accuracy:
 - Image Recognition and Object Recognition:
 - Large Scale Visual Recognition Challenge 2017 (ILSVRC2017)



Deep Learning x GPUs

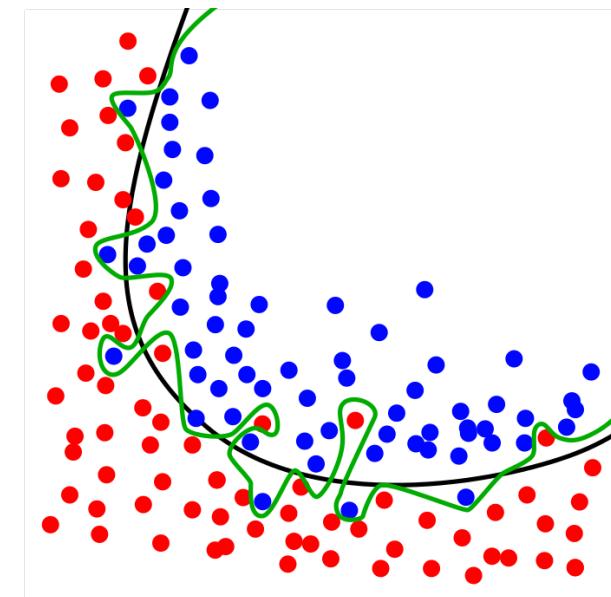
- What is the relation between DL and GPUs?

Deep Learning & GPU:

- Matrix operations
- Parallel computing paradigm

Problems with DNN (Deep MLP)

- **Overfitting:**
 - The more layers you have, the more degrees of freedom you have.
 - DNNs model rare dependencies in the training data.
- **Diffusion of Gradient:** error attenuates as it propagates to early layers.
 - **Early layers never learn!**



Main Deep Learning Architectures

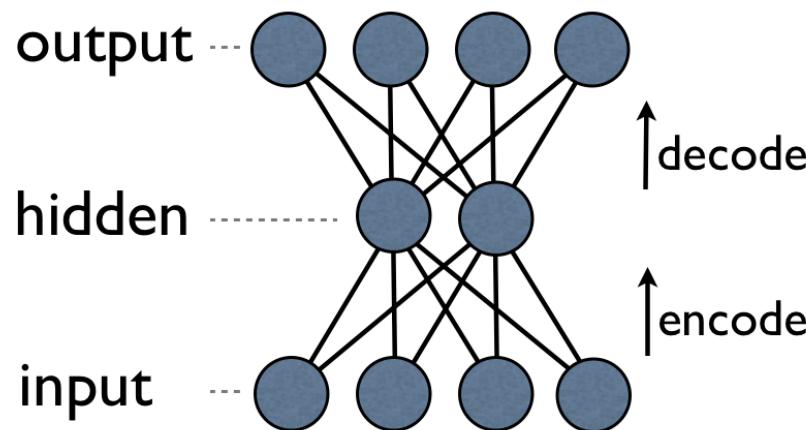
- Deep Belief Networks / Autoencoders
 - Greedy layer-wise pretraining, by Hinton *et al.*, 2006
- Deep Convolutional Neural Networks
 - La Net, by Le Cun et al, 1998.
- Deep Recurrent Networks
 - Long short-term Memory, by Sepp&Jurgen, 97.

Greedy layer-wise pretraining,
by Hinton *et al.*, 2006.

AUTOENCODERS

Autoencoders

- Autoencoders have just **one layer**.
- The aim of an autoencoder is to **learn a representation (encoding)** for a set of data, typically for the purpose of **dimensionality reduction**.



Autoencoders

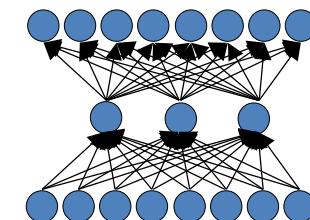
- An auto-encoder is trained, with an absolutely standard weight-adjustment algorithm to **reproduce the input.**
- By making this happen with (many) fewer units than the inputs, this forces the ‘hidden layer’ units to become **good feature detectors.**

Representation learning (hidden layer)

Entrada		Saída hs		Saída
10000000	→	.89 .04 .08	→	10000000
01000000	→	.15 .99 .99	→	01000000
00100000	→	.01 .97 .27	→	00100000
00010000	→	.99 .97 .71	→	00010000
00001000	→	.03 .05 .02	→	00001000
00000100	→	.01 .11 .88	→	00000100
00000010	→	.80 .01 .98	→	00000010
00000001	→	.60 .94 .01	→	00000001

1 0 0
0 1 1

Intermediate representation: "discover" binary code !!!

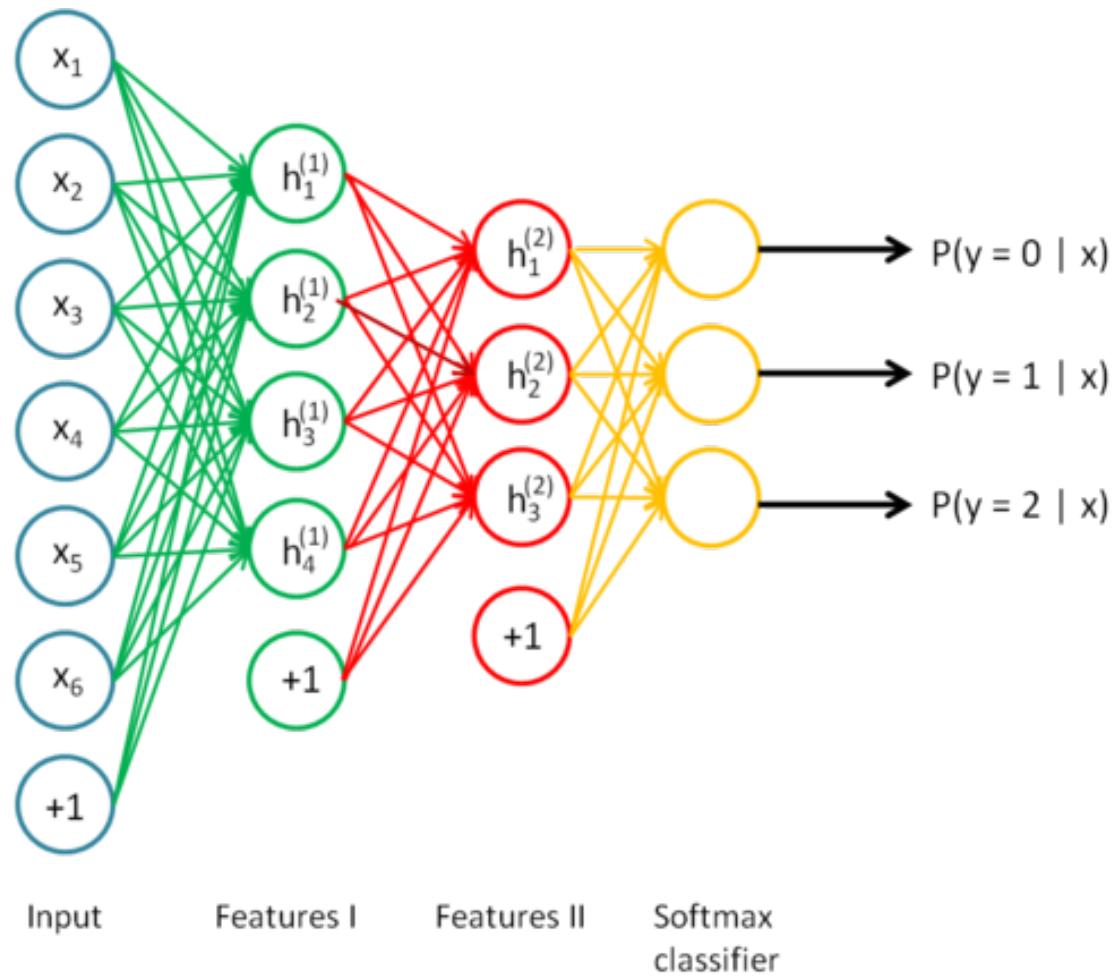


Greedy Layer-Wise Training

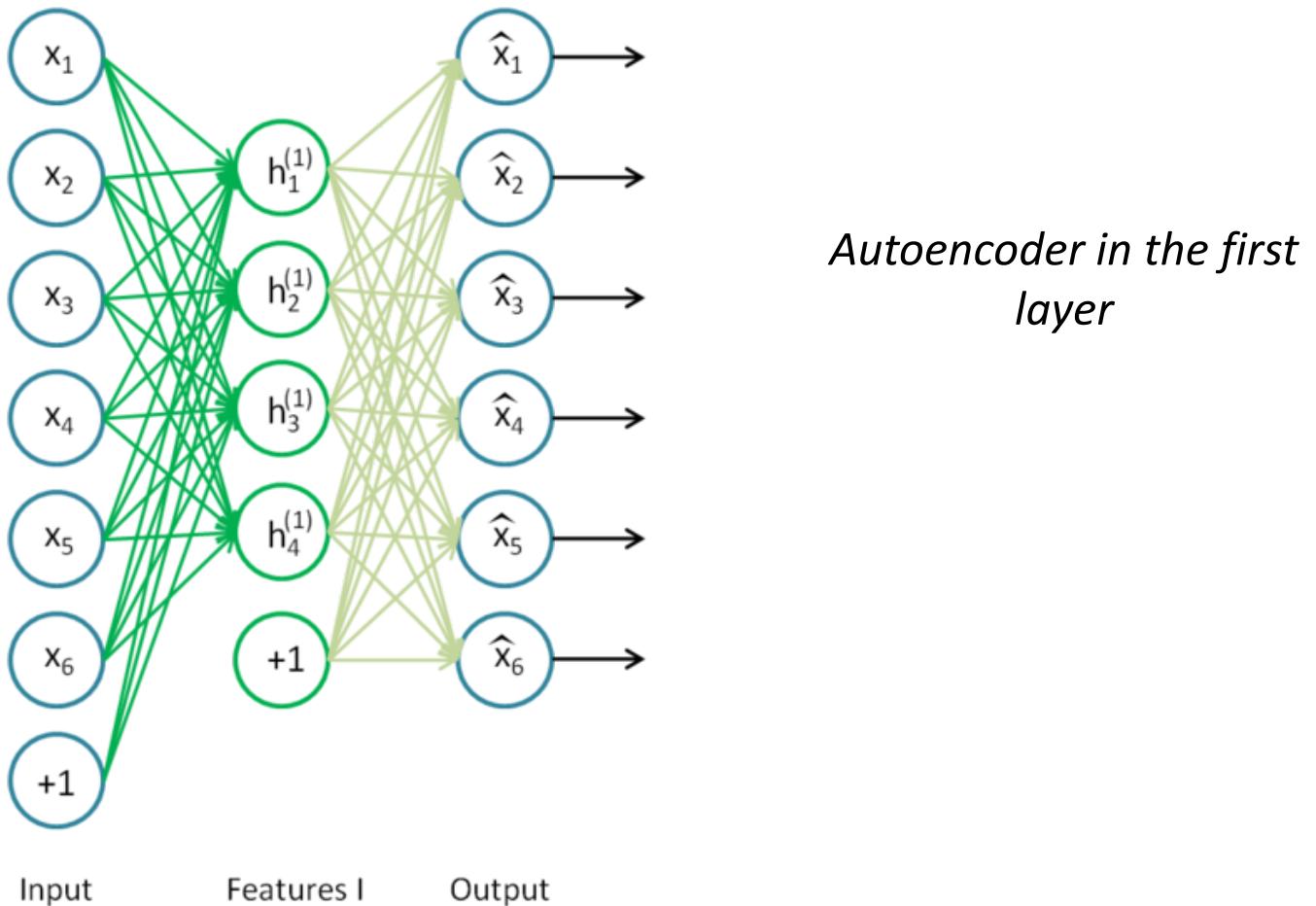
Geoffrey E. Hinton and Simon Osindero and Yee-Whye Teh. A fast learning algorithm for deep belief nets, 2006

1. Train first layer using your data without the labels.
2. Then freeze the first layer parameters and start training the second layer using the output of the first layer as the input to the second layer.
3. Repeat this for as many layers as desired:
 - This builds our set of robust features
4. Use the outputs of the final layer as inputs to a supervised layer/model and train the last supervised layer(s) (leave early weights frozen)
5. Unfreeze all weights and fine tune the full network by training with a supervised approach, given the pre-processed weight settings.

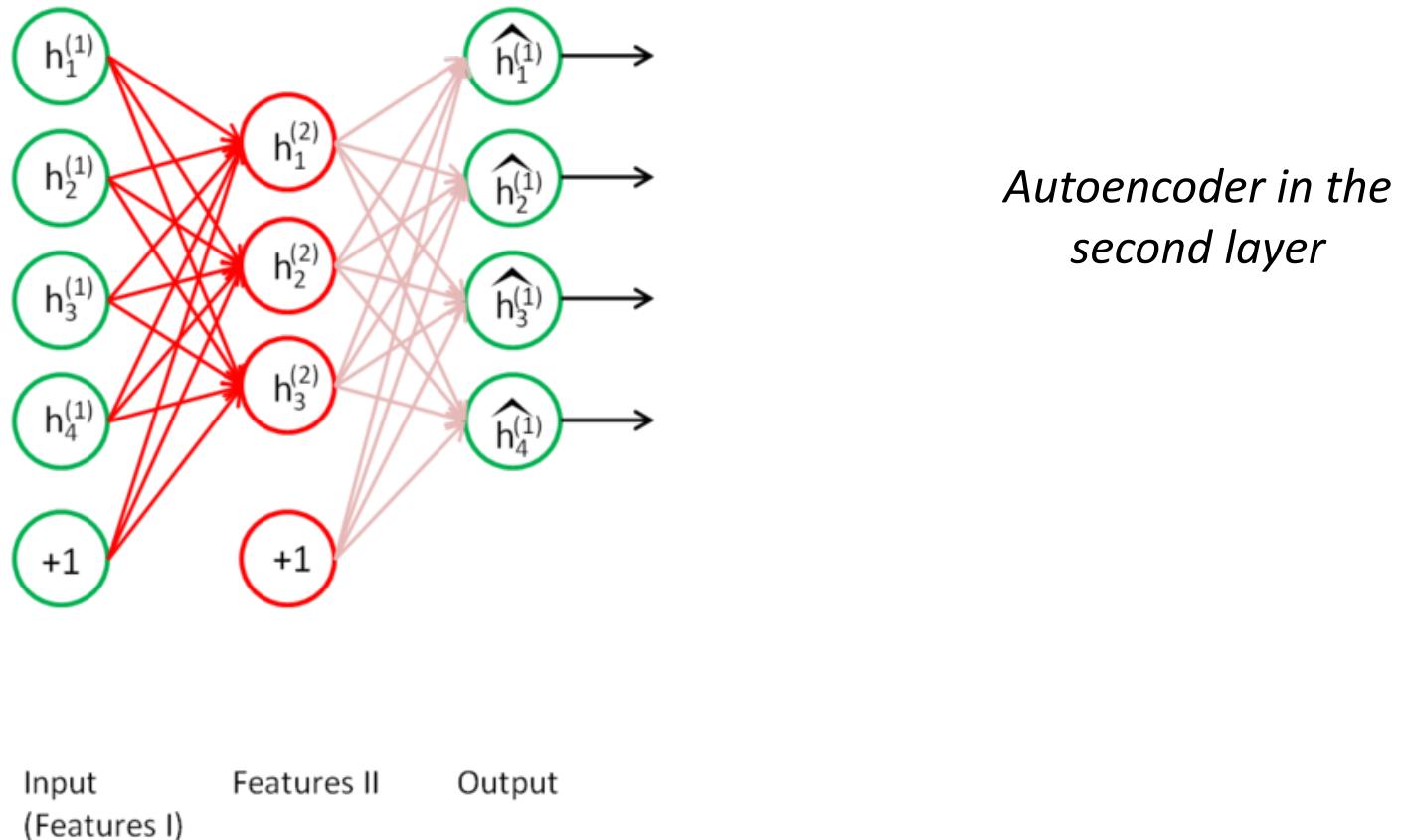
How to train this network?



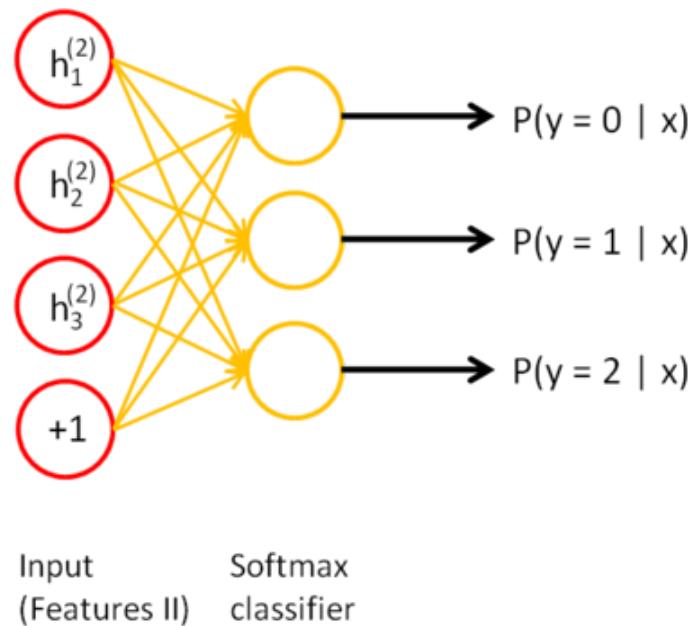
How to train this network?



How to train this network?

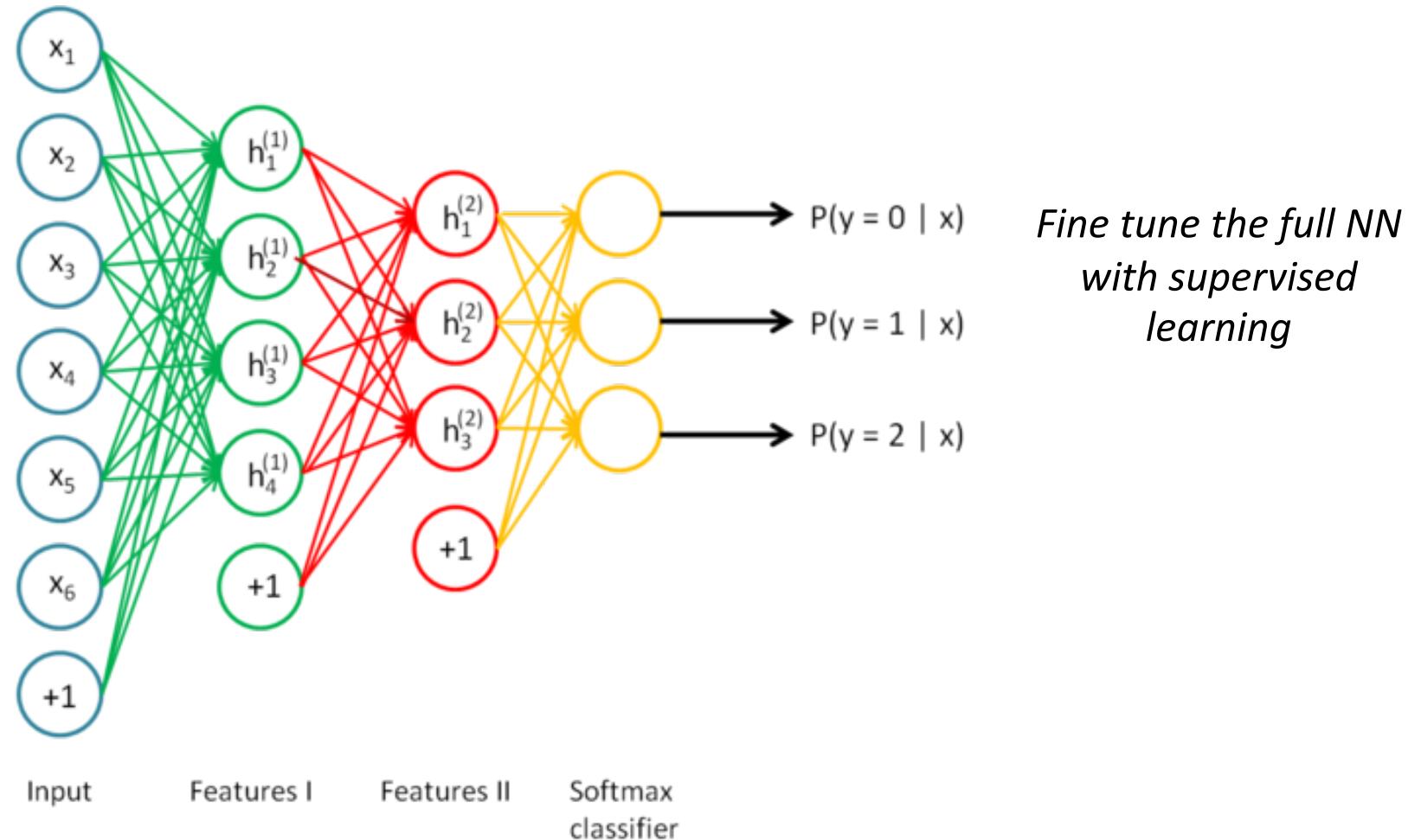


How to train this network?



*Supervised learning in
the last layer*

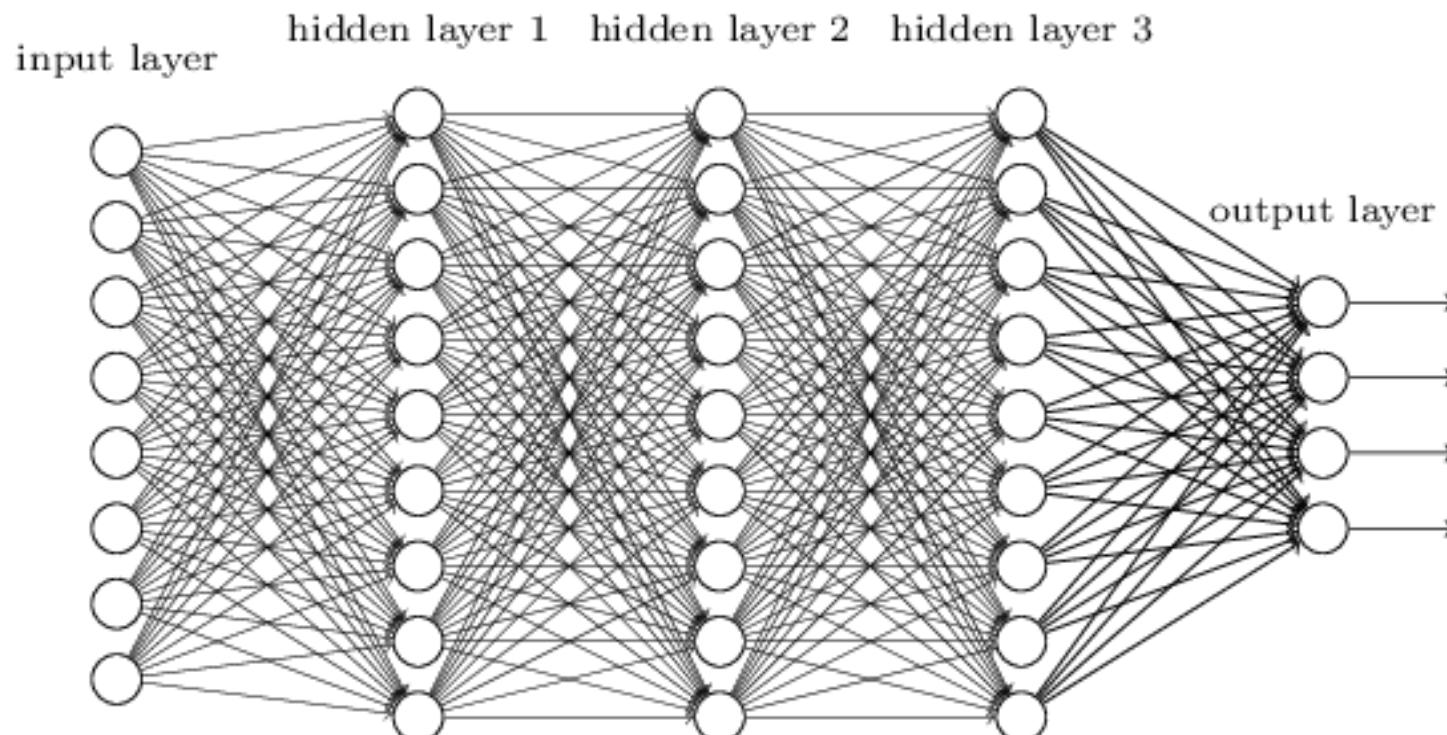
How to train this network?



CONVOLUTIONAL NEURAL NETWORKS

CNN

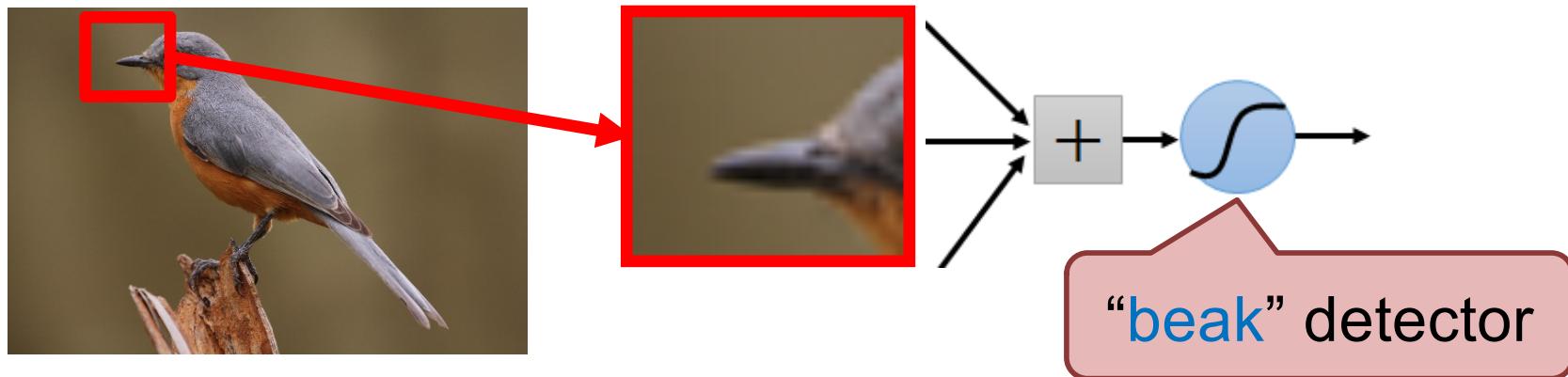
- We know it is good to learn a small model.
- From this fully connected model, do we really need all the edges?
- Can some of these be shared?



Consider learning an image:

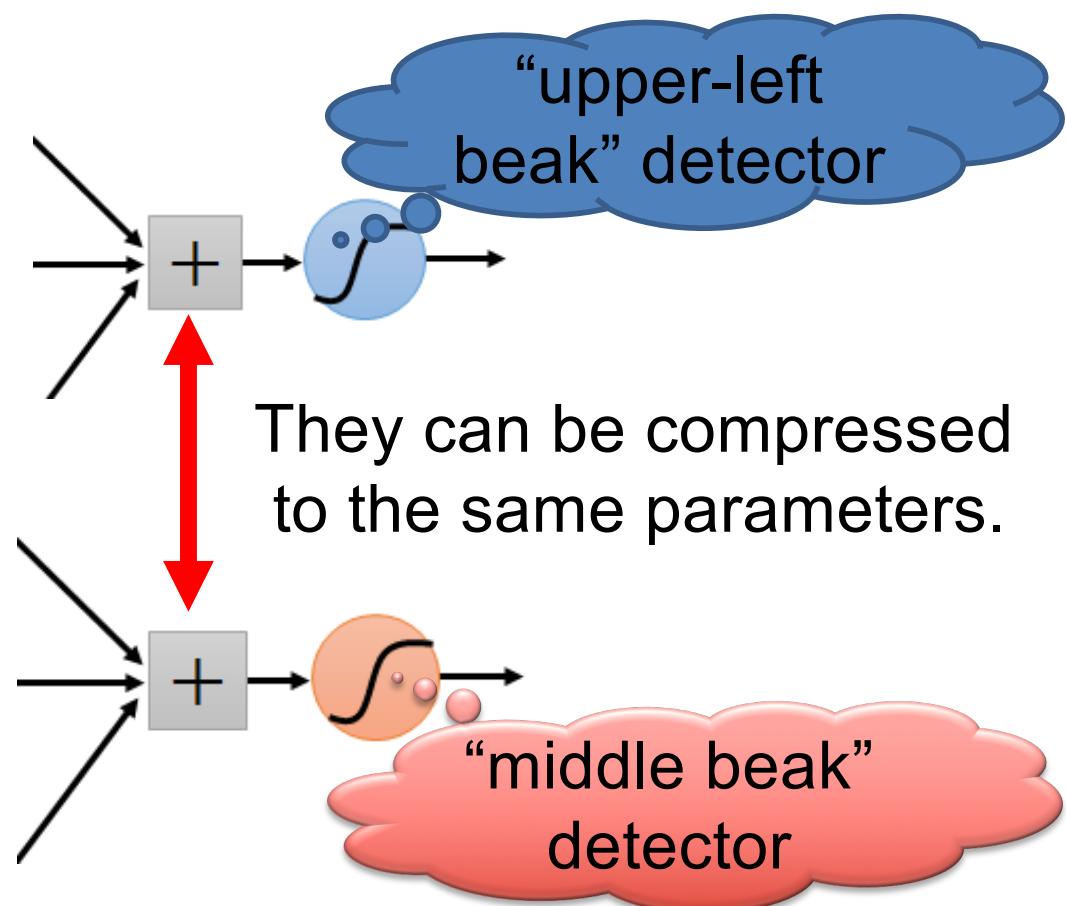
- Some patterns are much smaller than the whole image

Can represent a small region with fewer parameters



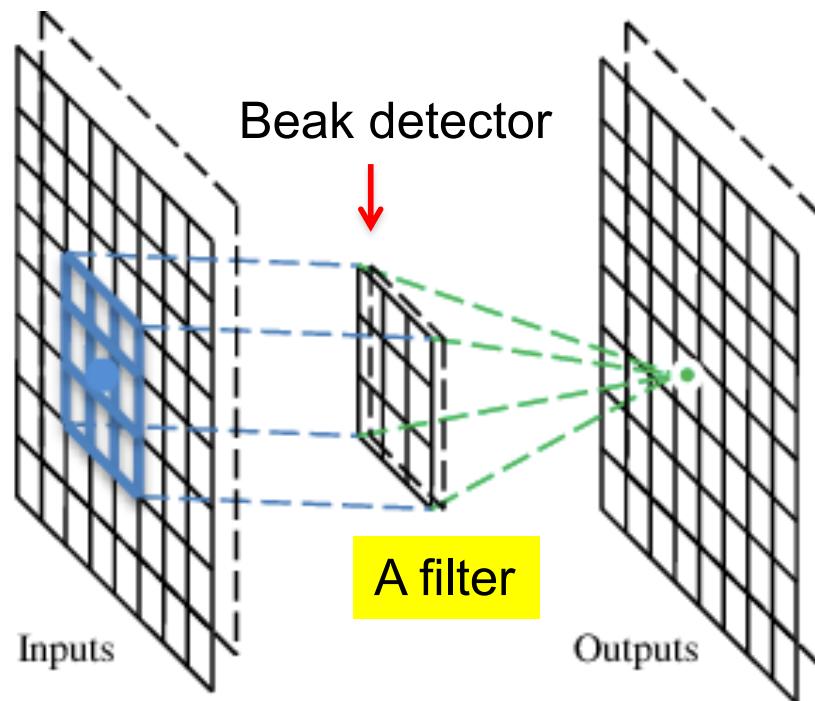
Same pattern appears in different places:
They can be compressed!

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

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).

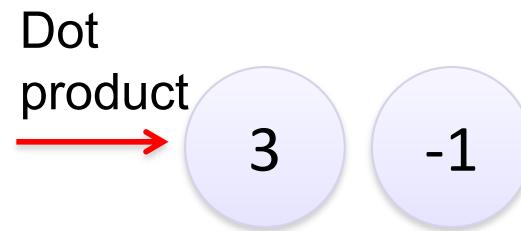
Convolution

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

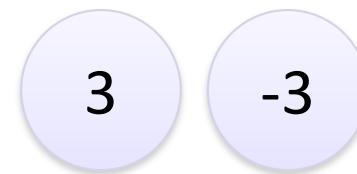
Convolution

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

Convolution

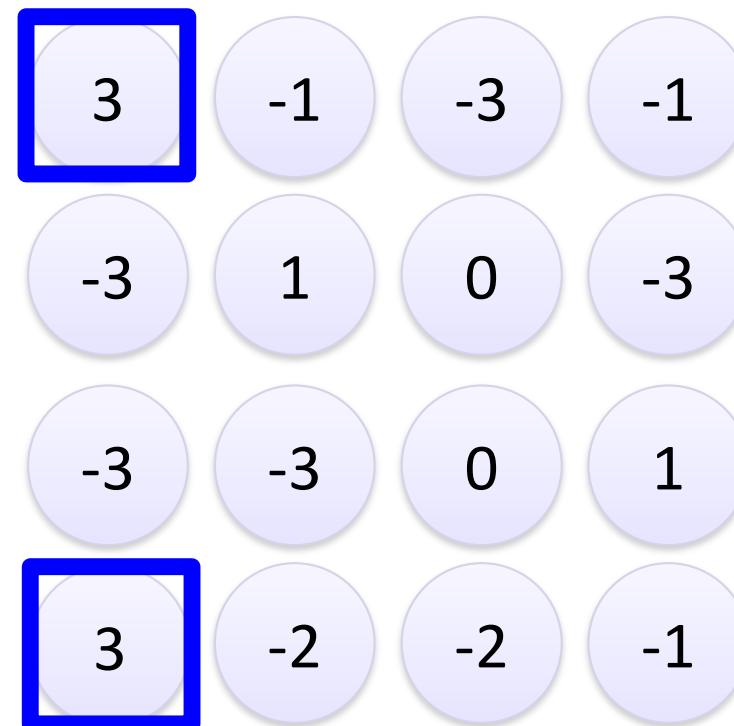
stride=1

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

6 x 6 image

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

Filter 1



Convolution

stride=1

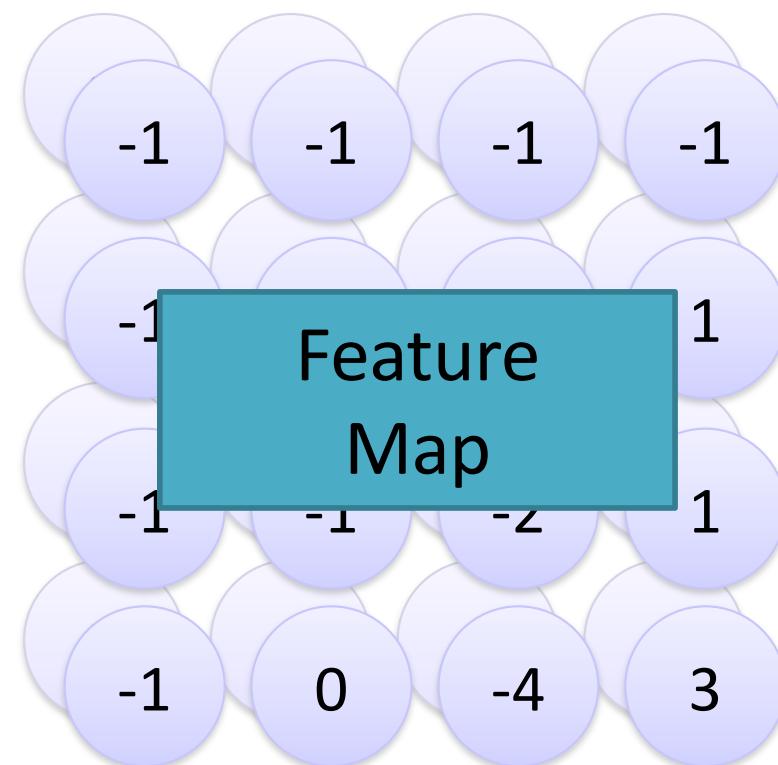
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

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

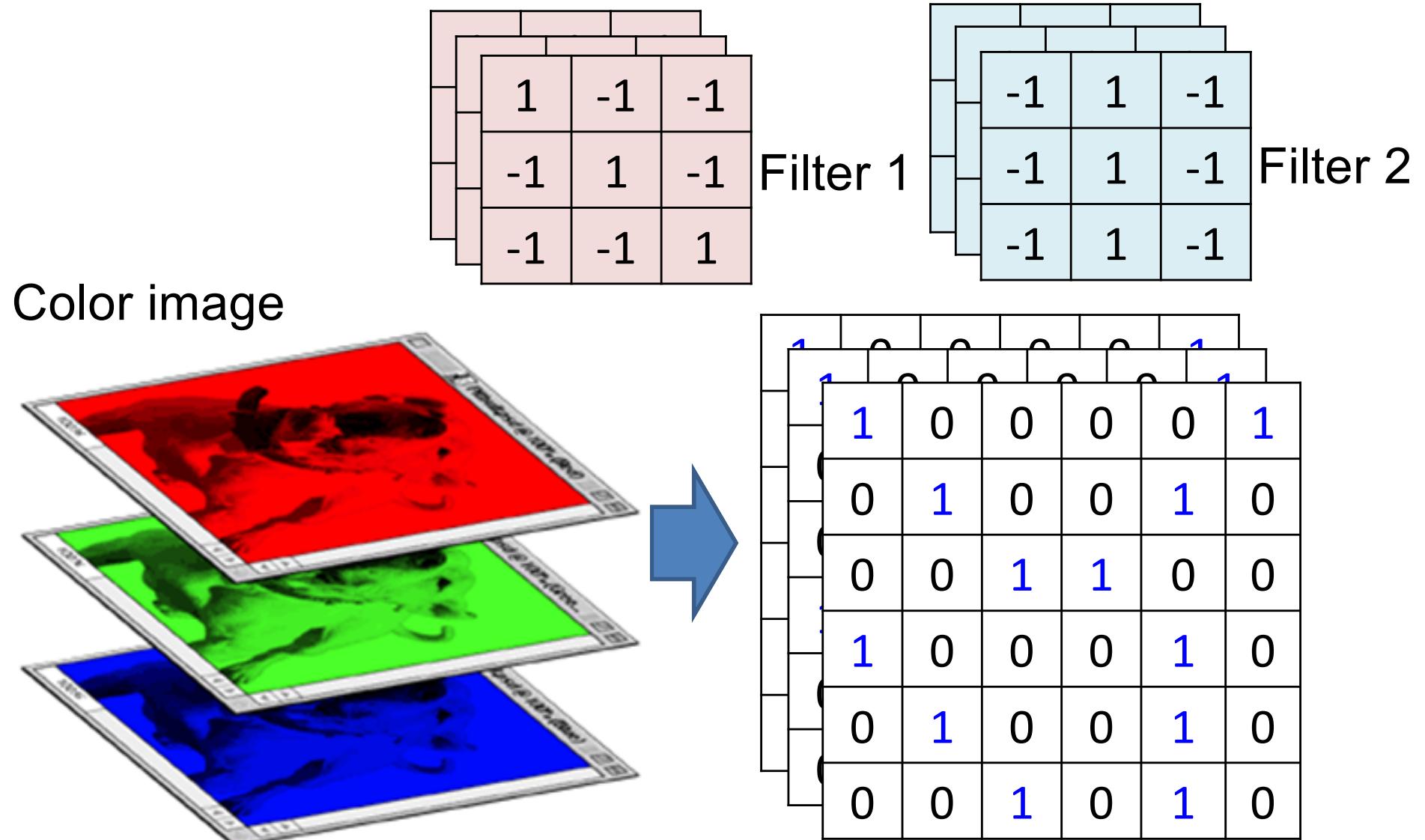
Filter 2

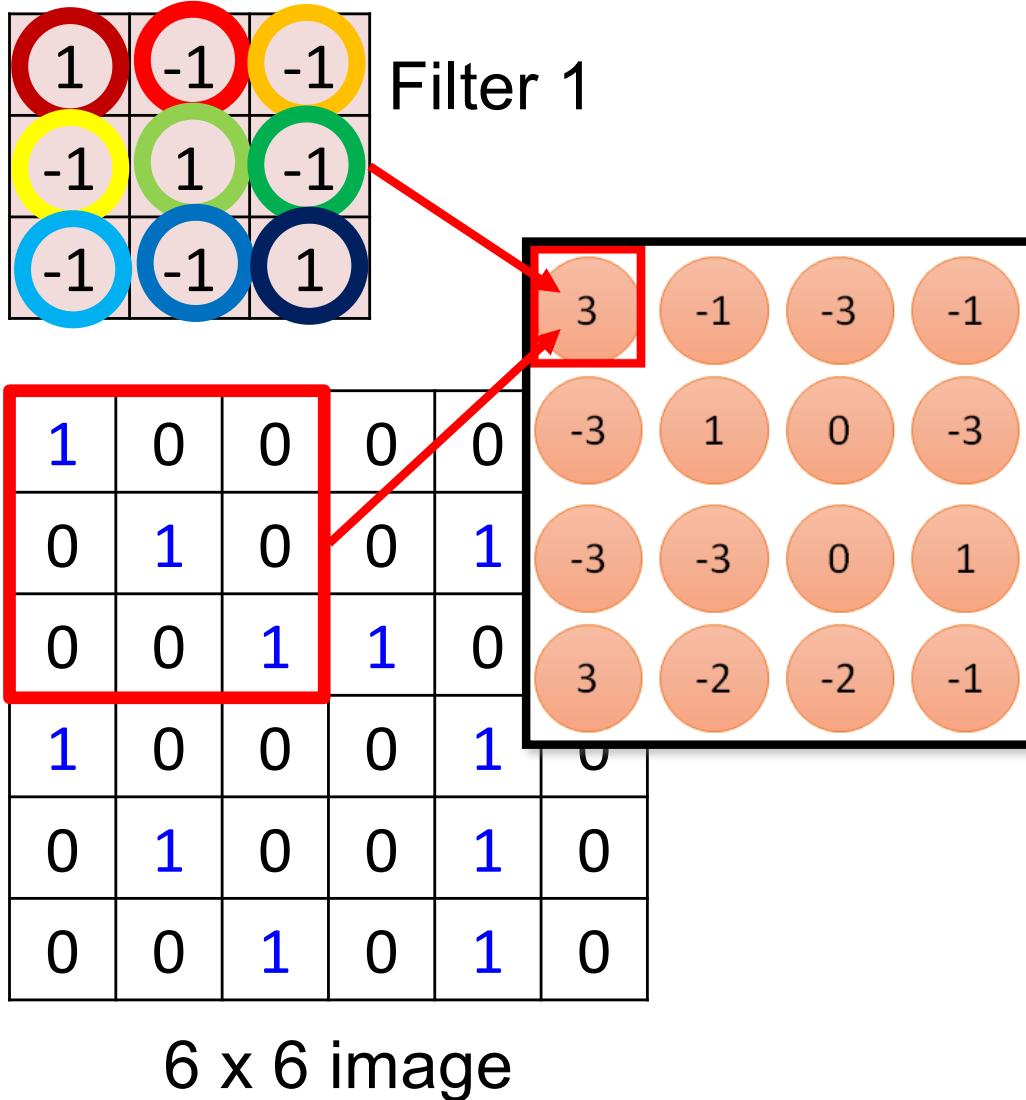
Repeat this for each filter



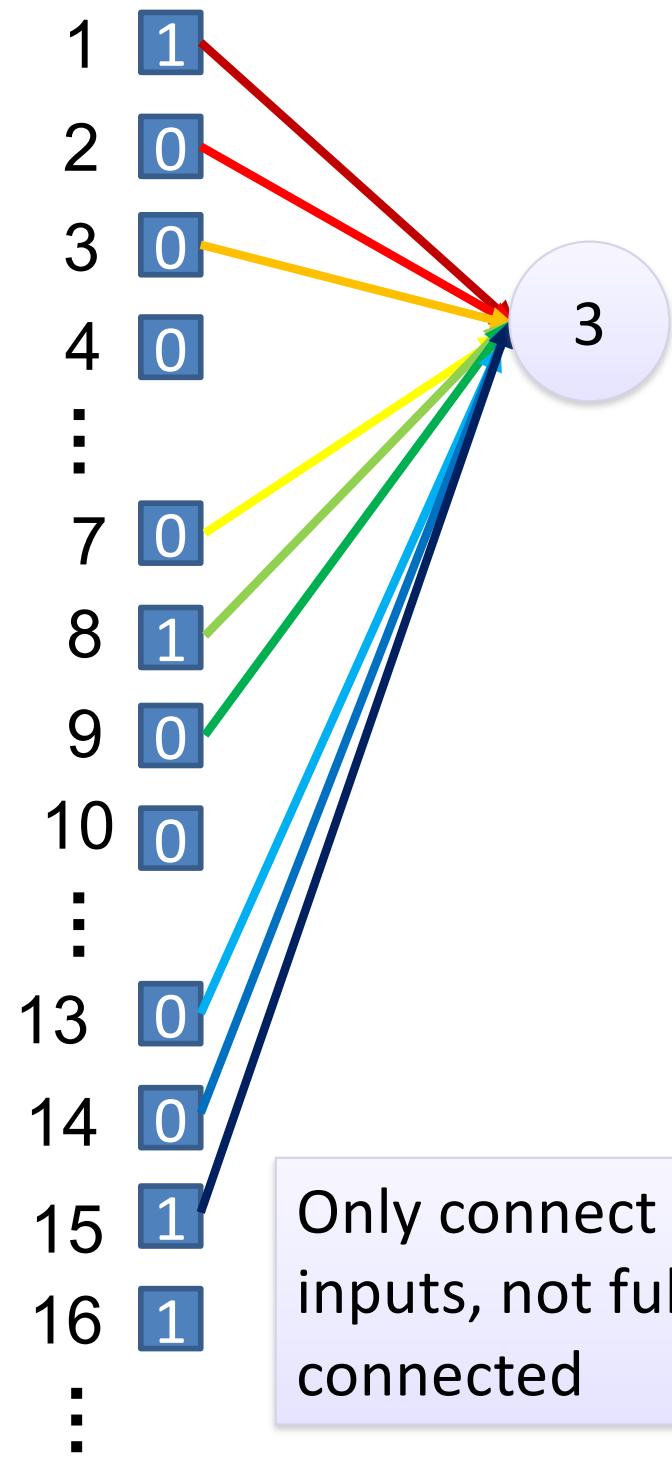
Two 4 x 4 images
Forming 2 x 4 x 4 matrix

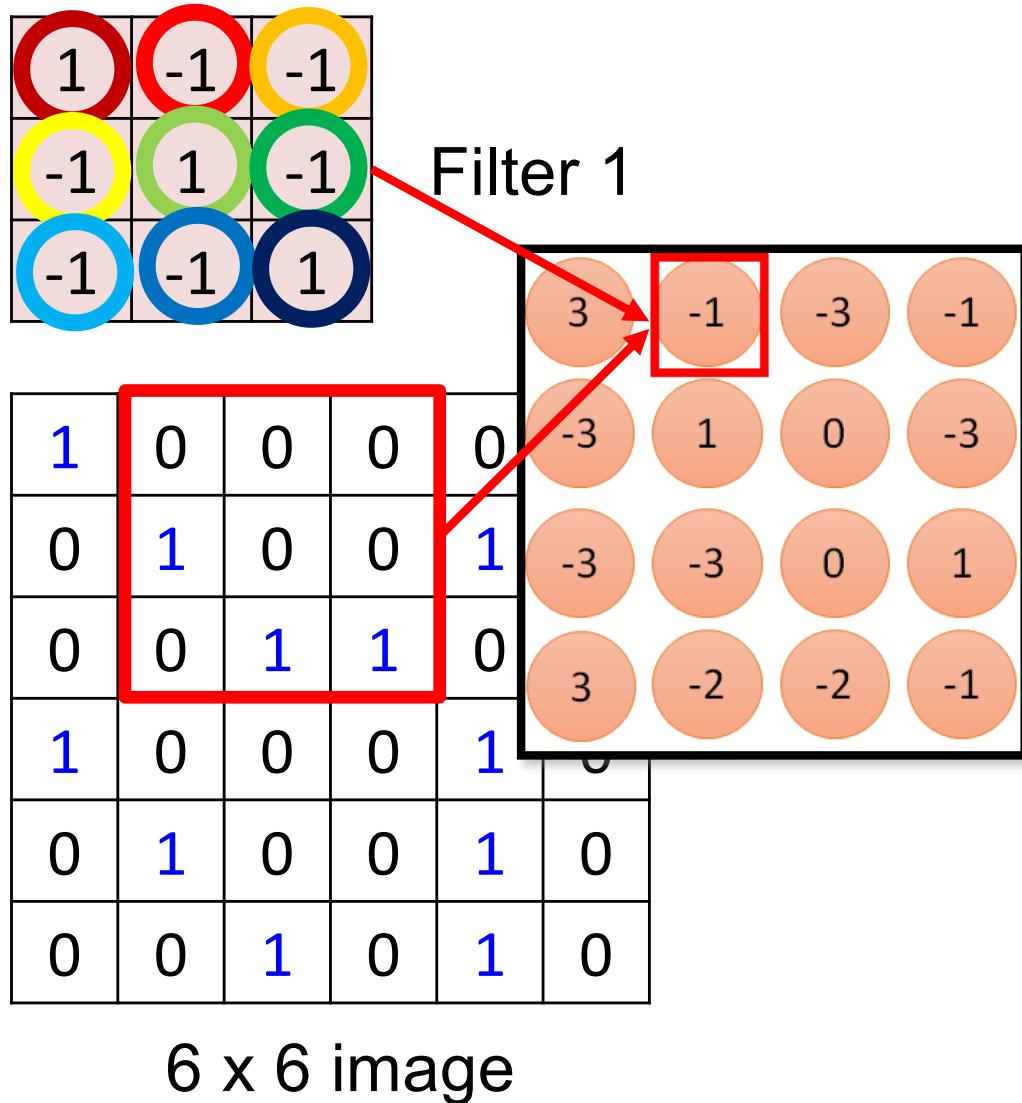
Color image: RGB 3 channels





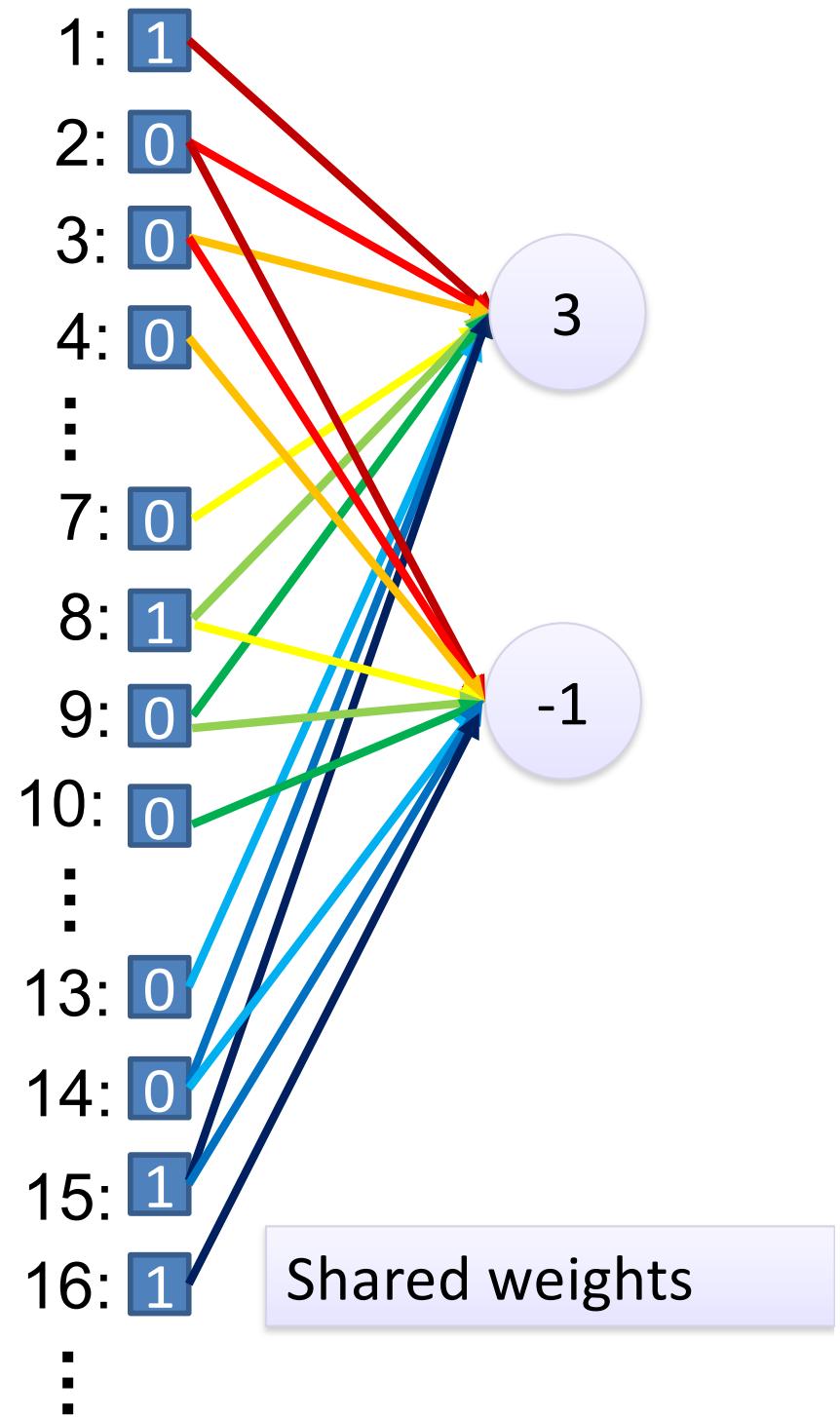
fewer parameters!





Fewer parameters

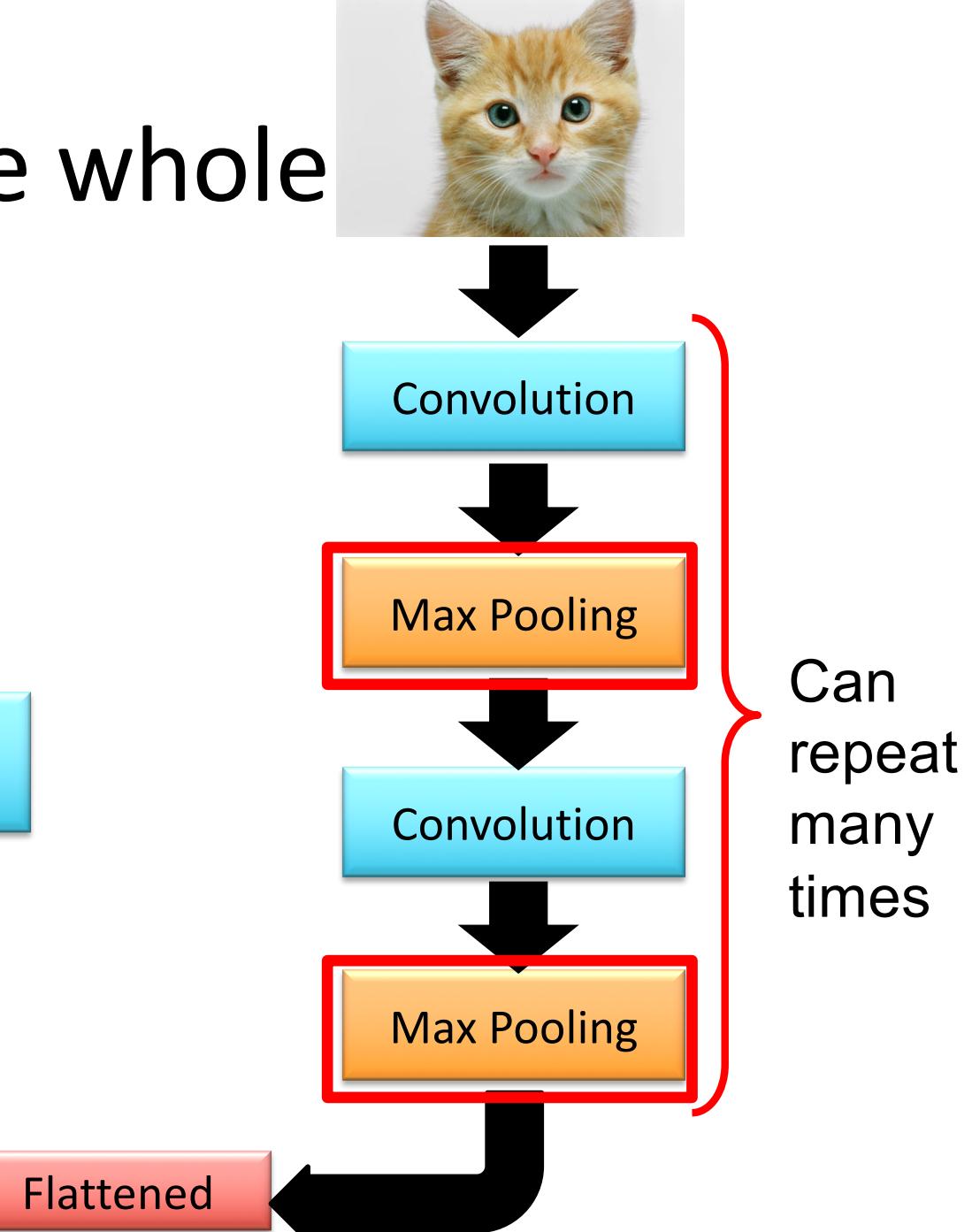
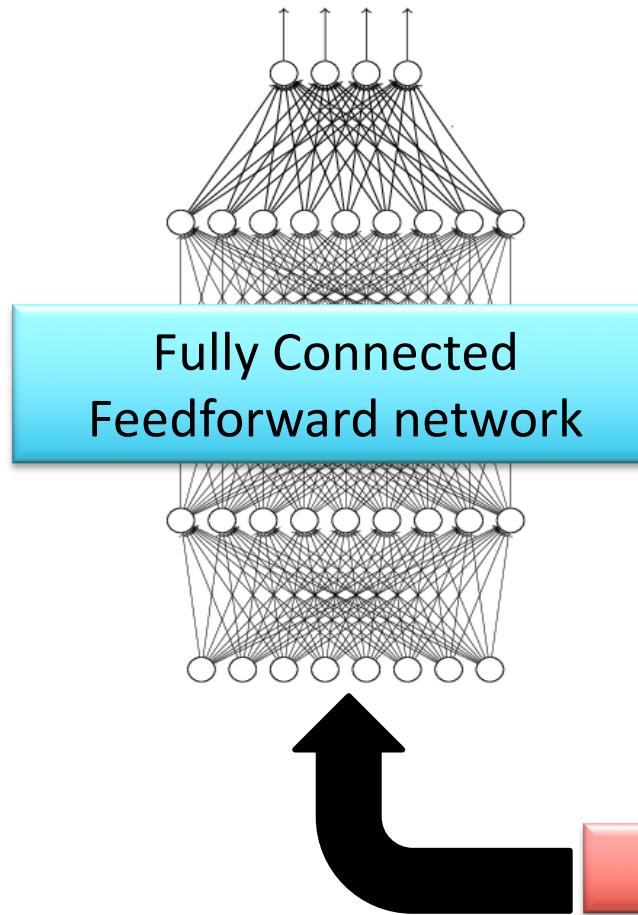
Even fewer parameters



The whole



cat dog



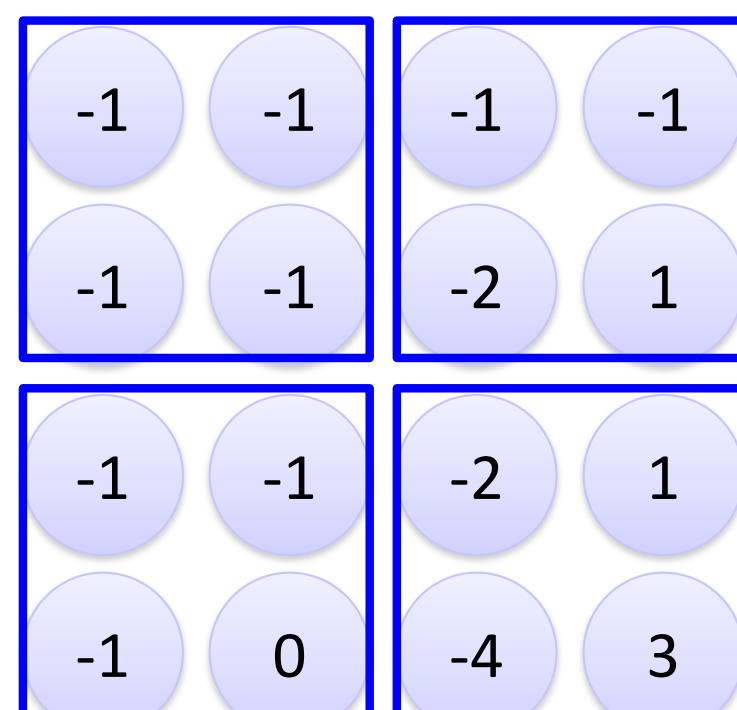
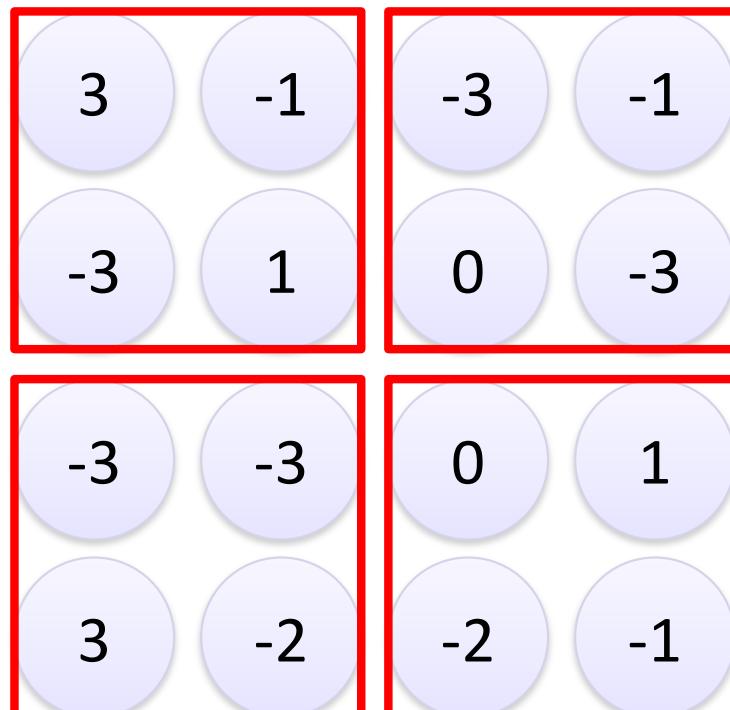
Max Pooling

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



bird

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

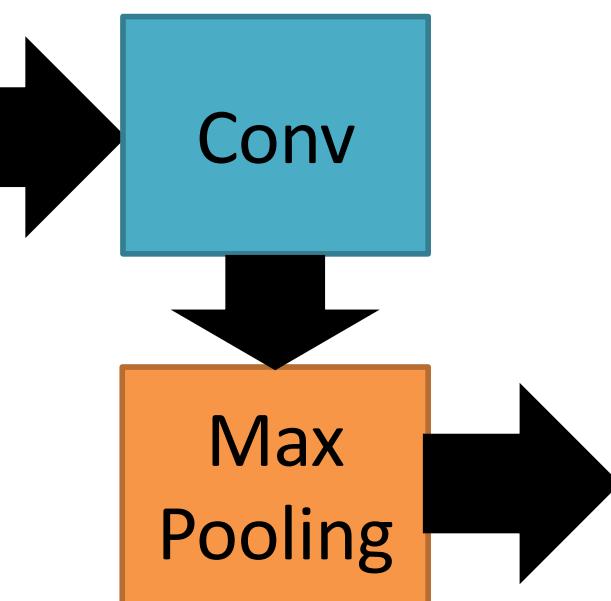
A CNN compresses a fully connected network in two ways:

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

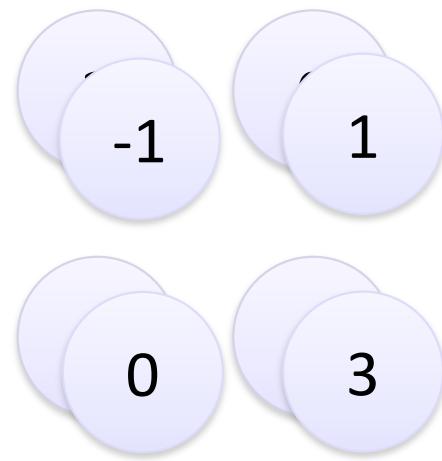
Max Pooling

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



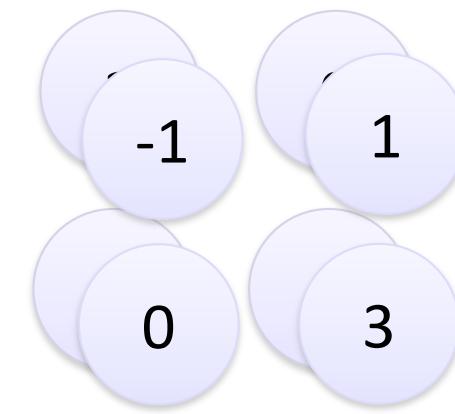
New image
but smaller



2 x 2 image

Each filter
is a channel

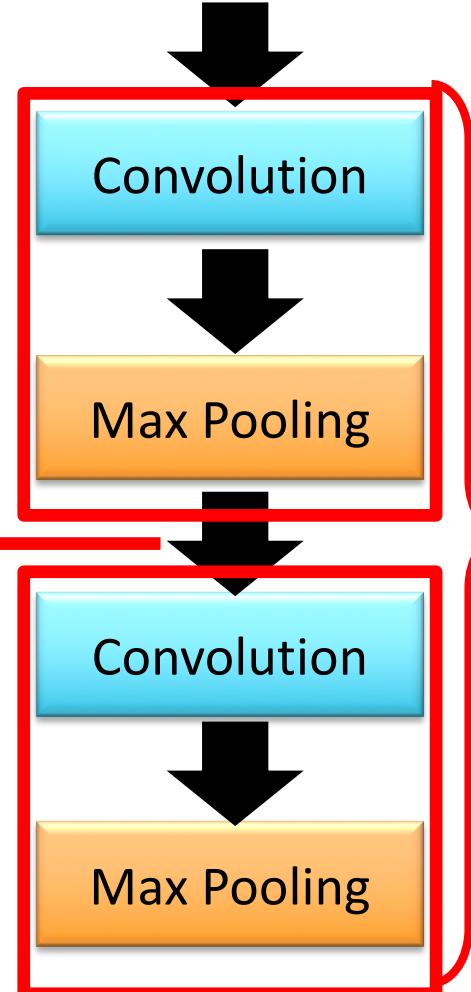
The whole



A new image

Smaller than the original image

The number of channels
is the number of filters

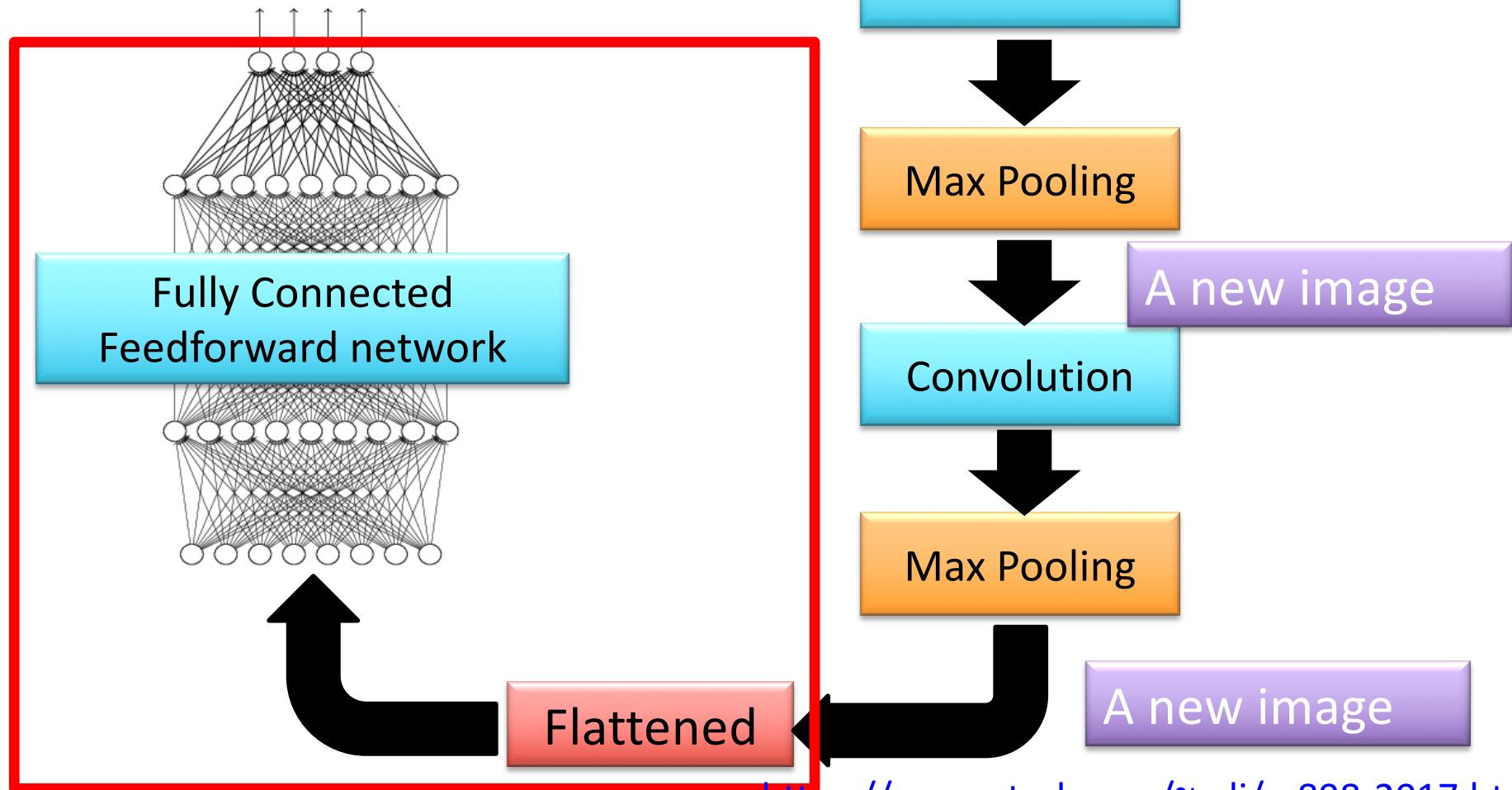


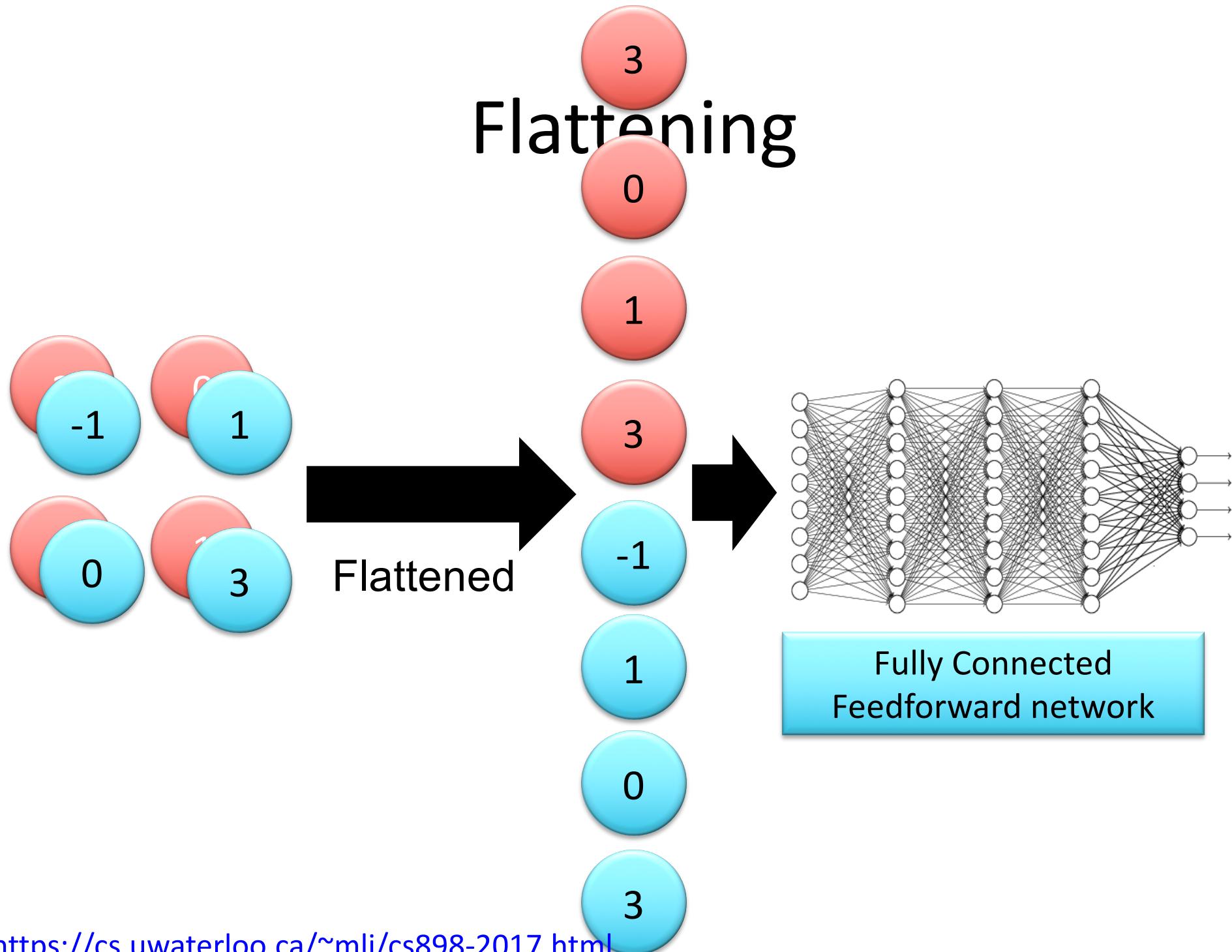
Can
repeat
many
times

The whole

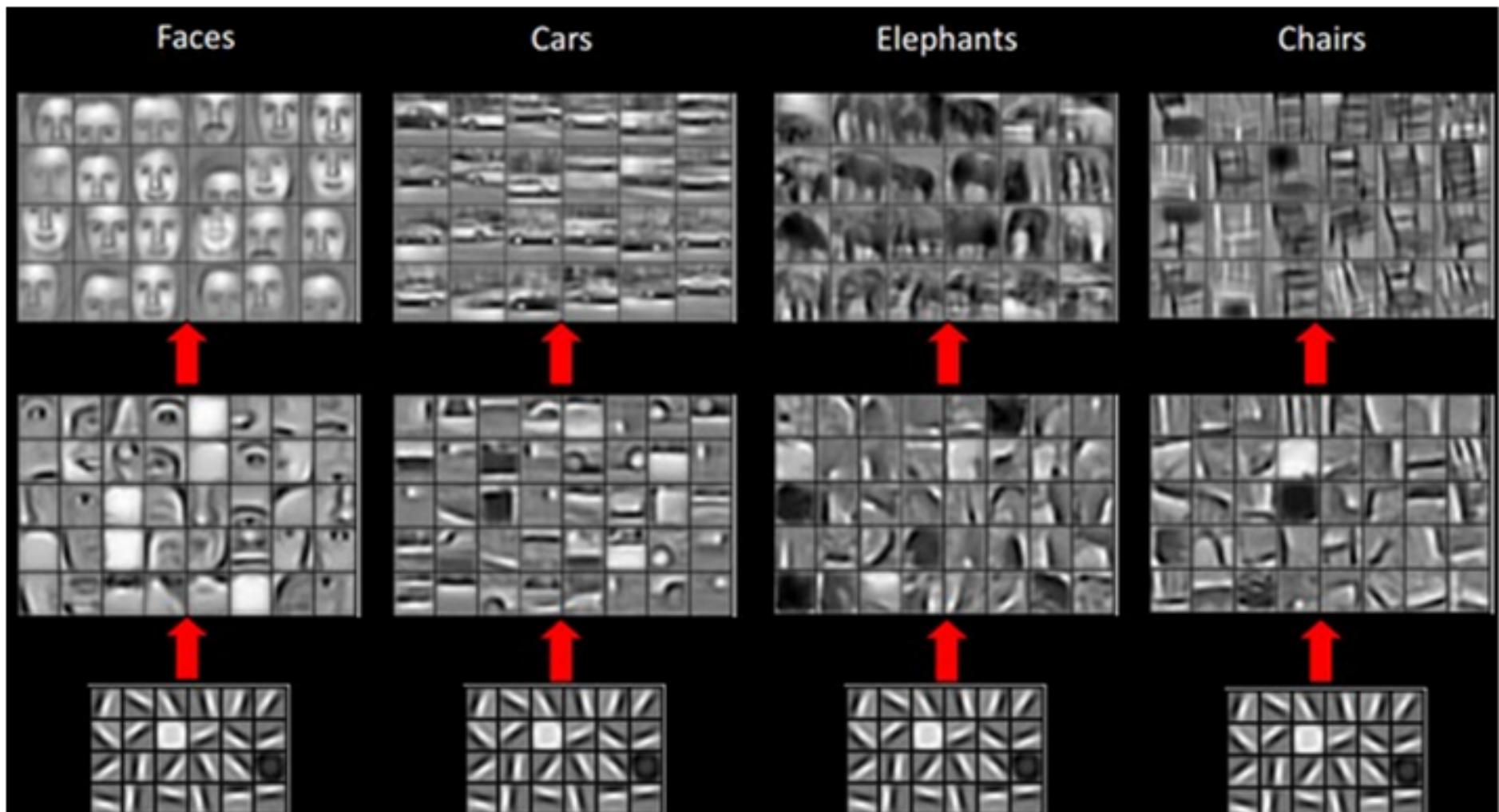


cat dog





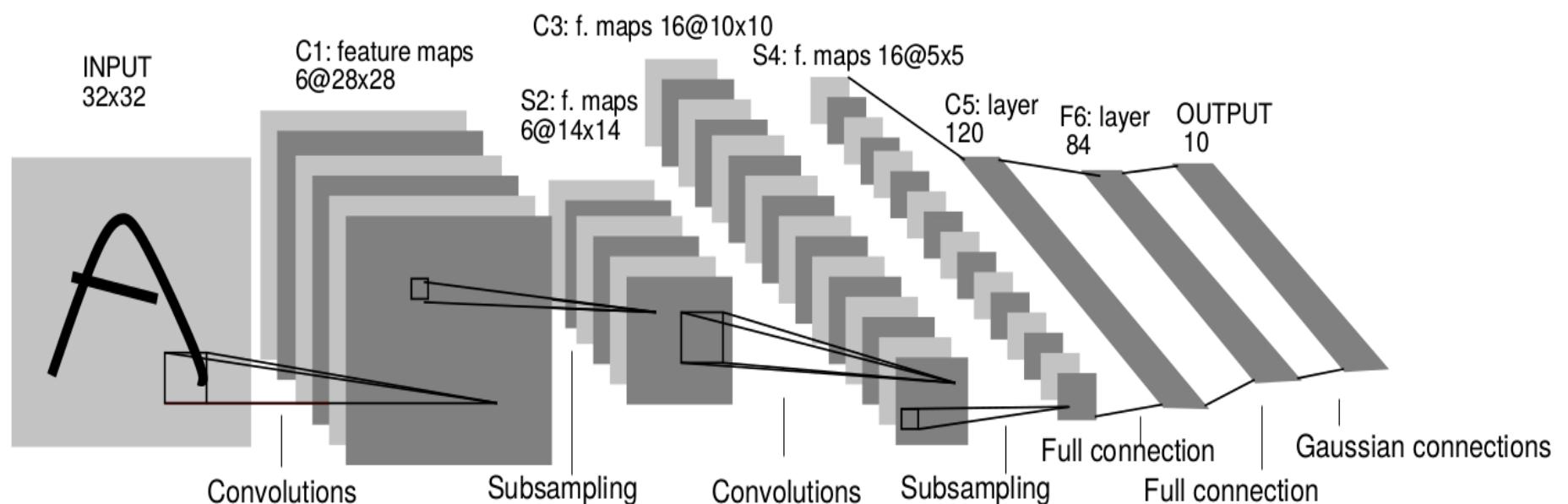
CNNs



<http://stats.stackexchange.com/questions/146413/why-convolutional-neural-networks-belong-to-deep-learning>

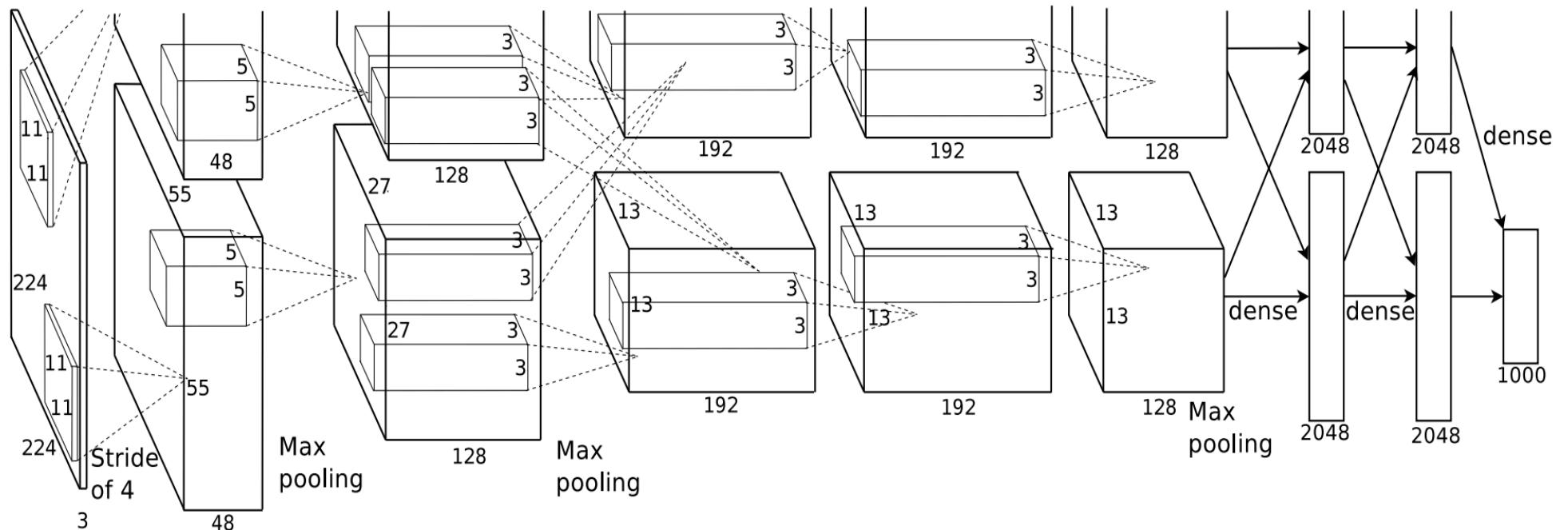
LeNet-5

LeCun, Bottou, Bengio & Haffner, 1998



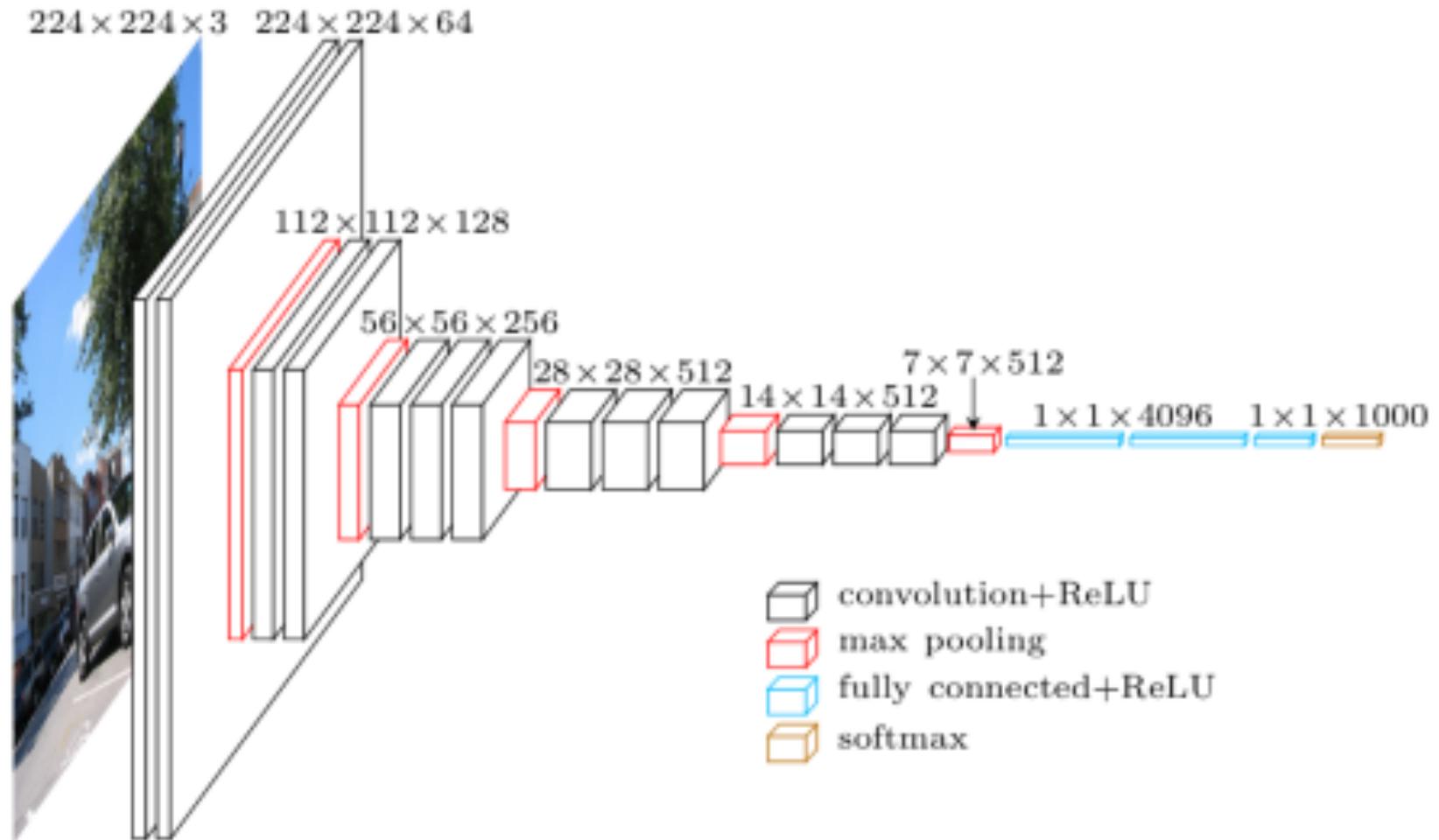
Alexnet

Krizhevsky, Sutskever & Hinton, 2012



VGG

Simonyan & Zisserman, 2014



DeepFace

Taigman, Yang, Ranzato & Wolf, 2014

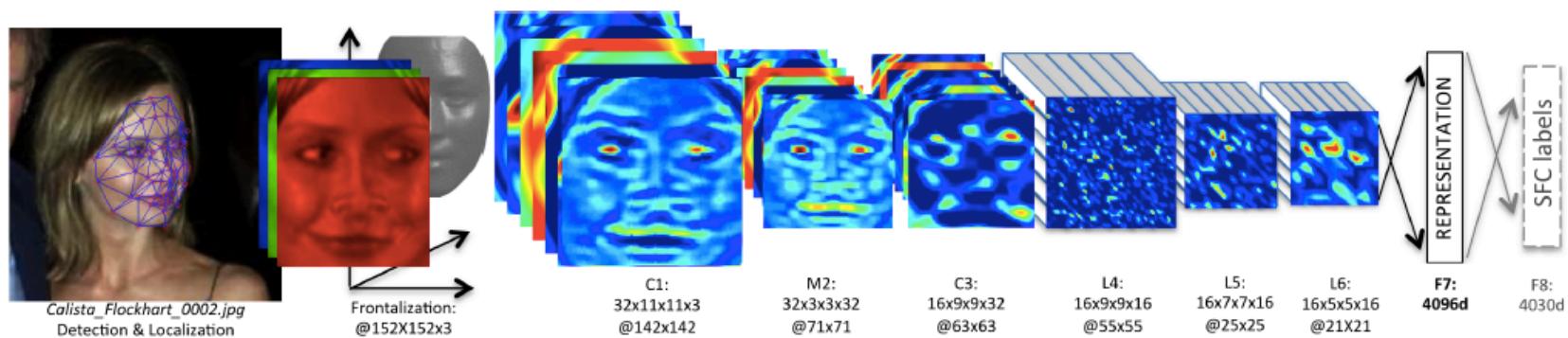
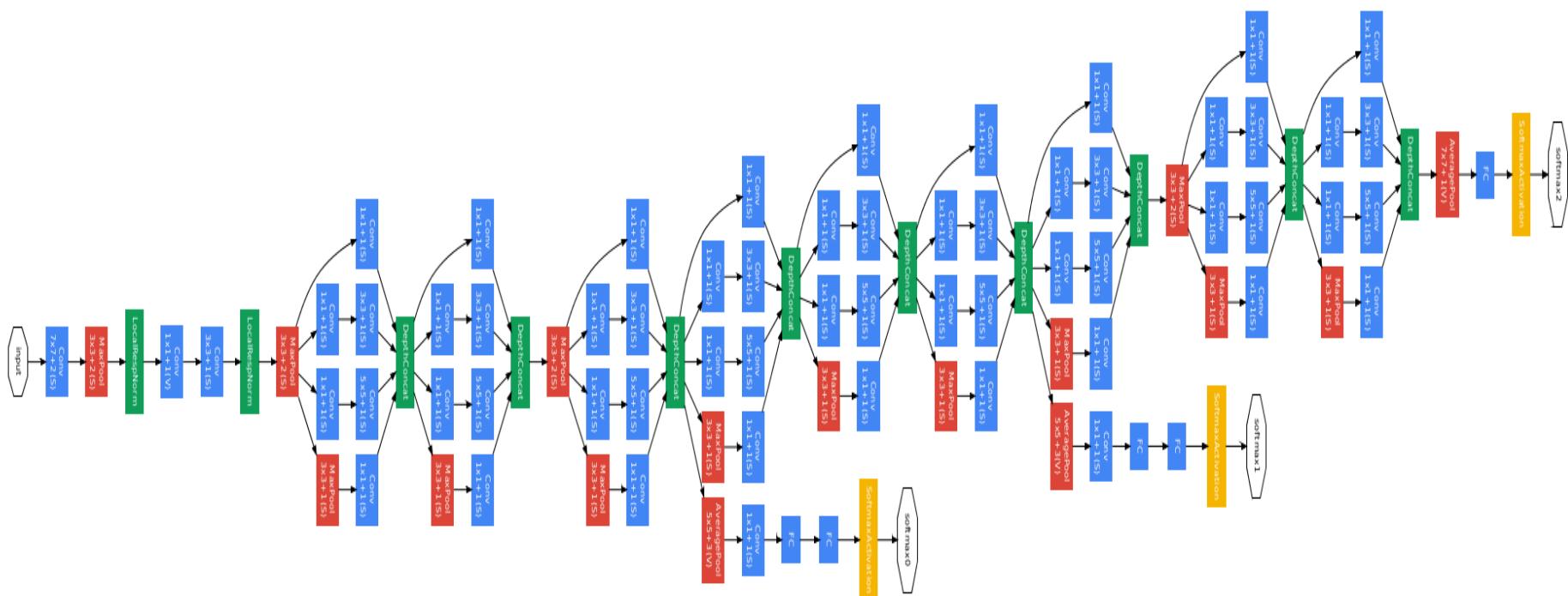


Figure 2. Outline of the **DeepFace** architecture. A front-end of a single convolution-pooling-convolution filtering on the rectified input, followed by three locally-connected layers and two fully-connected layers. Colors illustrate outputs for each layer. The net includes more than 120 million parameters, where more than 95% come from the local and fully connected layers.

DeepFace: Closing the Gap to Human-Level Performance in Face Verification
<https://research.fb.com/publications/deepface-closing-the-gap-to-human-level-performance-in-face-verification/>

GoogleLeNet

Szegedy et al, 2015.



Conclusion

- CNN: Special purpose net – Just for images or problems with strong **grid-like local spatial/temporal correlation**
- Once trained on one problem could use same net (often tuned) for a new similar problem – general creator of vision features
- Autoencoder could be used to find initial parameters
- **Lots of hand crafting and tuning** to find the right recipe of receptive fields, layer interconnections, etc.
 - Lots more Hyperparameters than standard nets, and even than other deep networks, since the structures of CNNs are more handcrafted
 - CNNs getting wider and deeper with speed-up techniques (e.g. GPU, ReLU, etc.) and lots of current research, excitement, and success

Fully Supervised Deep Learning

- Much recent success in doing fully supervised deep learning with extensions which diminish the effect of early learning difficulties (unstable gradient, etc.)
- **Patience** (now that we know it may be worth it), **faster computers**, and use of **GPUs**
- More efficient activation functions (e.g. ReLUs) in terms of both computation and avoiding $f(\text{net})$ saturation

Open problems

- A More Scientific Approach is Needed, not Just Building Better Systems...
 - Geoff Hinton, Yoshua Bengio & Yann LeCun, NIPS 2015