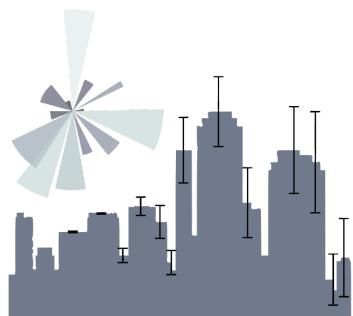


# principles of Urban Science 13



## Neural Networks

*dr.federica bianco*

| [fbb.space](http://fbb.space) | [fedhere](#) |  [fedhere](#)

*this slide deck:*

*[http://slides.com/federicabianco/pus2020\\_13](http://slides.com/federicabianco/pus2020_13)*



# MLTSA: *Recap*

# Machine Learning

**Data driven models for exploration of structure, prediction that learn parameters from data.**

# Machine Learning

**Data driven models for exploration of structure, prediction that learn parameters from data.**

**unupervised**

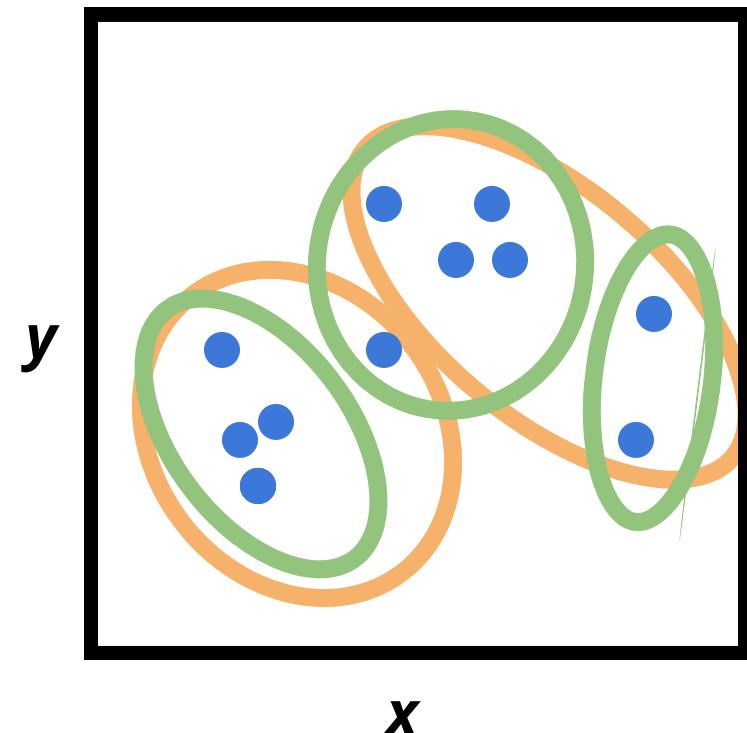
-----

set up: All features known for all observations

Goal: explore structure in the data

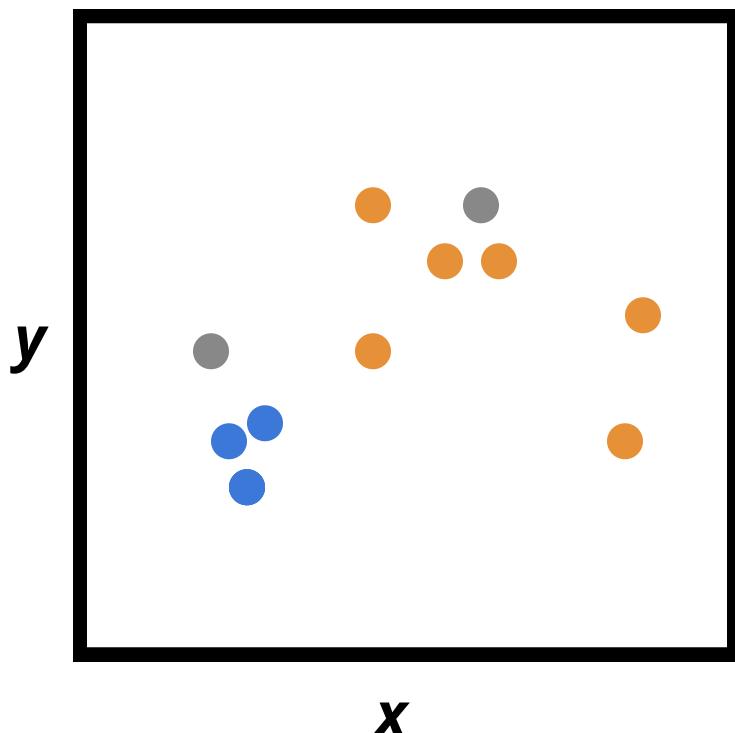
- data compression
- understanding structure

Algorithms: ***Clustering, (...)***



# Machine Learning

**Data driven models for exploration of structure, prediction that learn parameters from data.**



**supervised**

**set up:** All features known for a subset of the data; one feature cannot be observed for the rest of the data

**Goal:** predicting missing feature

- classification
- regression

**Algorithms:** *regression, SVM, tree methods, k-nearest neighbors, neural networks, (...)*

# Machine Learning

## unupervised

**set up:** All features known for all observations

**Goal:** explore structure in the data

- data compression
- understanding structure

**Algorithms:** *k-means clustering,*  
*agglomerative clustering,*  
*density based clustering, (...)*

## supervised

**set up:** All features known for a sunbset  
of the data; one feature cannot be  
observed for the rest of the data

**Goal:** predicting missing feature

- classification
- regression

**Algorithms:** *regression, SVM, tree*  
*methods, k-nearest neighbors,*  
*neural networks, (...)*

# Machine Learning

Learning relies on the definition of a ***loss function***

model parameters are learned by calculating a loss function for different parameter sets and trying to minimize loss  
(or a target function and trying to maximize)

e.g.

$$L1 = |target - prediction|$$

# Machine Learning

Learning relies on the definition of a ***loss function***

learning type	loss / target
unsupervised	intra-cluster variance / inter cluster distance
supervised	distance between prediction and truth

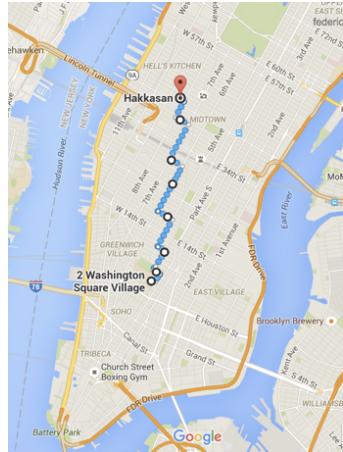
# Machine Learning

The definition of a loss function requires the definition of *distance* or *similarity*

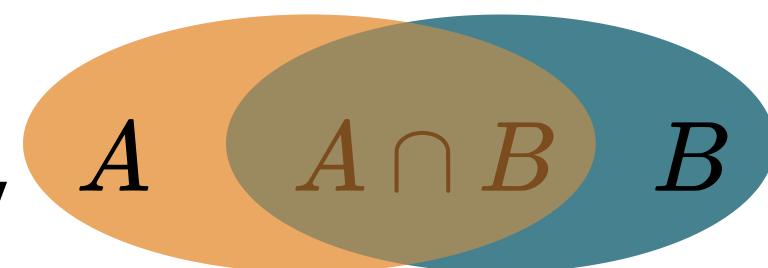
# Machine Learning

The definition of a loss function requires the definition of *distance* or *similarity*

**Minkowski distance**



**Jaccard similarity**



**Great circle distance**



NN:  
1  
*Neural Networks*



A LOGICAL CALCULUS OF THE  
IDEAS IMMANENT IN NERVOUS ACTIVITY

WARREN S. McCULLOCH and WALTER H. PITTS

Because of the “all-or-none” character of nervous activity, neural events and the relations among them can be treated by means of propositional logic. It is found that the behavior of every net can be described in these terms, with the addition of more complicated logical means for nets containing circles; and that for any logical expression satisfying certain conditions, one can find a net behaving in the fashion it describes. It is shown that many particular choices among possible neurophysiological assumptions are equivalent, in the sense that for every net behaving under one assumption, there exists another net which behaves under the other and gives the same results, although perhaps not in the same time. Various applications of the calculus are discussed.

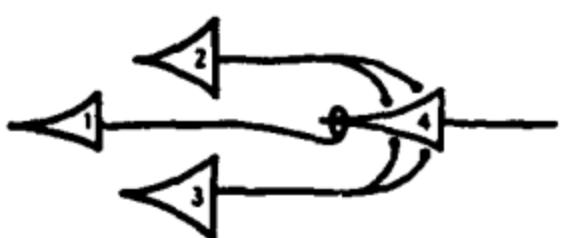
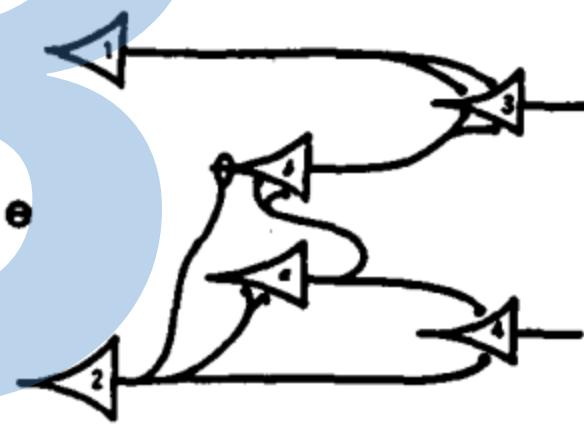
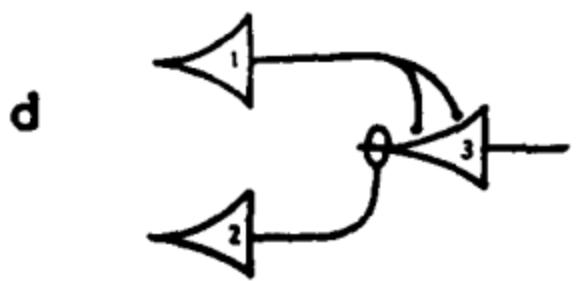
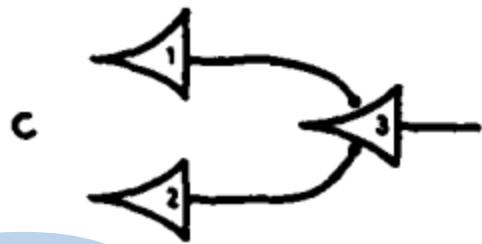
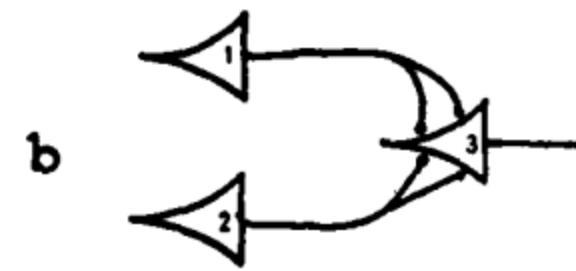
1943

THE THEORY: NETS WITHOUT CIRCLES

We shall make the following physical assumptions for our calculus.

1. The activity of the neuron is an “all-or-none” process.
5. The structure of the net does not change with time.

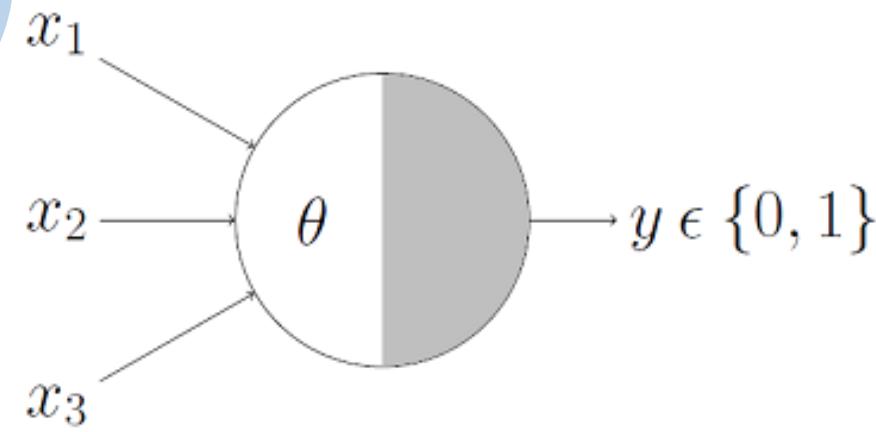
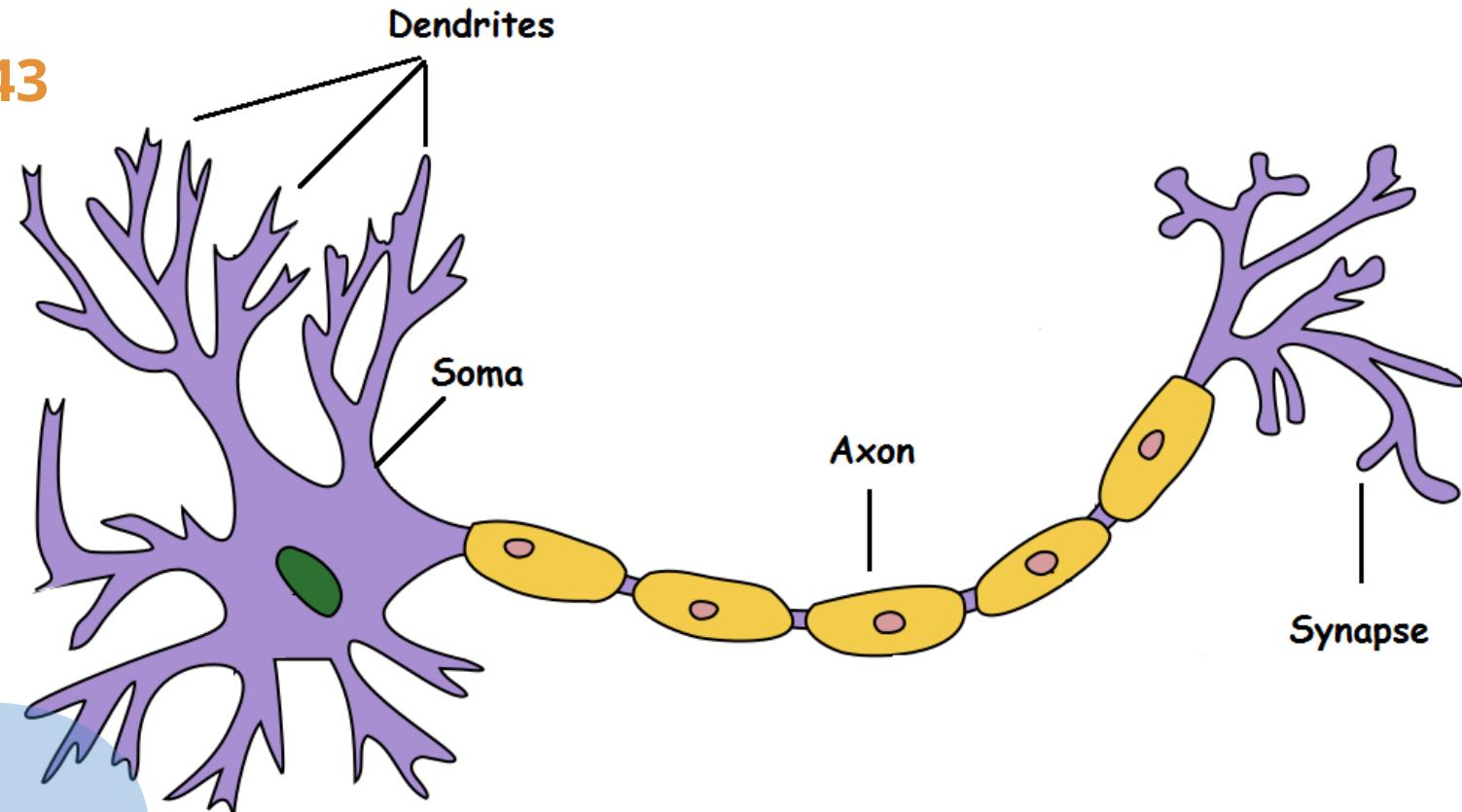
1943



M-P Neuron McCulloch & Pitts 1943

# M-P Neuron

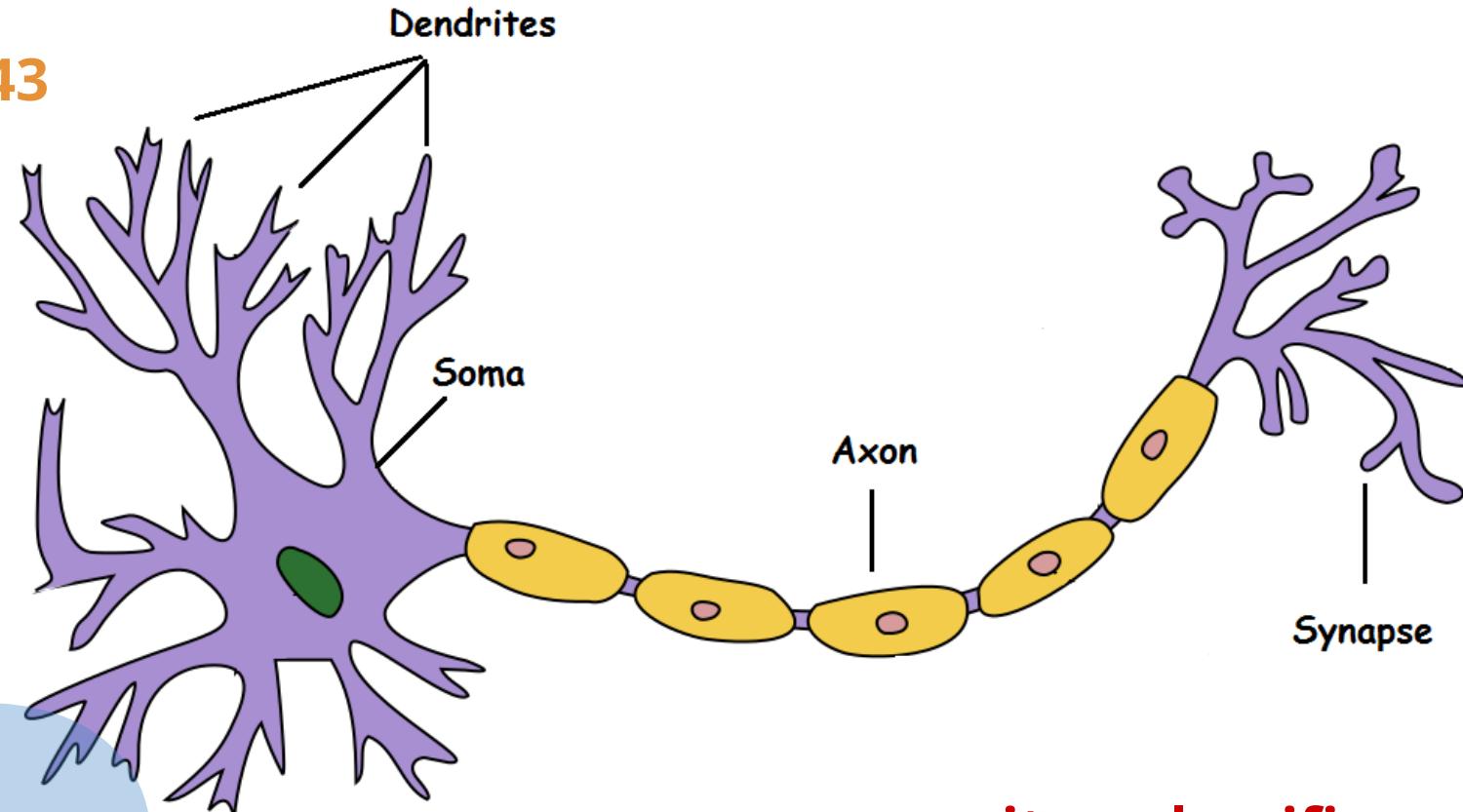
1943



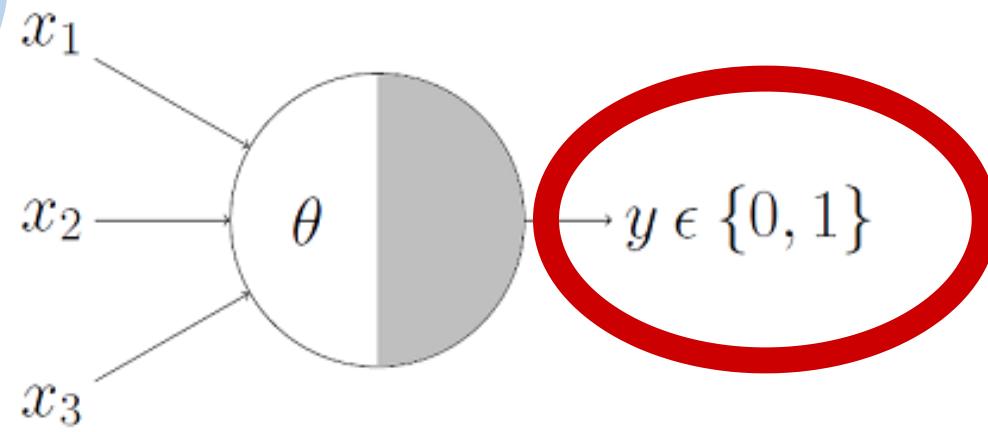
M-P Neuron McCulloch & Pitts 1943

# M-P Neuron

1943



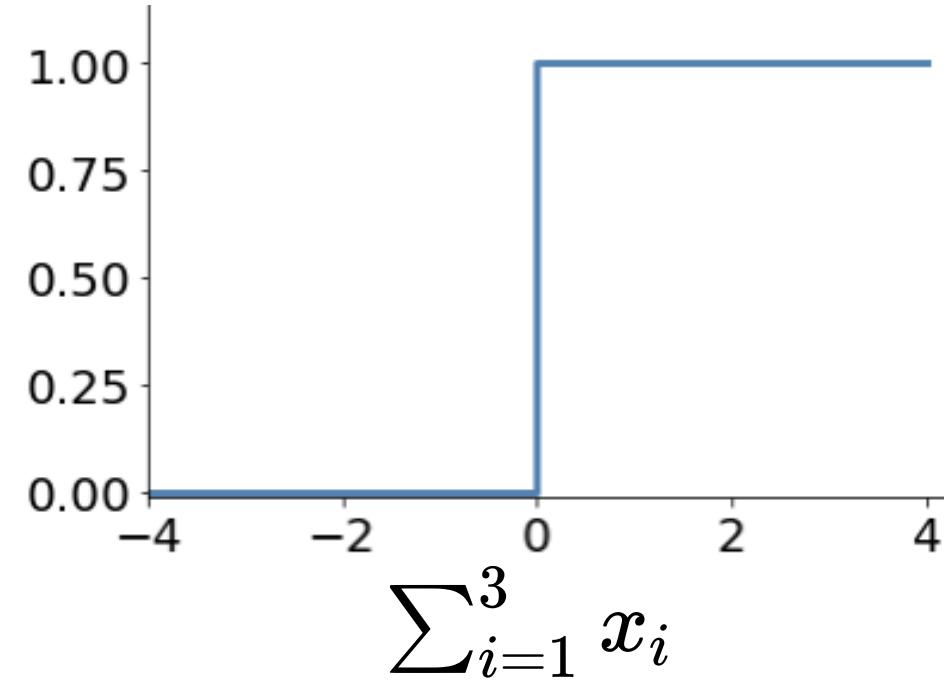
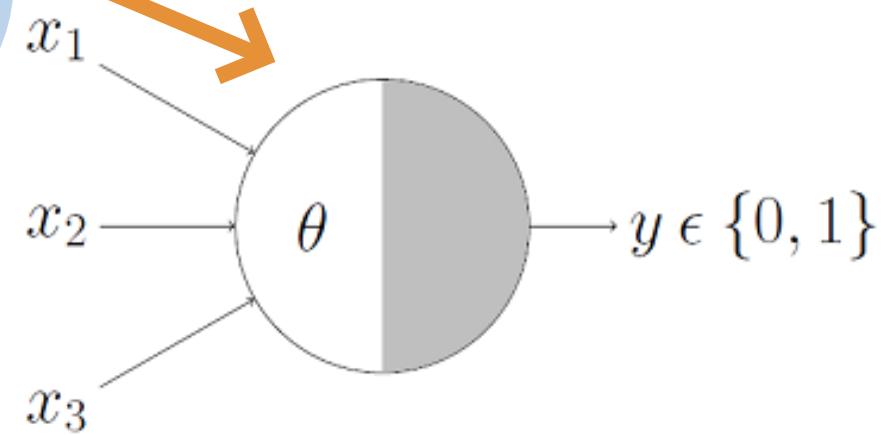
its a classifier



## M-P Neuron

1 if  $\sum_{i=1}^3 x_i \geq \theta$  else 0

1943



$$1 \text{ if } \sum_{i=1}^3 x_i \geq \theta \text{ else } 0$$

1943

if  $x_i$  is Bool (True/False)  
what value of  $\theta$   
corresponds to logical  
AND?

# The perceptron algorithm : 1958, Frank Rosenblatt

# Perceptron

*Psychological Review*  
Vol. 65, No. 6, 1958

## THE PERCEPTRON: A PROBABILISTIC MODEL FOR INFORMATION STORAGE AND ORGANIZATION IN THE BRAIN<sup>1</sup>

F. ROSENBLATT

*Cornell Aeronautical Laboratory*

If we are eventually to understand the capability of higher organisms for perceptual recognition, generalization, recall, and thinking, we must first have answers to three fundamental questions:

1. How is information about the physical world sensed, or detected, by the biological system?
2. In what form is information stored, or remembered?
3. How does information contained in storage, or in memory, influence recognition and behavior?

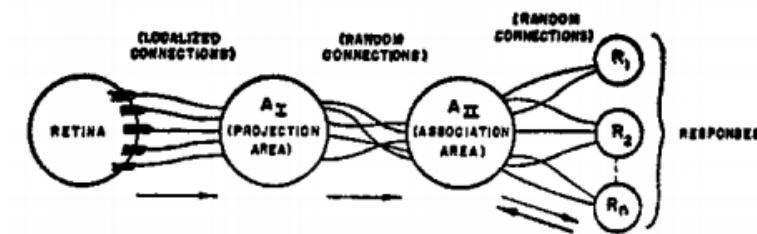


FIG. 1. Organization of a perceptron.

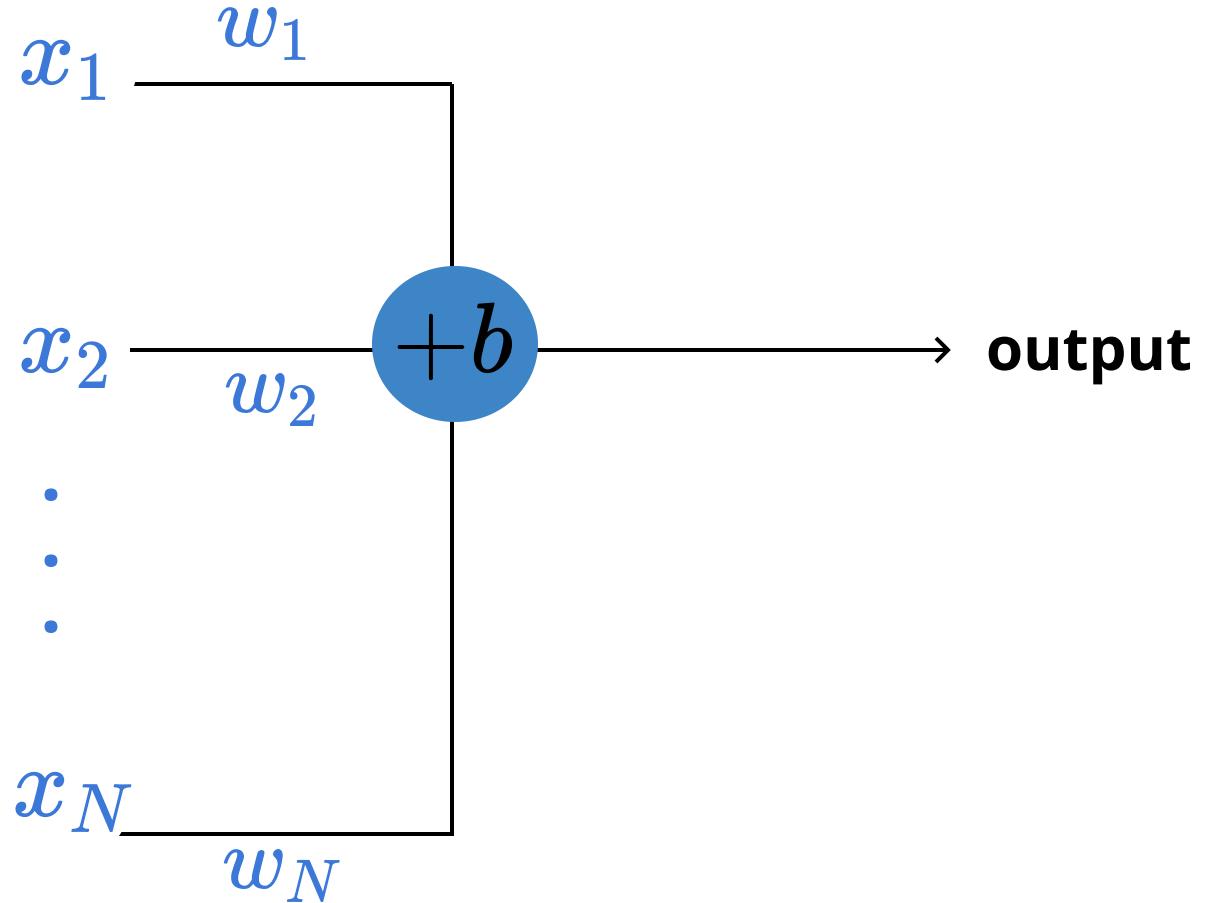
The perceptron algorithm : 1958, Frank Rosenblatt

# Perceptron

1 if  $\sum_{i=1}^N w_i x_i \geq \theta$  else 0

1958

linear regression:  
 $w_i$  weights  
 $b$  bias



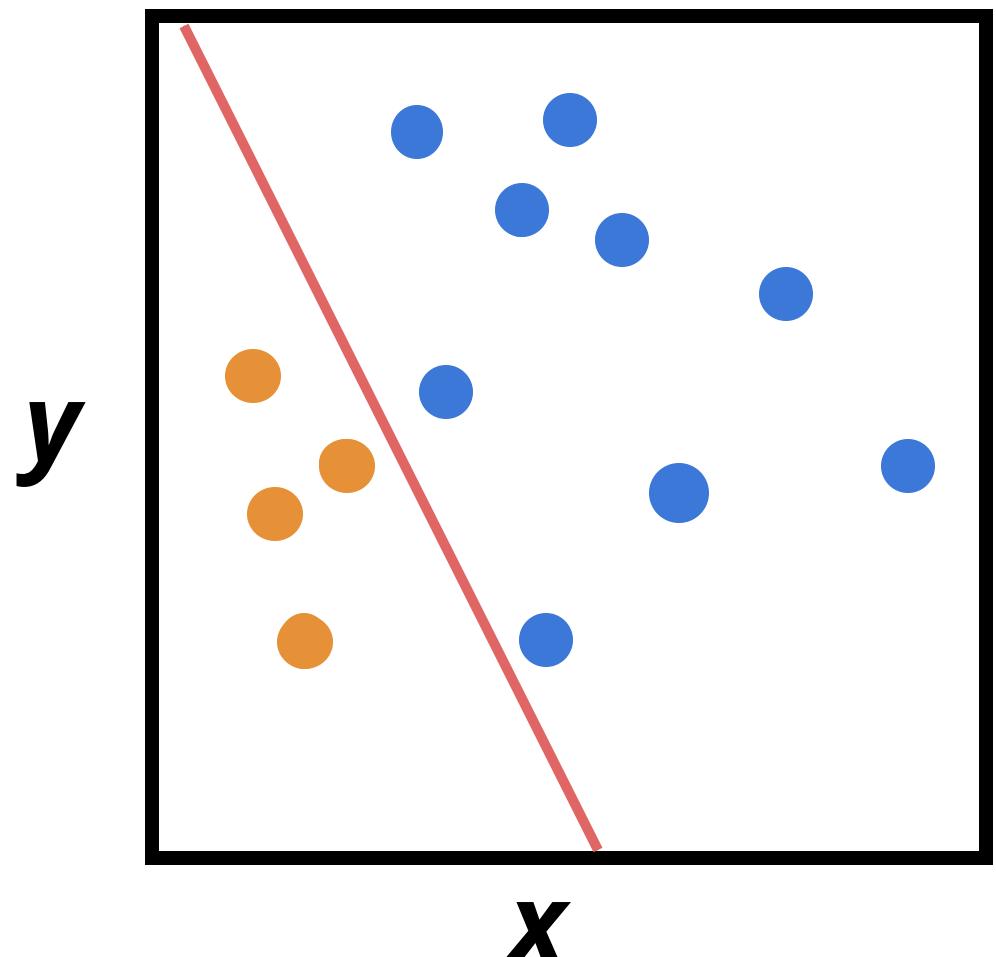
## The perceptron algorithm : 1958, Frank Rosenblatt

Perceptrons are **linear classifiers**: makes its predictions based on a linear predictor function

*combining a set of weights (=parameters) with the feature vector.*

$$y = \sum_i w_i x_i + b$$

1958

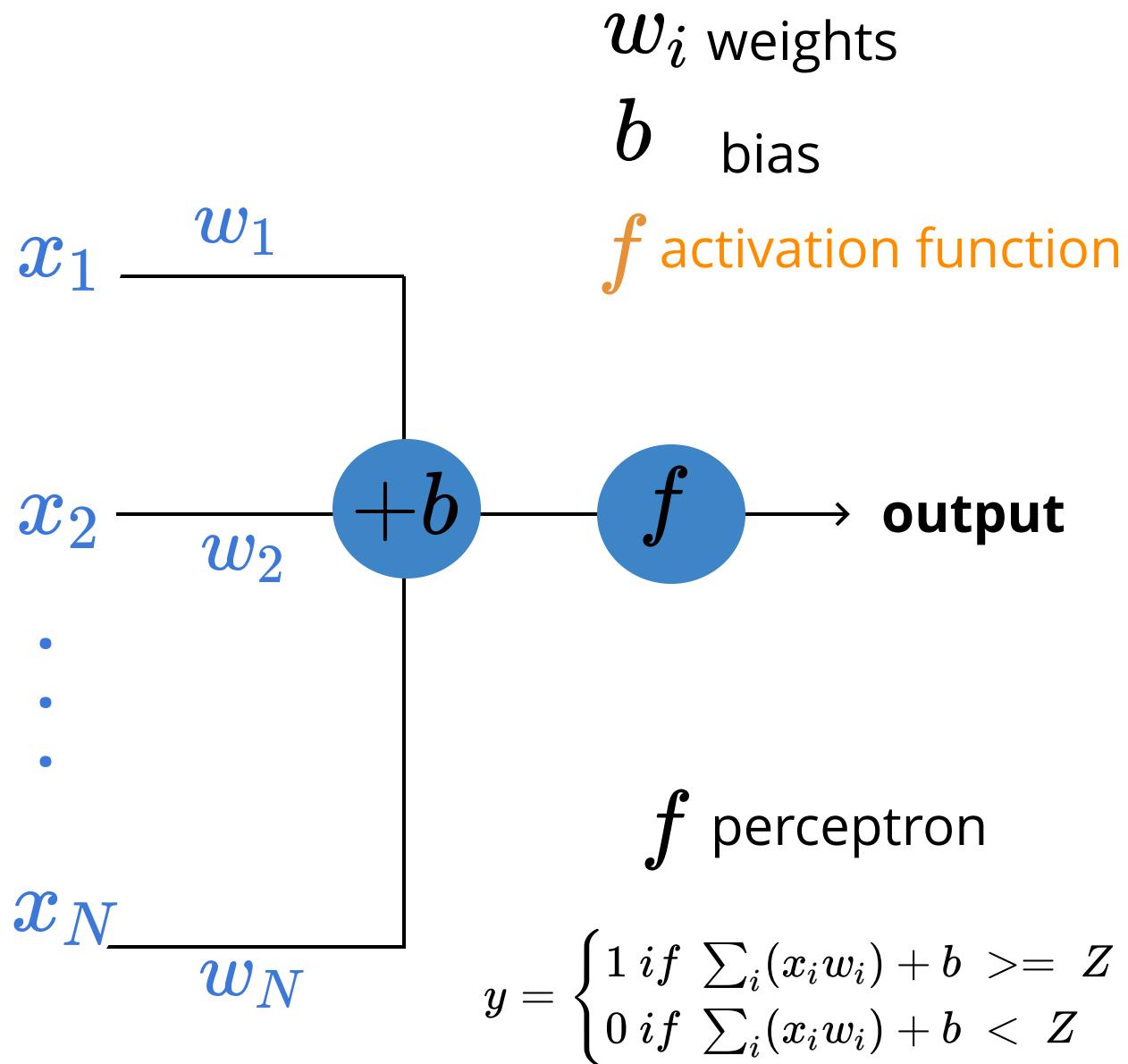
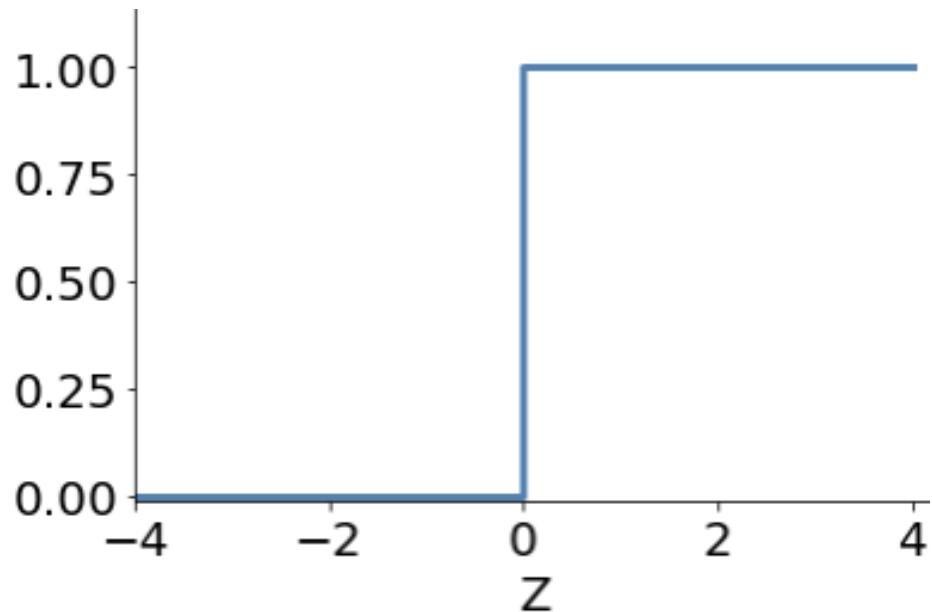


## The perceptron algorithm : 1958, Frank Rosenblatt

Perceptrons are **linear classifiers**: makes its predictions based on a linear predictor function

combining a set of weights (=parameters) with the feature vector.

$$y = f(\sum_i w_i x_i + b)$$

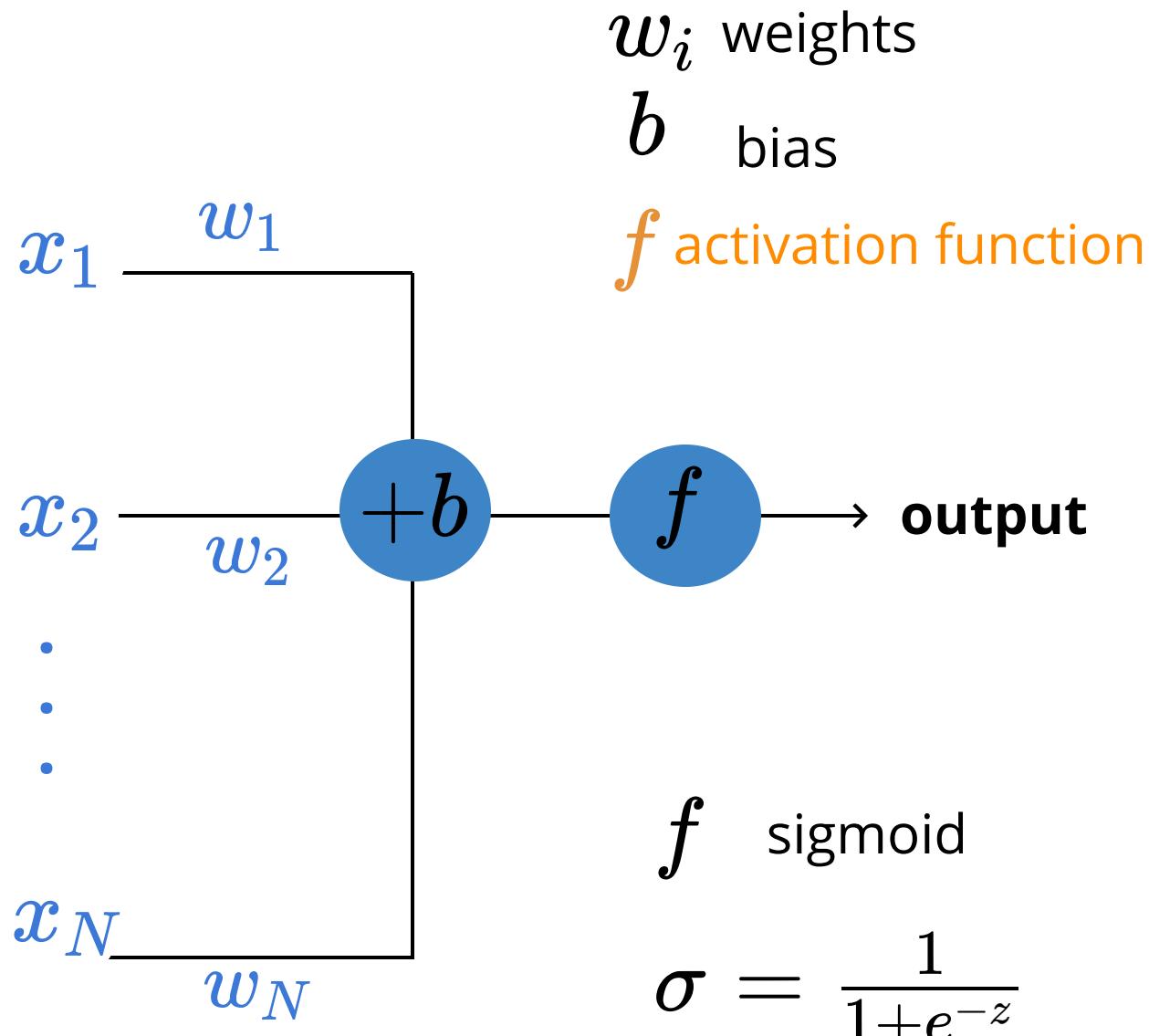
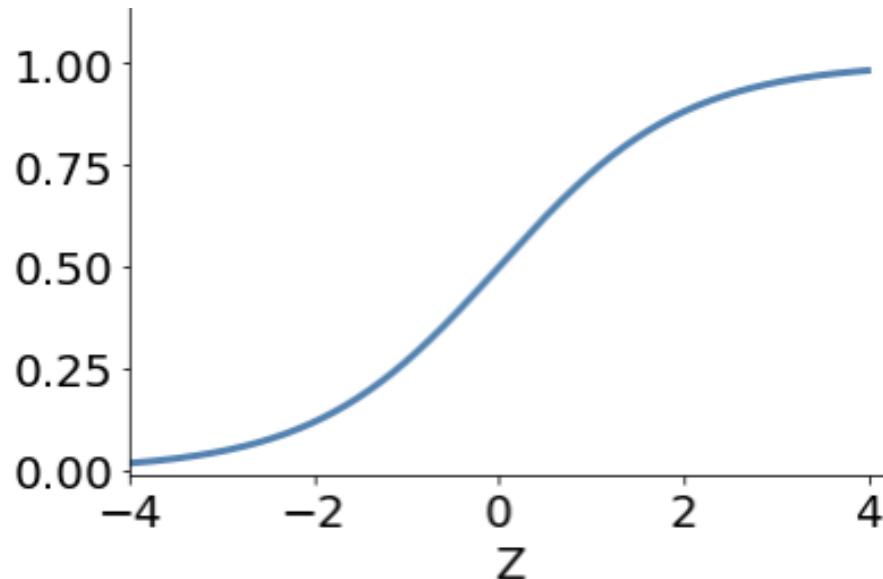


## The perceptron algorithm : 1958, Frank Rosenblatt

Perceptrons are **linear classifiers**: makes its predictions based on a linear predictor function

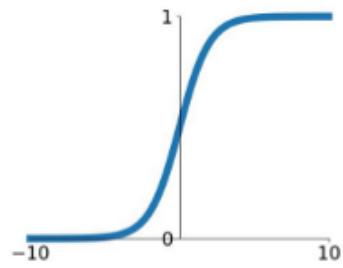
combining a set of weights (=parameters) with the feature vector.

$$y = f(\sum_i w_i x_i + b)$$



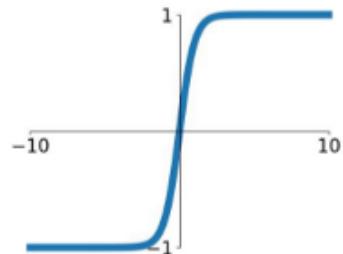
## Sigmoid

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



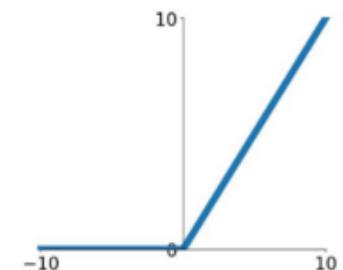
## tanh

$$\tanh(x)$$



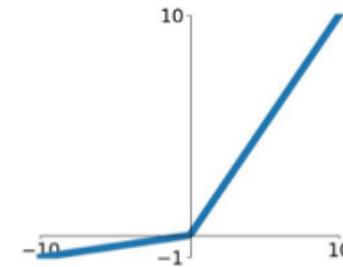
## ReLU

$$\max(0, x)$$



## Leaky ReLU

$$\max(0.1x, x)$$

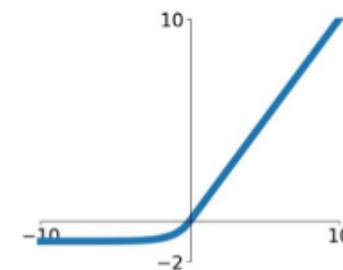


## Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

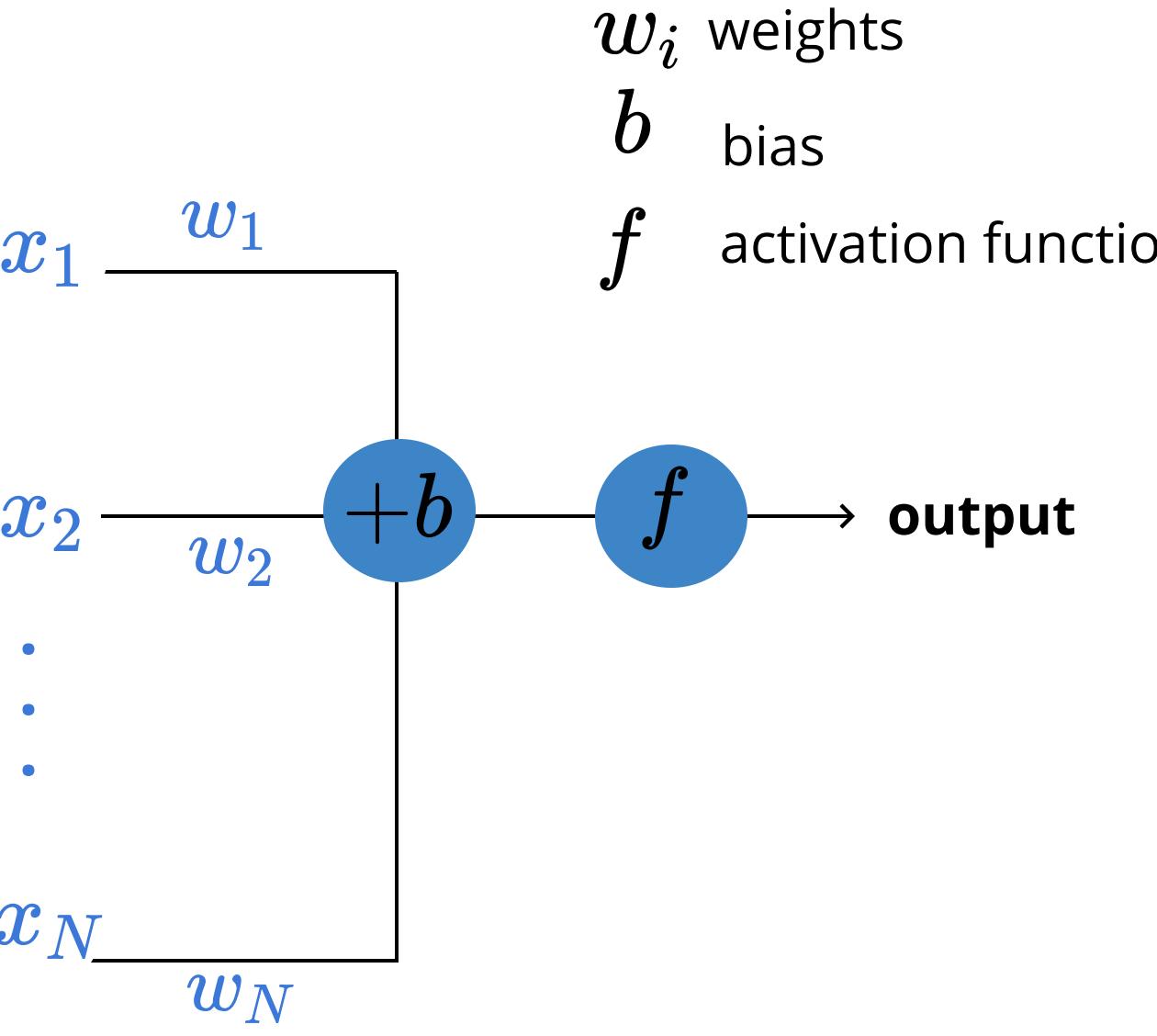
## ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



# The perceptron algorithm : 1958, Frank Rosenblatt

# Perceptron



# The perceptron algorithm : 1958, Frank Rosenblatt

## Perceptron



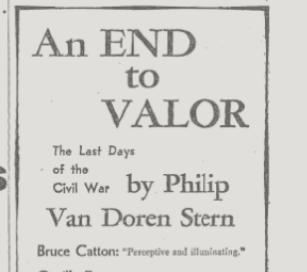
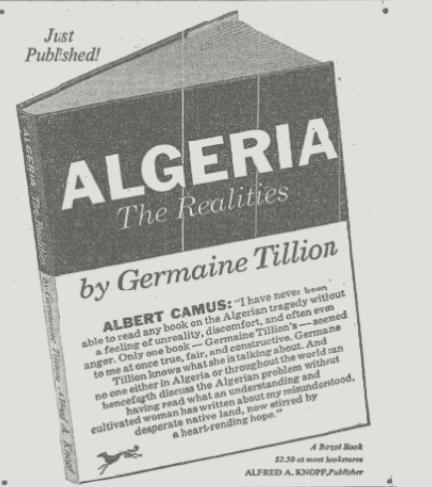
# The New York Times

July 8, 1958

## NEW NAVY DEVICE LEARNS BY DOING; Psychologist Shows Embryo of Computer Designed to Read and Grow Wiser

*The Navy revealed the embryo of an electronic computer today that it expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence.*

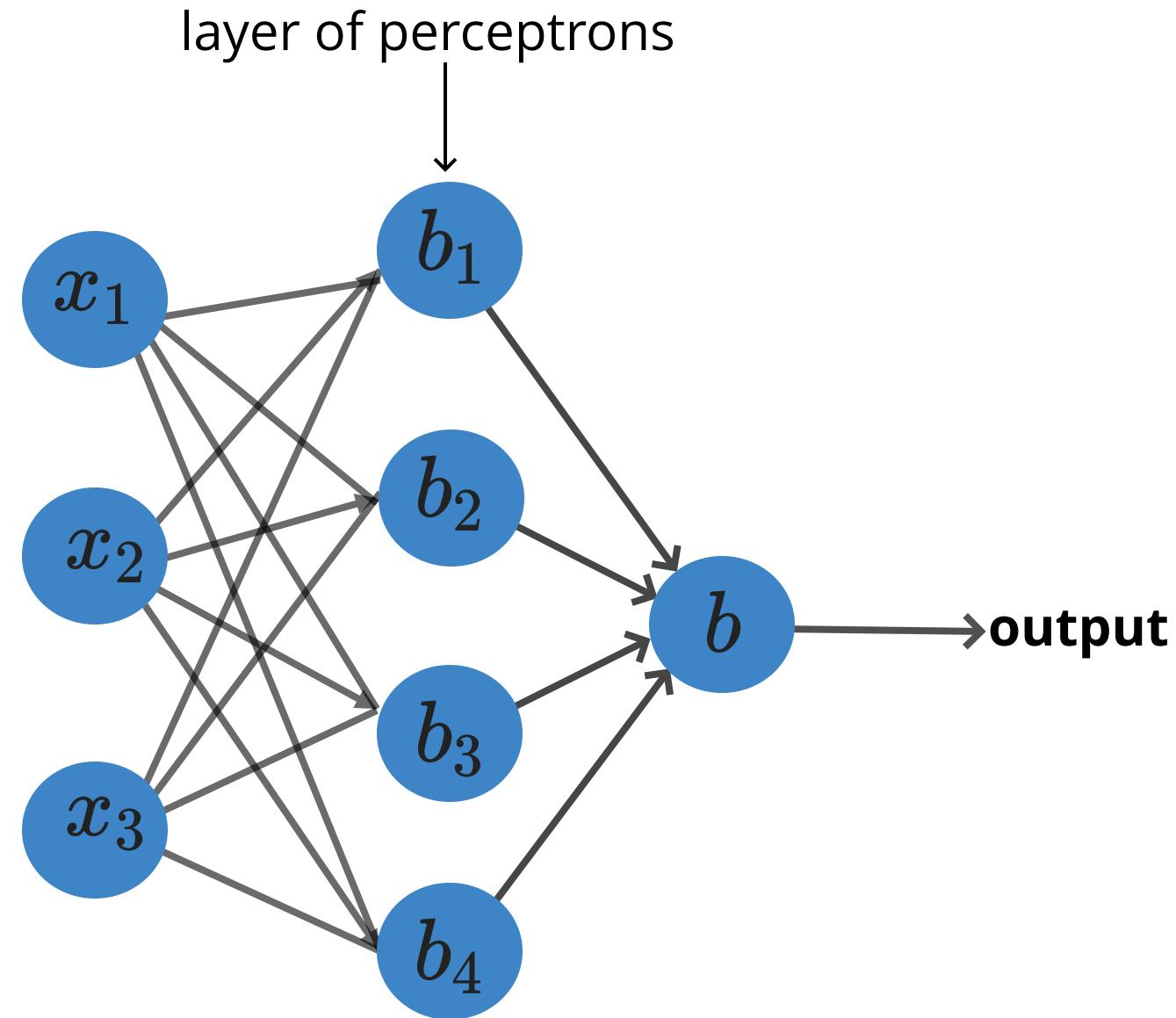
*The embryo - the Weather Bureau's \$2,000,000 "704" computer - learned to differentiate between left and right after 50 attempts in the Navy demonstration*



A large, light blue graphic of a question mark is positioned above the text. It has a thick, rounded top curve and a sharp, triangular point at the bottom. The interior of the question mark is white, and it sits on a horizontal base that tapers to the right.

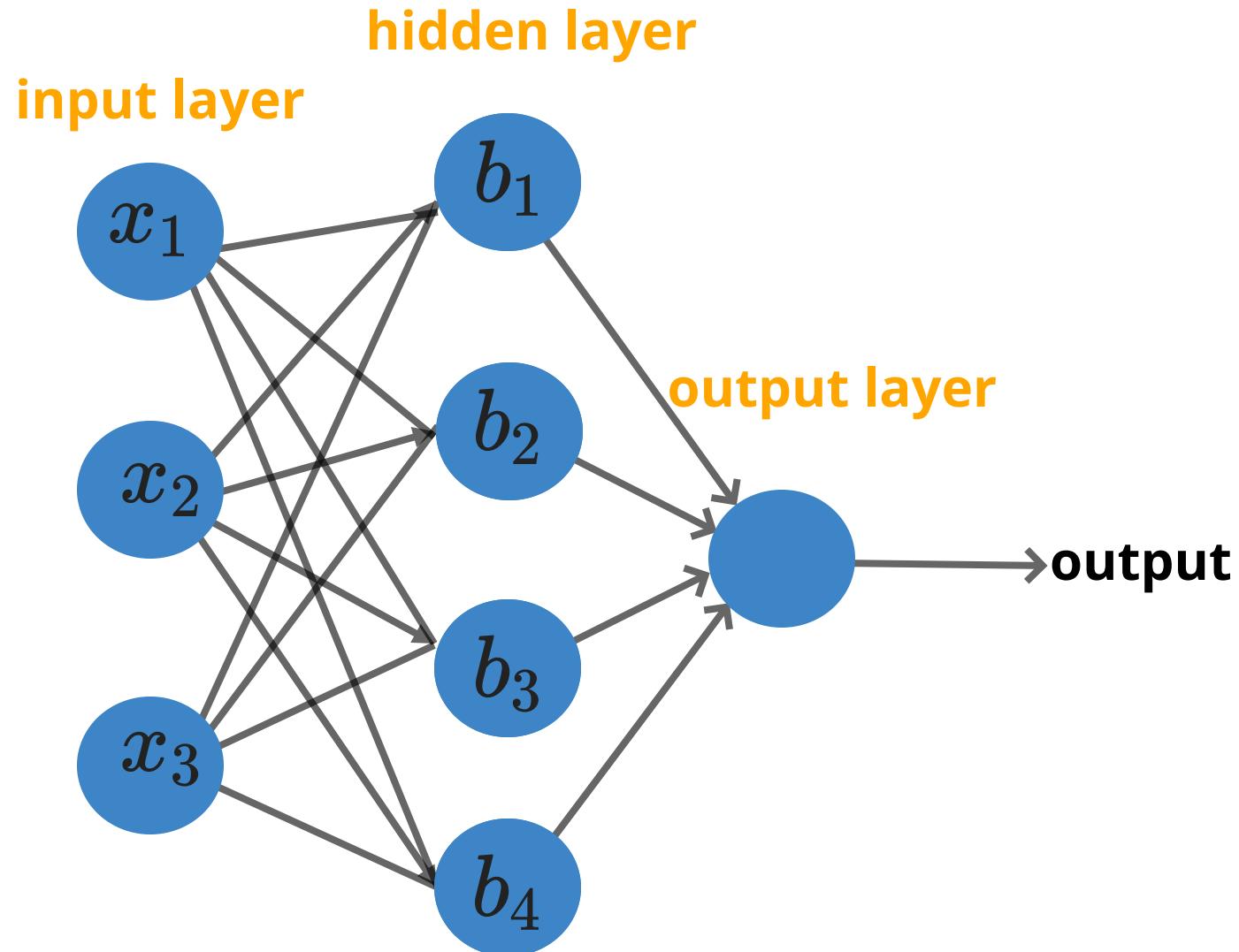
DNN:  
*Deep Learning*

# multilayer perceptron



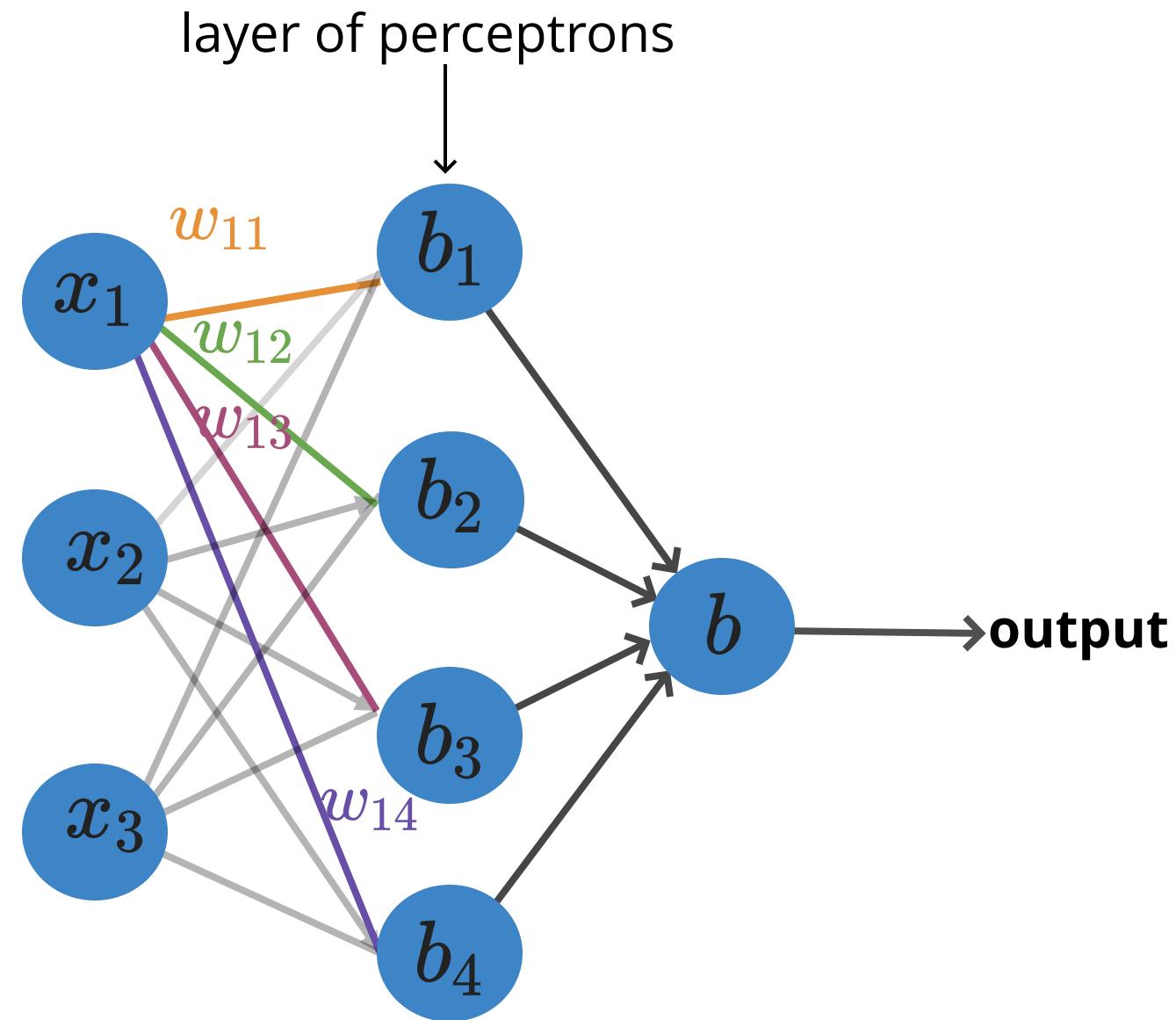
# multilayer perceptron

1970: multilayer  
perceptron architecture

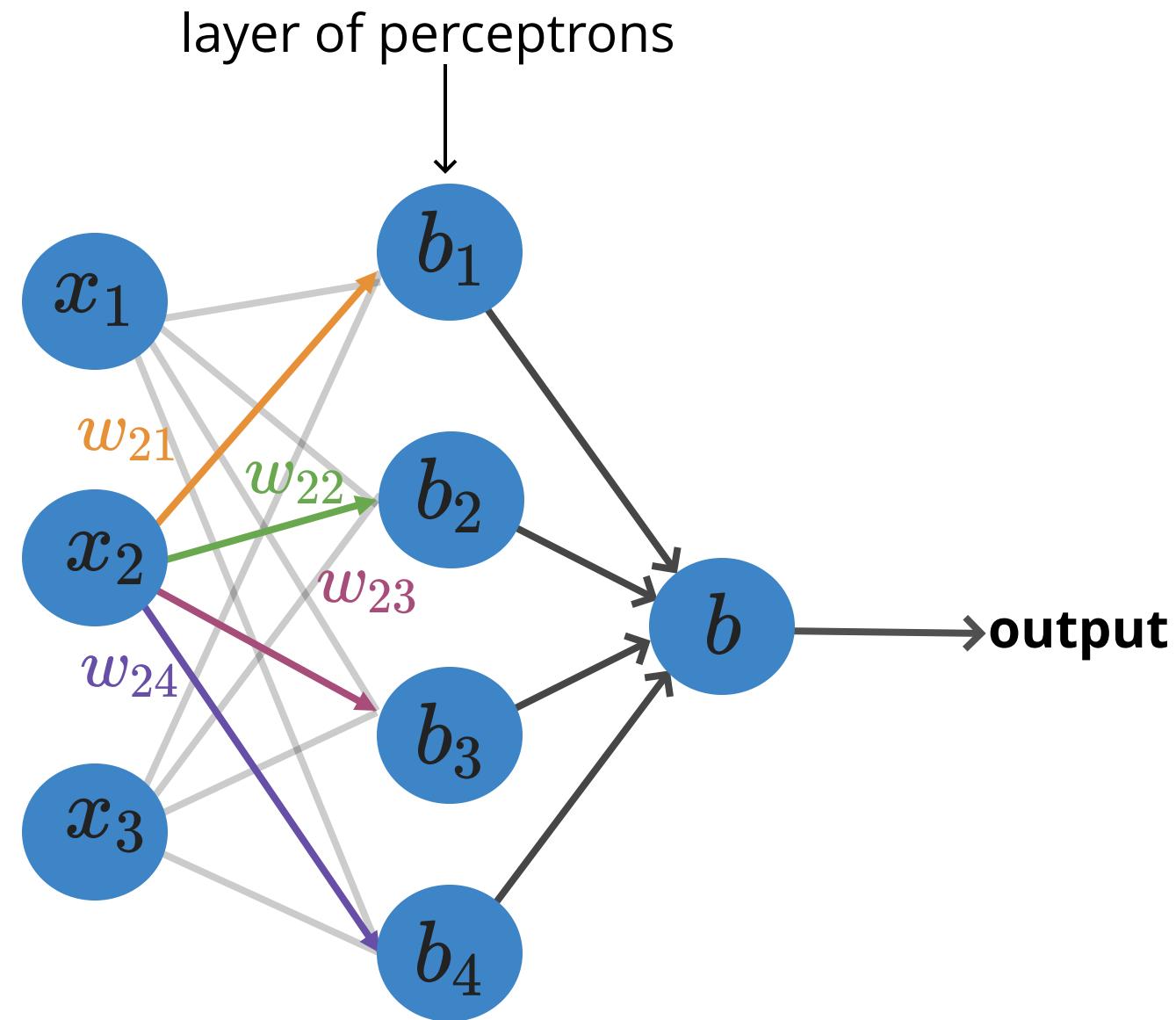


**Fully connected:** all nodes go to  
all nodes of the next layer.

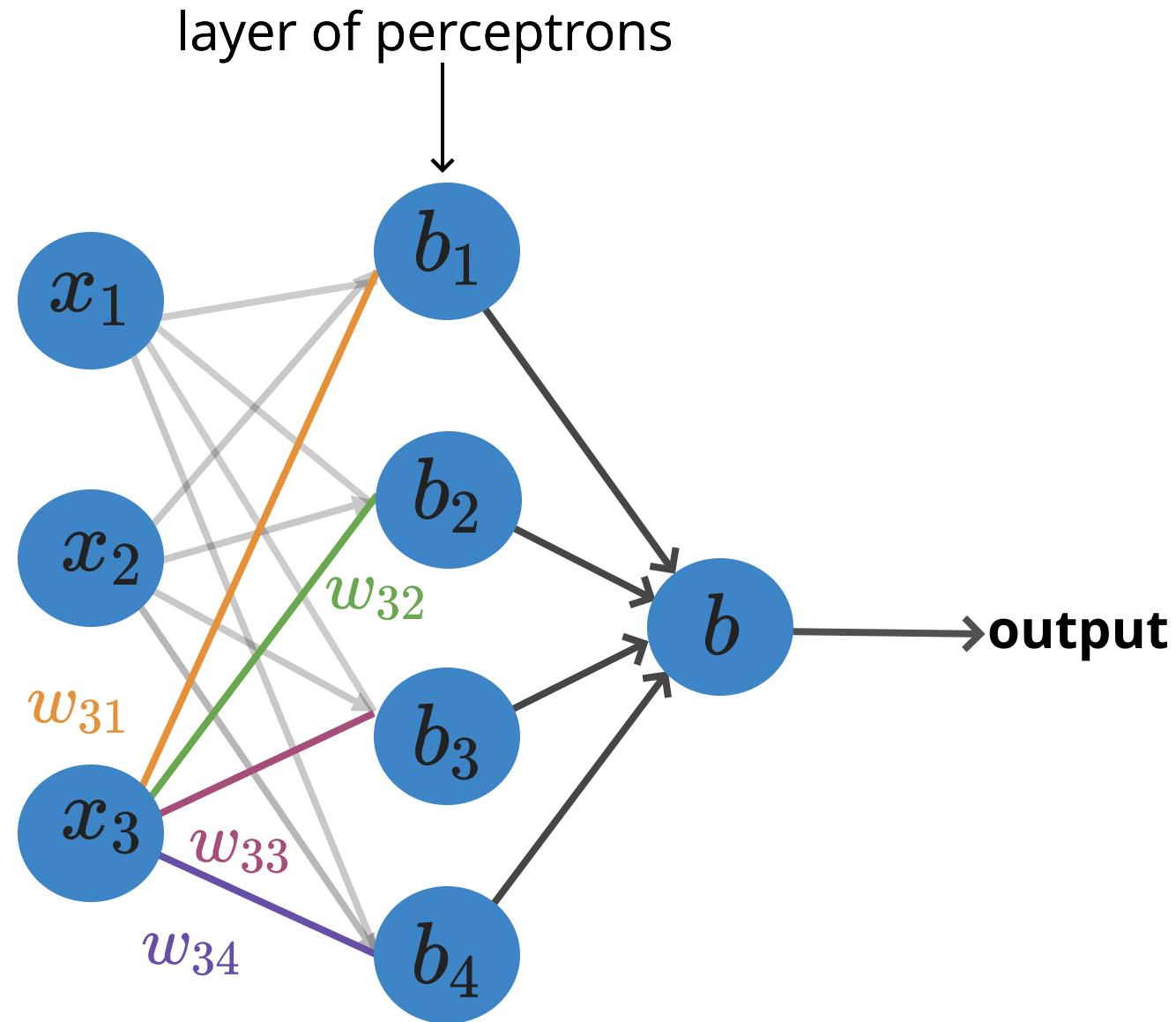
# multilayer perceptron



# multilayer perceptron

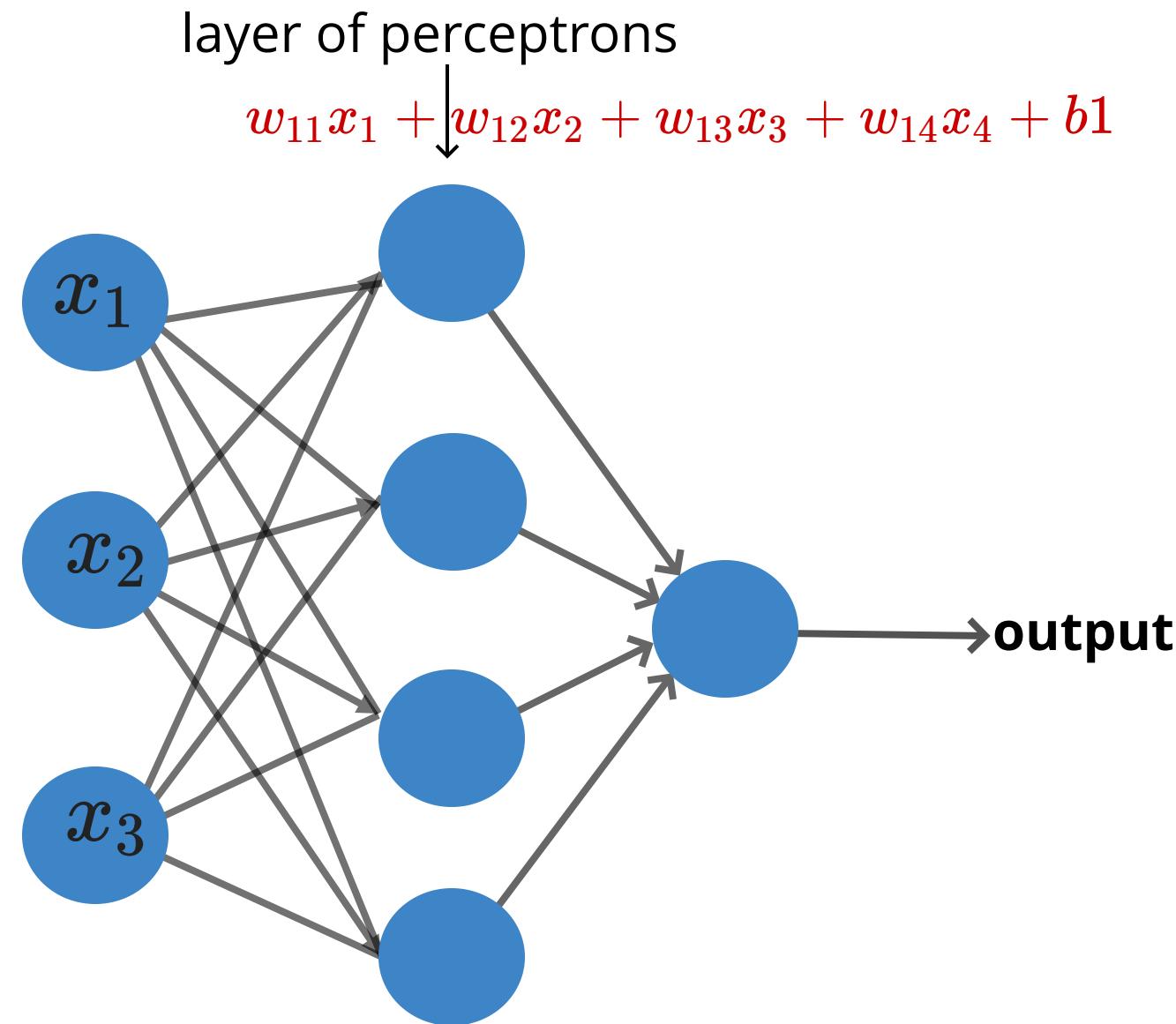


# multilayer perceptron



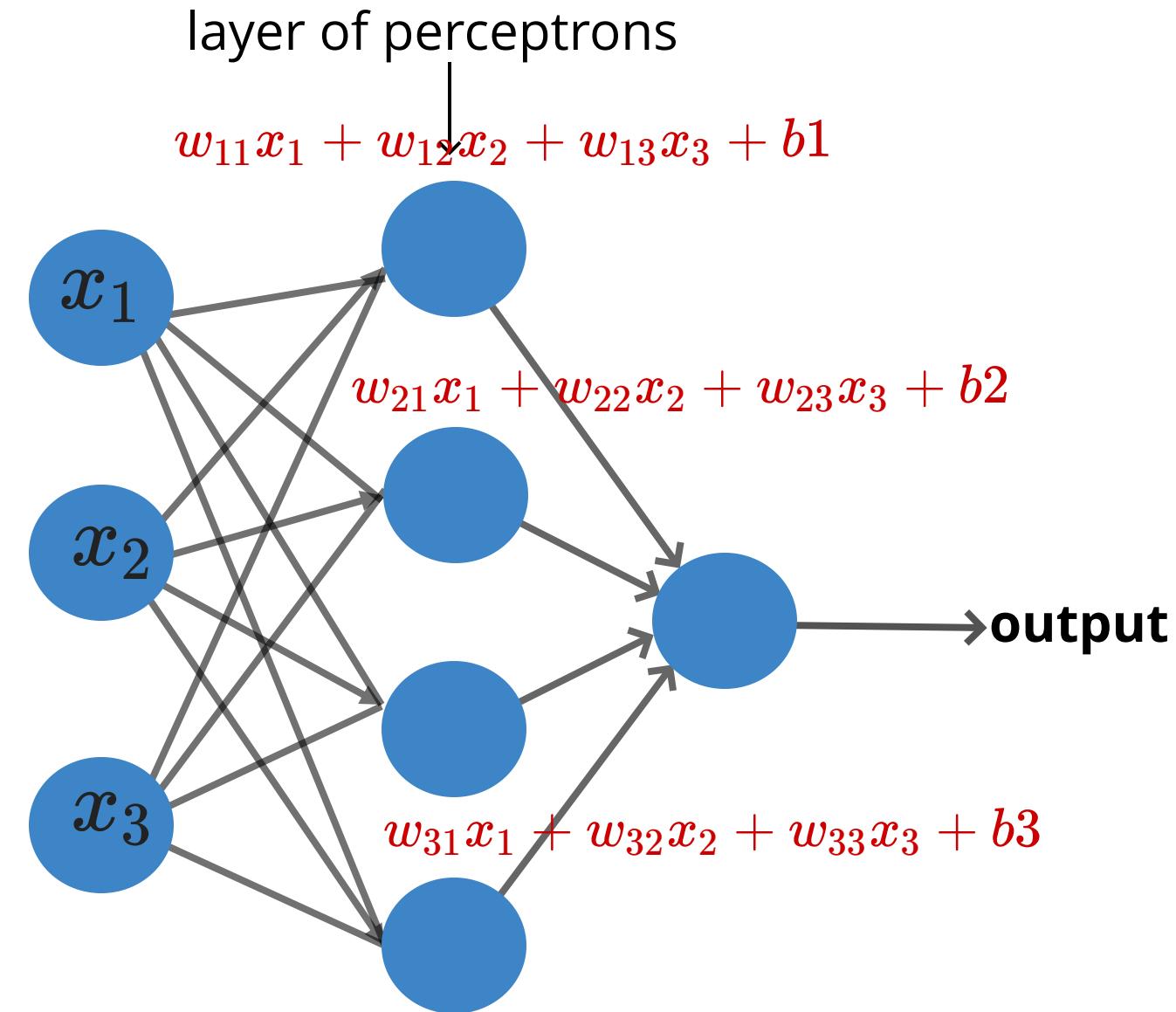
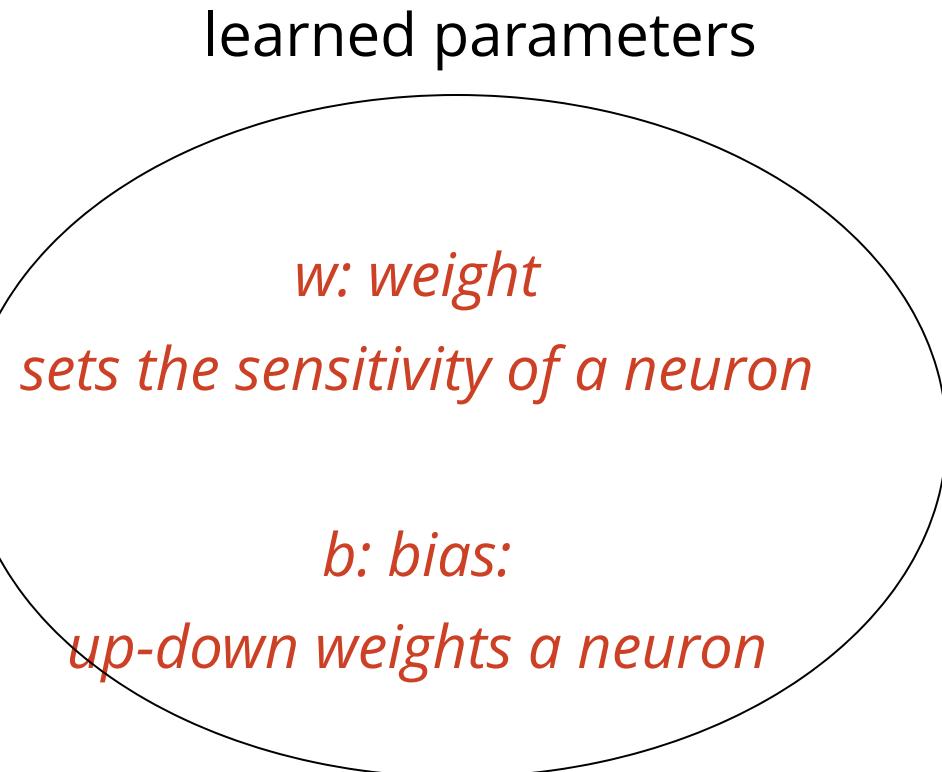
# multilayer perceptron

Fully connected: all nodes go to all nodes of the next layer.



# multilayer perceptron

Fully connected: all nodes go to all nodes of the next layer.



# multilayer perceptron

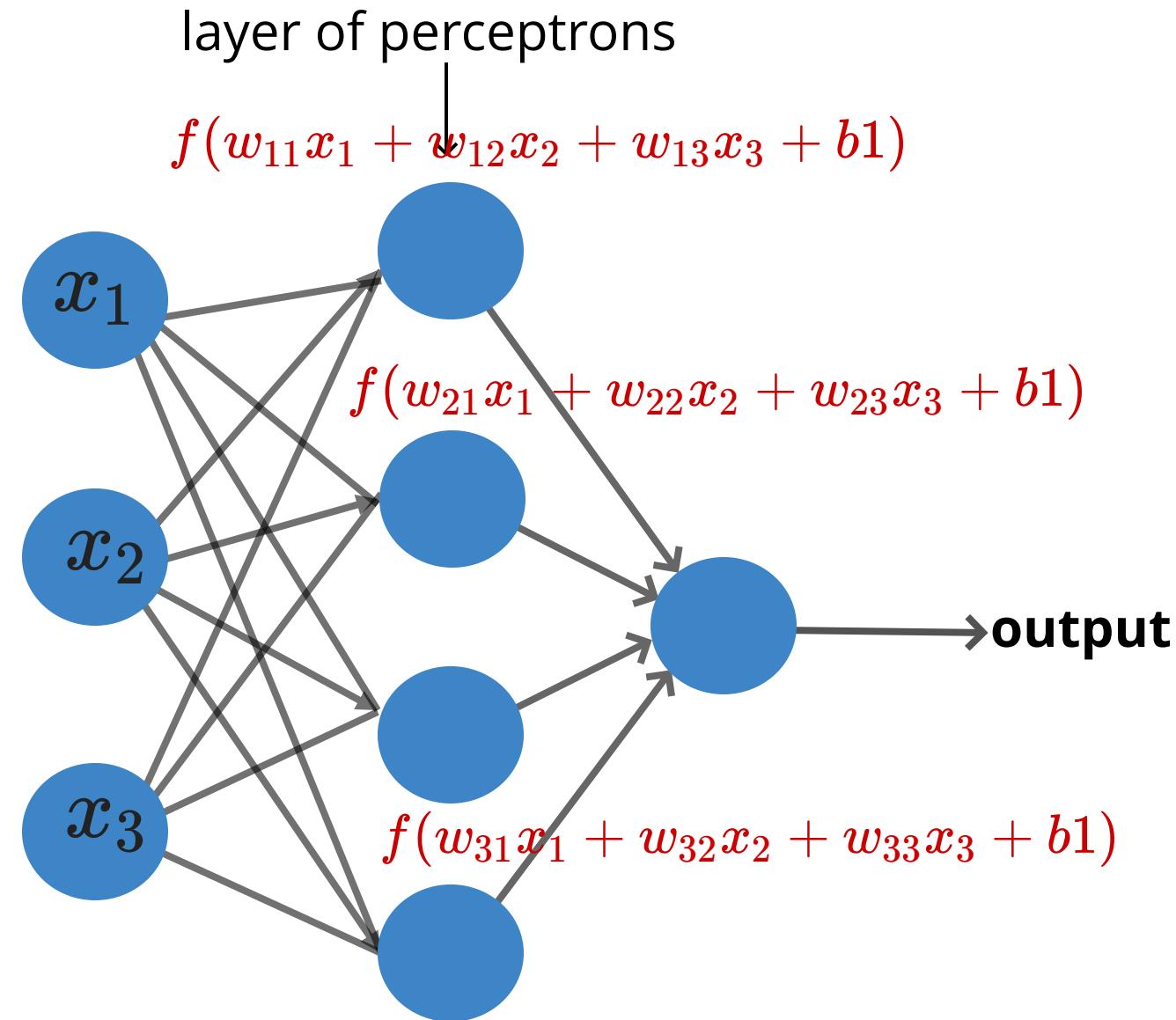
Fully connected: all nodes go to all nodes of the next layer.

*f: activation function:*  
turns neurons on-off

*w: weight*

*sets the sensitivity of a neuron*

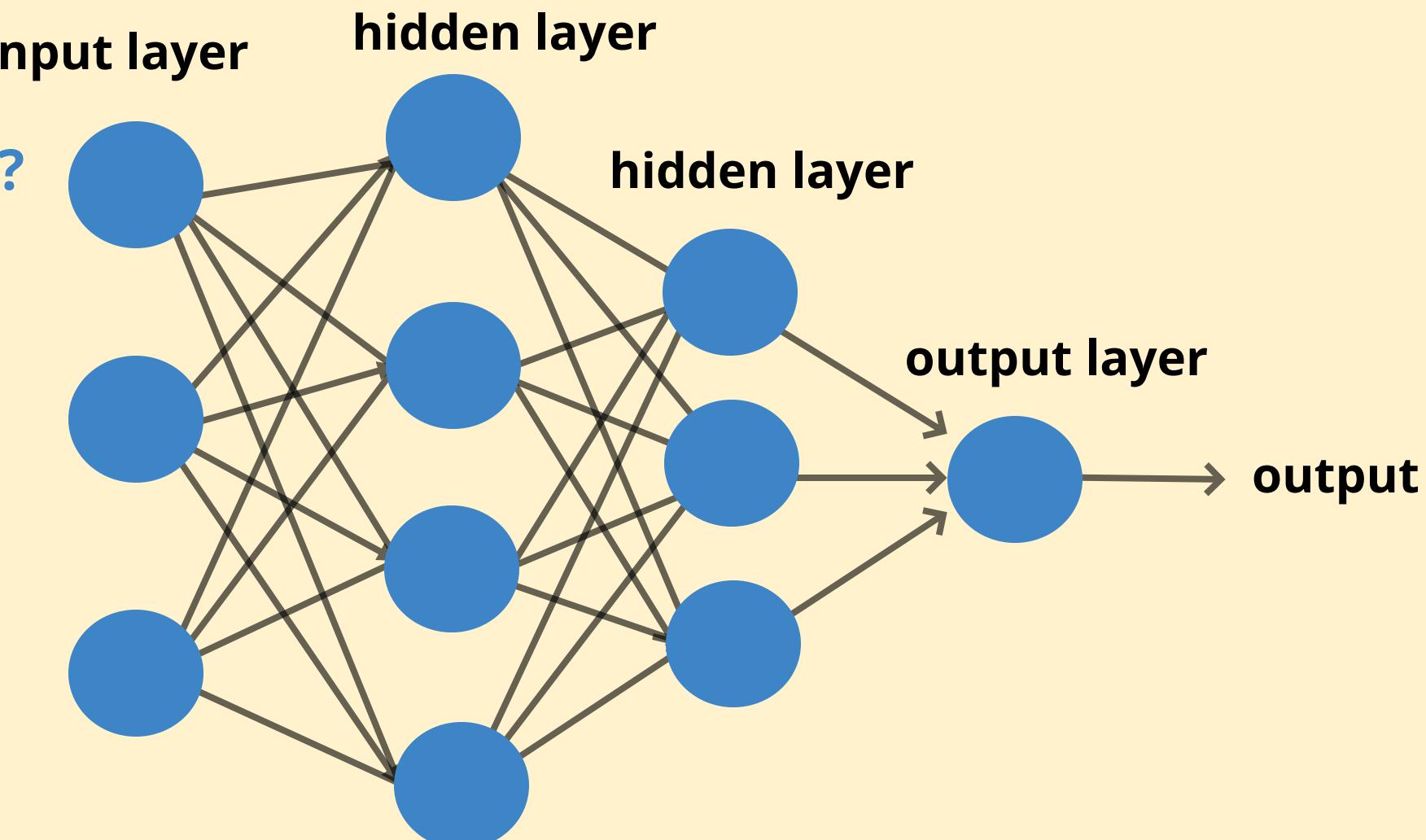
*b: bias:*  
*up-down weights a neuron*



# DNN: *hyperparameters* of DNN

# EXERCISE

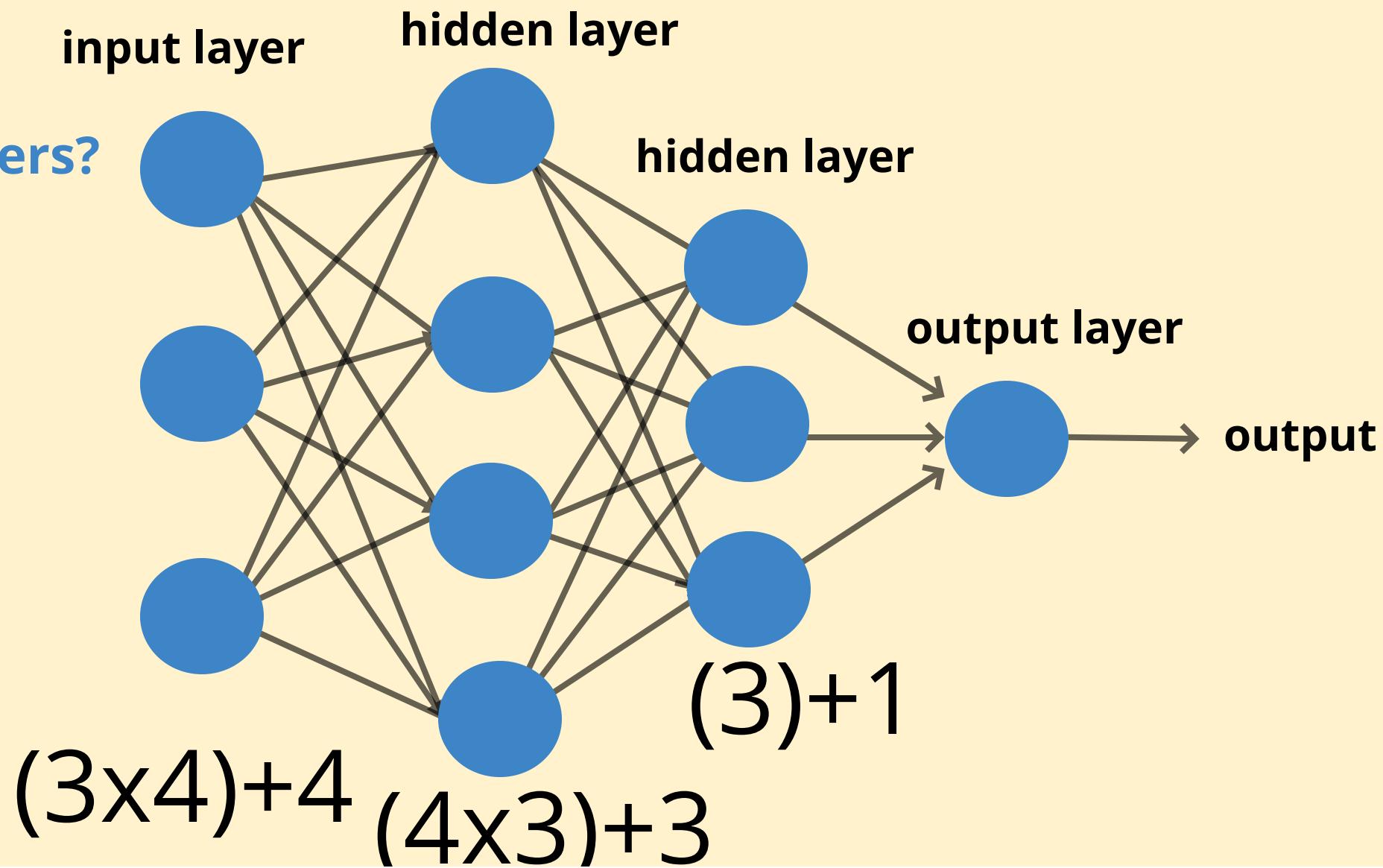
how many parameters?



# EXERCISE

how many parameters?

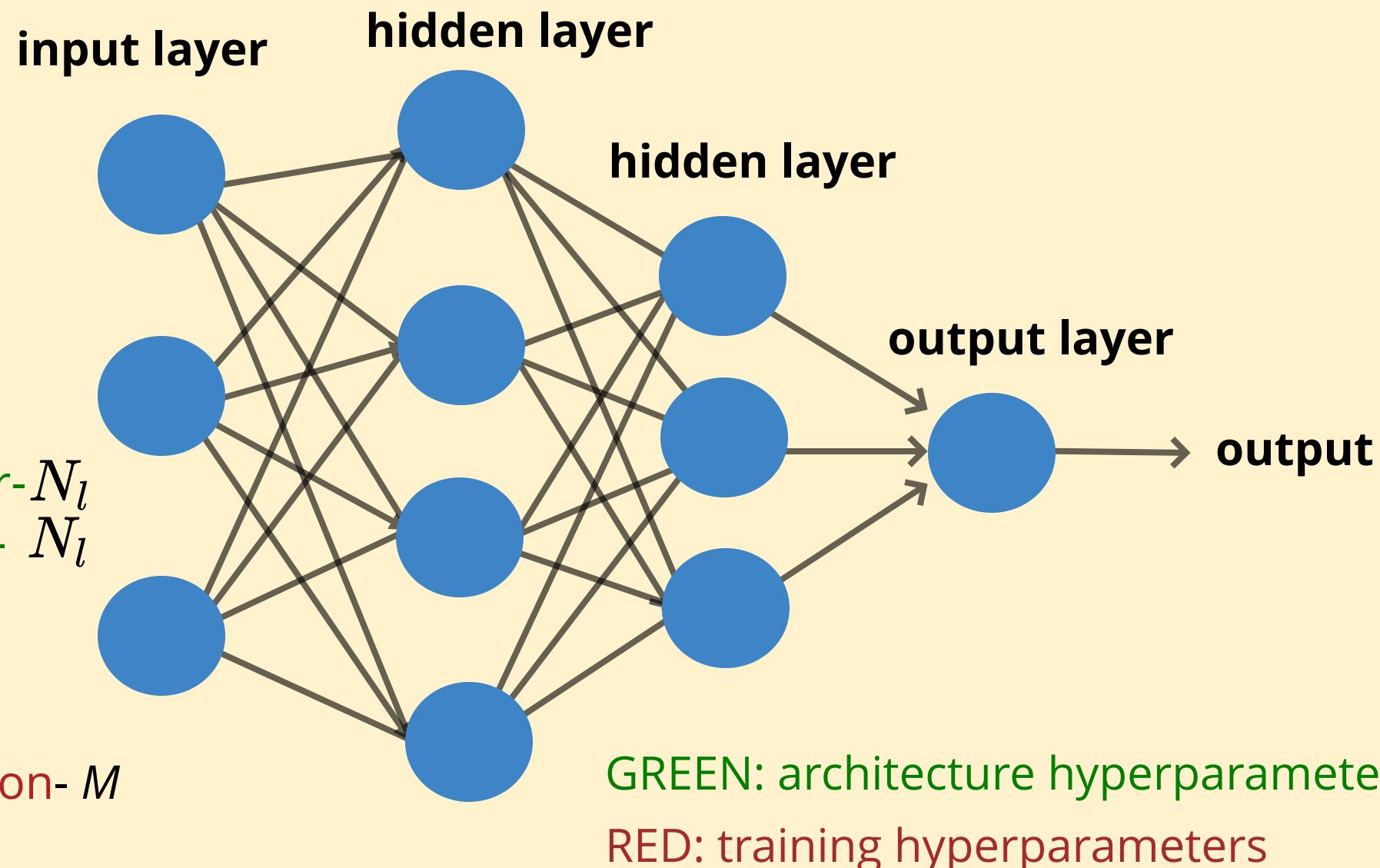
35



# EXERCISE

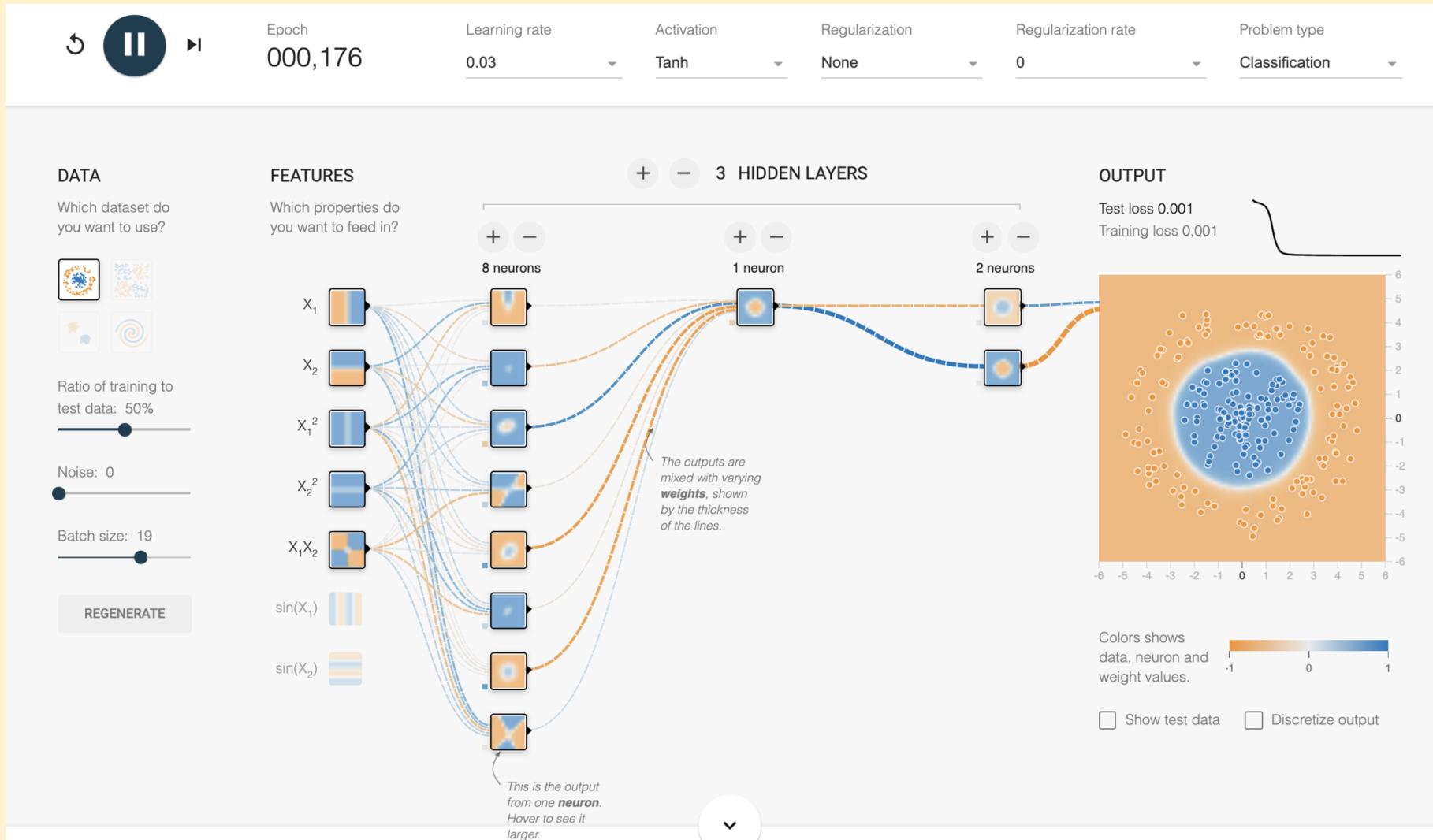
how many  
hyperparameters?

1. number of layers- 1
2. number of neurons/layer-  $N_l$
3. activation function/layer-  $N_l$
4. layer connectivity-  $N_l$ ??
5. optimization metric - 1
6. optimization method - 1
7. parameters in optimization-  $M$



# EXERCISE

<http://playground.tensorflow.org/>





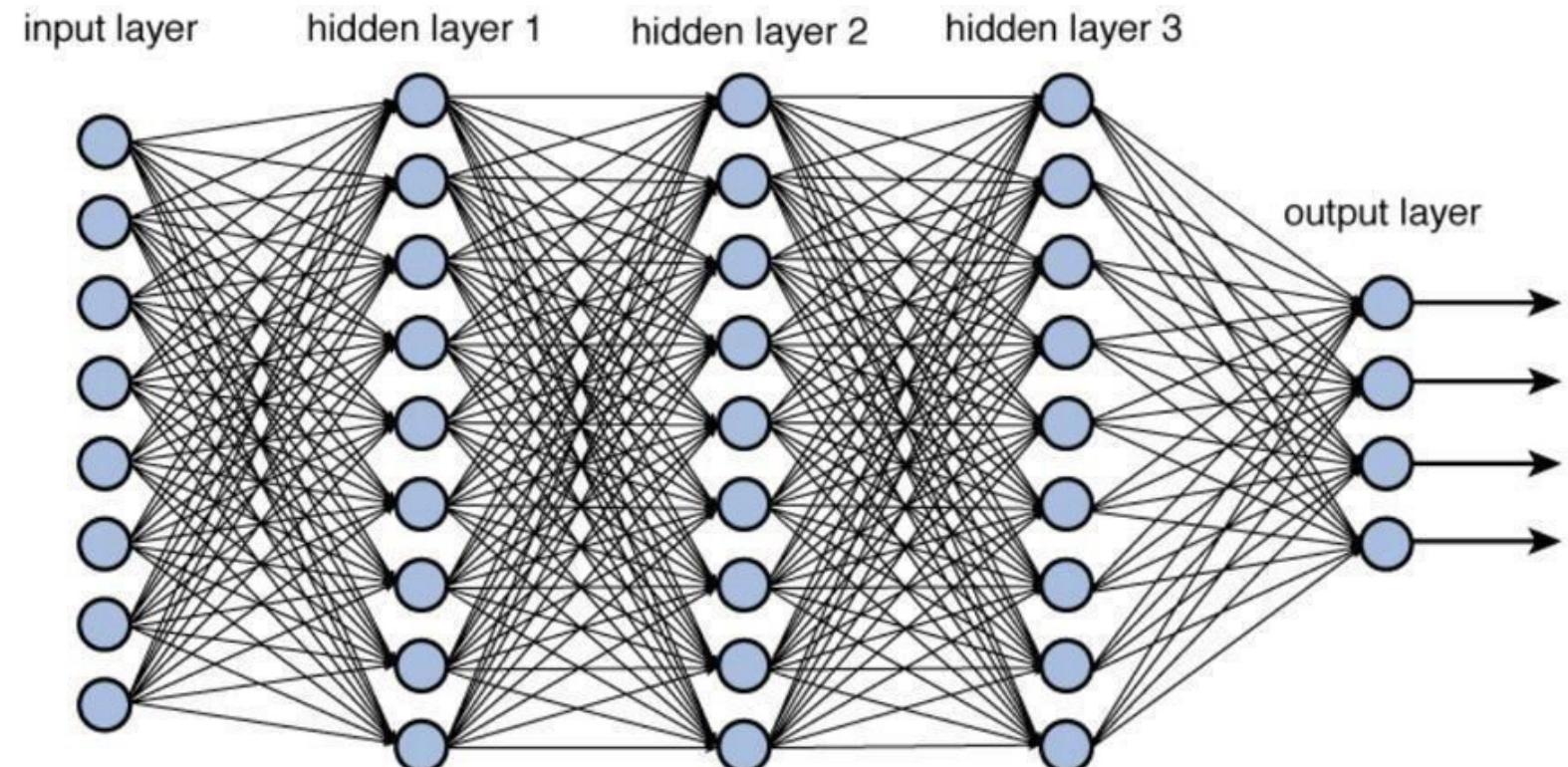
**DNN:**

*training DNN*

[https://colab.research.google.com/drive/13c9uJ\\_fPGjszgsyEuYWafR2F4\\_n-IXeZ](https://colab.research.google.com/drive/13c9uJ_fPGjszgsyEuYWafR2F4_n-IXeZ)

# deep neural net

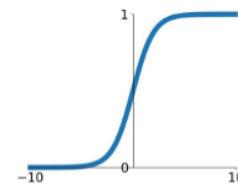
1986: Deep Neural Nets



*f: activation function:*  
turns neurons on-off

**Sigmoid**

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



*w: weight*

*sets the sensitivity of a neuron*

*b: bias:*

*up-down weights a neuron*

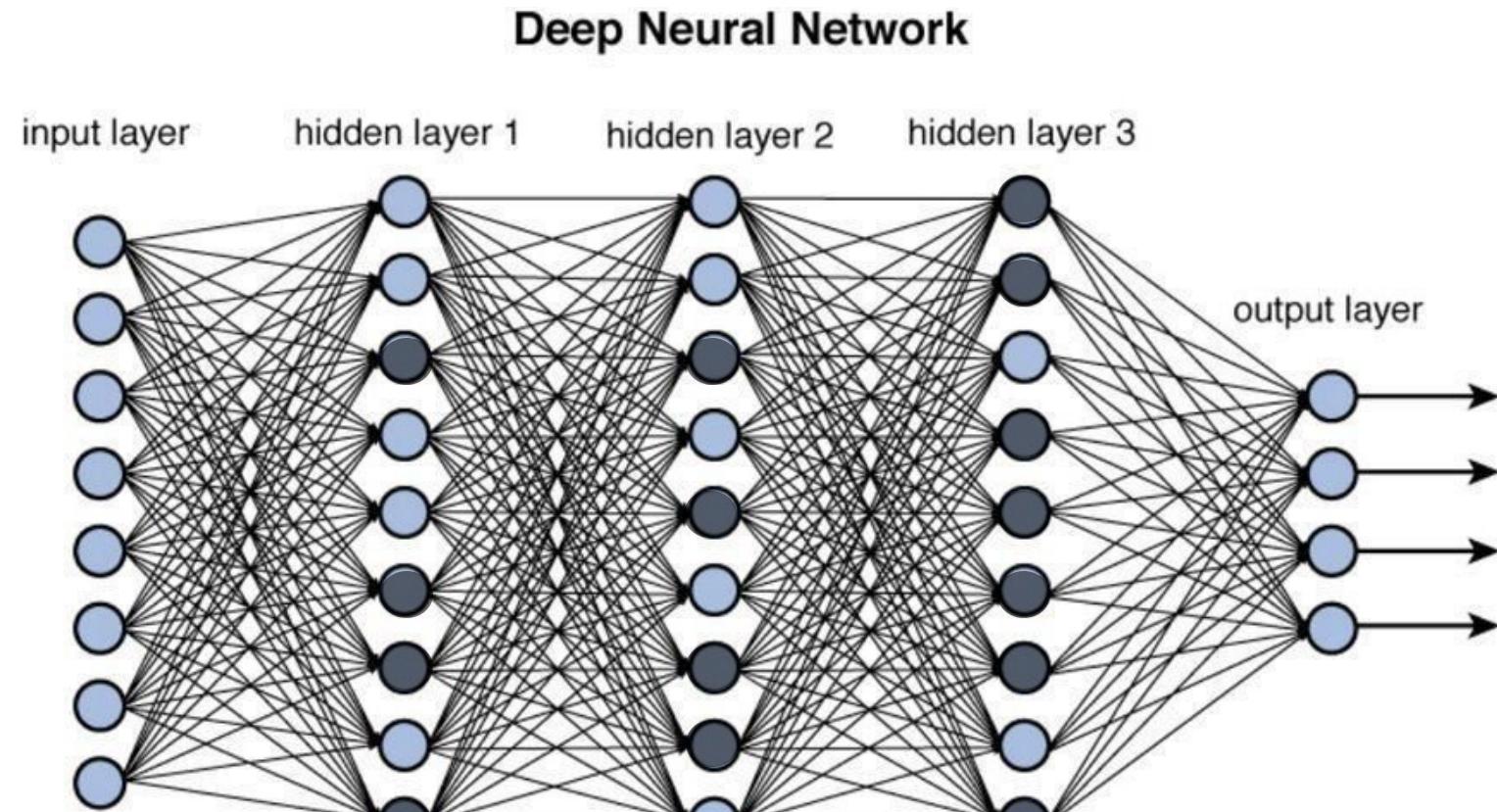


Figure 12.2 Deep network architecture with multiple layers.

In a CNN these layers would not be fully connected  
except the last one

$$\vec{y} = f_N(\dots(f_1(\vec{x}W_i + b_1\dots W_N + b_N)))$$

# Seminal paper

## Y. LeCun 1998

# Gradient-Based Learning Applied to Document Recognition

Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner

## *Abstract—*

Multilayer Neural Networks trained with the backpropagation algorithm constitute the best example of a successful Gradient-Based Learning technique. Given an appropriate network architecture, Gradient-Based Learning algorithms can be used to synthesize a complex decision surface that can classify high-dimensional patterns such as handwritten characters, with minimal preprocessing. This paper reviews various methods applied to handwritten character recognition and compares them on a standard handwritten digit recognition task. Convolutional Neural Networks, that are specifically designed to deal with the variability of 2D shapes, are shown to outperform all other techniques.

Real-life document recognition systems are composed of multiple modules including field extraction, segmentation, recognition, and language modeling. A new learning paradigm, called Graph Transformer Networks (GTN), allows such multi-module systems to be trained globally using Gradient-Based methods so as to minimize an overall performance measure.

Two systems for on-line handwriting recognition are described. Experiments demonstrate the advantage of global training, and the flexibility of Graph Transformer Networks.

A Graph Transformer Network for reading bank check is also described. It uses Convolutional Neural Network character recognizers combined with global training techniques to provides record accuracy on business and personal checks. It is deployed commercially and reads several million checks per day.

## I. INTRODUCTION

Over the last several years, machine learning techniques, particularly when applied to neural networks, have played an increasingly important role in the design of pattern recognition systems. In fact, it could be argued that the availability of learning techniques has been a crucial factor in the recent success of pattern recognition applications such as continuous speech recognition and handwriting recognition.

The main message of this paper is that better pattern recognition systems can be built by relying more on automatic learning, and less on hand-designed heuristics. This is made possible by recent progress in machine learning and computer technology. Using character recognition as a case study, we show that hand-crafted feature extraction can be advantageously replaced by carefully designed learning machines that operate directly on pixel images. Using document understanding as a case study, we show that the traditional way of building recognition systems by manually integrating individually designed modules can be replaced by a unified and well-principled design paradigm, called *Graph Transformer Networks*, that allows training

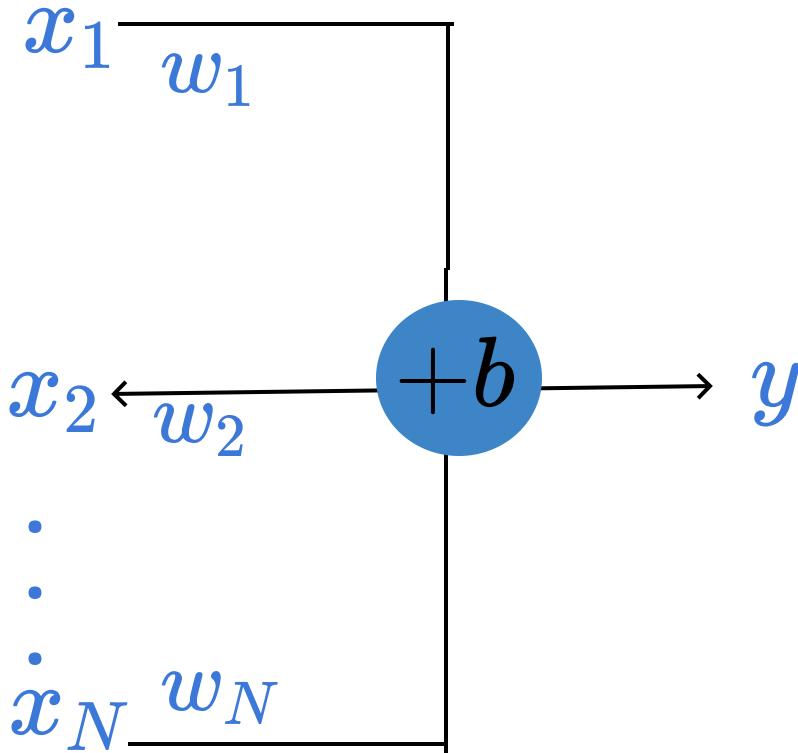
<http://yann.lecun.com/exdb/publis/pdf/lecun-01a.pdf>

# back-propagation

$y$ : prediction

Find the best parameters by finding the minimum of the L2 hyperplane

Any linear model:



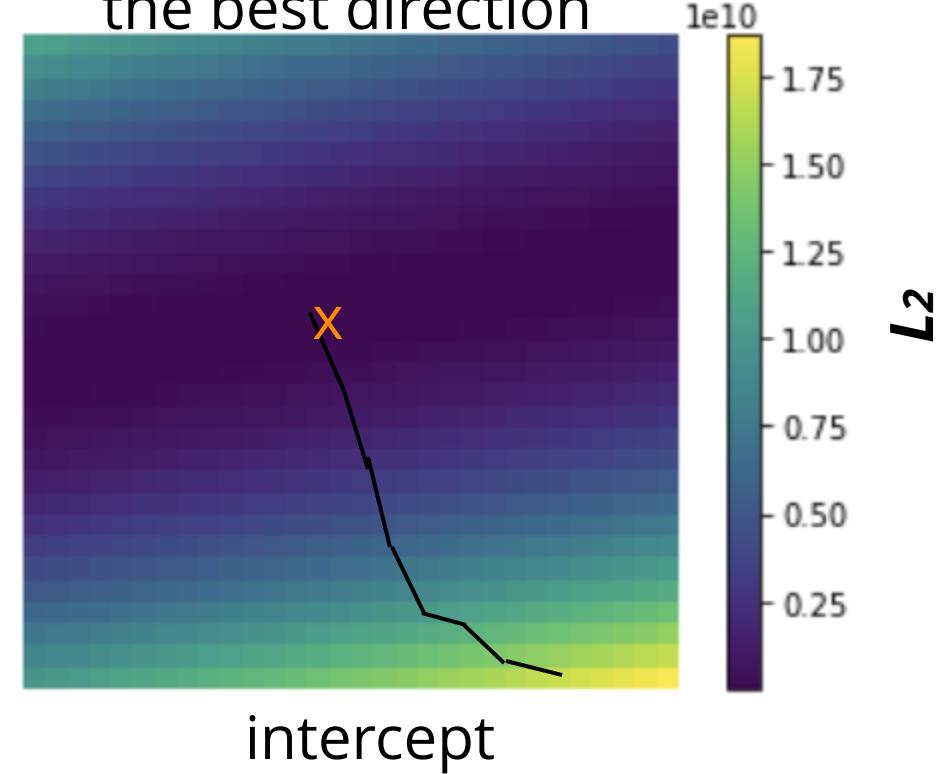
Error: e.g.

$$L_2 = (y - y_{\text{true}})^2$$

$y_{\text{true}}$  : target

at every step look around and choose the best direction

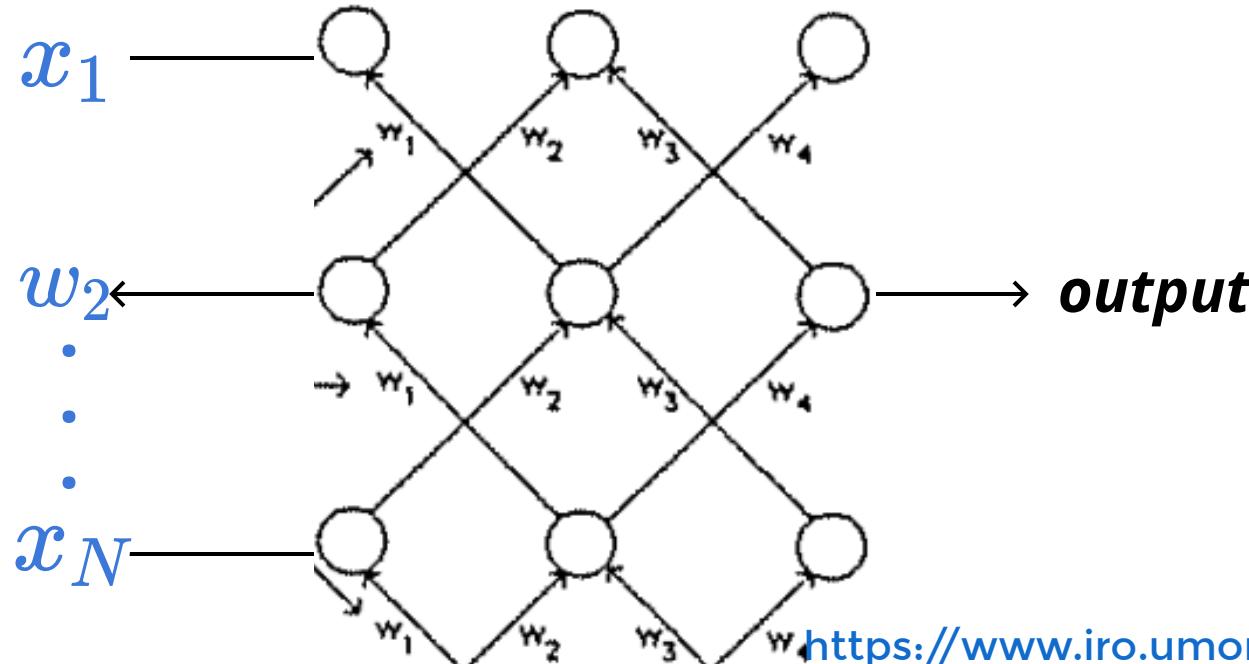
slope



# back-propagation

how does linear descent look when you have a whole network structure with hundreds of weights and biases to optimize??

$$x_j = \sum_i y_i w_{ji} \quad y_j = \frac{1}{1+e^{-x_j}}$$



nature

## Learning representations by back-propagating errors

David E. Rumelhart, Geoffrey E. Hinton & Ronald J. Williams

Nature 323, 533–536(1986) | Cite this article

22k Accesses | 7872 Citations | 167 Altmetric | Metrics

### Abstract

We describe a new learning procedure, back-propagation, for networks of neurone-like units. The procedure repeatedly adjusts the weights of the connections in the network so as to minimize a measure of the difference between the actual output vector of the net and the desired output vector. As a result of the weight adjustments, internal ‘hidden’ units which are not part of the input or output come to represent important features of the task domain, and the regularities in the task are captured by the interactions of these units. The ability to create useful new features distinguishes back-propagation from earlier, simpler methods such as the perceptron-convergence procedure<sup>1</sup>.

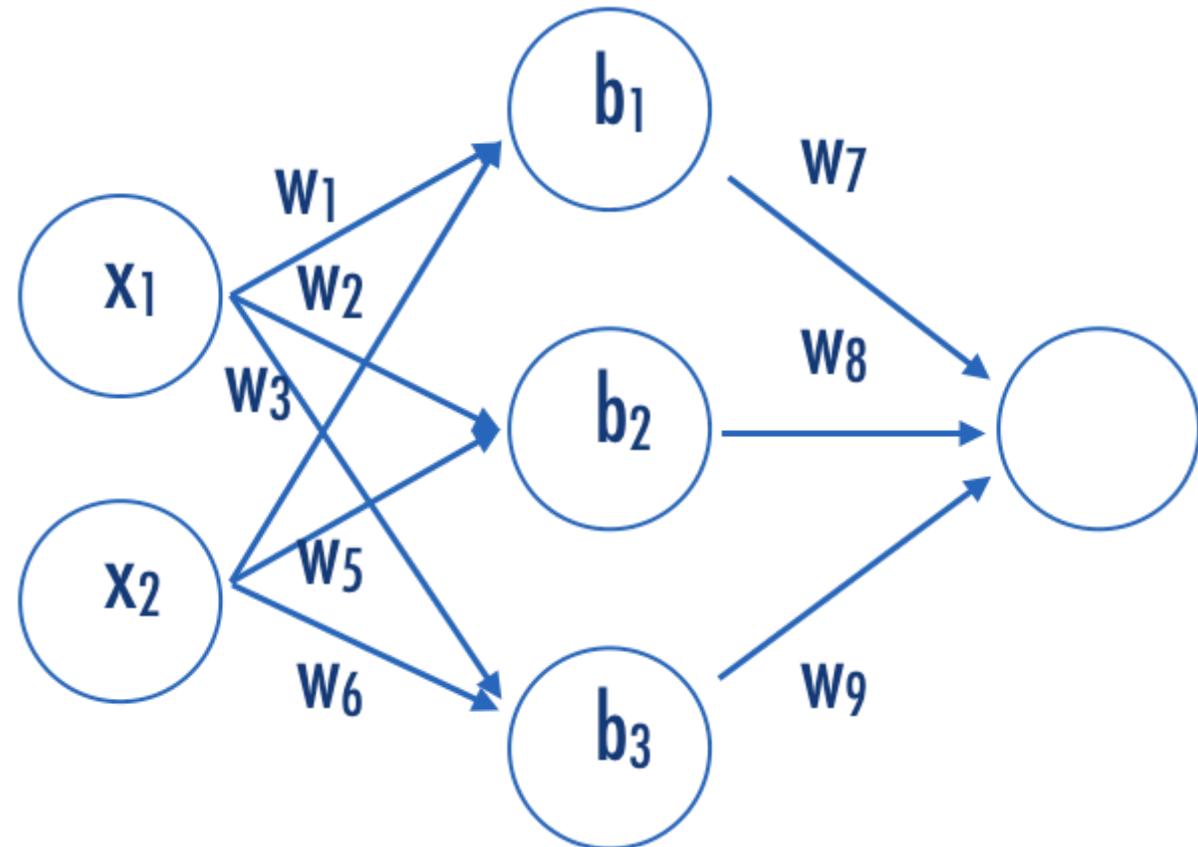
Training models with this many parameters requires a lot of care:

- . defining the metric
- . optimization schemes
- . training/validation/testing sets

But just like our simple linear regression case, the fact that small changes in the parameters leads to small changes in the output for the right activation functions.

define a cost function, e.g.

$$C = \frac{1}{2} |y - a^L|^2 = \frac{1}{2} \sum_j (y_j - a_j^L)^2$$



$$\text{output} = \frac{1}{1 + e^{-\frac{w_7}{1 + e^{-w_1 x_1 - w_4 x_2 - b_1}} - \frac{w_8}{1 + e^{-w_2 x_1 - w_5 x_2 - b_2}} - \frac{w_9}{1 + e^{-w_3 x_1 - w_6 x_2 - b_3}} - b_4}}$$

$$\vec{y} = f_N(\dots(f_1(\vec{x}W_i + b_1\dots W_N + b_N)))$$

# Training a DNN

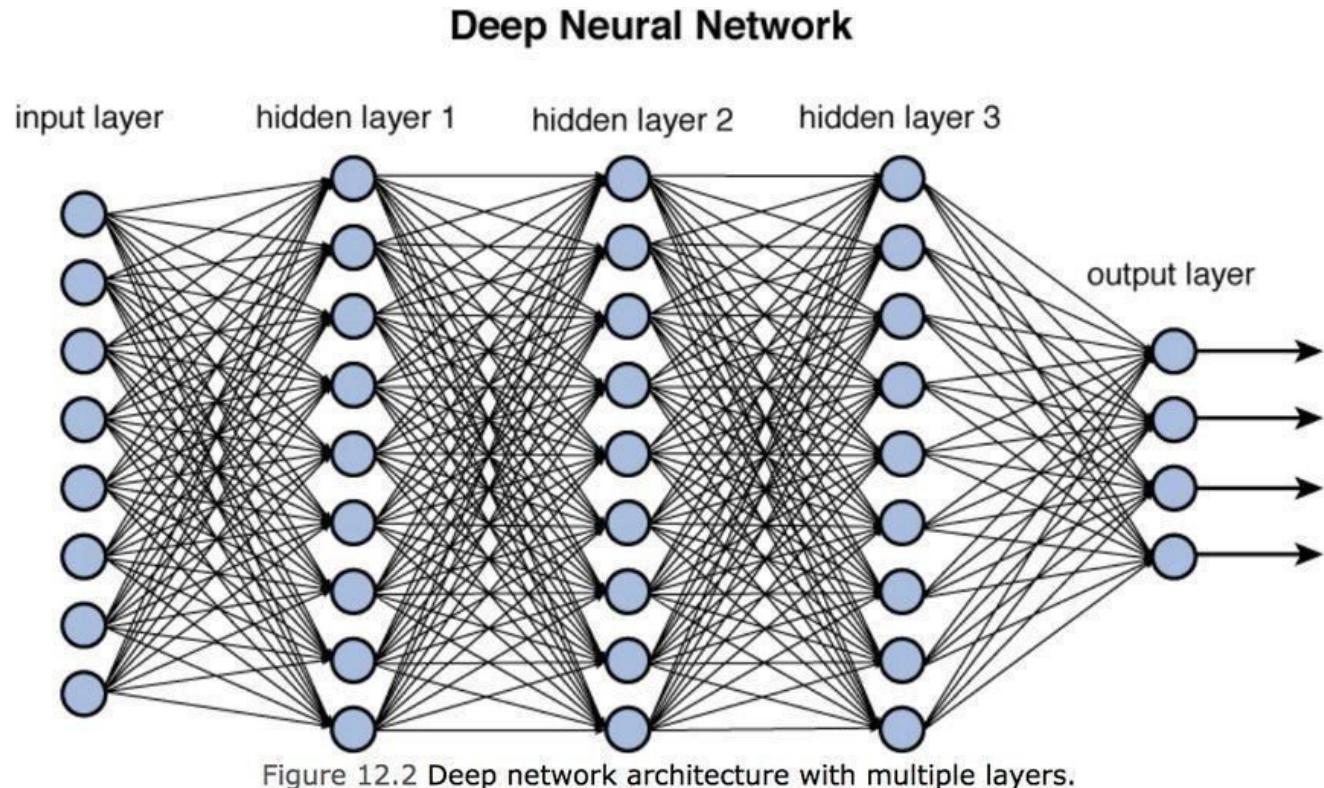
Training models with this many parameters requires a lot of care:

- . defining the metric
- . optimization schemes
- . training/validation/testing sets

But just like our simple linear regression case, the fact that small changes in the parameters leads to small changes in the output for the right activation functions.

define a cost function, e.g.

$$C = \frac{1}{2} |y - a^L|^2 = \frac{1}{2} \sum_j (y_j - a_j^L)^2$$



feed data forward through network and calculate cost metric

for each layer, calculate effect of small changes on next layer

$$\vec{y} = f_N(\dots(f_1(\vec{x}W_i + b_1\dots W_N + b_N)))$$

# Training a DNN back-propagation

how does linear descent look when you have a whole network structure with hundreds of weights and biases to optimize??

think of applying just gradient to a function of a function of a function... use:

- 1) partial derivatives, 2) chain rule

define a cost function, e.g.  $C = \frac{1}{2}|y - a^L|^2 = \frac{1}{2} \sum_j (y_j - a_j^L)^2$

**An equation for the error in the output layer,  $\delta^L$ :** The components of  $\delta^L$  are given by

$$\delta_j^L = \frac{\partial C}{\partial a_j^L} \sigma'(z_j^L). \quad (\text{BP1})$$

matrix-based form, as

$$\delta^L = \nabla_a C \odot \sigma'(z^L). \quad (\text{BP1a})$$

Here,  $\nabla_a C$  is defined to be a vector whose components are the partial derivatives  $\partial C / \partial a_j^L$ . You can think of  $\nabla_a C$  as expressing the rate of change of  $C$  with respect to the output activations

The backpropagation equations provide us with a way of computing the gradient of the cost function. Let's explicitly write this out in the form of an algorithm:

1. **Input  $x$ :** Set the corresponding activation  $a^1$  for the input layer.
2. **Feedforward:** For each  $l = 2, 3, \dots, L$  compute  $z^l = w^l a^{l-1} + b^l$  and  $a^l = \sigma(z^l)$ .
3. **Output error  $\delta^L$ :** Compute the vector  $\delta^L = \nabla_a C \odot \sigma'(z^L)$ .
4. **Backpropagate the error:** For each  $l = L-1, L-2, \dots, 2$  compute  $\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$ .
5. **Output:** The gradient of the cost function is given by

$$\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l \text{ and } \frac{\partial C}{\partial b_j^l} = \delta_j^l.$$

Examining the algorithm you can see why it's called *backpropagation*. We compute the error vectors  $\delta^l$  backward, starting from the final layer. It may seem peculiar that we're going

<http://neuralnetworksanddeeplearning.com/chap2.html>

$$\vec{y} = f_N(\dots(f_1(\vec{x}W_i + b_1\dots W_N + b_N)))$$

# *Punch Line*

Deep Neural Net are not  
some fancy-pants  
methods, they are just  
linear models with a  
bunch of parameters

# *Black Box?*

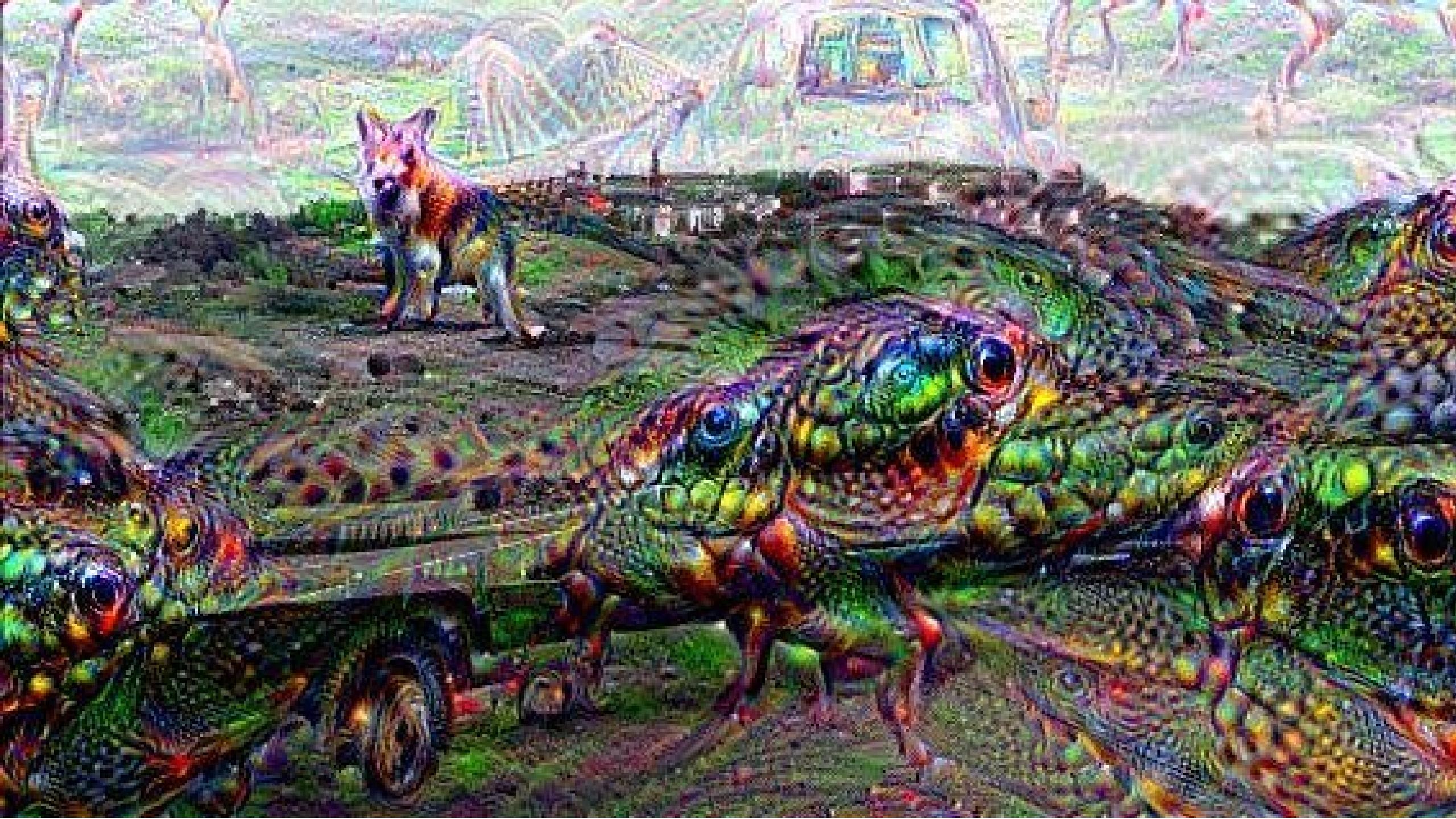
Because they have many  
parameters they are  
difficult to "interpret" (no  
easy feature extraction)

that is ok because they are  
prediction machines

A deep dream image of a dog's face, where the features are heavily stylized and colored with a rainbow-like palette. The dog has large, expressive eyes and a wide, open mouth. The background is a blurred, colorful landscape of green fields and distant buildings, creating a dreamlike atmosphere.

*deep dreams*





# what is happening in DeepDream?

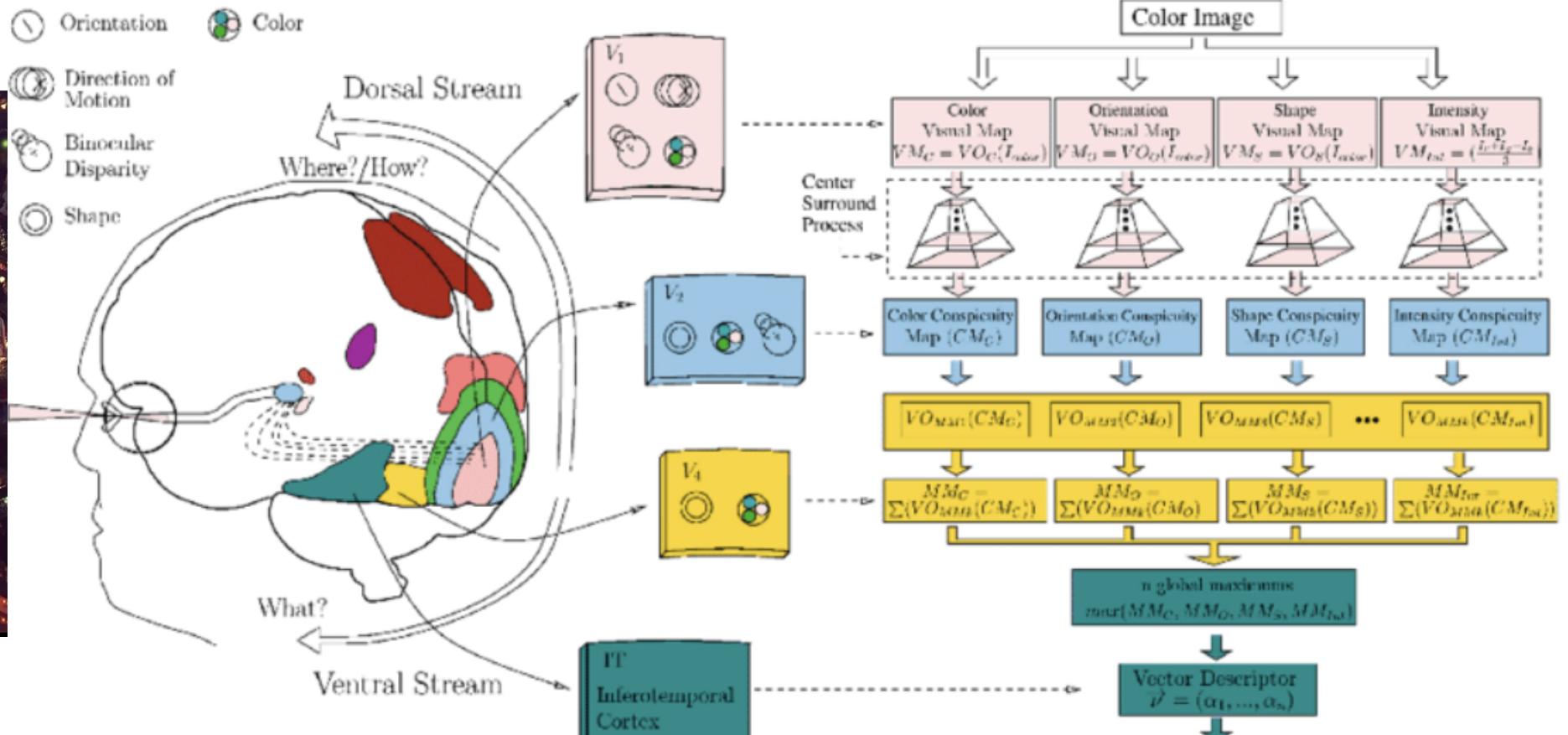
Deep Dream (DD) is a google software, a pre-trained NN (originally created on the Cafe architecture, now imported on many other platforms including tensorflow).

The high level idea relies on training a convolutional NN to recognize common objects, e.g. dogs, cats, cars, in images. As the network learns to recognize those objects it develops its layers to pick out "features" of the NN, like lines at a certain orientations, circles, etc.

The DD software runs this NN on an image you give it, and it loops on some layers, thus "manifesting" the things it knows how to recognize in the image.

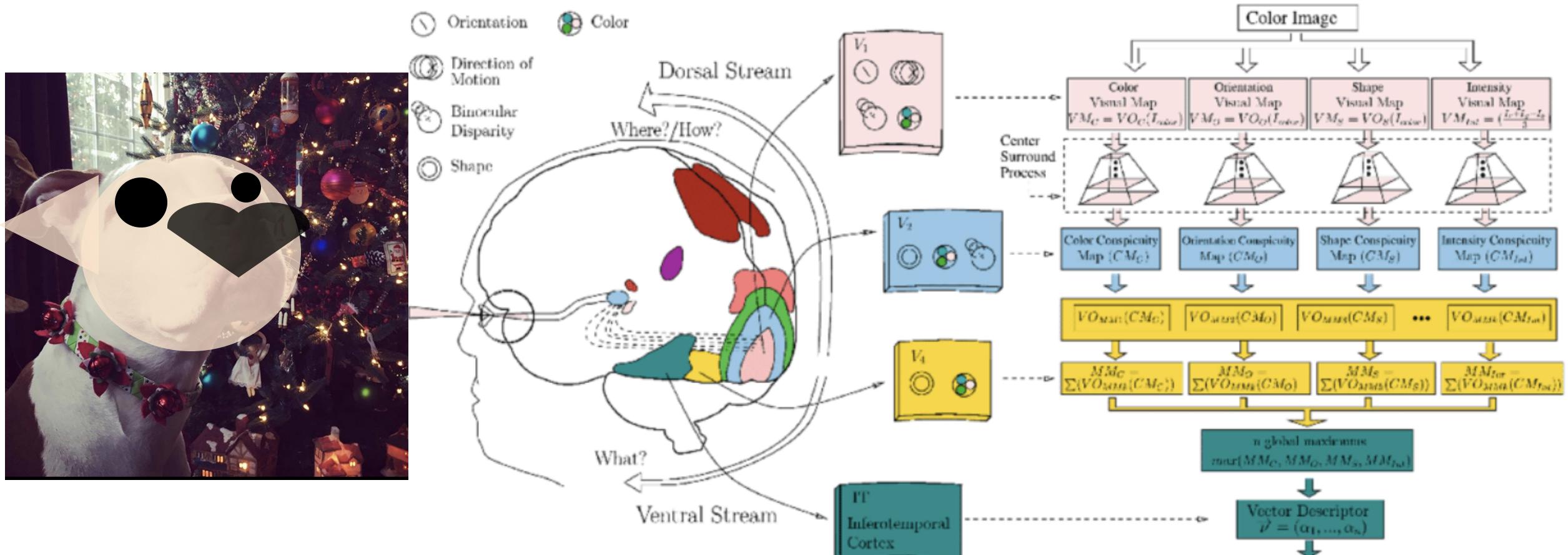


@akumadog

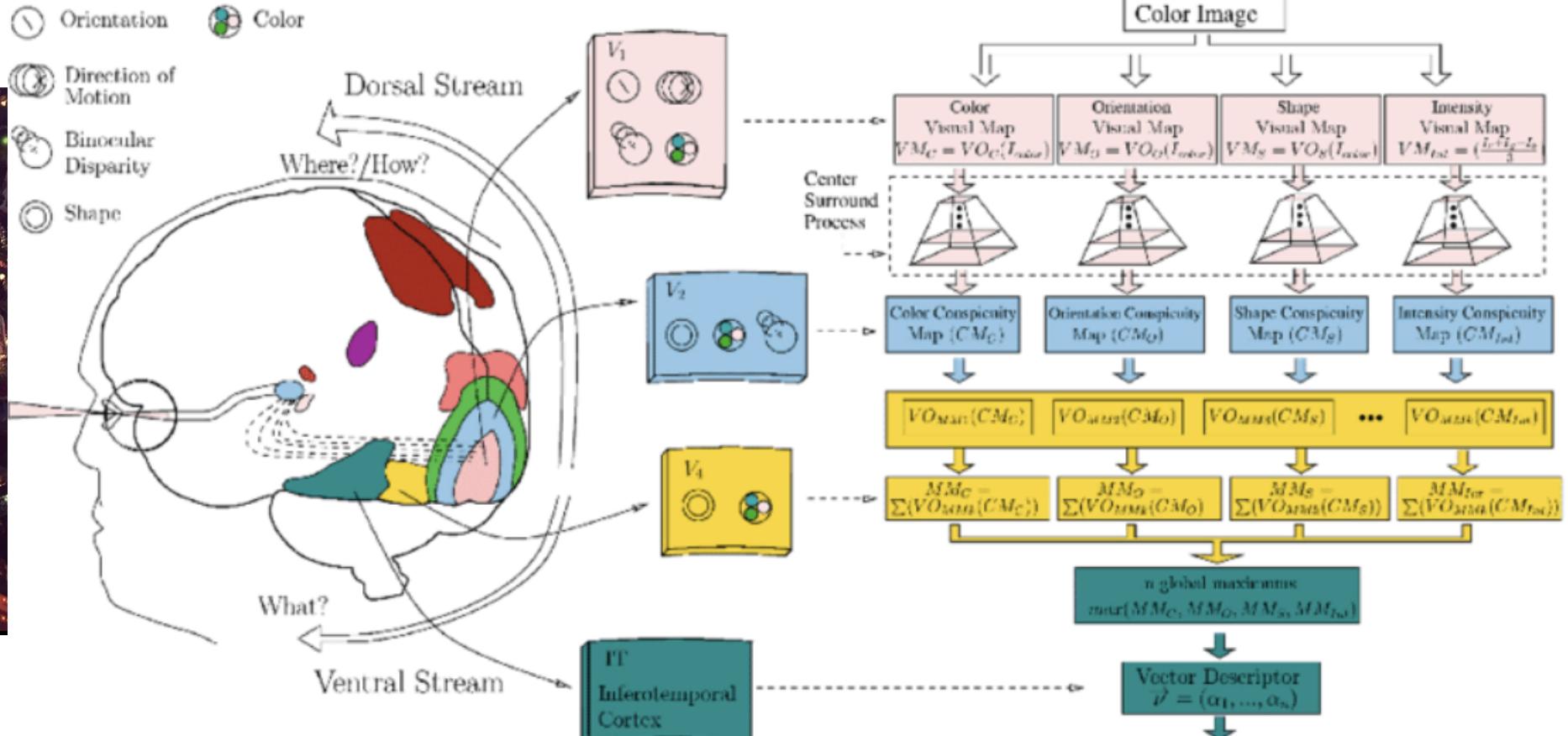


Brain Programming and the Random Search in Object Categorization Olague et al 2017

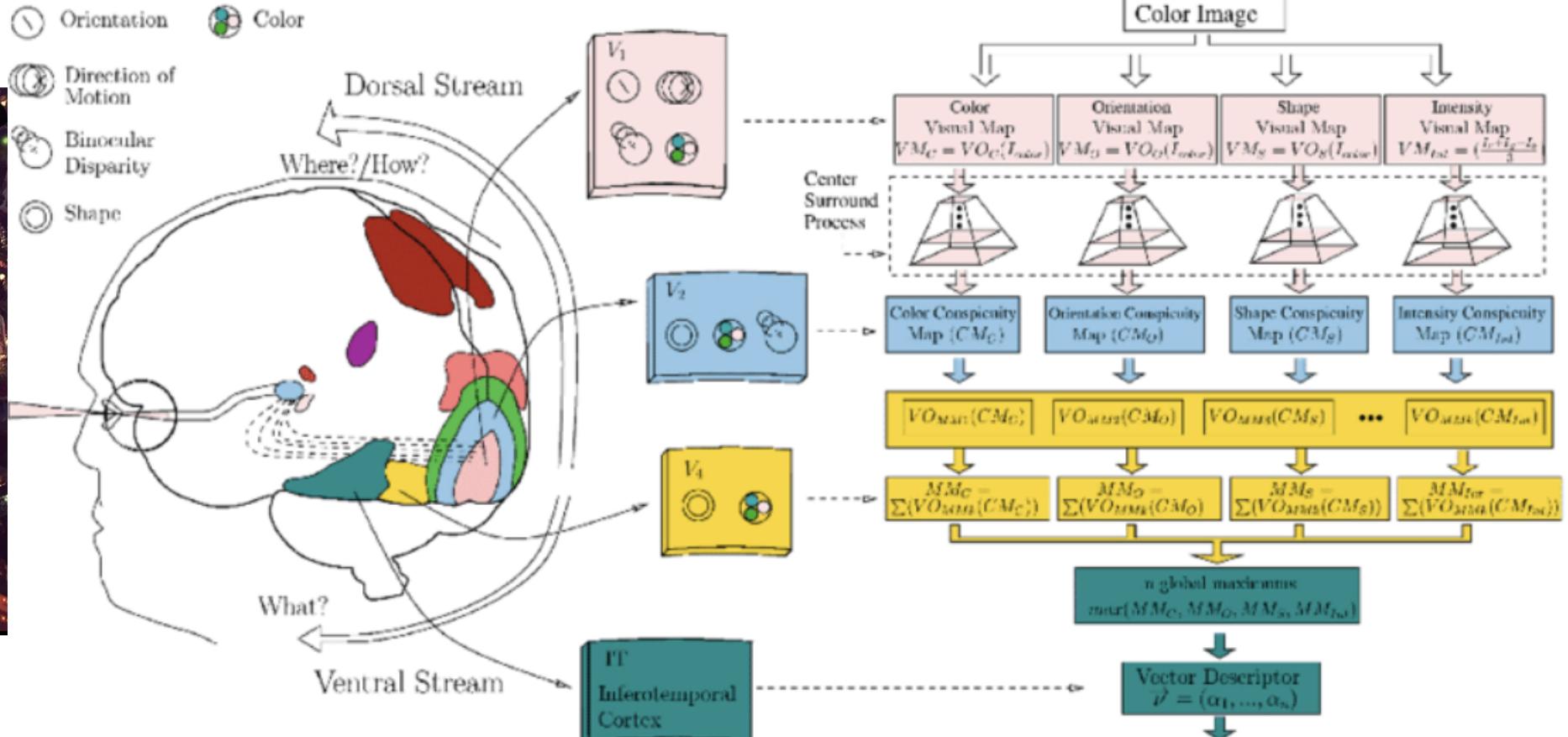
*The visual cortex learns hierarchically: first detects simple features, then more complex features and ensembles of features*



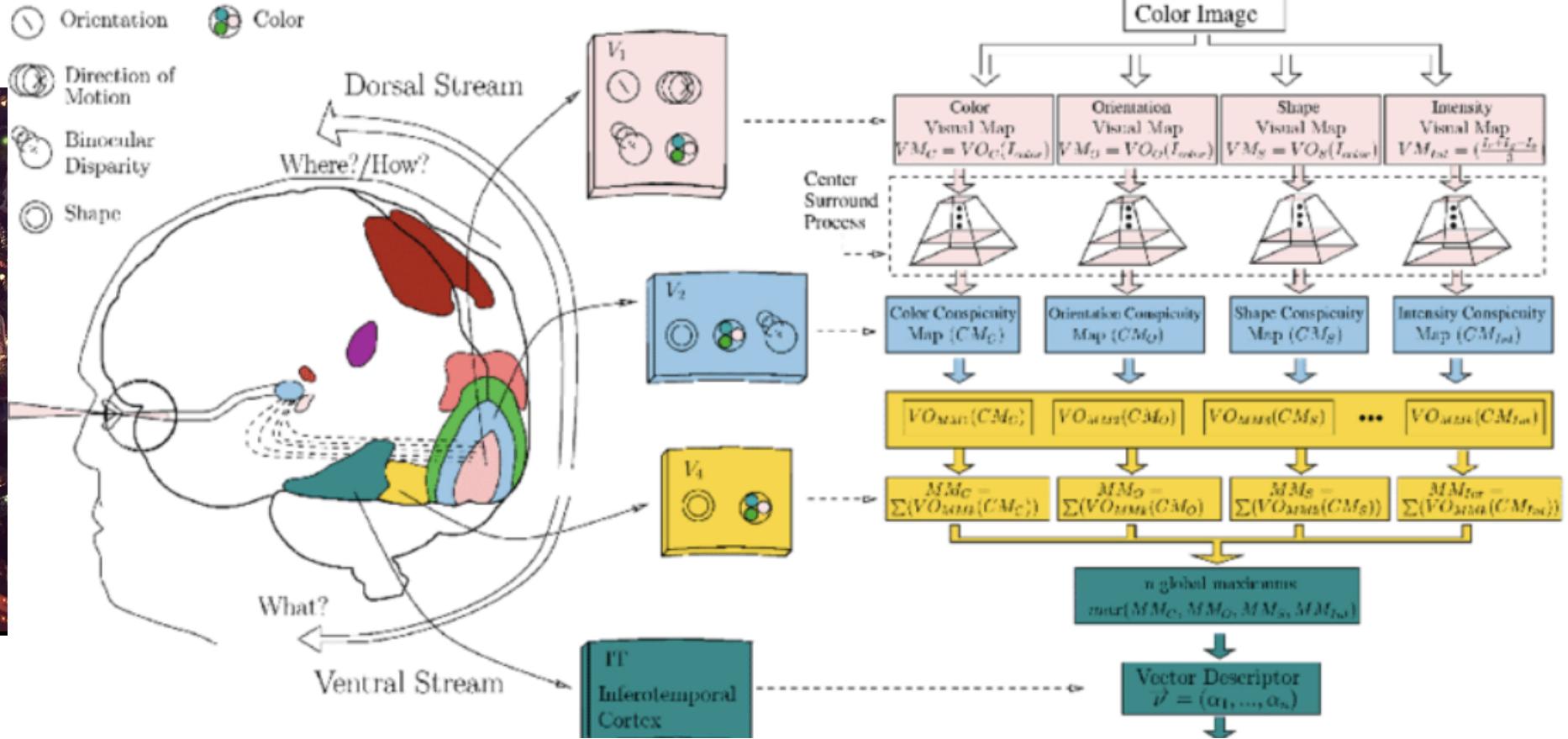
*The visual cortex learns hierarchically: first detects simple features, then more complex features and ensembles of features*



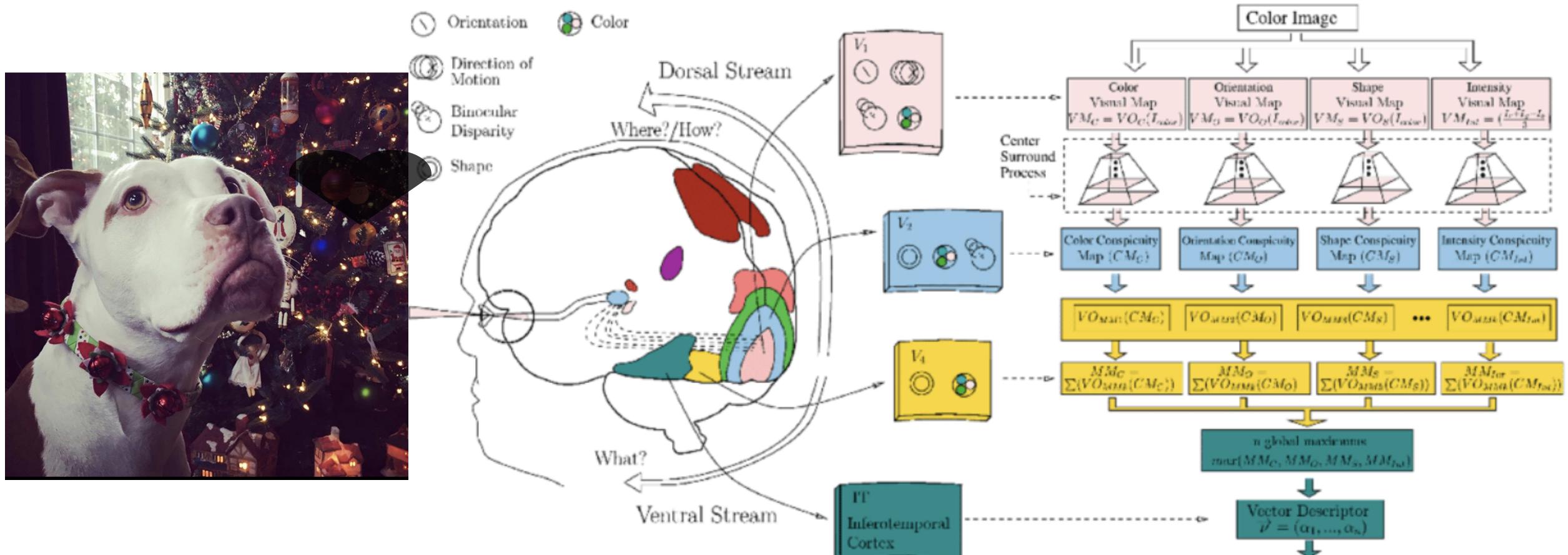
*The visual cortex learns hierarchically: first detects simple features, then more complex features and ensembles of features*



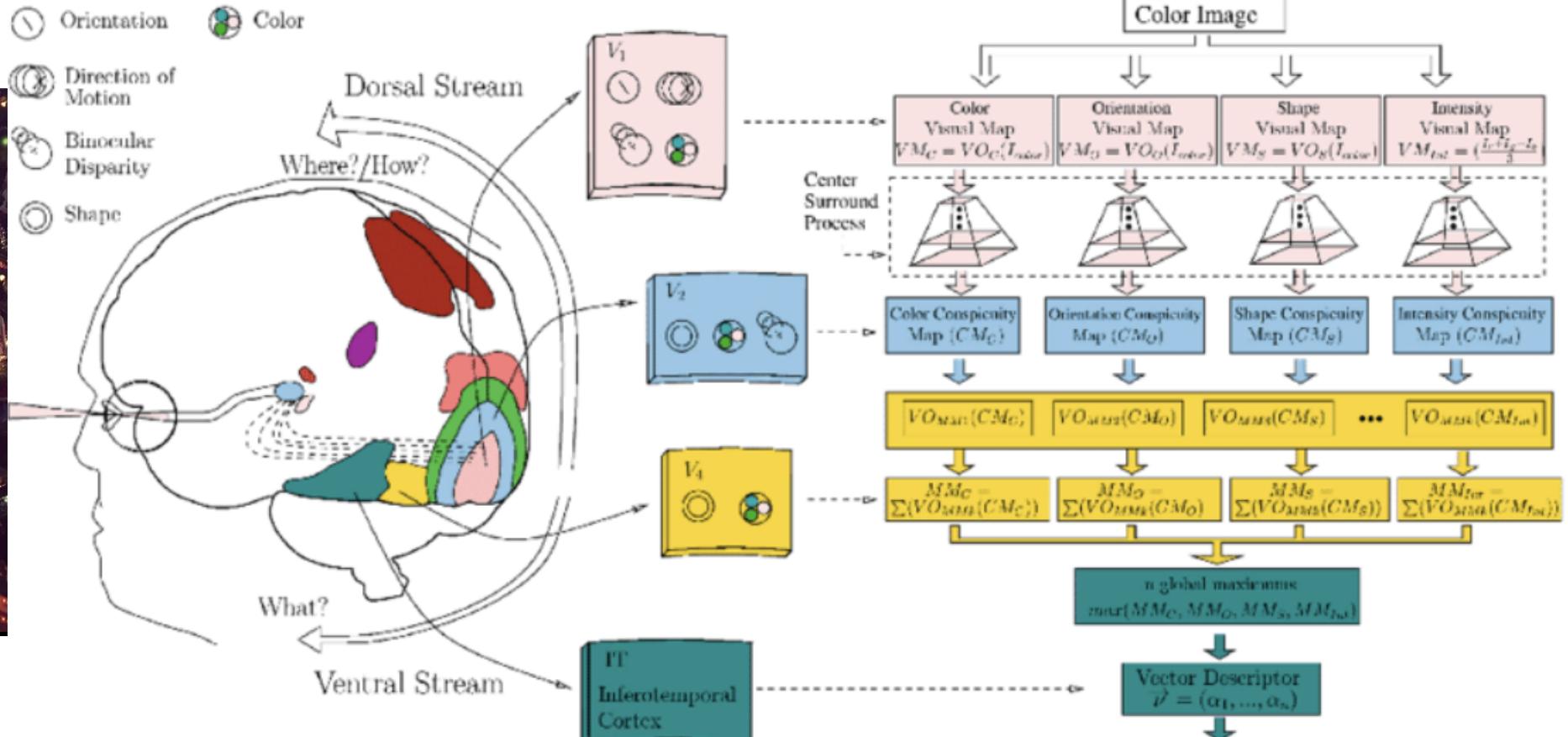
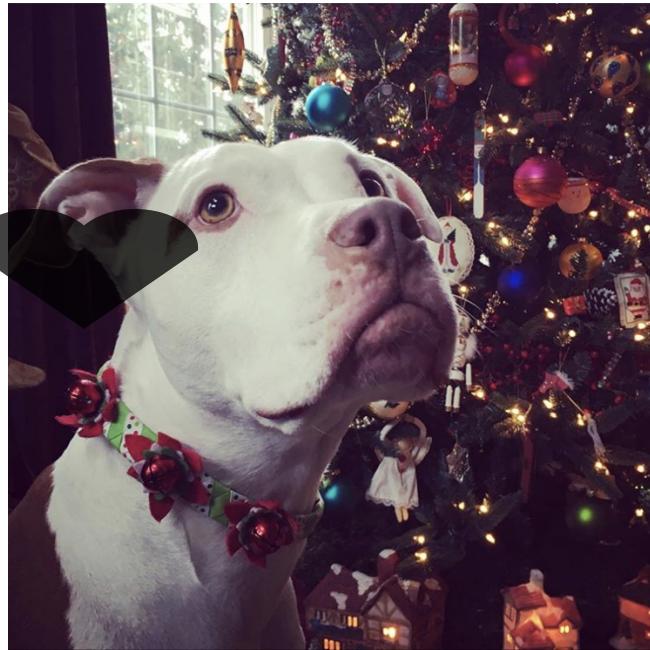
*The visual cortex learns hierarchically: first detects simple features, then more complex features and ensembles of features*



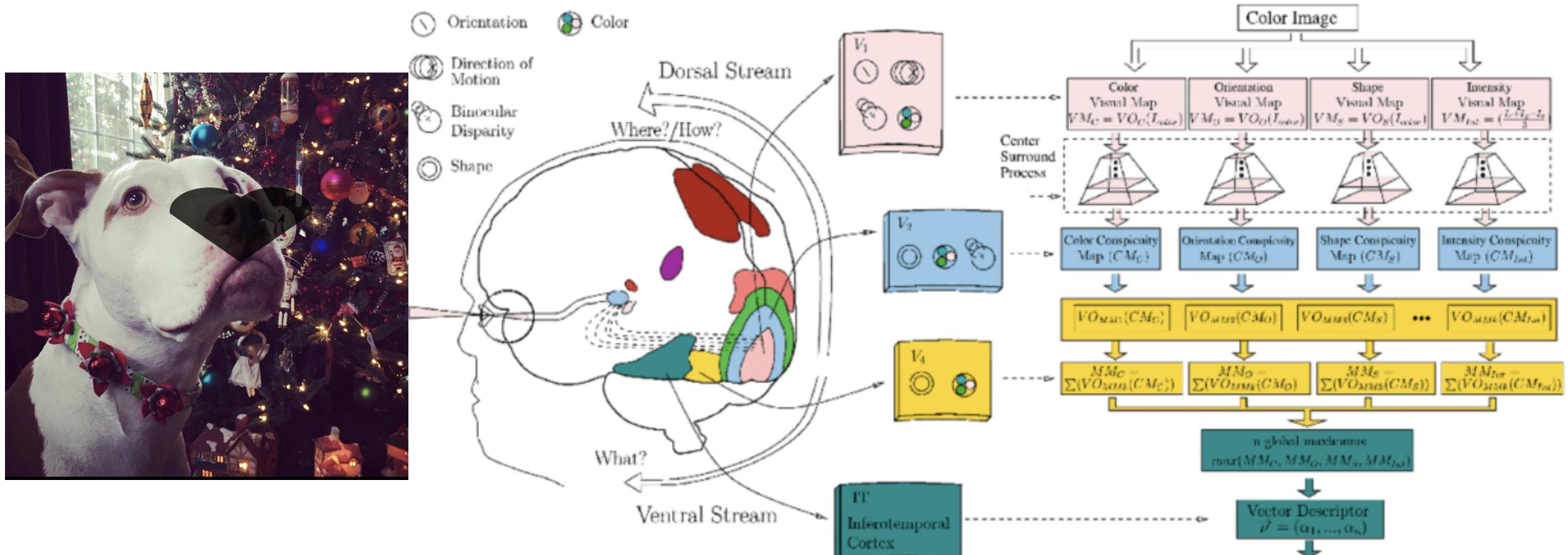
*The visual cortex learns hierarchically: first detects simple features, then more complex features and ensembles of features*



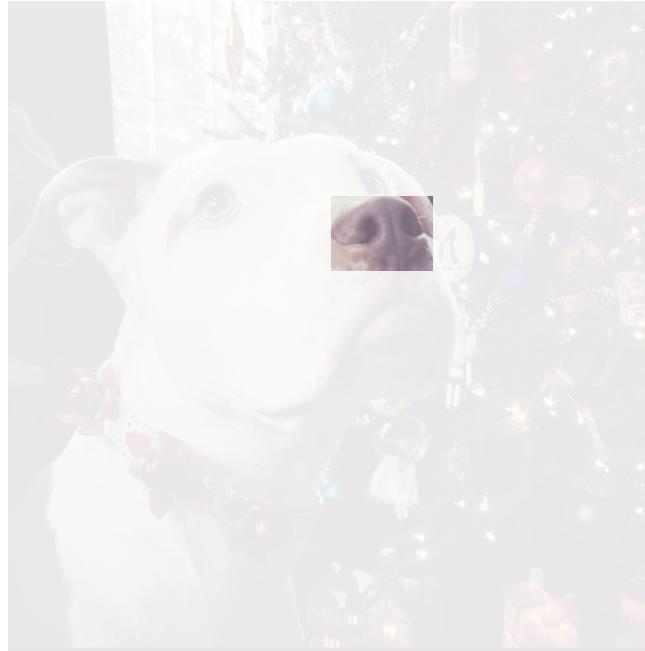
*The visual cortex learns hierarchically: first detects simple features, then more complex features and ensembles of features*



*The visual cortex learns hierarchically: first detects simple features, then more complex features and ensembles of features*



*The visual cortex learns hierarchically: first detects simple features, then more complex features and ensembles of features*



Orientation

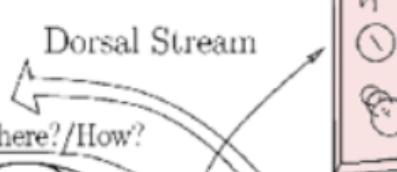
Direction of Motion

Binocular Disparity

Shape

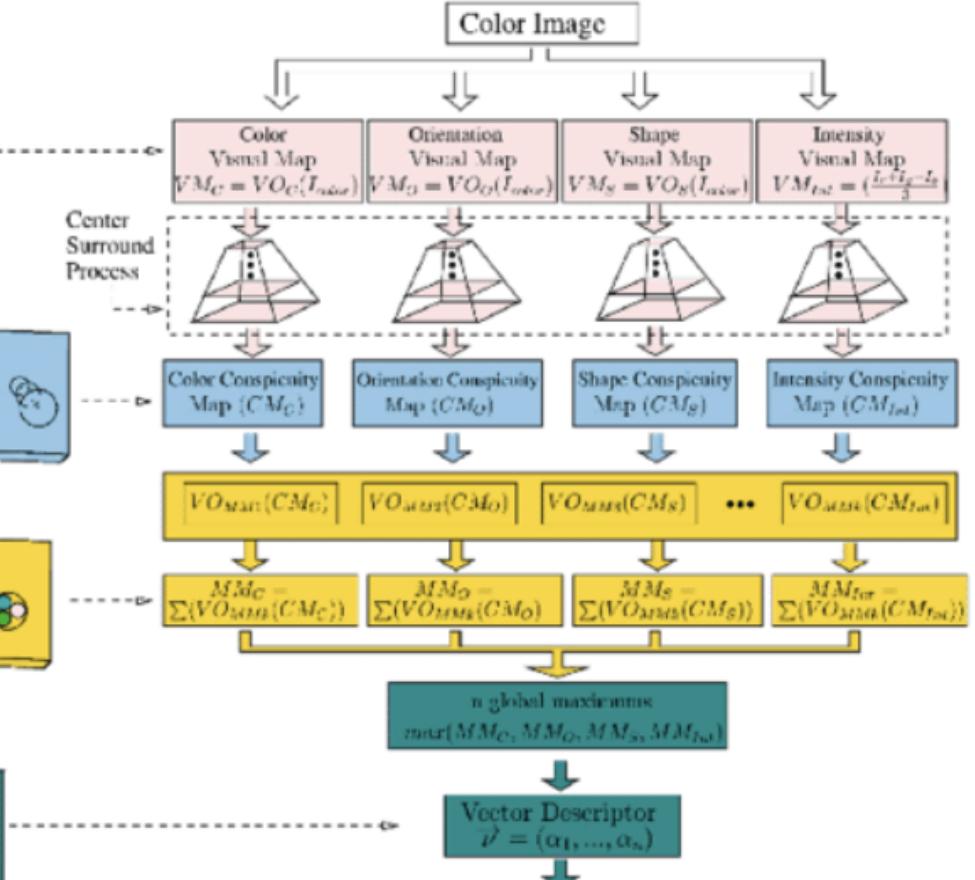
Color

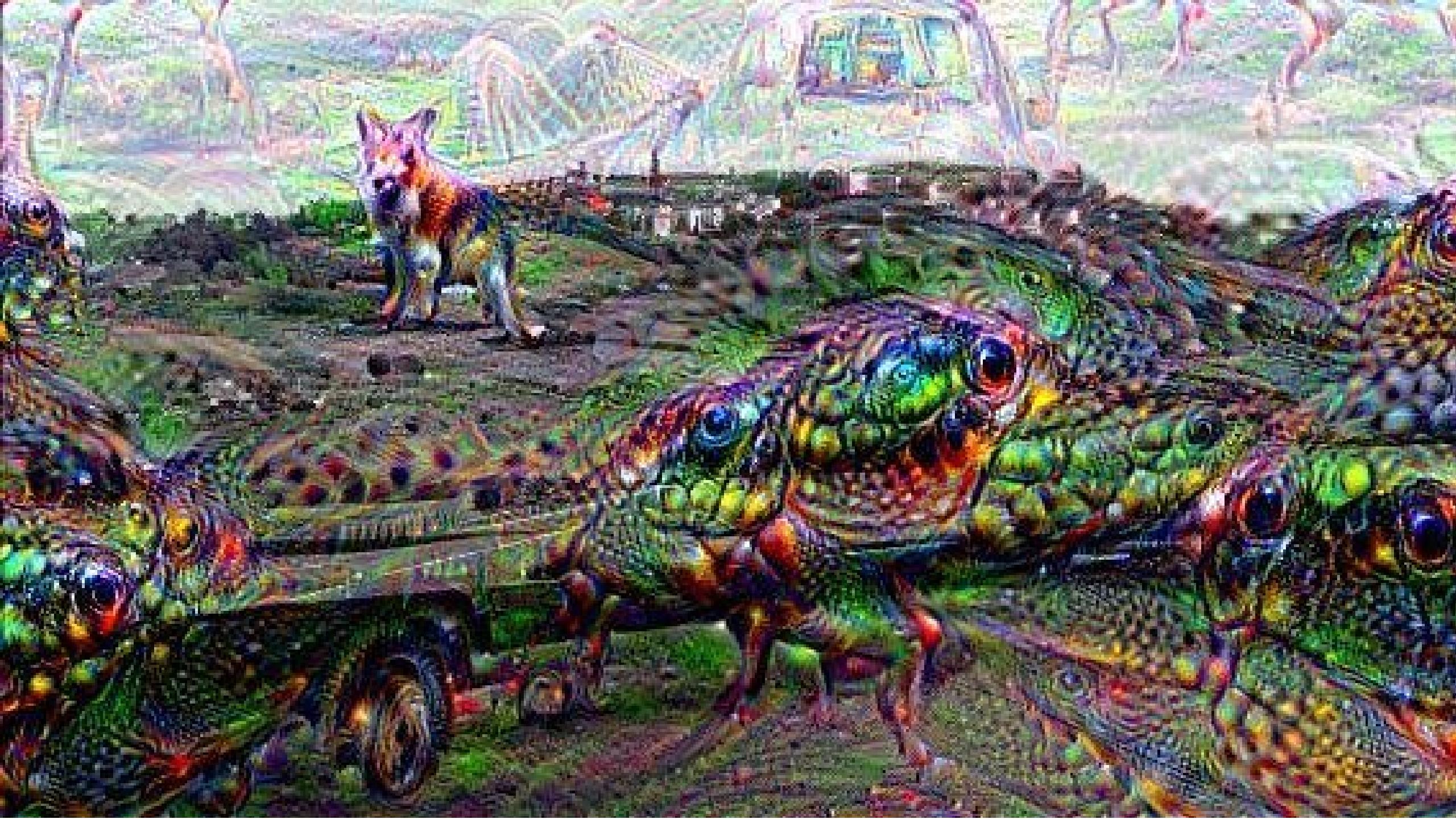
Dorsal Stream  
Where?/How?



Ventral Stream  
What?

IT  
Inferotemporal Cortex





Neural Network and Deep Learning

an excellent and free book on NN and DL

<http://neuralnetworksanddeeplearning.com/index.html>

History of NN

<https://cs.stanford.edu/people/eroberts/courses/soco/projects/neural-networks/History/history2.html>

resources

# Inceptionism: Going Deeper into Neural Networks

Wednesday, June 17, 2015

<https://ai.googleblog.com/2015/06/inceptionism-going-deeper-into-neural.html>

[https://github.com/fedhere/PUS2020\\_FBianco/blob/master/HW13/Instructions\\_PUS2020\\_deeplab.ipynb](https://github.com/fedhere/PUS2020_FBianco/blob/master/HW13/Instructions_PUS2020_deeplab.ipynb)

homework