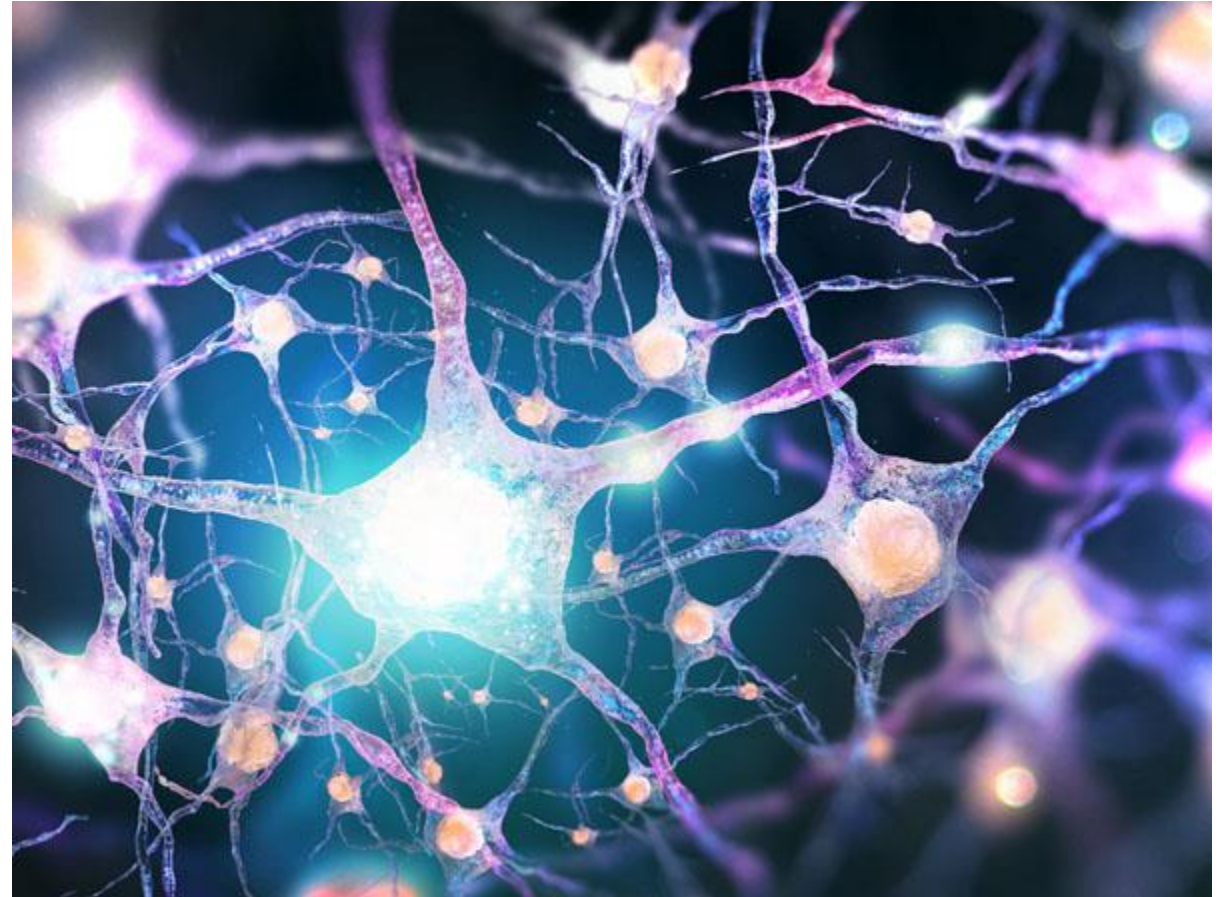# Lecture 3: Introduction to Reinforcement Learning

GEARS Program
Summer 2024
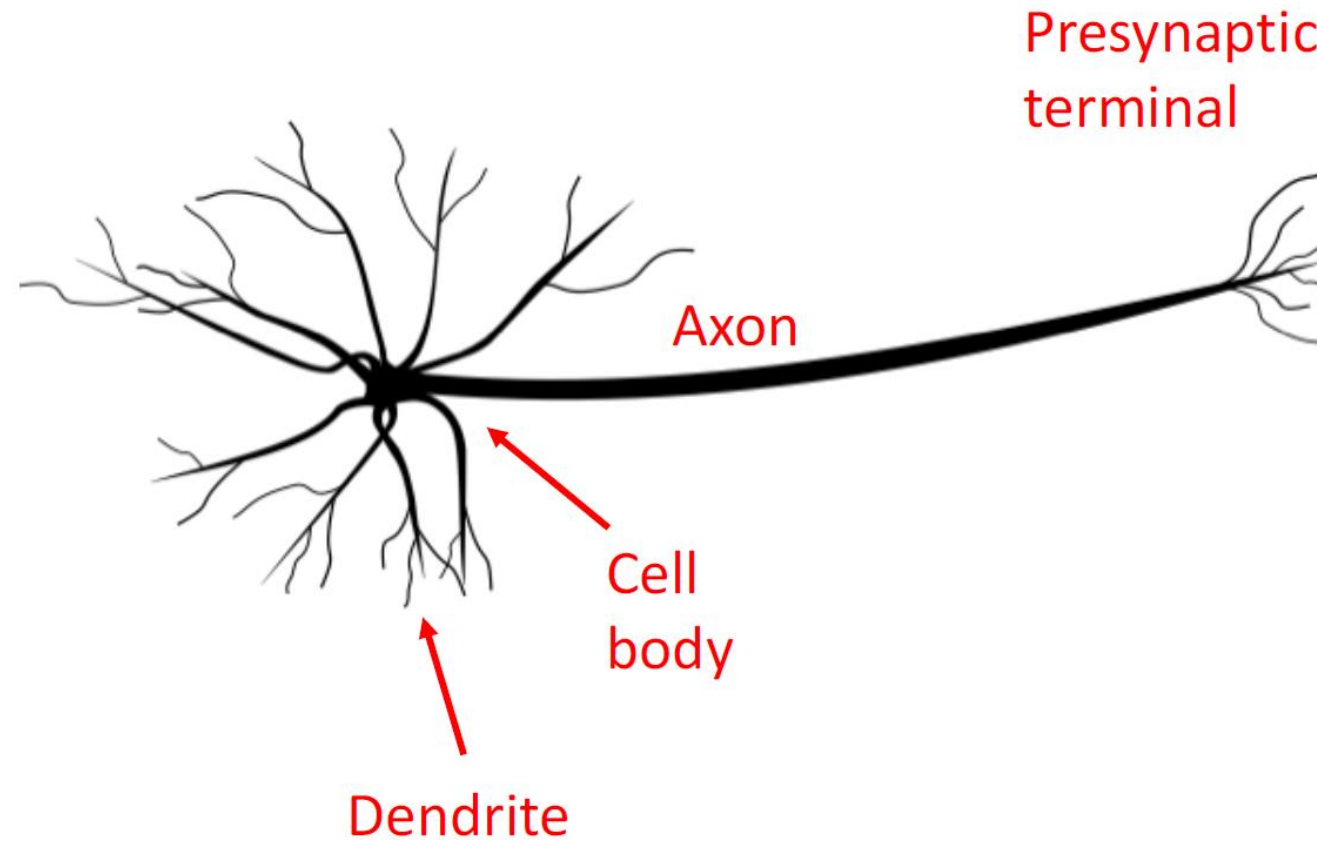
Ivan Lopez-Sanchez, PhD

# The Biological Neural Network
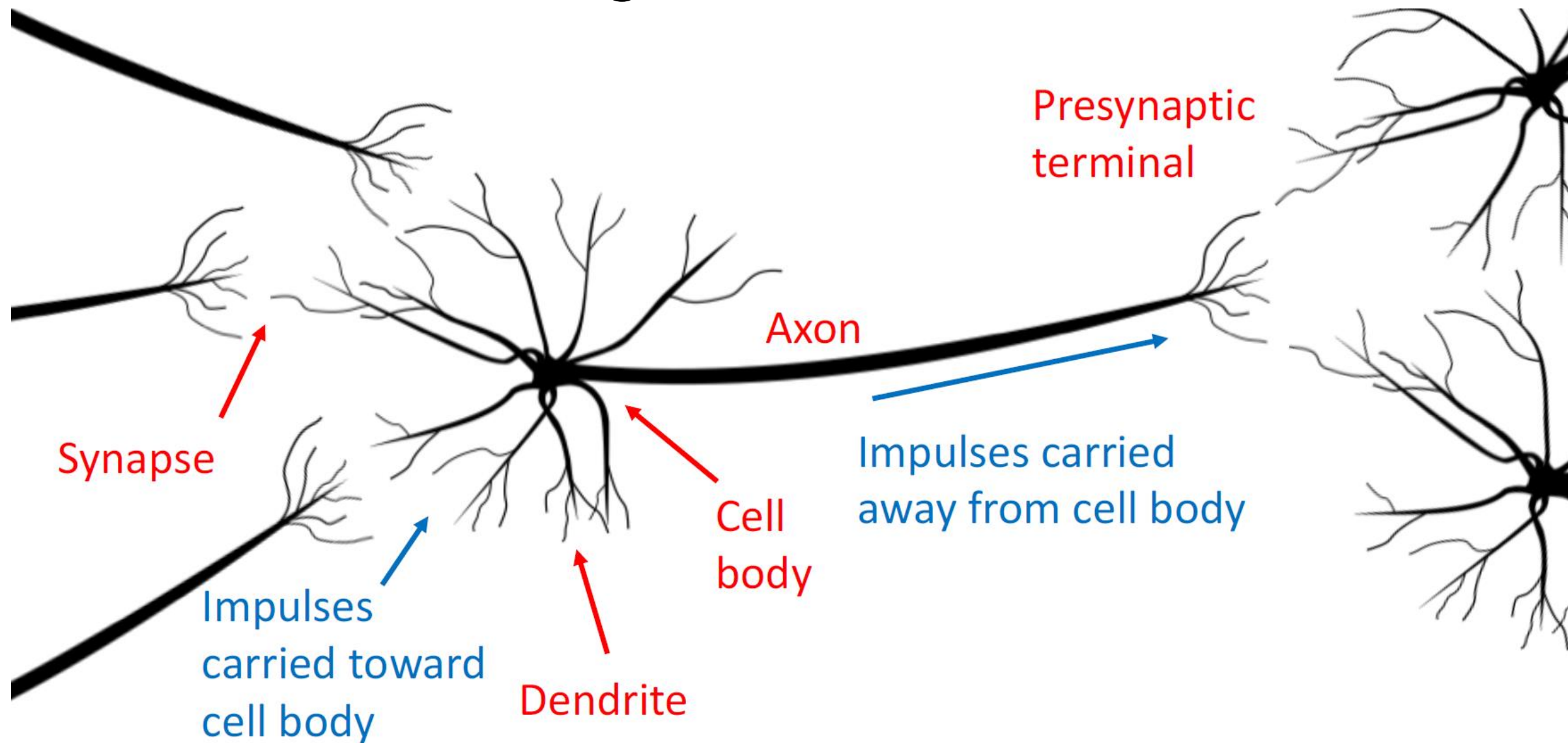
- Our brain is made of huge amounts of interconnected neurons.
- These neurons are our processing units.

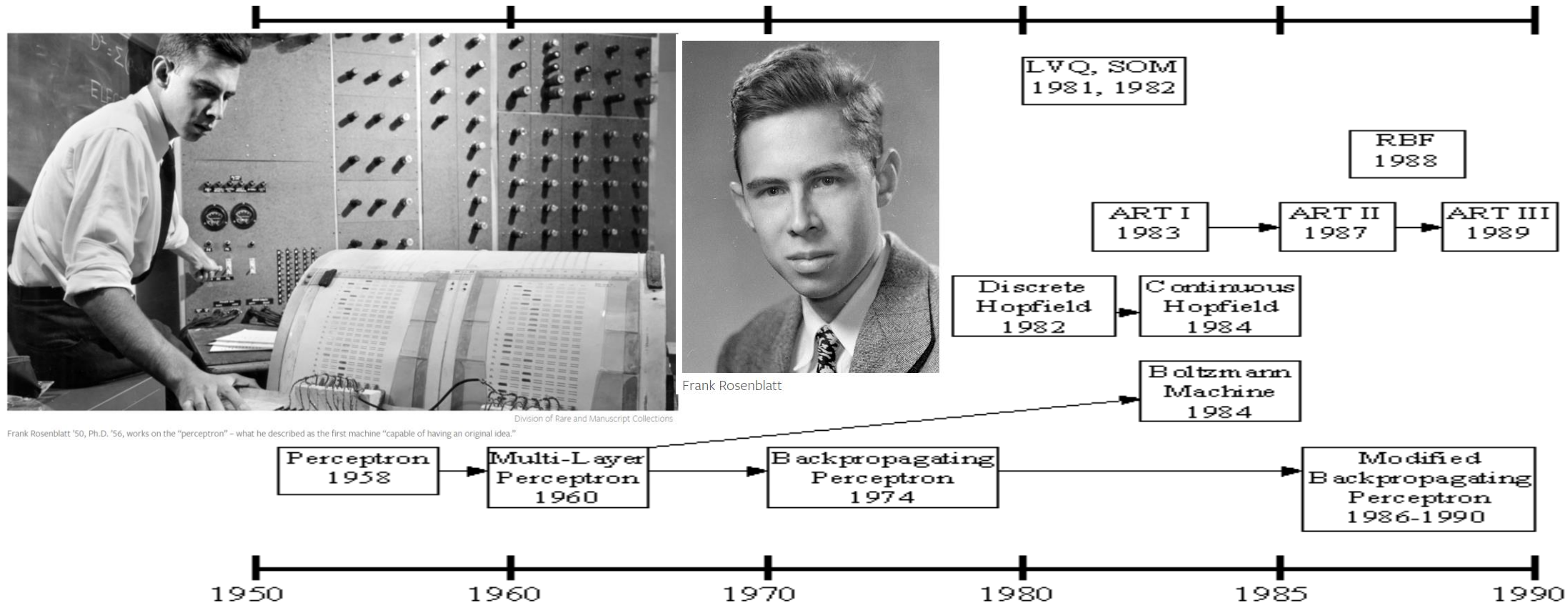# The Biological Neural Network

# The Biological Neural Network



Presynaptic terminal

Axon

Impulses carried away from cell body

Synapse

Cell body

Impulses carried toward cell body

Dendrite

Biomechatronics and Intelligent Robotics (BIRO) Lab

# The first Artificial Neural Networks (ANNs)

- History of the ANNs stems from the **1940s**, the decade of the first electronic computer.
- However, the first important step took place in **1958** when **Rosenblatt introduced the first concrete neural model**, the perceptron. Rosenblatt also took part in constructing the first successful neurocomputer, the Mark I Perceptron.
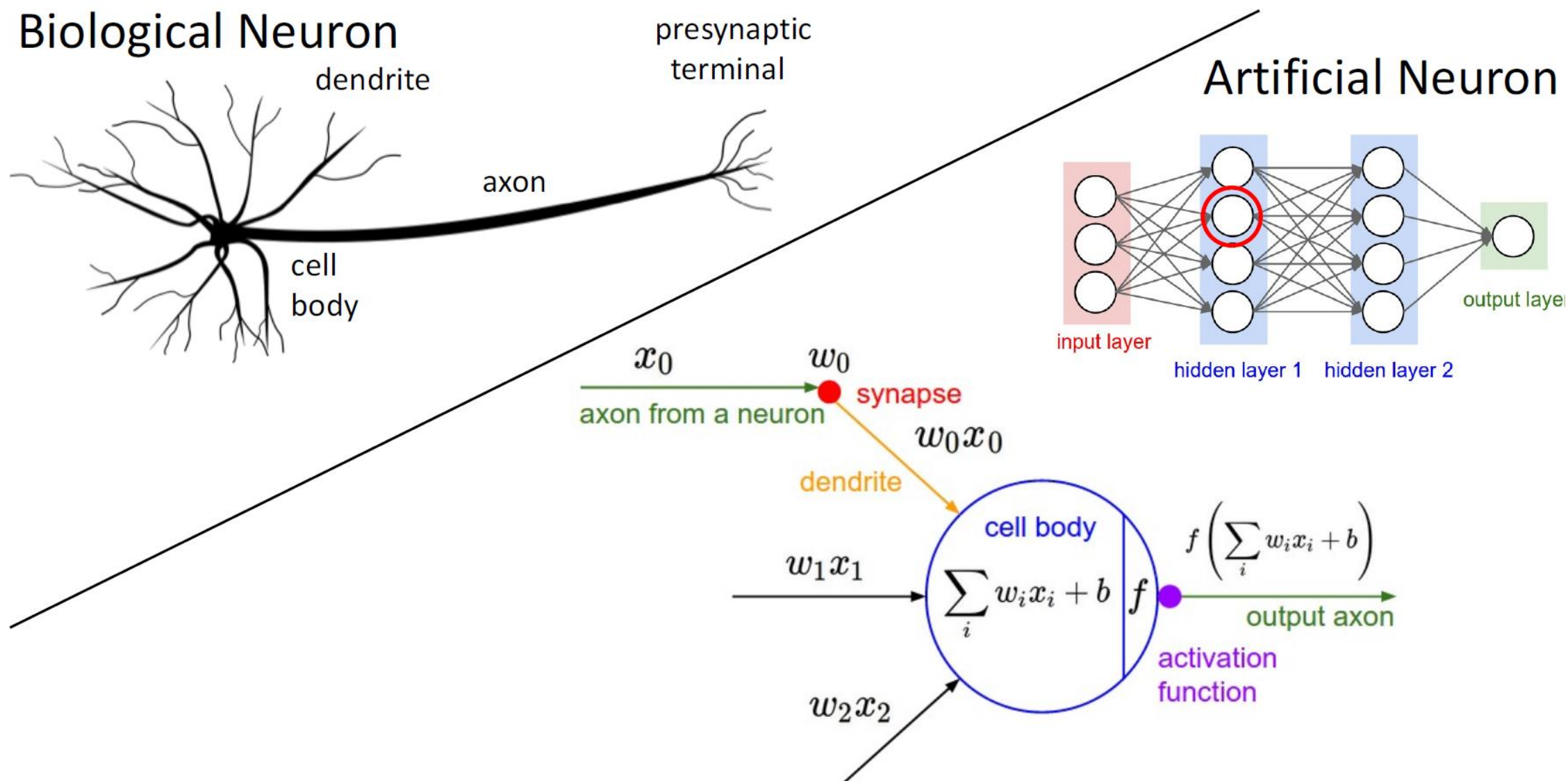


Division of Rare and Manuscript Collections

Frank Rosenblatt '50, Ph.D. '56, works on the "perceptron" – what he described as the first machine "capable of having an original idea."

Frank Rosenblatt

LVQ, SOM
1981, 1982

RBF
1988

ART I
1983 → ART II
1987 → ART III
1989

Discrete Hopfield
1982 → Continuous Hopfield
1984

Boltzmann Machine
1984

Perceptron
1958 → Multi-Layer Perceptron
1960 → Backpropagating Perceptron
1974 → Modified Backpropagating Perceptron
1986-1990

1950        1960        1970        1980        1985        1990
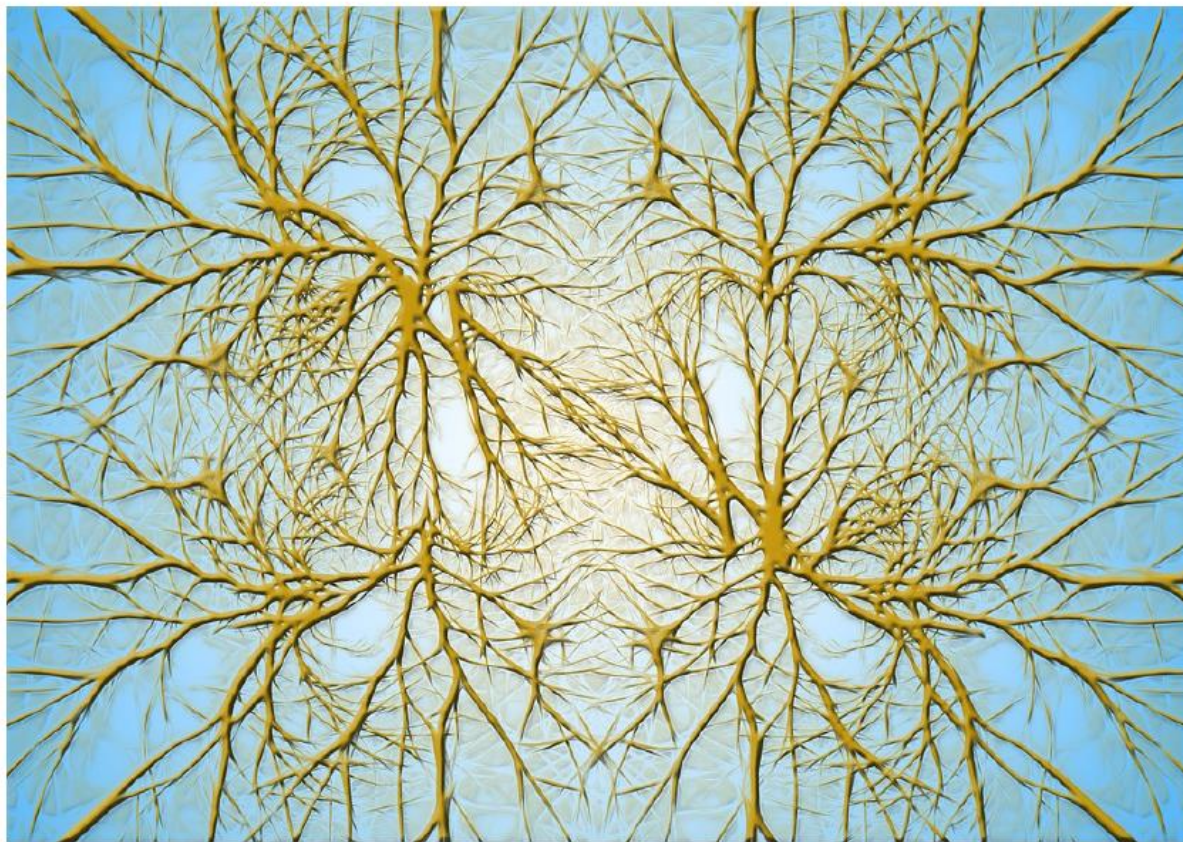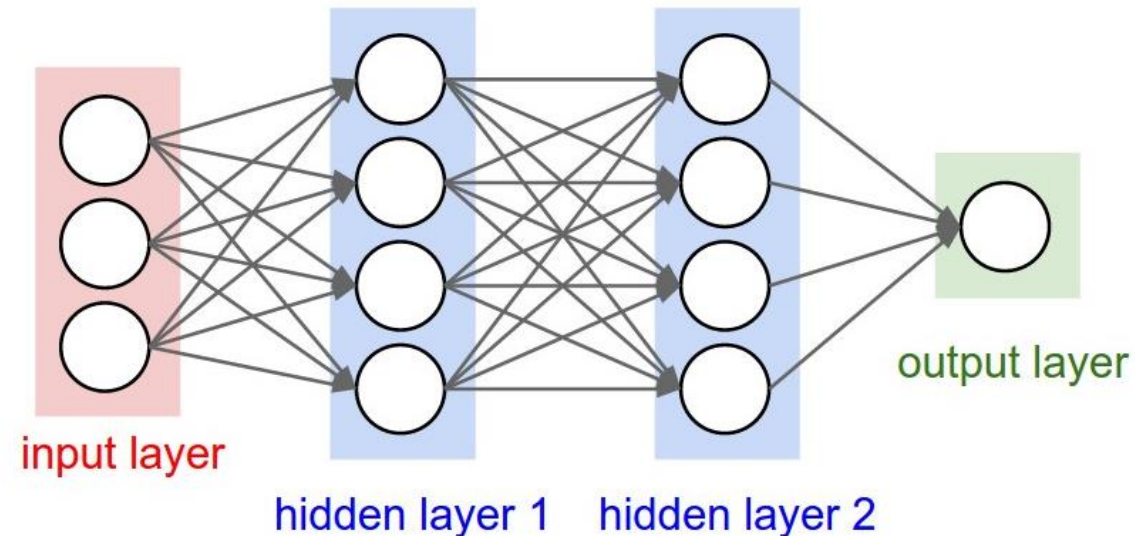
# The relation between…

# The relation between…

**Biological Neurons:**
Complex connectivity patterns



**Neurons in a neural network:**
Organized into regular layers for computational efficiency



input layer

hidden layer 1    hidden layer 2

output layer

NC STATE UNIVERSITY

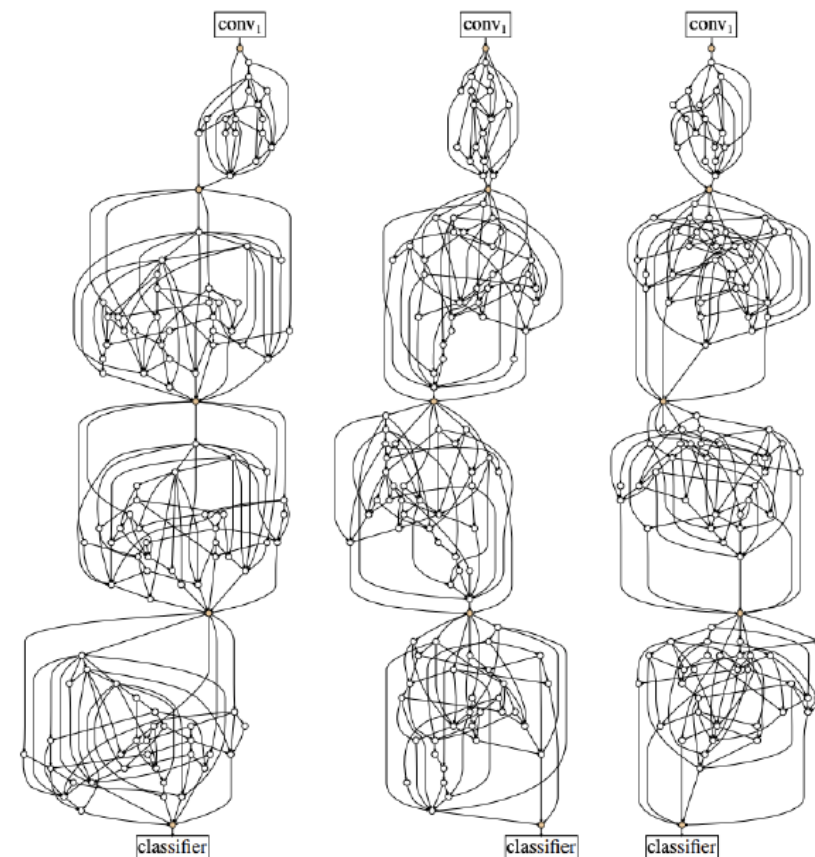Biomechatronics and Intelligent Robotics (BIRO) Lab

# The relation between…

## Biological Neurons: Complex connectivity patterns



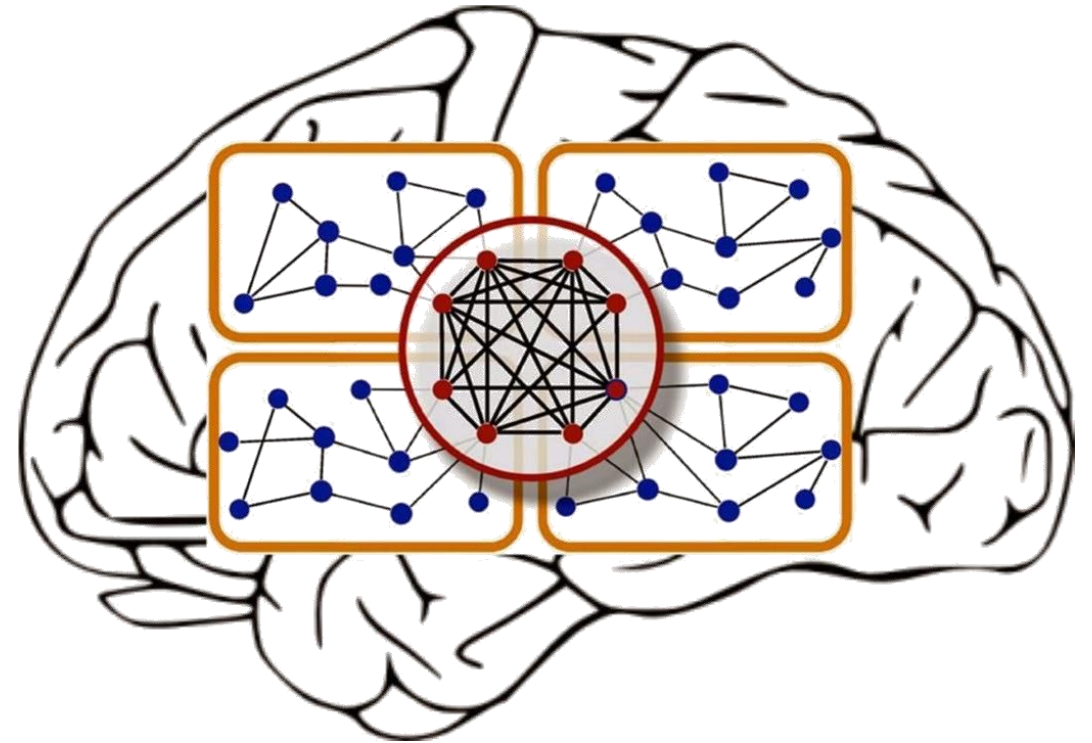## But neural networks with random connections can work too!



Xie, Saining, et al. "Exploring randomly wired neural networks for image recognition." *Proceedings of the IEEE/CVF International Conference on Computer Vision.* 2019.

**NC STATE UNIVERSITY**

**BIRO** Biomechatronics and Intelligent Robotics (BIRO) Lab

# The Artificial Neural Network

NC STATE UNIVERSITY

Biomechatronics and Intelligent Robotics (BIRO) Lab

# Why ANN?

There are two basic reasons why we are interested in building artificial neural networks (ANNs):

- **Technical viewpoint:** Some problems such as character recognition or the prediction of future states of a system require massively parallel and adaptive processing.

- **Biological viewpoint:** ANNs can be used to replicate and simulate components of the human (or animal) brain, thereby giving us insight into natural information processing.



**NC STATE UNIVERSITY**

**BIRO**

**Biomechatronics and Intelligent Robotics (BIRO) Lab**
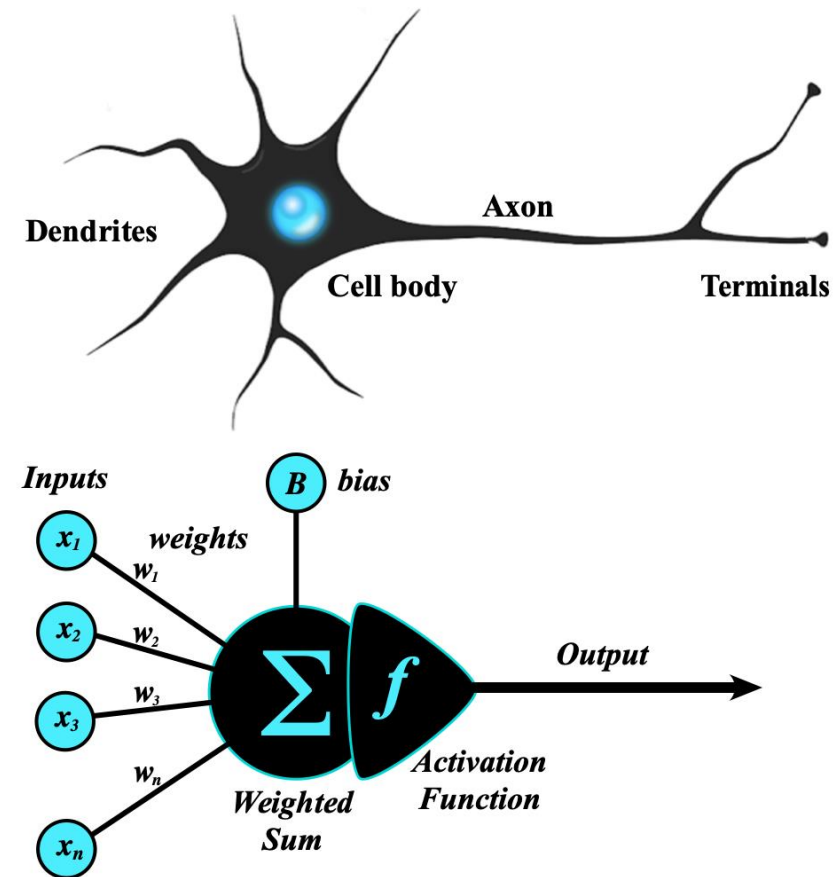
# The Artificial Neuron: Perceptron

The "building blocks" of neural networks are the **neurons**.

- Artificial Neurons are also called units or nodes. However, in the ANNs field are better known as **perceptron**.

Basically, each perceptron:

- Receives input from many other neurons.

- Changes its internal state (activation) based on the current input.

- Sends one output signal to many other neurons, possibly including its input neurons (recurrent network).

$$y = \left( \sum_{i=1}^{n} w_i x_i + b \right) = \varphi(\mathbf{w}^T \mathbf{x} + b)$$

# The single layer ANN

- Is the simplest NN and it is know as Functional Link Neural Network.
- It consist of a single layer of perceptrons.

$$W_\eta^\top \boldsymbol{\sigma}_\eta \left( V_\eta^\top \boldsymbol{\gamma}_\eta \right)$$
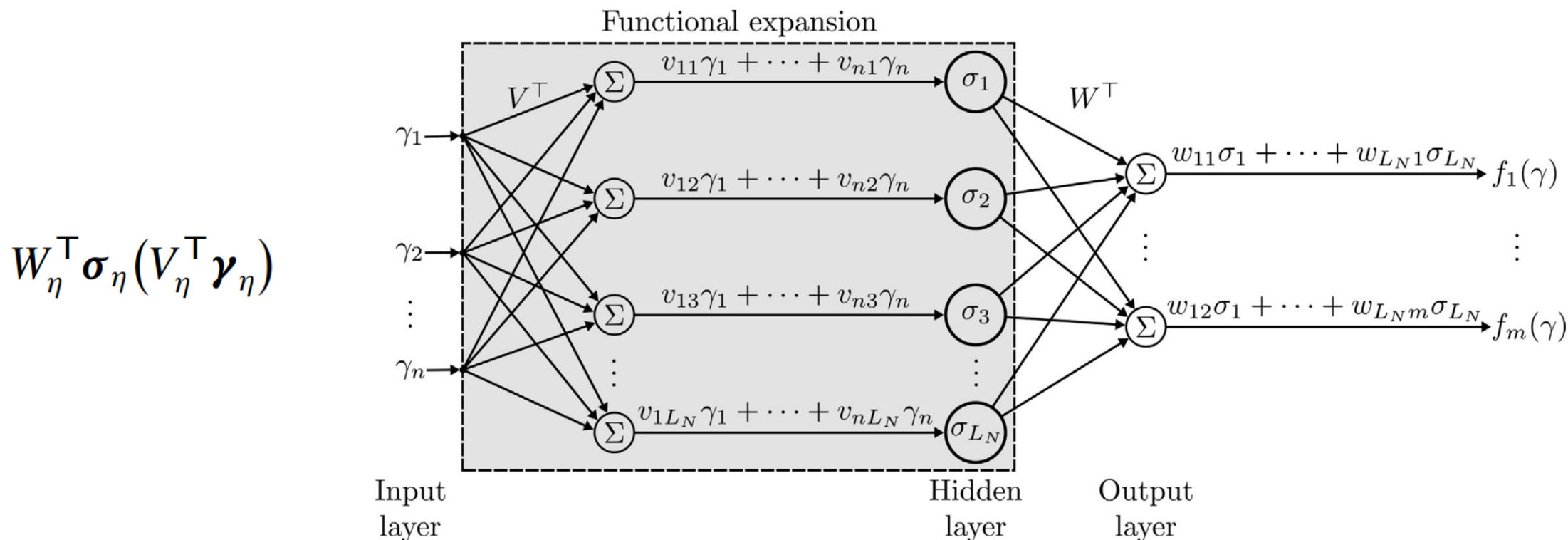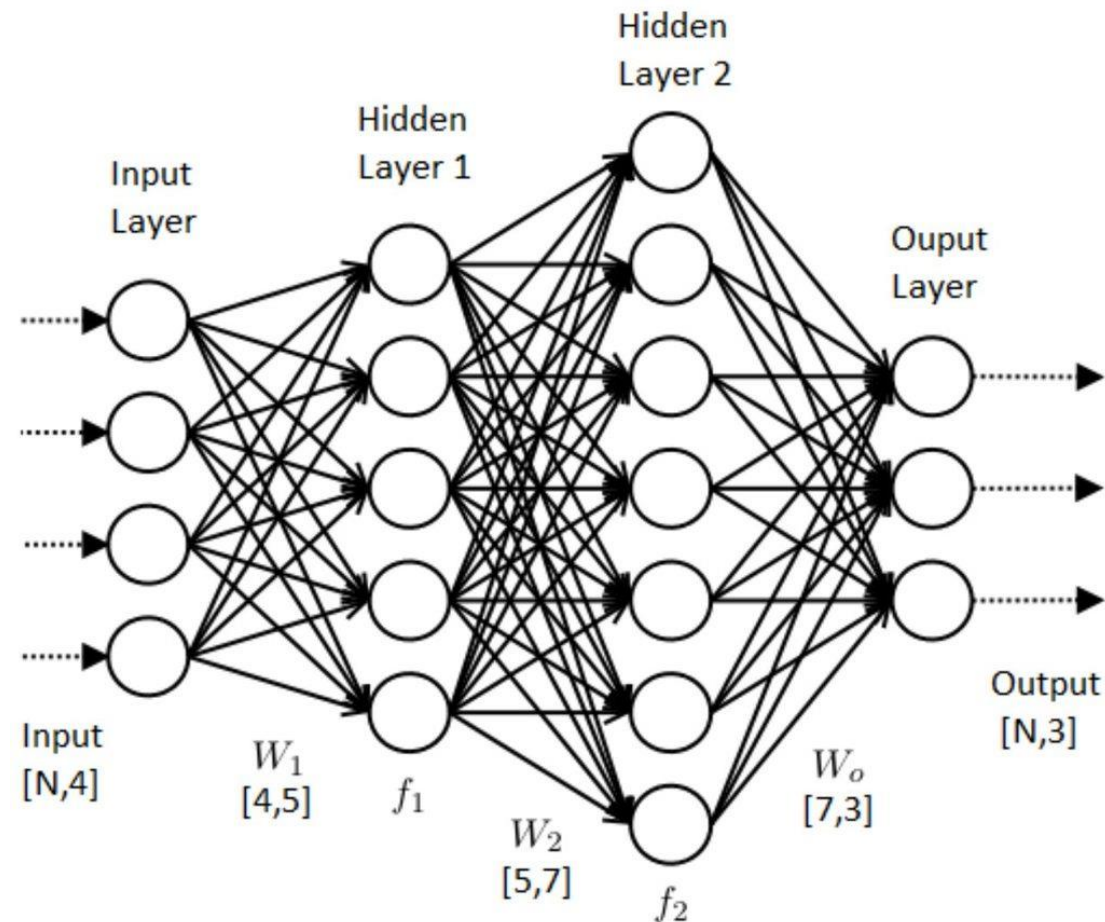


Fig. 3. Diagram of a FLNN where $\boldsymbol{\gamma} \in \mathbb{R}^n$ is the NN's extended input vector, $V \in \mathbb{R}^{n \times L_N}$ and $W \in \mathbb{R}^{L_N \times m}$ are the input and output weights matrices, respectively, $\boldsymbol{\sigma} \in \mathbb{R}^{L_N}$ is the vector of activation functions, $\boldsymbol{f}(\boldsymbol{\gamma}) \in \mathbb{R}^m$ NN's output vector, and $L_N \in \mathbb{R}$ is the number of neurons in the hidden layer.

Lopez-Sanchez, Ivan, Ricardo Pérez-Alcocer, and Javier Moreno-Valenzuela. "Trajectory tracking double two-loop adaptive neural network control for a Quadrotor." Journal of the Franklin Institute 360.5 (2023): 3770-3799.

**NC STATE UNIVERSITY**

**Biomechatronics and Intelligent Robotics (BIRO) Lab**
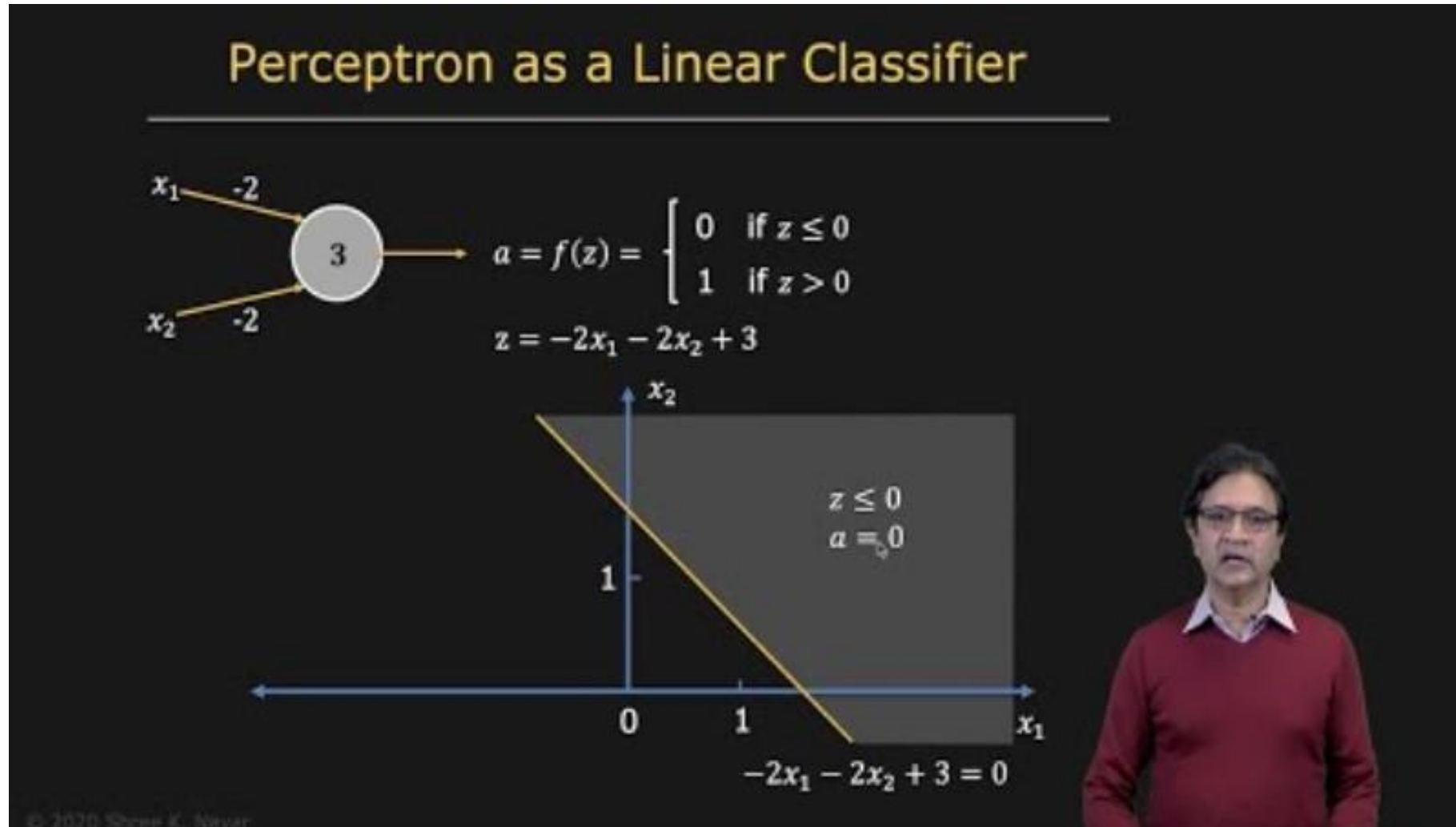
# Then, an ANN

- Consists of a pool of simple processing units which communicate by sending signals to each other over a large number of weighted connections.
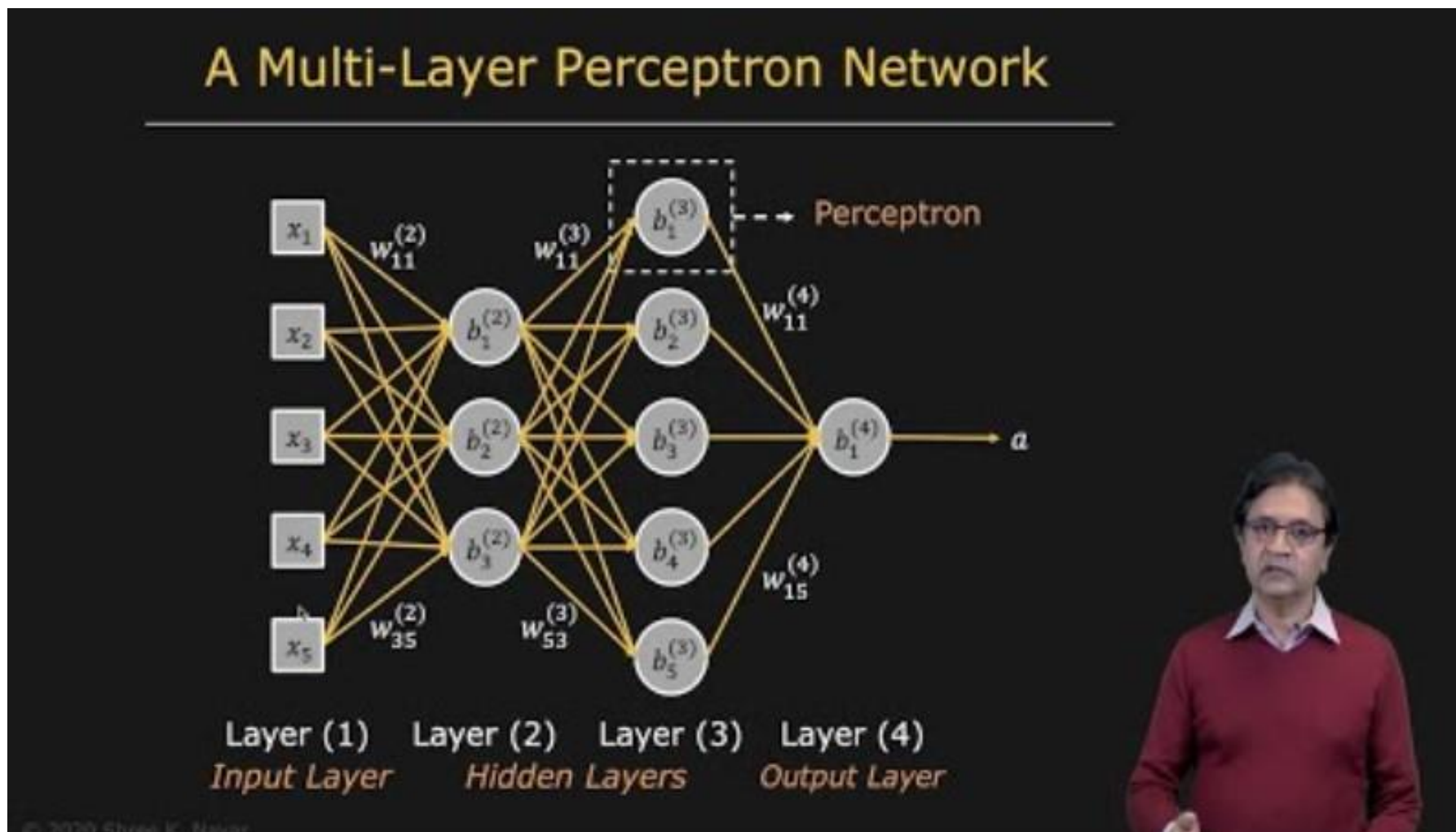
# The Artificial Neuron: Perceptron

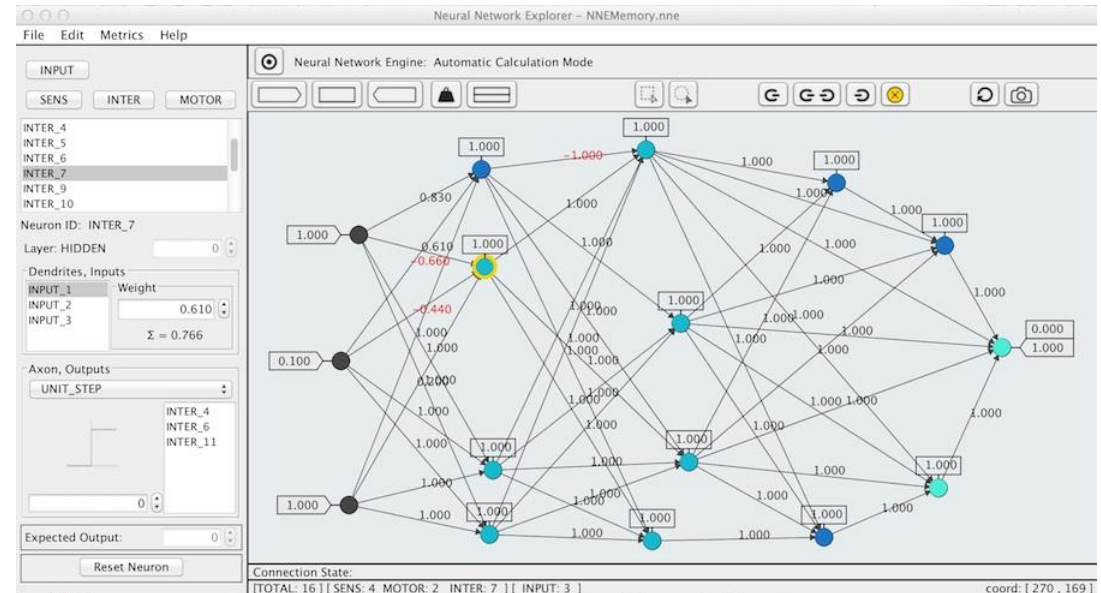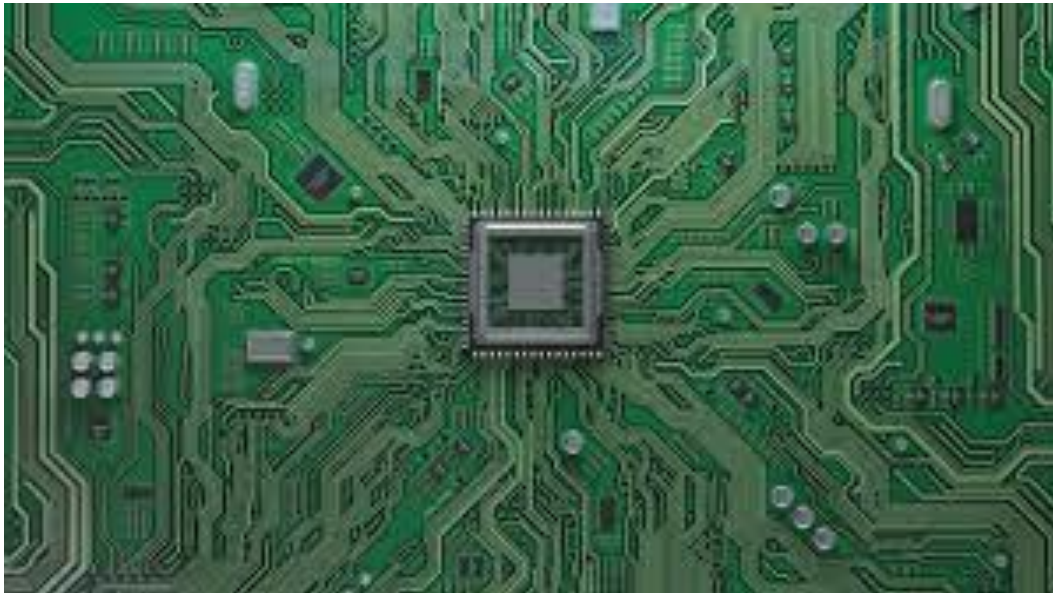The "building blocks" of neural networks are the **neurons.** https://www.youtube.com/watch?v=OFbnpY_k7js

# Multilayer Perceptron (MLP) NN

The base of artificial neural networks are the MLPs. https://www.youtube.com/watch?v=torNuKNLwBE

# In general words

- An ANN is either a hardware implementation or a computer program which strives to simulate the information processing capabilities of its biological exemplar.

- ANNs are typically composed of a great number of interconnected artificial neurons.

- The artificial neurons are simplified models of their biological counterparts.

- ANN is a technique for solving problems by constructing software that works like our brains.

Biomechatronics and
Intelligent Robotics (BIRO) Lab

# Reinforcement learning



ARTIFICIAL INTELLIGENCE

MACHINE LEARNING

>SUPERVISED LEARNING
>UNSUPERVISED LEARNING
>REINFORCEMENT  LEARNING

DATA SCIENCE

DEEPL LEARNING

>ANN
>CNN
>RNN

**NC STATE**
UNIVERSITY

BIRO

**Biomechatronics and**
**Intelligent Robotics (BIRO) Lab**

# General interaction in RL



$$V(s) = max_a \left( \gamma V(s') \right)$$

$$a = \{ \uparrow, \rightarrow, \downarrow, \leftarrow \}$$

# Agent

- Entity that learns to make decisions by interacting with an environment.
  It is form by:
  - Policy
  - Value functions
  - Model of the environment (optional)

  Its goal is to optimize its behavior and <u>maximize cumulative rewards</u> over time.

  It makes decisions based on its <u>observations of the environment</u> and receives <u>feedback</u> in the form of <u>rewards</u>.

  The learning process involves exploring the environment, exploiting known information, and <u>updating the policy based on feedback</u> received from the environment.

**NC STATE UNIVERSITY**

**Biomechatronics and Intelligent Robotics (BIRO) Lab**

# Policy

- Policy ($\pi$)
  - The policy is a strategy used by the agent to decide which action to take in a given state. It can be deterministic (always selecting the same action for a state) or stochastic (selecting actions according to a probability distribution).

- The policy can be represented as a mapping from states to actions, $\pi:S{\to}A$, or as a probability distribution over actions, $\pi(a|s)$.

# Value function

- The value function estimates the expected cumulative reward that can be obtained from a state (or state-action pair).

- State Value Function ($V$): The expected cumulative reward starting from state $s$ and following the policy $\pi$.

$$V^{\pi}(s) = E_{\pi}\left[\sum_{t=0}^{\infty} \gamma^t R_{t+1} | S_t = s\right]$$

- Action Value Function (Q): The expected cumulative reward starting from state $s$, taking action $a$, and following the policy $\pi$.

$$Q^{\pi}(s,a) = E_{\pi}\left[\sum_{t=0}^{\infty} \gamma^t R_{t+1} | S_t = s, A_t = a\right]$$

# Value functions

```
Inputs
  │
  ▼
Value        Determines        Policy      That         Total
Function  ──────────────▶              maximizes    reward
  │                              ──────────────▶
  ▼
Outputs
```

State value function

state ──▶ V(s) ──▶ number

State-action value function

state
action ──▶ Q(s,a) ──▶ number
                      q-value

- State value function let us know how good is to be in certain state.

- State-action value functions let us know how good is to take certain action given the actual state.

**NC STATE UNIVERSITY**

**Biomechatronics and Intelligent Robotics (BIRO) Lab**

# Interaction of the elements
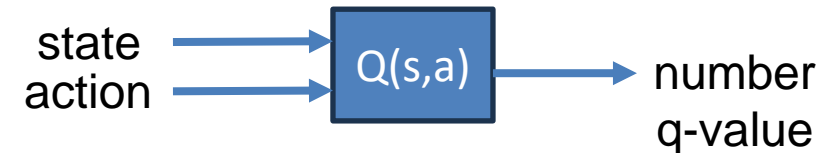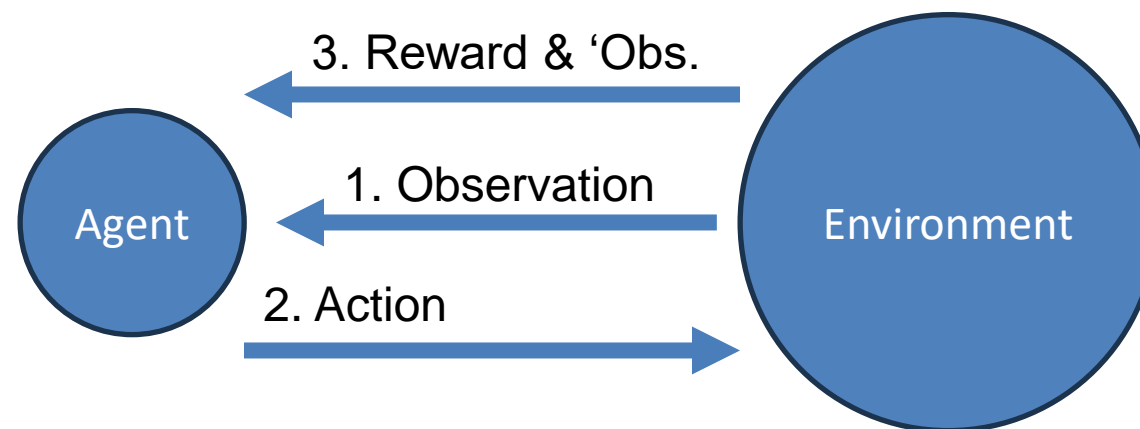
3. Reward & 'Obs.

1. Observation

**Agent**

**Environment**

2. Action

Bellman's equation example

$$V(s) = max_a \left( R(s,a) + V(s') \right) \qquad V(s) = max_a \left( \gamma V(s') \right)$$

Value | States

$$V(s) = max_a \left( R(s,a) + \gamma V(s') \right)$$

Actions

Reward
-1

Discount factor
0.9

| 1 |

| -3 | 2.6 → 4 |

| 2 |

$4 \cdot 0.9 - 1 = 2.6$

| 2 ← 1 → 2 → 3 → 4 → +5 |
| 3 ← 2 1 → 2 → 3 → 4 |
| +4 0 → 1 → 2 → 3 |
| 3 -1 |
| 2 ← 1 ← 0 ← -1 ← -2 ← -2 |
| 1 ← 0 ← -1 ← -2 ← -2 -1 |

https://youtu.be/SgC6AZss478?si=1JzdH4ZsYAvHJ4CN

**BIRO**

**Biomechatronics and Intelligent Robotics (BIRO) Lab**

# Q-Learning

- It is about the total reward

- The Q-function:
  - States
  - Actions $$R_T = Q(s, a)$$
  - Total cumulative reward
- The policy is the function that maps from states to actions $a = \pi(s)$
- Putting them together $R_T = Q_\pi(s, a)$

- The reward function $R_T = R_1 + R_2 + R_3 + \cdots + R_n$
  - Discount factor

$$R_T = R_1 + \gamma R_2 + \gamma^2 R_3 + \cdots + \gamma^n R_{n+1} = \sum_{i=0}^{n} \gamma^i R_{i+1}$$

- The Bellman's Equation (for this purpose)

$$Q(s, a) \leftarrow Q(s, a) + \alpha(r(s', a') + \gamma \max_a Q(s', a) - Q(s, a))$$

**NC STATE**
UNIVERSITY

**Biomechatronics and**
**Intelligent Robotics (BIRO) Lab**

# Deep Q-Network (DQN)

Deep reinforcement learning algorithm.

- Uses a neural network to approximate the Q-values. The network takes the state as input and outputs Q-values for all possible actions.

- Uses experience replay, where experiences (state, action, reward, next state) are stored in a replay buffer. During training, random mini-batches are sampled from this buffer to break correlation between consecutive experiences and stabilize training.

- Employs a target network (a copy of the main network), which is updated less frequently. This helps to stabilize training by providing consistent target Q-values.

- The neural network is trained using a gradient descent algorithm, minimizing the mean squared error between the predicted Q-values and the target Q-values:

$$L = (r(s, a) + \gamma \max_{a'} Q^*(s', a') - Q(s, a))^2$$

- It is suitable for environments with large or continuous state spaces, where the Q-table would be impractical. The use of neural networks allows it to generalize across similar states and actions.
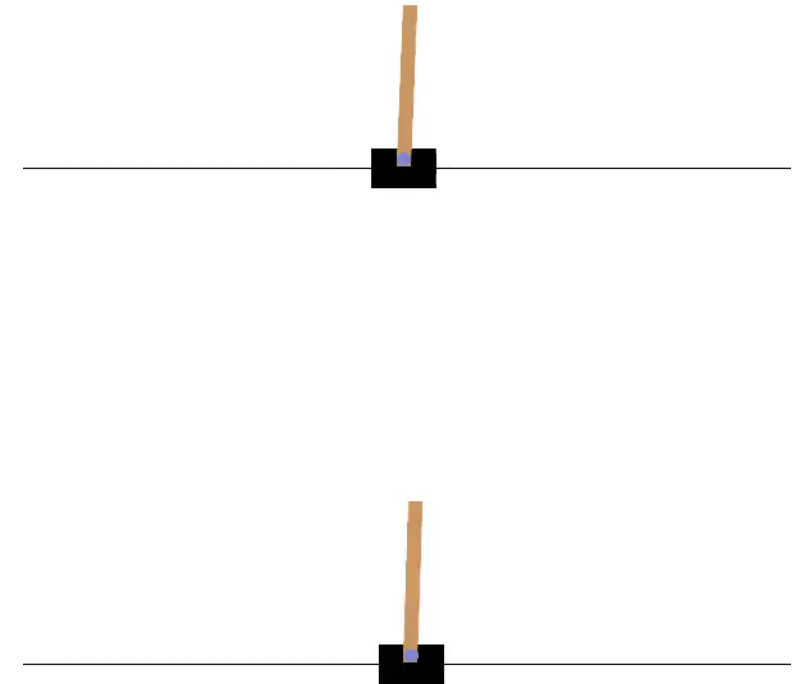
# Q-Learning VS DQN

| | Q-learning | Deep Q-Network |
|---|---|---|
| Q-values representation | Q-table | Neural Network |
| Scalability | Small and discrete state-action spaces | Large or continuous state spaces |
| Learning stability | Stable only on small state action spaces | Stable training given its experience replay and target network feature |
| Update mechanism | Bellman's equation | Gradient descent and loss minimization (other algorithms can be employed) |
| Computational Requirements | Low computational power and memory consumption | High computational resources needed |
| When to use it | When dealing with small, discrete state-action spaces | When dealing with large or continuous state-action spaces |

**Biomechatronics and Intelligent Robotics (BIRO) Lab**

# Cart-pole example

**Objective**

1. Understand the basic concepts in reinforcement learning:
   1. Environment
   2. Agent
   3. States
   4. Observation
   5. Action
   6. Reward
   7. Policy

2. Compare the performance of the trained policy regarding
   1. Number of neurons
   2. Number of network hidden layers
   3. Training episodes

Using Google Colab and OpenAI Gym

# Cart-pole Environment



- **Cart and Pole Dynamics:**
  - **Cart: A small cart that can move left or right along a frictionless track.**
  - **Pole: A pole attached to the cart by an unactuated rotational joint. The pole starts upright and can fall to either side.**

- **State Space (four-dimensional vector):**
  - **cart_position: The position of the cart on the track.**
  - **cart_velocity: The velocity of the cart.**
  - **pole_angle: The angle of the pole from the vertical.**
  - **pole_velocity_at_tip: The velocity of the tip of the pole.**

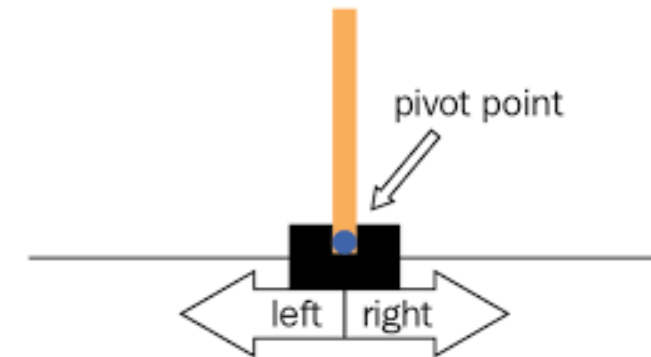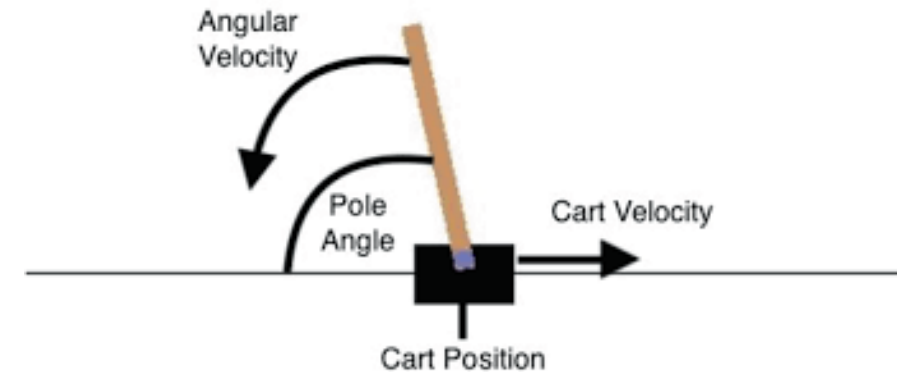- **Action Space (discrete with two possible actions):**
  - **0: Apply a force to the cart to move it to the left.**
  - **1: Apply a force to the cart to move it to the right.**

- **Reward:**
  - **The agent receives a reward of +1 for every time step the pole remains upright and within the allowed boundaries.**

- **Termination Conditions (episode):**
  - **The pole angle exceeds ±12 degrees from the vertical.**
  - **The cart position exceeds ±2.4 units from the center.**
  - **The episode length reaches a maximum of 200 steps (500 steps in our example).**



**NC STATE UNIVERSITY**

**Biomechatronics and Intelligent Robotics (BIRO) Lab**

# Homework 1

**Based on the [Examples](#) obtained in the [Google Colab Notebook](#).**

1. Compare the performance and reward values per episode for different models employing the same agent (Deep Q-Network). Please describe your findings in text and attach figures and videos, if necessary.

    1. Training episodes: Model with 2 hidden layers and 24 neurons per layer with relu activation funcions for the hidden layers and a linear activation function for the output layer.

        1. 5,000 training episodes.
        2. 10,000 training episodes.
        3. 20,000 training episodes.

    2. Number of neurons per layer: Model with 2 hidden layers with relu activation functions for the hidden layers and a linear activation function for the output layer, 5000 training episodes.

        1. 24 neurons per layer
        2. 32 neurons per layer
        3. 40 neurons per layer

    3. Number of hidden layers: Model with 24 neurons per layer with relu activation funcions for the hidden layers and a linear activation function for the output layer, 5,000 training episodes.

        1. 2 hidden layers
        2. 4 hidden layers
        3. 6 hidden layers

    HINT: you can use several Google Colab notebooks simultaneously to accelerate the process of training each model.

**NC STATE UNIVERSITY**

**BIRO** Biomechatronics and Intelligent Robotics (BIRO) Lab

# Homework 1

**Based on the [Examples](#) obtained in the [Google Colab Notebook](#).**

1. Submit a brief answer to the questions below using both text and figures.
   1. Diagram of each of the neural network architectures (without repeating them) specifying:
      1. Number of inputs
      2. Number of outputs
      3. Weights between each layer (trainable parameters)
      4. Biases (per layer)
      5. Total of trainable parameters

HINT: You can use any software to make the diagrams (or can be done by hand).