# Introduction to deep learning in computational biology

**María Rodríguez Martínez**

Technical Lead Systems Biology, IBM, Zürich Research Lab
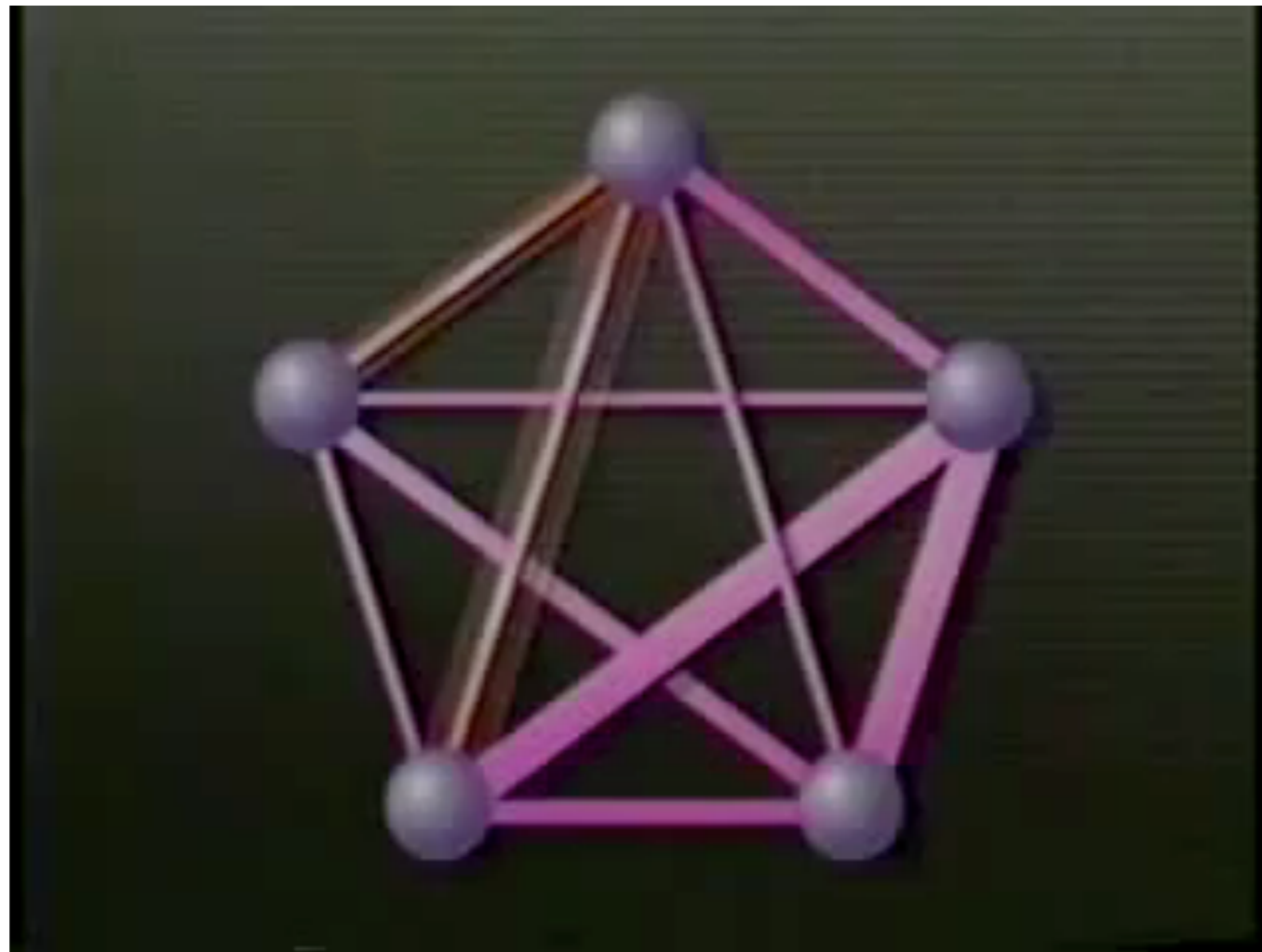
# Overview

- Introduction to deep learning
  - History and motivation
  - Activations functions
  - Cost functions
  - Backpropagation
  - Regularization
  - Optimization
- Multi-Layer Perceptron (MLP)
- Auto-enconders (AE)
- Convolutional Neural Networks (CNN)
- Recurrent Neural Networks (RNN)

# Overview

- Introduction to deep learning
  - History and motivation
  - Activations functions
  - Cost functions
  - Backpropagation
  - Regularization
  - Optimization
- Multi-Layer Perceptron (MLP)
- Auto-enconders (AE)
- Convolutional Neural Networks (CNN)
- Recurrent Neural Networks (RNN)
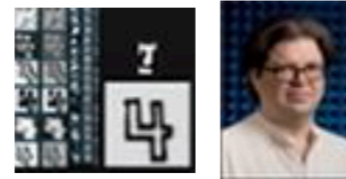
# The beginnings: perceptron.

# A brief history of deep learning
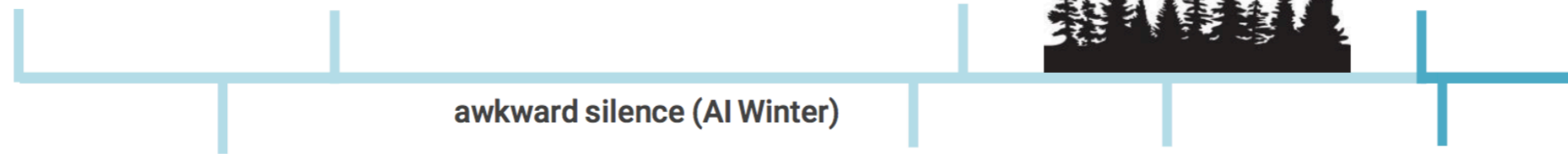


**1958** Perceptron     **1974** Backpropagation

Convolution Neural Networks for Handwritten Recognition
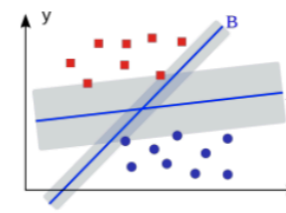**1998**

Google Brain Project on 16k Cores
**2012**

awkward silence (AI Winter)

**1969**
Perceptron criticized

**1995**
SVM reigns

**2006**
Restricted Boltzmann Machine

**2012**
AlexNet wins ImageNet
IMGENET

https://www.slideshare.net/LuMa921/deep-learning-a-visual-introduction
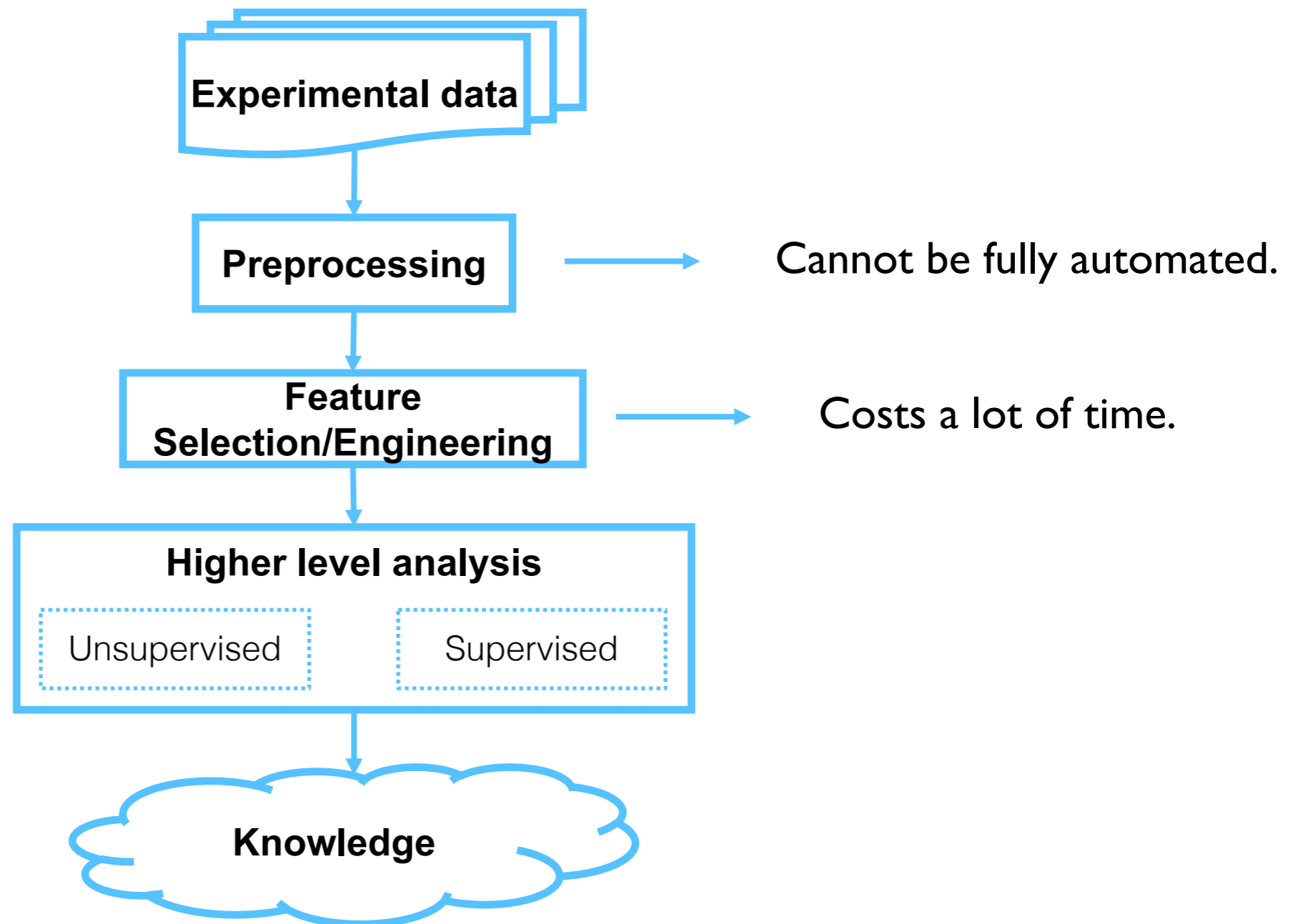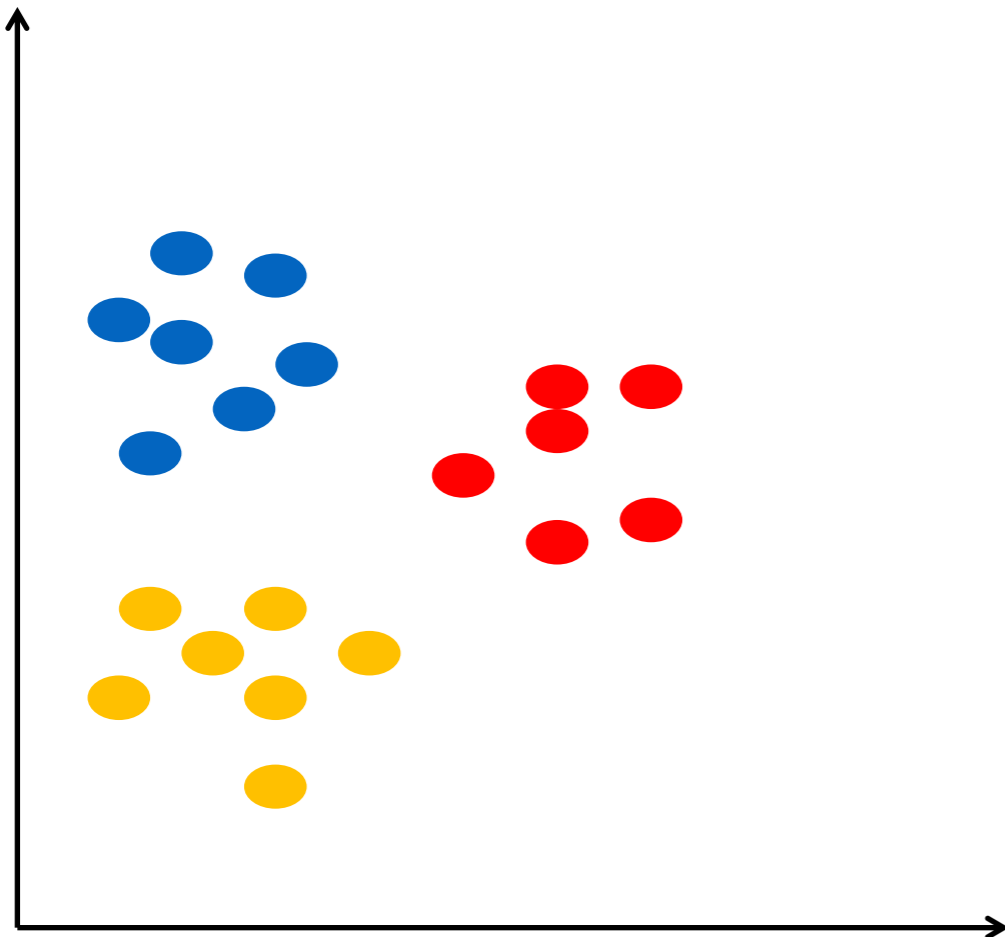
# Machine learning

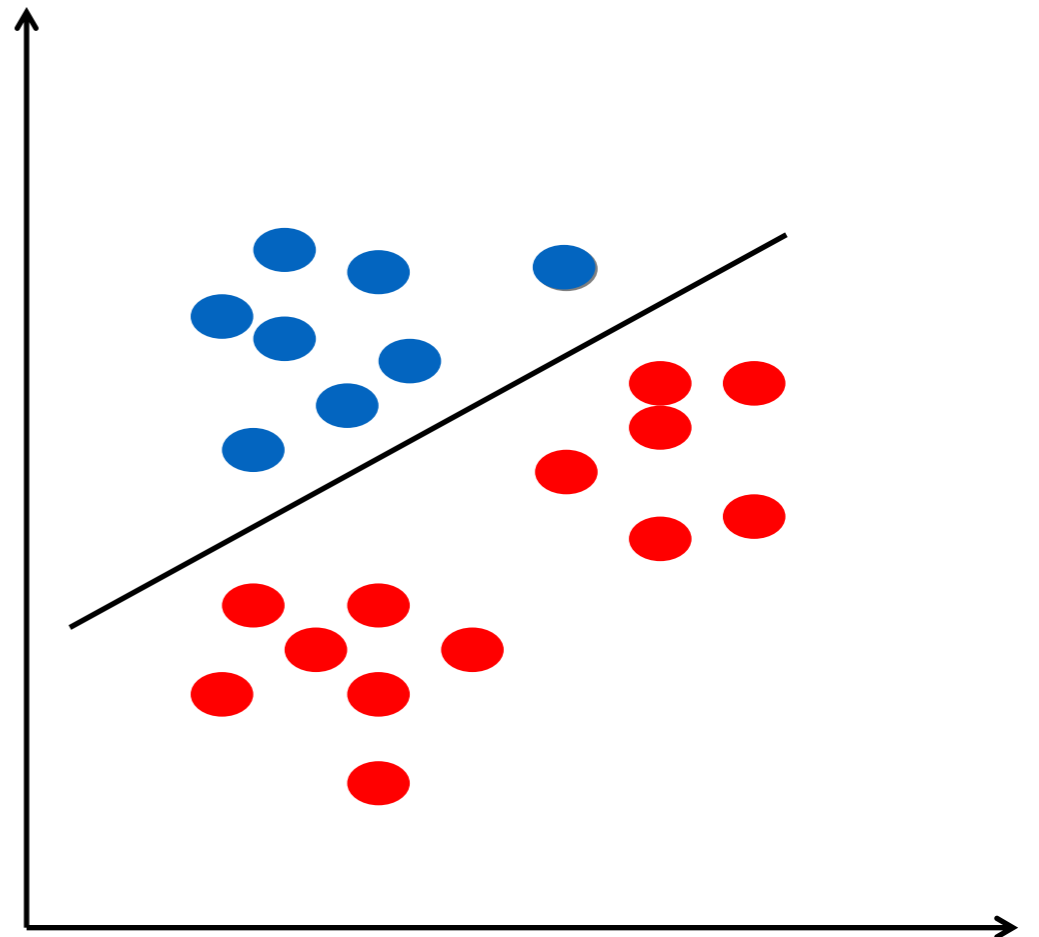Machine Learning is a type of Artificial Intelligence that provides computers with the ability to learn without being explicitly programmed.



**Experimental data**

**Preprocessing** → Cannot be fully automated.

**Feature Selection/Engineering** → Costs a lot of time.

**Higher level analysis**

Unsupervised    Supervised

**Knowledge**

# Learning approaches



Dimensionality reduction: e.g. PCA, tSNE
Clustering: e.g. Phenograph, FlowSOM

Classification: SVMs, Random Forests

**Supervised Learning**: Learning with a labeled training set.
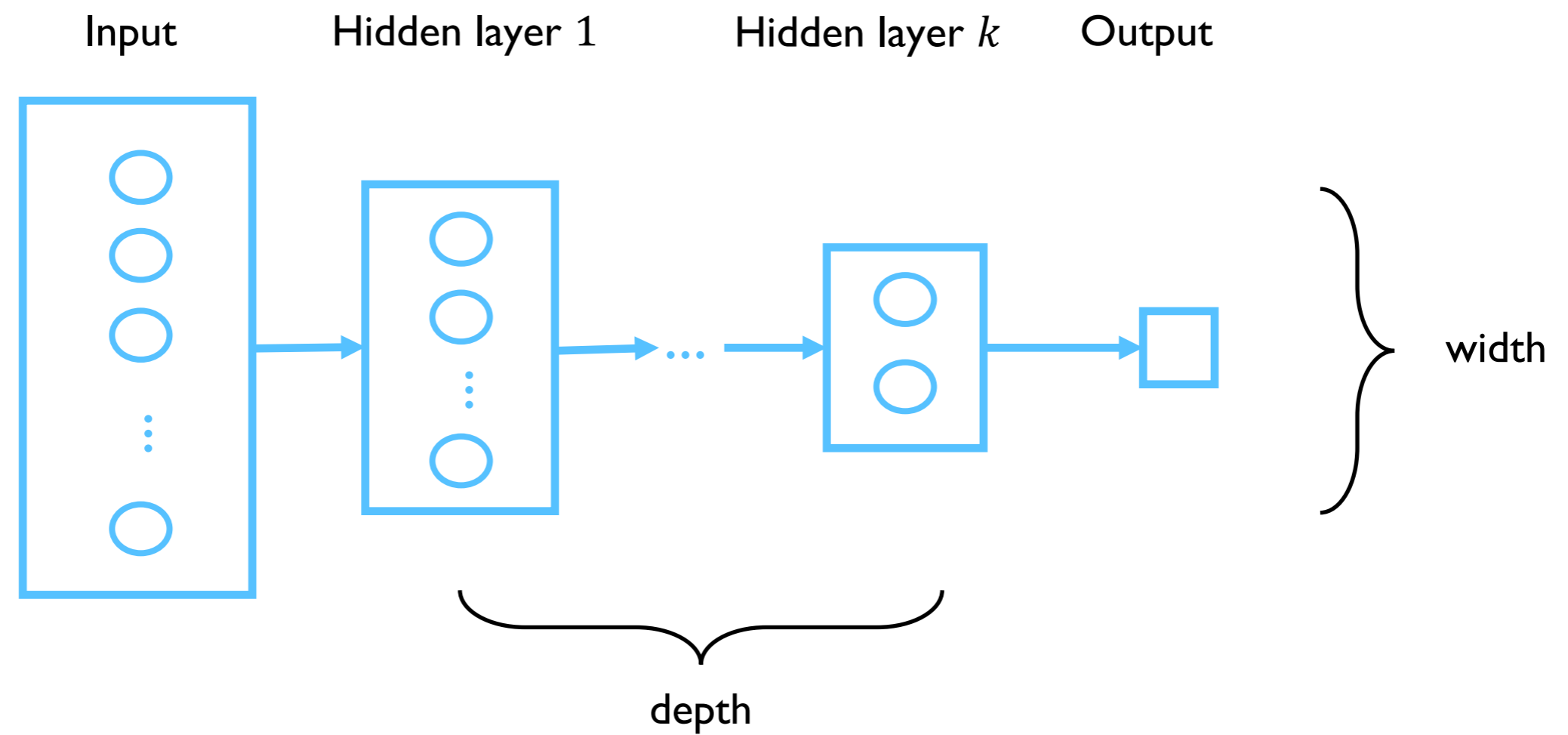 E.g. email spam detector with training set of already labeled emails.
**Unsupervised Learning**: Discovering patterns in unlabeled data.
 E.g. cluster similar documents based on the text content .
**Reinforcement Learning**: learning based on feedback or reward.
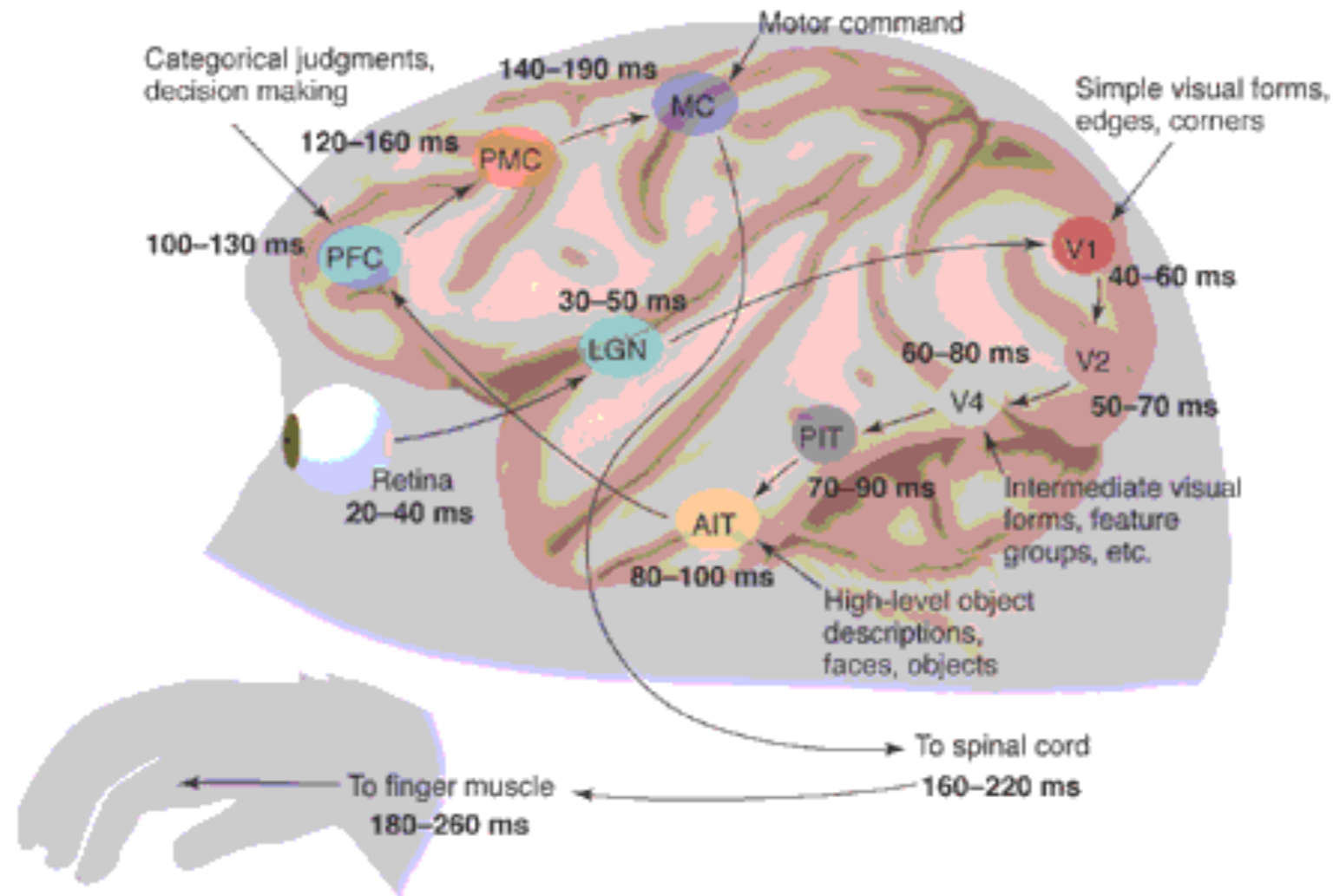 E.g. learn to play chess by winning or losing.

# Neural networks



- Learn data representations. Exceptional effective at learning patterns.
- Use a hierarchy of layers that mimic the neural networks of our brain.
- Can learn highly complex patterns if sufficient data is available for training.

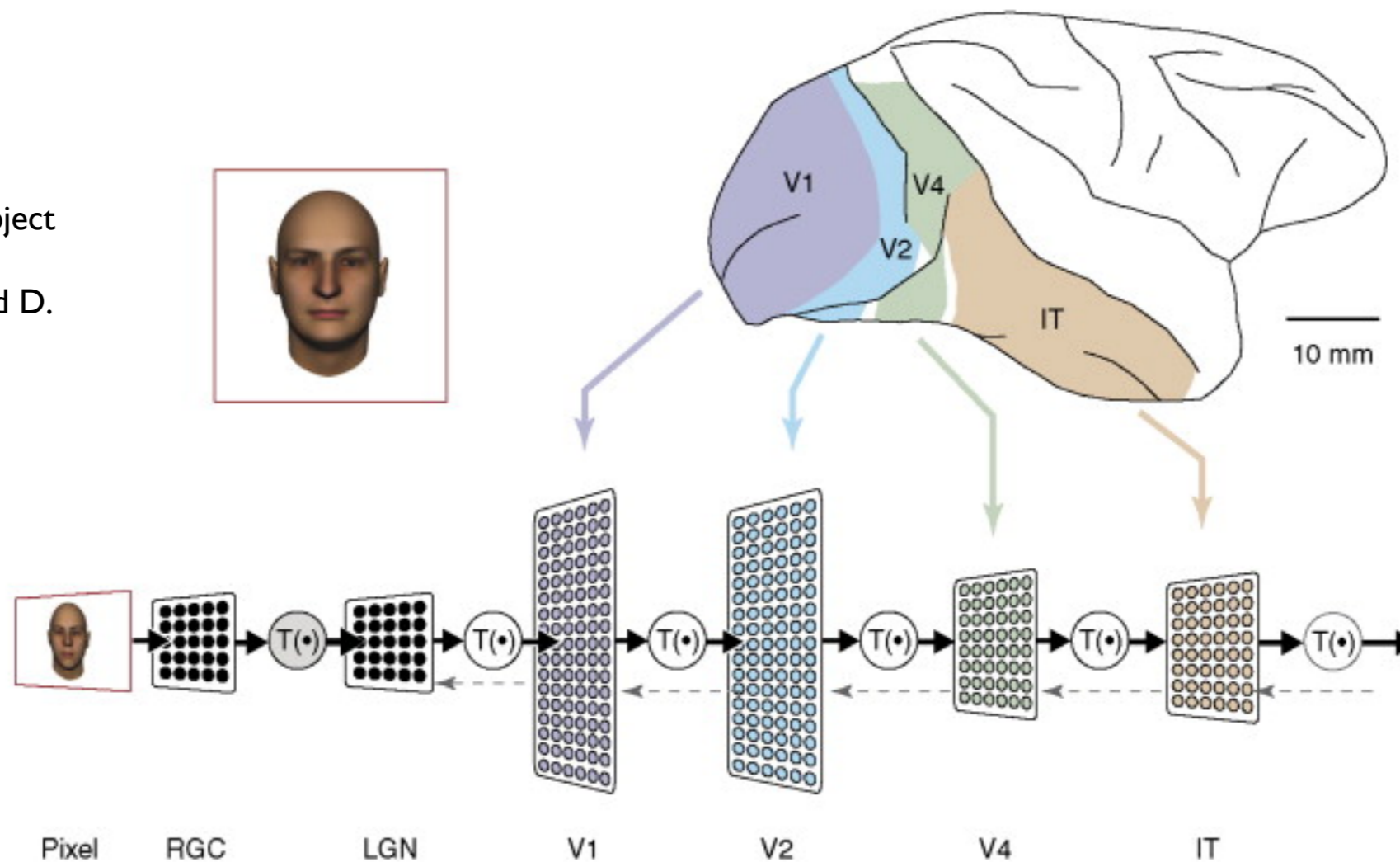# The mammalian visual cortex is hierarchical



Simon J. Thorpe, Michèle Fabre-Thorpe, *Science* 2001

- First hierarchy of neurons are sensitive to edges.
- Brain regions further down the visual pipeline are sensitive to more complex structures (e.g. faces).
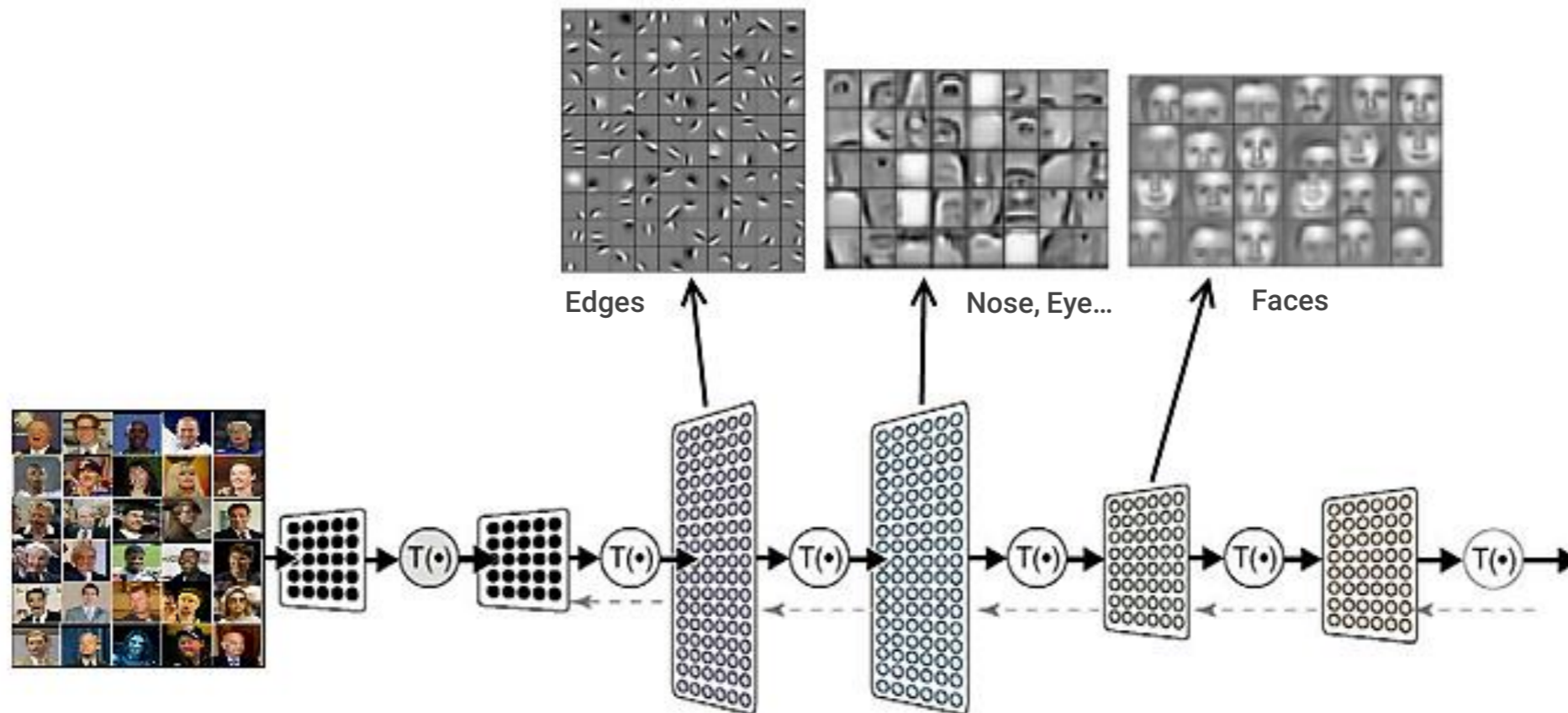- The strength of the connections between neurons represents long term knowledge.

# DNNs mimic the neuronal hierarchical connectivity.

Untangling invariant object recognition
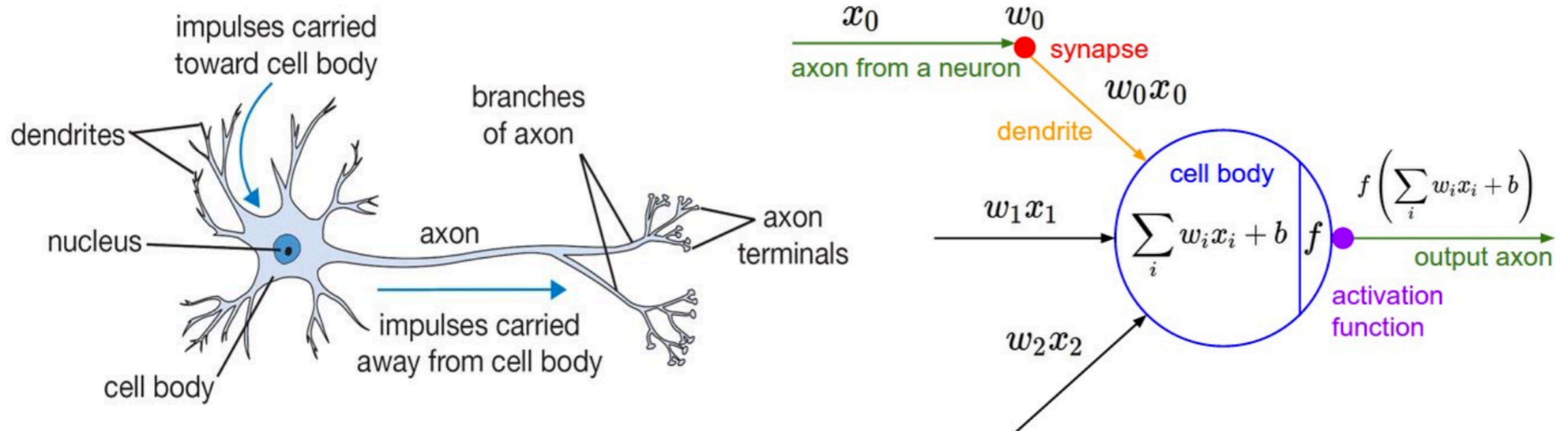James J. DiCarlo, David D. Cox, 2007



- Deep neural networks (DNNs) consists of a hierarchy of layers.
- Each layer transforms the input data into more abstract representations:
  e.g. edge -> nose -> face.
- The output layer combines those features to make predictions.

# DNNs mimic the neuronal hierarchical connectivity.



Edges    Nose, Eye...    Faces

- Deep neural networks (DNNs) consists of a hierarchy of layers.
- Each layer transforms the input data into more abstract representations:
    e.g. edge -> nose -> face.
- The output layer combines those features to make predictions.
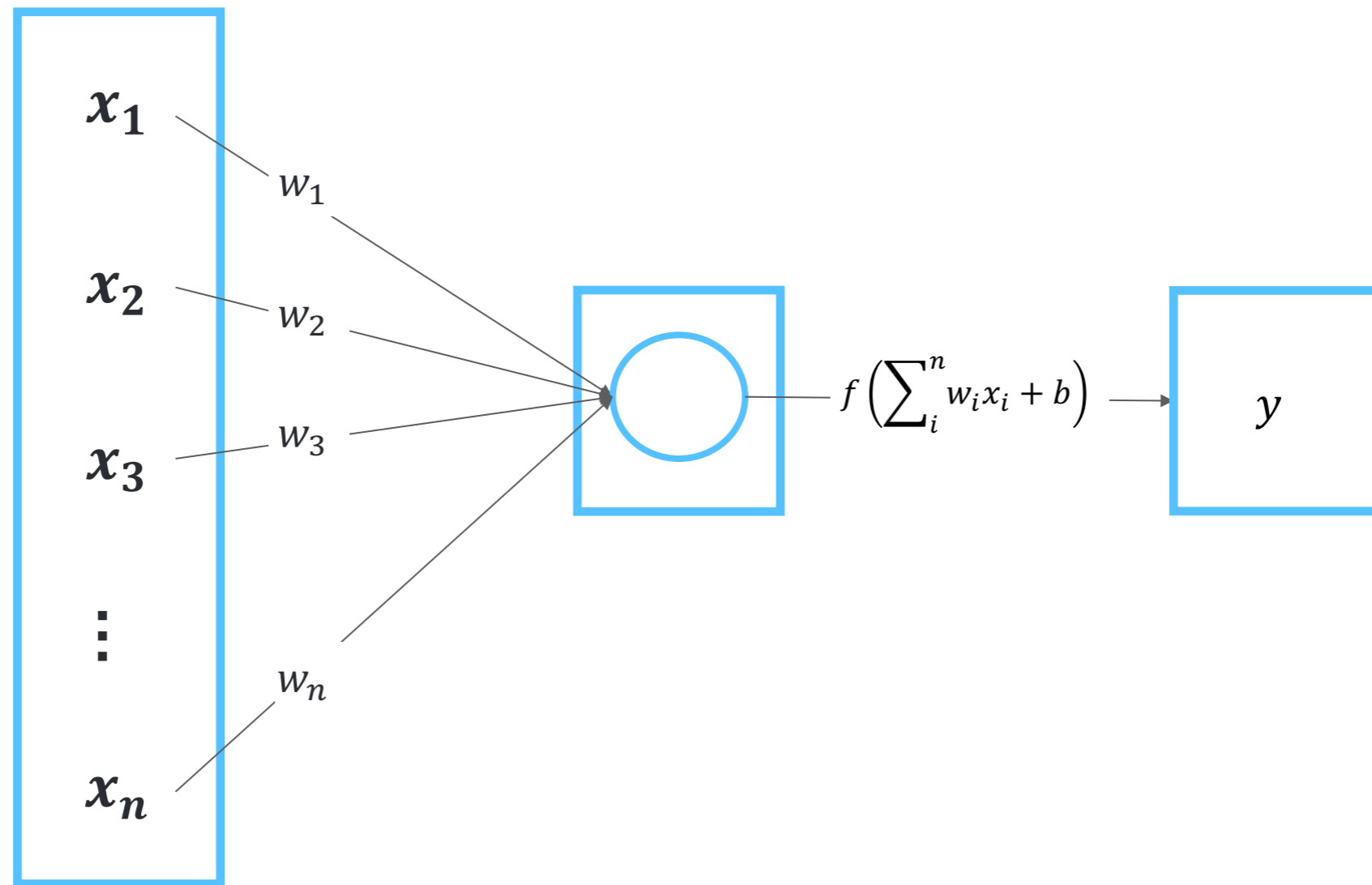
# Biological vs. artificial neurons
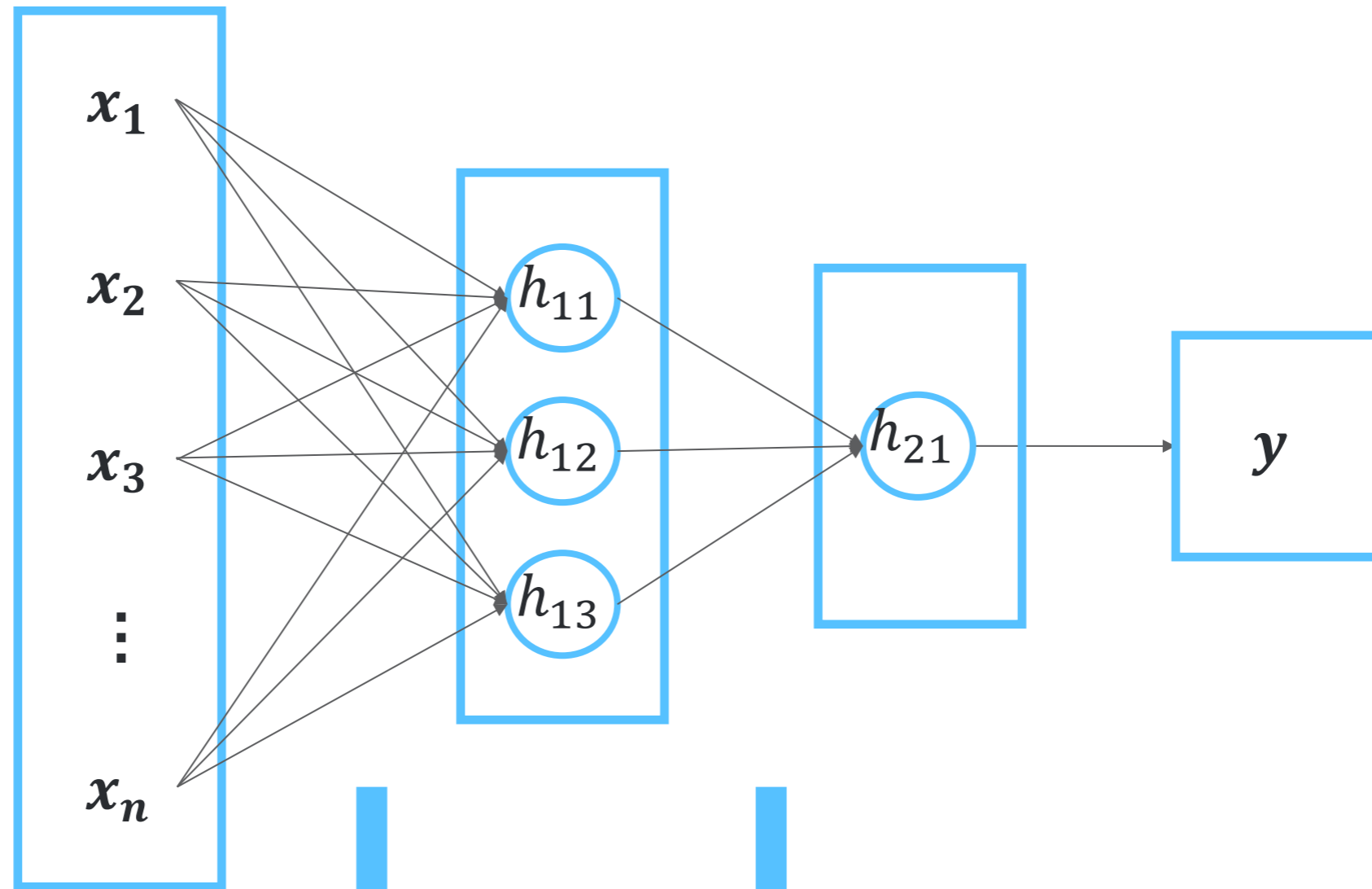


http://cs231n.github.io/neural-networks-1/

- Neurons filter and detect specific features or patterns (e.g. edge, nose) by receiving a weighted input, transforming it with the activation function and passing it to the outgoing connections.
  - Each neuron performs a dot product with the input and its weights, adds the bias and applies the activation function.
- Artificial neurons mimic brain neurons.

# Simplest neural network



$$f\left(\sum_{i}^{n} w_i x_i + b\right)$$

- Weights and biases are the learnable parameters.
- **Weight**: controls the strength of the connection. Weights near zero mean changing this input will not change the output.
- **Bias**: measure of how easy it is to get a node to fire. A node with a large bias will tend produce large positive outputs.
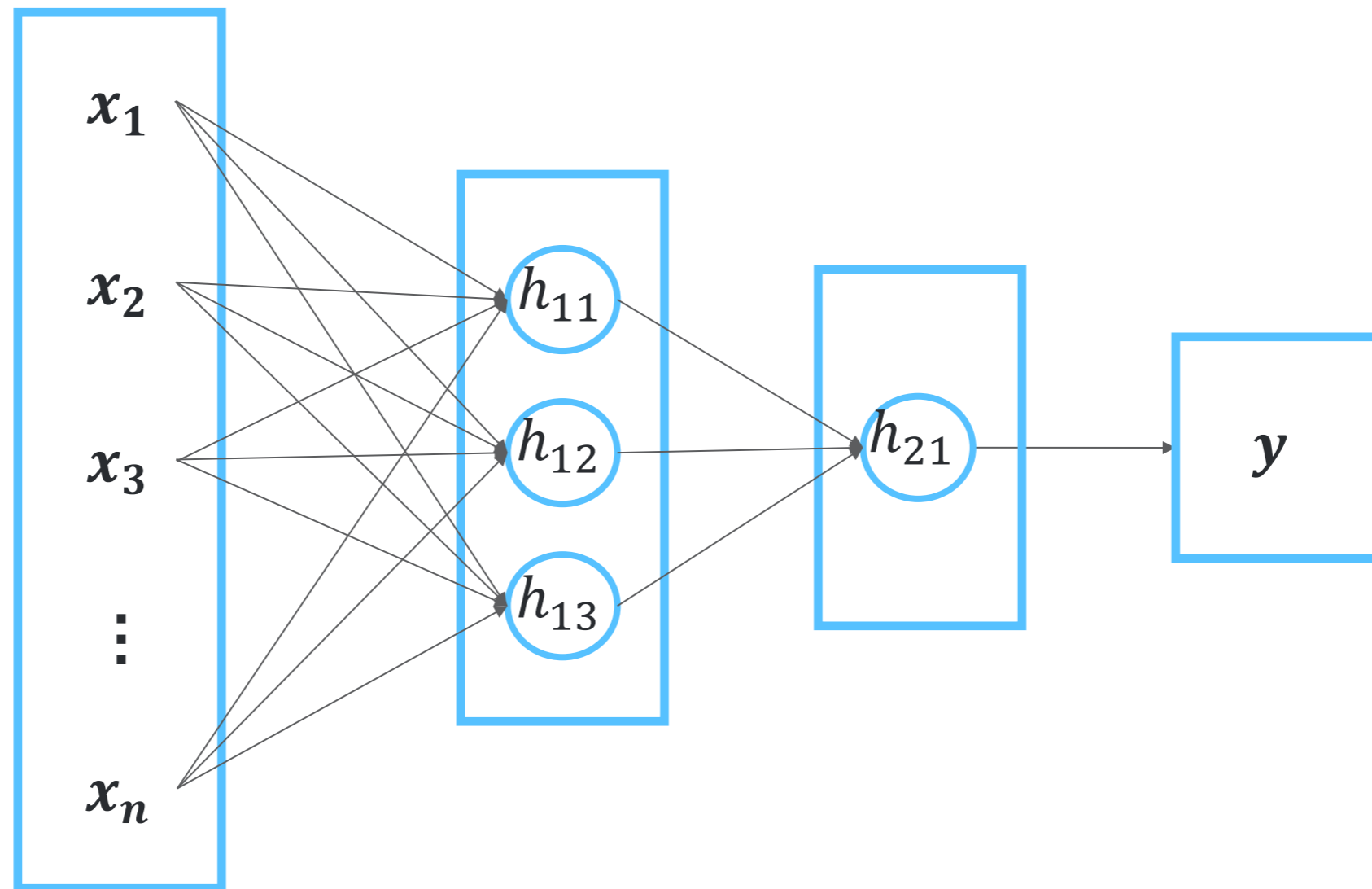
# A more realistic example



$$W^1 = \begin{bmatrix} w^1_{11} & w^1_{12} & w^1_{13} \\ w^1_{21} & w^1_{22} & w^1_{23} \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ w^1_{n1} & w^1_{n2} & w^1_{n3} \end{bmatrix}, b = \begin{bmatri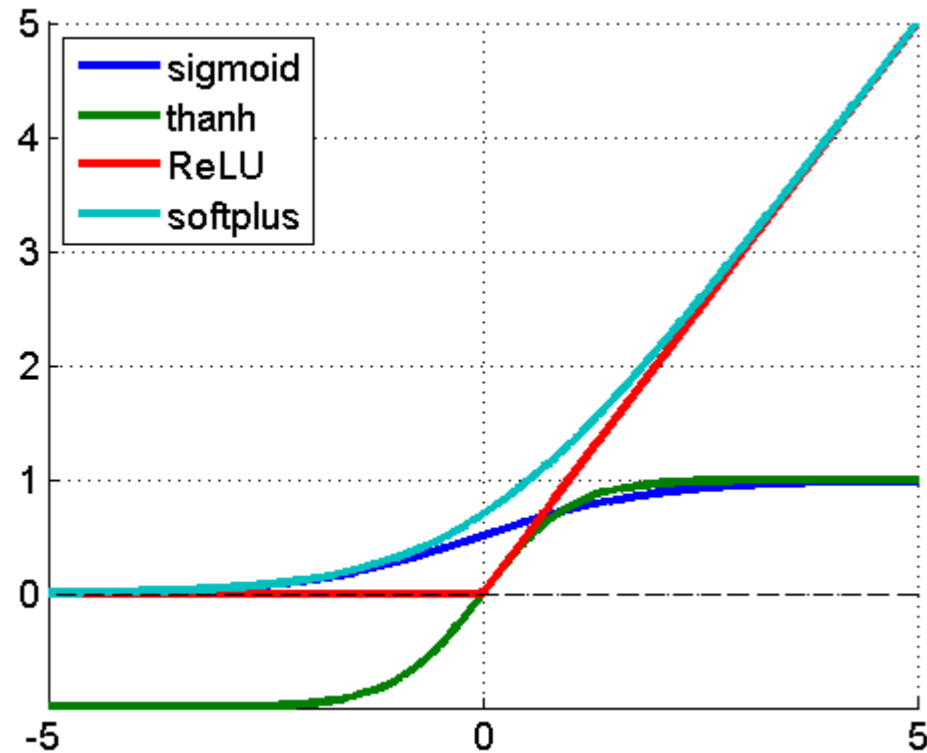x} b^1_1 \\ b^1_2 \\ b^1_3 \end{bmatrix} \qquad W^2 = \begin{bmatrix} w^2_{11} \\ w^2_{12} \\ w^2_{13} \end{bmatrix}, b = [b^2_1]$$

# A more realistic example



- Each DNN consists of one input, one output and multiple fully-connected hidden layers in between.
- Each layer is represented as a series of neurons that progressively extract higher-level features of the input until the final layer makes a decision about what the input shows.
- The more layers the network has, the more abstract features it can learn.

# Commonly used activation functions

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

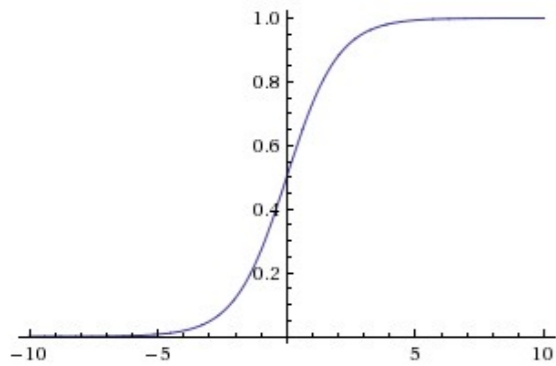$$Tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$ReLU(x) = max(0, x)$$

$$Softplus(x) = log(1 + e^x)$$

- Activations functions are non-linear. Non-linearity is needed to learn complex representations of data, otherwise the DNN would be just a linear function (analogous to PCA).
- Most deep networks use ReLU in hidden layers:
    - it trains much faster (constant derivative),
    - improve discriminative performance,
    - prevents the gradient vanishing problem.
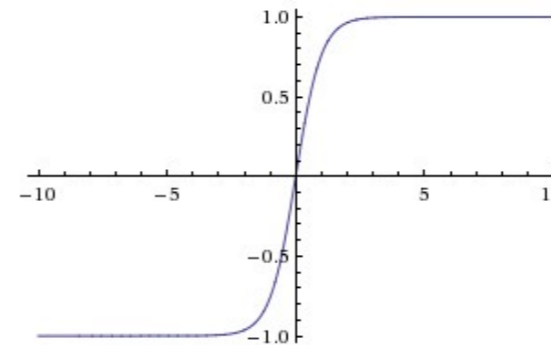
# Activation functions

## Sigmoid

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

+     Resembles neuronal firing
−     Saturation → zero gradients
−     Sigmoid outputs are not zero-centered

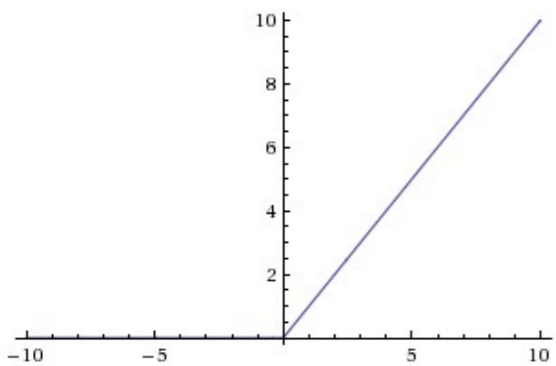## Hyperbolic Tangent (tanh)

$$\tanh(x) = \frac{1 - e^{-2x}}{1 + e^{-2x}} = 2\sigma(2x) - 1$$

+     Outputs are zero-centered
−     Saturation → zero gradients

## Rectified Linear Unit

$$\mathbf{ReLU}(x) = \max(0, x)$$

+     Cheap operation
+     Accelerates convergence
−     Large gradients → Dying ReLUs

## Additional activation functions

- **Leaky ReLUs**: solves the dying ReLUs issue, results are not consistent
- **Maxout:** generalization of ReLUs, no saturation, no dying, very expensive to compute

# Examples of cost functions

- A cost function measures how well a neural network predicts the expected outputs given the training samples.
- A cost function is single valued function:

$$C_i(W, B, S_i, O_i)$$

Weights     Biases    Sample i      Expected output of sample i

- **Cost function requirements:**
  - The cost function C must be able to be written as an average over individual training samples:

$$C(W, B, S, O) = \frac{1}{n} \sum_{i=1}^{n} C_i(W, B, S_i, O_i)$$

  - The cost function C must not depend on any network activation value besides the activation value of the output layer, $a_j^L$.

- $C(W, B, S, O) \approx 0$ means the DNN is well trained.

# Cost functions

- **Mean square error** (aka maximum likelihood and sum squared error):

$$C(W, B, S, O) = \frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{m} \left(a_j^L - O_{ij}\right)^2$$

- **Cross-entropy** (aka Bernoulli negative log-likelihood and binary cross-entropy):

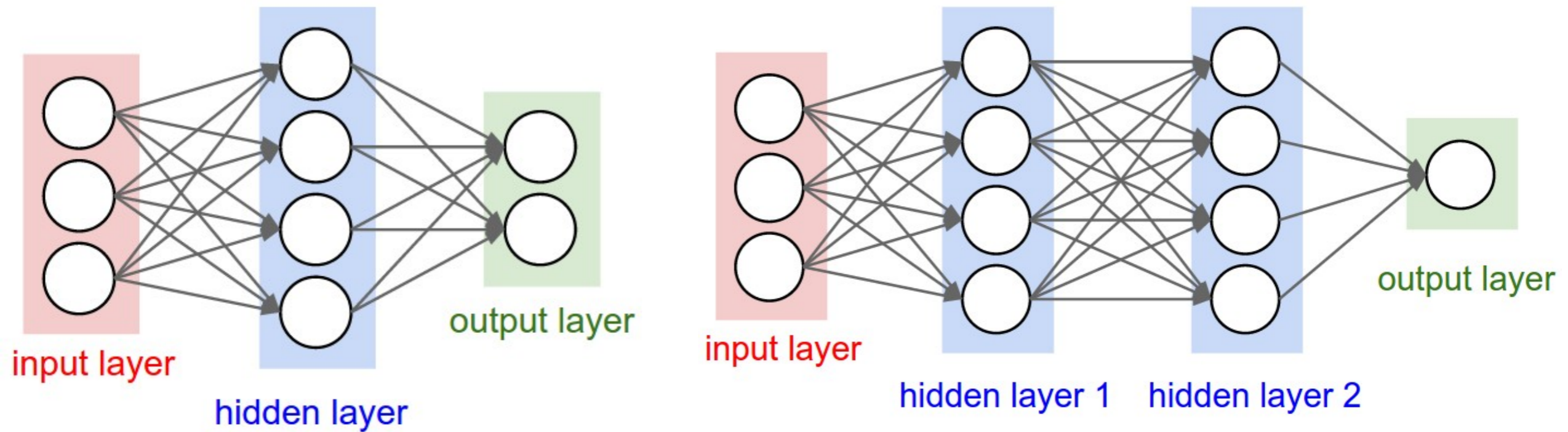$$C(W, B, S, O) = -\frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{m} \left[O_{ij} \log a_j^L + \left(1 - O_{ij}\right) \log\left(1 - a_j^L\right)\right]$$

- **Kullback–Leibler divergence** (aka information divergence, information gain and relative entropy):

$$C(W, B, S, O) = \frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{m} O_{ij} \log \frac{O_{ij}}{a_j^L}$$

- **Hellinguer distance**

$$C(W, B, S, O) = \frac{1}{n\sqrt{2}} \sum_{i=1}^{n} \sum_{j=1}^{m} \left(\sqrt{a_j^L} - \sqrt{O_{ij}}\right)^2$$

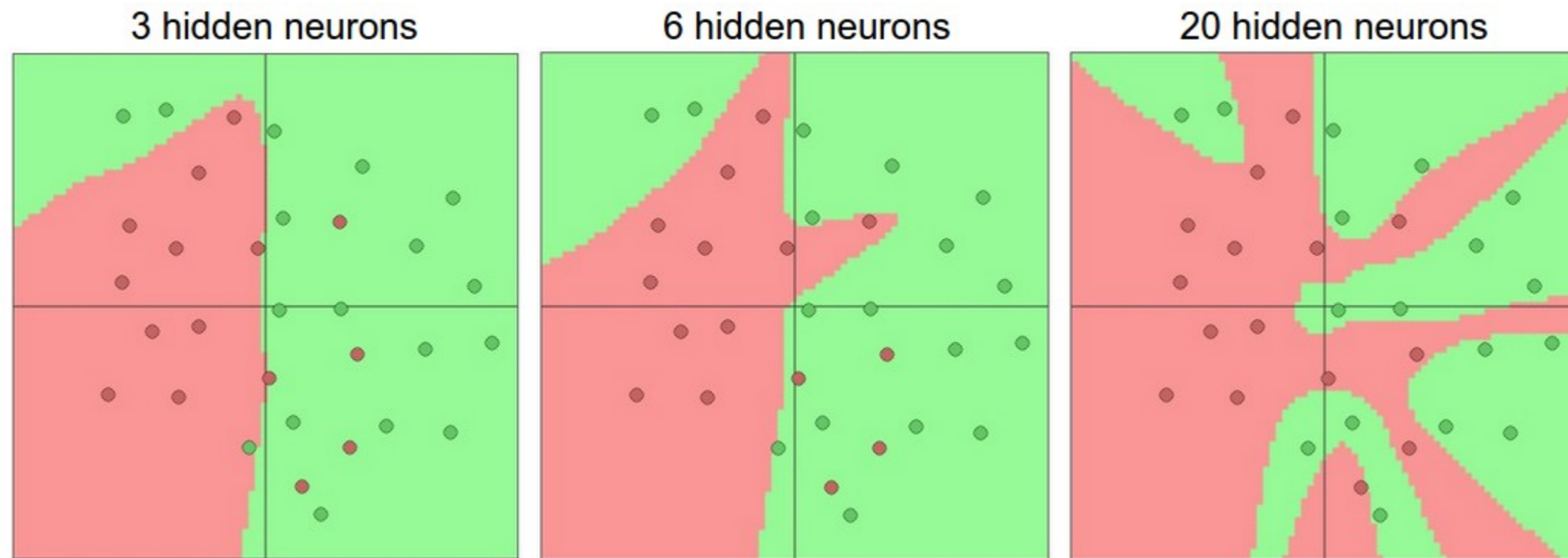# The importance of network architecture



- The capacity of a network can be increased with the number of layers and units per layer.
- As a rule of thumb, going deeper results in more expressive networks, while going wider may lead to overfitting
  - more layers lead to more nested functions and non-linearities that increase the abstraction power, while more units in the same layer usually add features of the same complexity, which might lead to redundancy.

# Overfitting



- **Overfitting** occurs when a model with high capacity fits the noise in the data instead of the (assumed) underlying relationship.

# Overfitting



3 hidden neurons     6 hidden neurons     20 hidden neurons

- Larger DNNS can represent more complicated functions. Should therefore we go always very deep?
    - No, DNNs with more neurons can express more complicated functions, however, large networks trained on scarce data might lead to overfitting.
- When data is scarce, it is essential to implement methods to prevent overfitting (L2 regularization, dropout, input noise, etc).
- In practice, it is always better to use methods to control overfitting instead of reducing the number of neurons.

http://cs231n.github.io/neural-networks-1/#bio

# How to prevent overfitting.

- Early stopping:
  - Stop training as so on as the error on the validation set is higher than it was the last time it was checked.

- Noise addition:
  - Dropout: dropping out units (both hidden and visible) in a neural network.
  - Add noise to data (e.g. denoising autoencoders): we train the network to reconstruct the input from a corrupted version of it.

- Regularization penalties :
  - Create weight penalties L1 and L2.

- Dataset augmentation:
  - Create fake data and add it to the training set.

# Dropout



- At each training iteration a dropout layer randomly removes some nodes in the network with probability p along with all of their incoming and outgoing connections.
- Dropout can be applied to hidden or input layer.
- Why it works:
  - Prevents co-adaptation between neurons.
  - Dropout is an example of ensemble technique, where multiple thinned networks with shared parameters are averaged out.

# Weight regularization

- L2 norm
  - penalizes the square value of the weight (p = 2).
  - tends to drive all the weights to smaller values.

- L1 norm

  - penalizes the absolute value of the weight (p = 1)
  - tends to drive some weights to exactly zero (introducing sparsity in the model), while allowing some weights to be big.

$$C(W, B, S, O)_{regularized} = C(W, B, S, O) + \lambda \sum_{i,j} |w_{ij}|^{p}$$

# The effect of regularization



$\lambda = 0.001$     $\lambda = 0.01$     $\lambda = 0.1$

- **The effects of regularization strength**:
  - Each neural network above has 20 hidden neurons, but increasing the regularization strength makes its final decision regions smoother.

# Backpropagation



$$h_{11} = f(\sum w_1^1 x + b_1)$$

$x_1$

$x_2$

$x_3$

$\vdots$

$x_n$

$h_{11}$

$h_{12}$

$h_{13}$

$h_{21}$

$\widehat{y}$

$X$

Forward propagation

# Backpropagation



$X$  Forward propagation  $L(\hat{y}, y)$

Loss function

# Backpropagation



$\delta_{11} = w_{11}^2 \delta_{21}$

$\delta_{21} = \hat{y} - y$

$x_1$

$x_2$

$x_3$

$\vdots$

$x_n$

$h_{11}$

$h_{12}$

$h_{13}$

$h_{21}$

$\hat{y}$

Loss function

$X$     Backward propagation     $L(\hat{y}, y)$

# Backpropagation

# Weight update

$$w_{11}^1(new) = w_{11}^1 + \eta\, \delta_{11} \frac{\partial}{\partial w_{ij}^k} J(\boldsymbol{w})\, x_1$$

learning rate

$w_{11}^1$

$x_1$

$x_2$

$x_3$

$\vdots$

$x_n$

$h_{11}$

$h_{12}$

$h_{13}$

$h_{21}$

$y$

# The effect of different learning rates



- Low learning rates decrease linearly (slow convergence).
- High learning rates initially decrease exponentially, but saturate at higher values: there is too much "energy" in the optimization and the parameters keep bouncing chaotically, unable to settle in a good minimum.

http://cs231n.github.io/neural-networks-3/

# Optimizers

**Stochastic gradient descent:**
1. Choose an initial vector of parameters
2. Repeat until convergence:
- Randomly shuffle training examples
- Move the weight vector towards the direction of steepest descent by learning rate $\eta$

**Advanced choices:**
- **Momentum:** save the update at each iteration, and determine the next update as a linear combination of the gradient and the previous update
- **Adaptive learning rate methods:** RMSprop, Adagrad, Adam

# Different optimizers achieve very different convergence rates

Images credit: Alec Radford.

# Hyperparameter tuning and reproducibility



DNNs can involve many hyperparameters. The most common include:
- initial learning rate
- momentum
- regularization strength (L2 penalty, dropout strength, etc)

A grid search exploration of all possible parameter combination might not be the most efficient way of tuning the DNN!

# Trained network

# Trained network

# Trained network

# Training and testing

# Overview

# Multi-layer perceptron (MLP)



- An MLP is a DNN that:
  - Consists at least of three layers of nodes (i.e. there is at least one hidden layer).
  - It is always feedforward (no loops are allowed).
  - Consecutive layers are fully connected.

- A single hidden layer is sufficient to make MLPs **a universal approximator**. However usually there are substantial benefits to using more than one hidden layer.

# Deep learning in genomics

# Deep learning frameworks

High-level frameworks make deep learning easier

**Deep Learning Frameworks:**

- Keras
- Lasagne
- Caffe

**Graph compilers:**

- Theano
- Tensor Flow

**Linear Algebra Libraries:**

- PyCuda (python)
- CUDAMat (python)
- JCuda (java)

# Thank you!



**CompSysBio team @ IBM Research**