



VLSI FOR MACHINE LEARNING

PEDRO JULIAN, DIEGO GIGENA, NICOLÁS
RODRIGUEZ

CAE 2023, CORDOBA





ORGANIZATION

- ☰ 1. Introduction
- ☰ 2. Basic elements of DNN
 - ☰ Neurons, layers, Convolution, Relu, MaxPool, etc.
- ☰ 3. Introduction to Pytorch
 - ☰ Tutorial: Mnist with 2 layers.
- ☰ 4. Case study: Tutorial and example
 - ☰ Task 1.1: You will have to design a multilayer NN using the tools (dataset provided, but your design)
- ☰ 5. VLSI realizations
 - ☰ Architectures, quantization
 - ☰ Basic elements (multipliers, adders, registers, etc)
 - ☰ Task 1.2: quantify size of designed multilayer NN
- ☰ 6. ML architectures (depth wise, Resnet, LSTM, transformers)
- ☰ 7. Review of case study
 - ☰ Task 1.3: Given an available area, e.g. 2mm², estimate resources you can build and the time to execute the multilayer NN

	Lunes	Martes	Mier.
9-10:40	Introduction - Pedro	Práctica	Práctica
11-13	Basic blocks of DNN -Nico	Caso de studio - Diego	Práctica
14-15:40	Intro Pytorch - Diego	Implementaciones en VLSI - Nico	Revisión de proyectos - ALL
16-18	Práctica	Práctica	Advanced architectures - ALL



1. INTRO

01

02

03

04

05

Machine Learning and methods

Biology and NN

Artificial NN

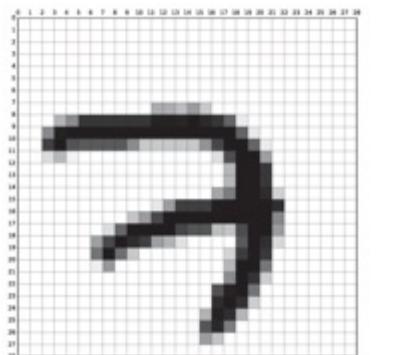
Latest results

Conclusions

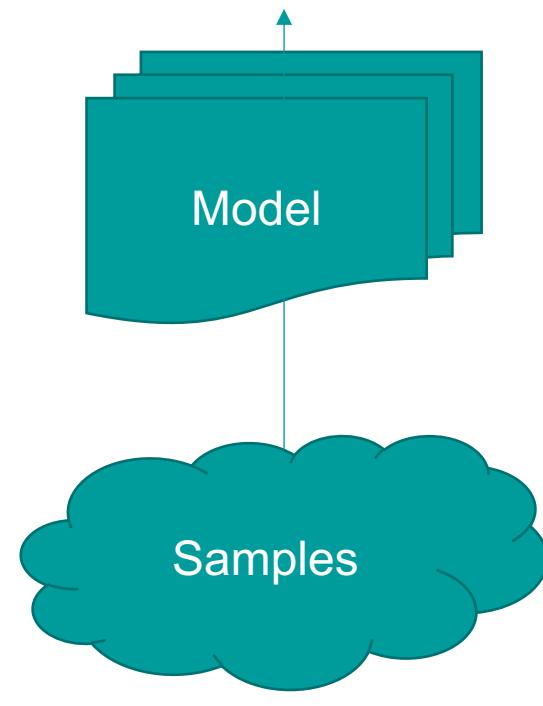


MACHINE LEARNING

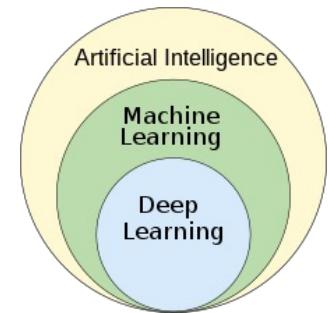
- Machine learning (ML) is a field of inquiry devoted to understanding and building methods that "learn" – that is, methods that leverage data to improve performance on some set of tasks.^[1] It is seen as a part of [artificial intelligence](#).



Input



0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9



METHODS



☰ Supervised Learning

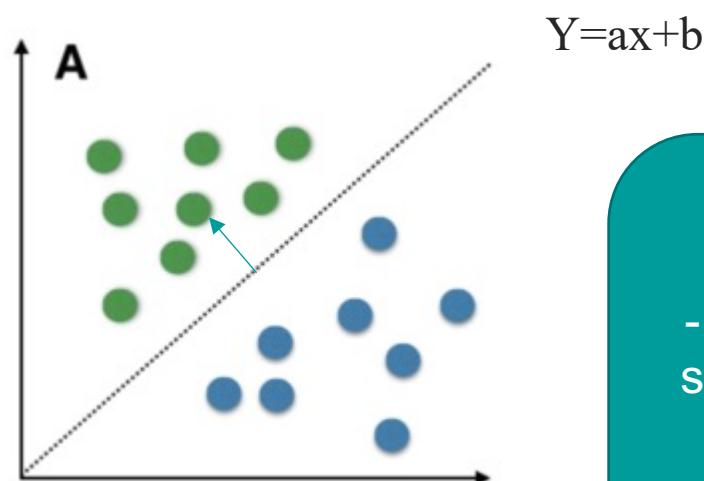
☰ Unsupervised Learning

☰ Reinforcement Learning



SUPERVISED LEARNING

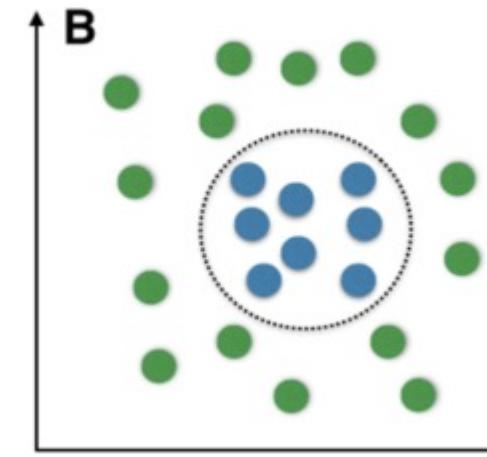
- ☰ Training set has a set of input data and the desired output
- ☰ Example: Support vector machine



$Y = ax + b$

Model

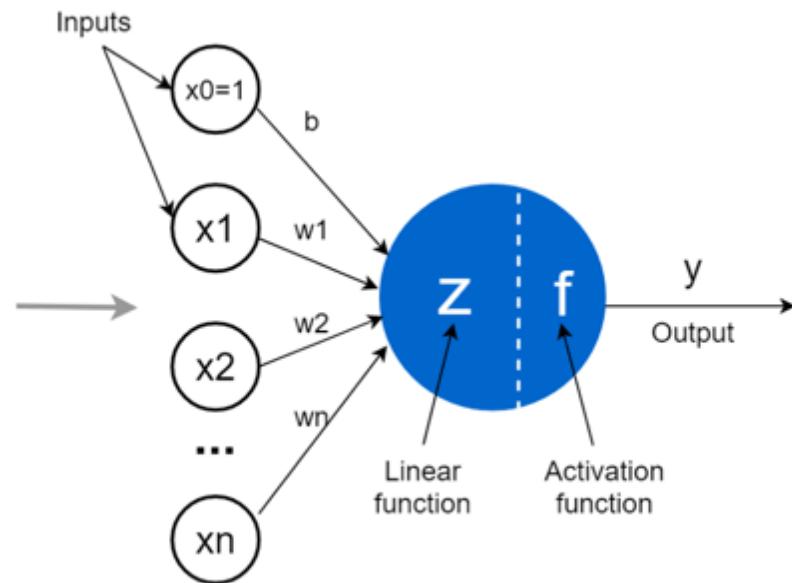
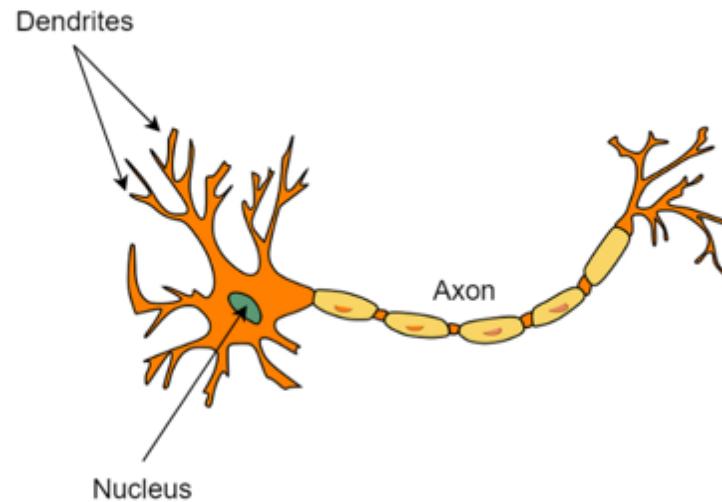
- Change a and b,
so that for all input
data, the output
classes are on the
right side of the line



Different models

SUPERVISED LEARNING: ANN

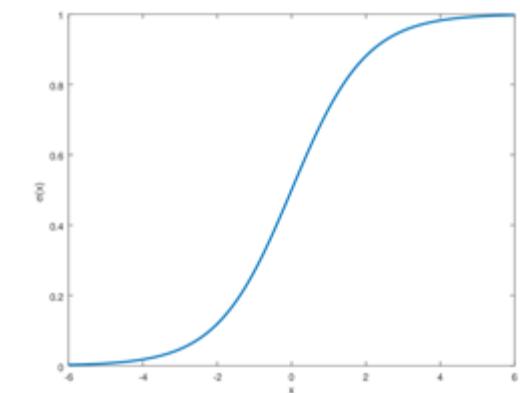
- ≡ Artificial Neural Networks:
- ≡ Single Neuron



Single Neuron:

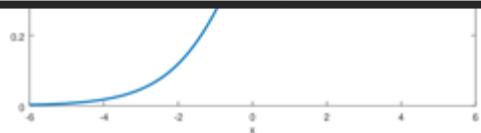
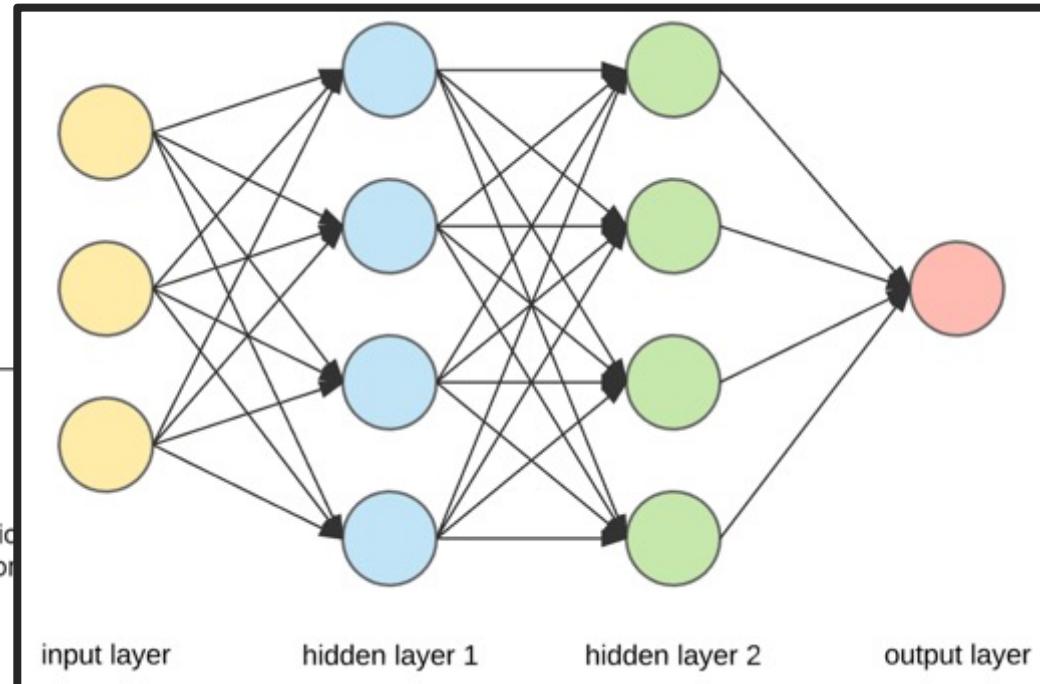
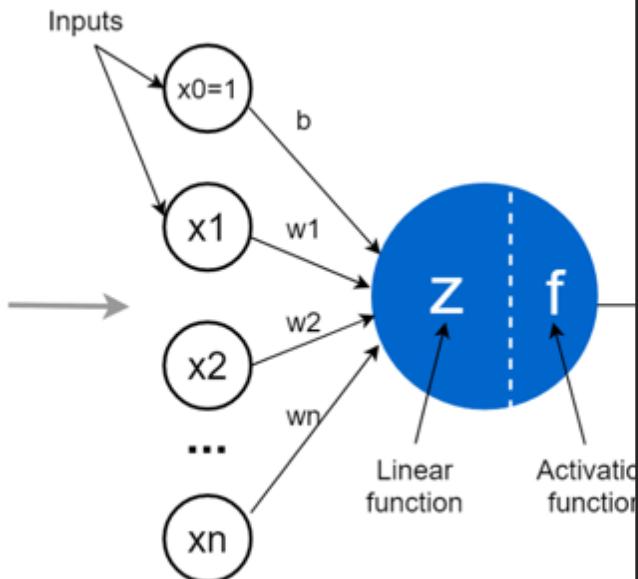
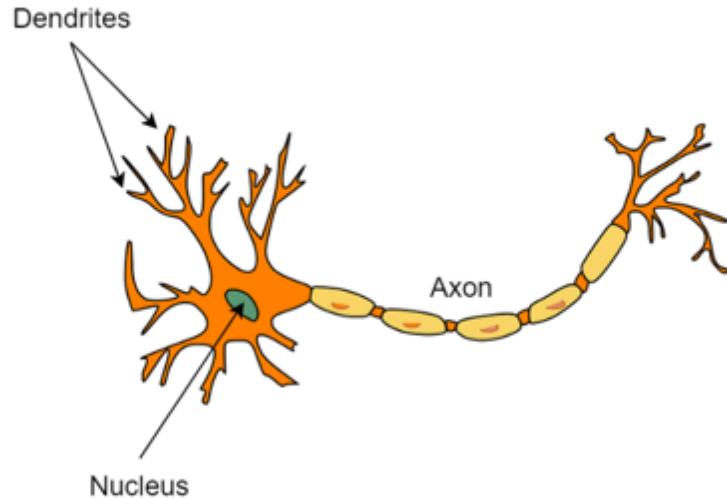
Model

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



SUPERVISED LEARNING: ANN

- ≡ Artificial Neural Networks:
- ≡ Single Neuron



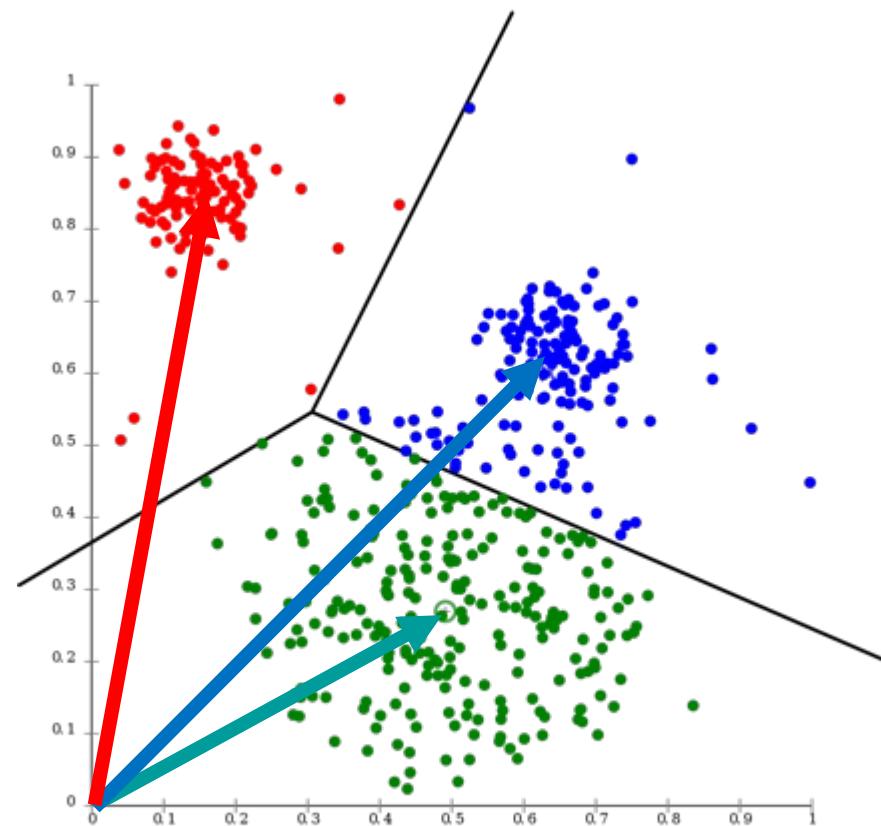
UNSUPERVISED LEARNING

- ☰ No tagged desired output:
- ☰ Model extracts patterns
- ☰ Example: Clustering

K-means clustering

Model:

- Assign k centroids
- Cluster elements so the distance is minimized



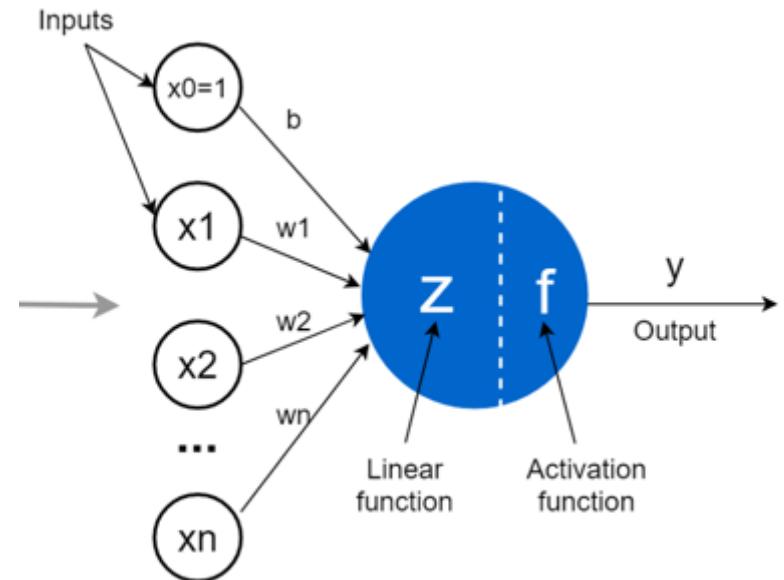
UNSUPERVISED LEARNING

- ☰ No tagged desired output:
- ☰ Model extracts patterns
- ☰ Example: Hebbian Learning

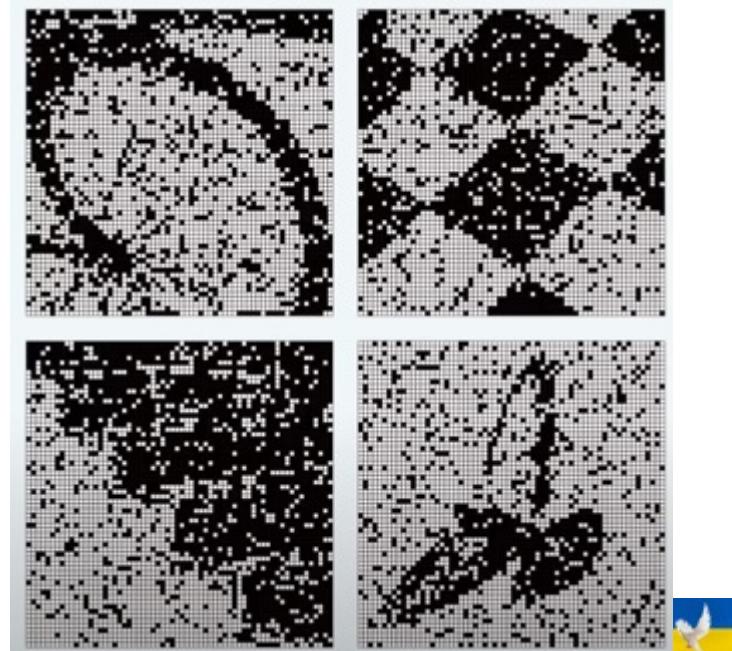
Hebbian Learning

Model:

- ANN
- Weight w_i is increased if input x_i and output are “active” together



Example: Hopfield Network



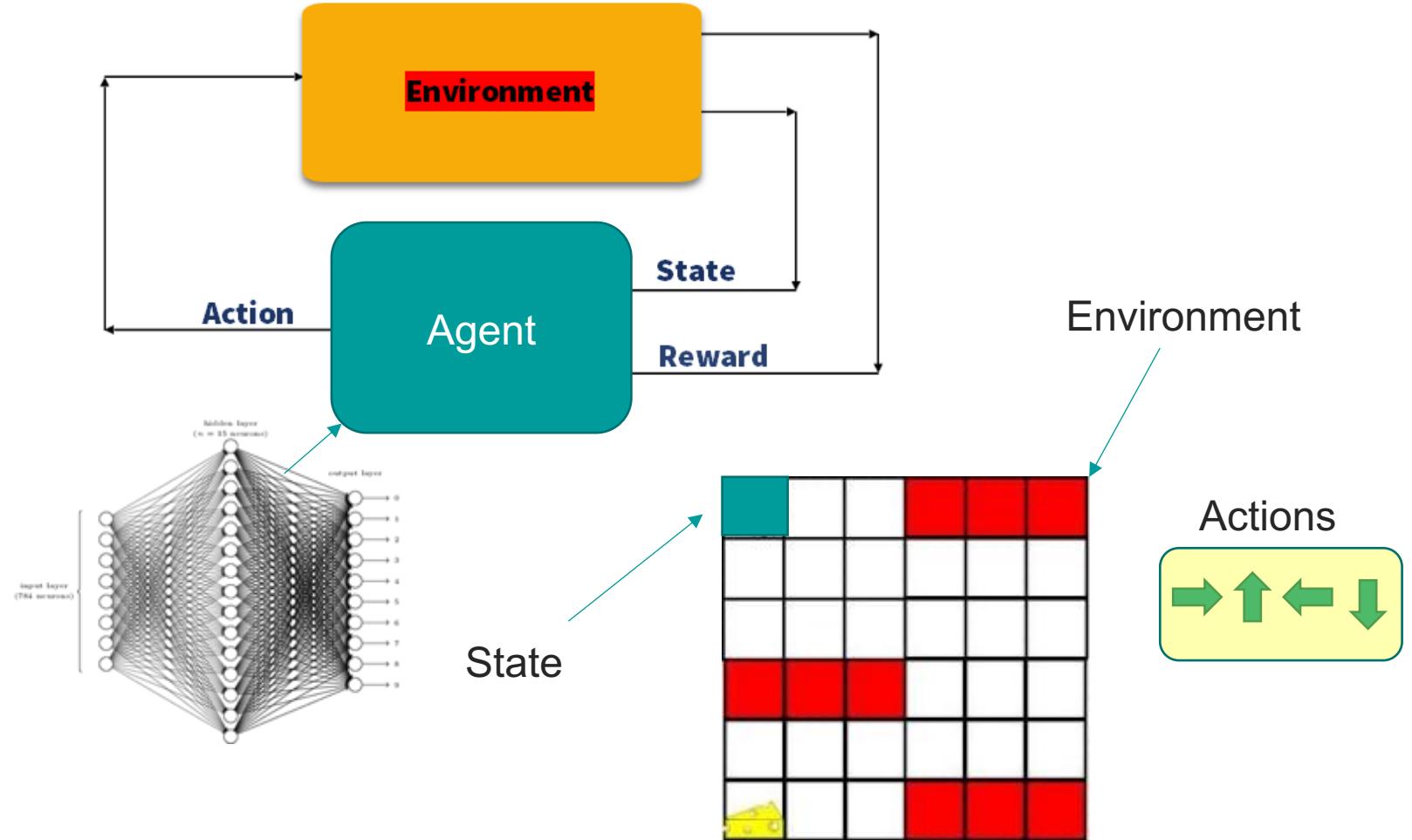
REINFORCEMENT LEARNING

- ≡ Agent interacts with environment
- ≡ Long term reward

Weak supervision

Model:

- agent learns to perform a task through repeated trial and error interactions with a dynamic environment
- Another example: chess



NEURONS

≡ How nature operates?



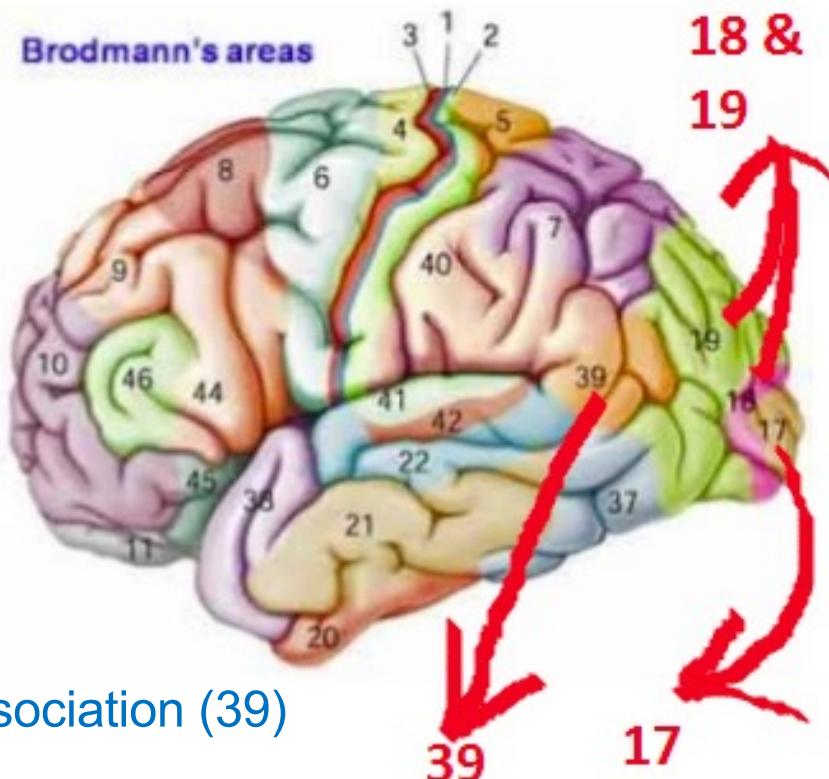
LESSONS FROM BIOLOGY



SAL
SILICON AUSTRIA LABS



Visual association (18&19)



- ≡ 100×10^9 neurons
- ≡ 1000×10^{12} connections
- ≡ Frequency: 10Hz
- ≡ Power: 20W
- ≡ Efficiency: 16mW/cm³



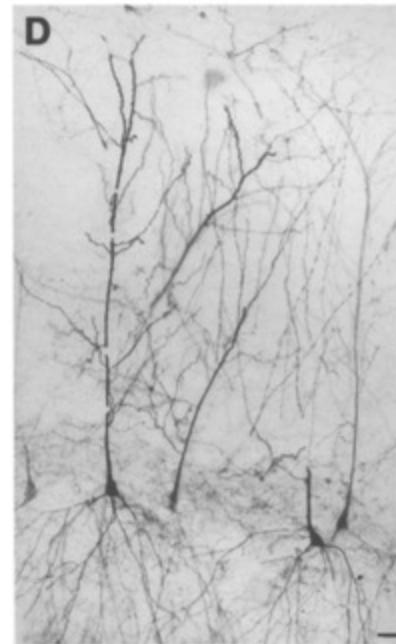
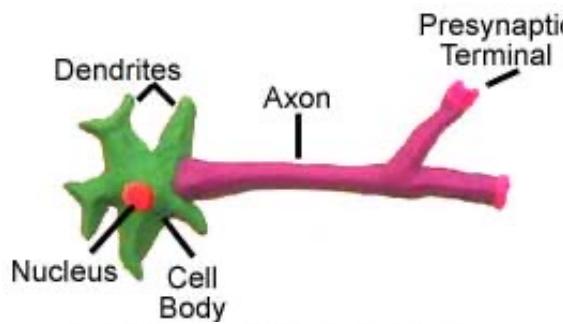
CELLULAR BIO-PROCESSING



- ≡ Local communication
- ≡ Large collection of cells (Massive parallel processing)
- ≡ Low speed
- ≡ Balance between complexity vs. area (volumen)
- ≡ 3D ...

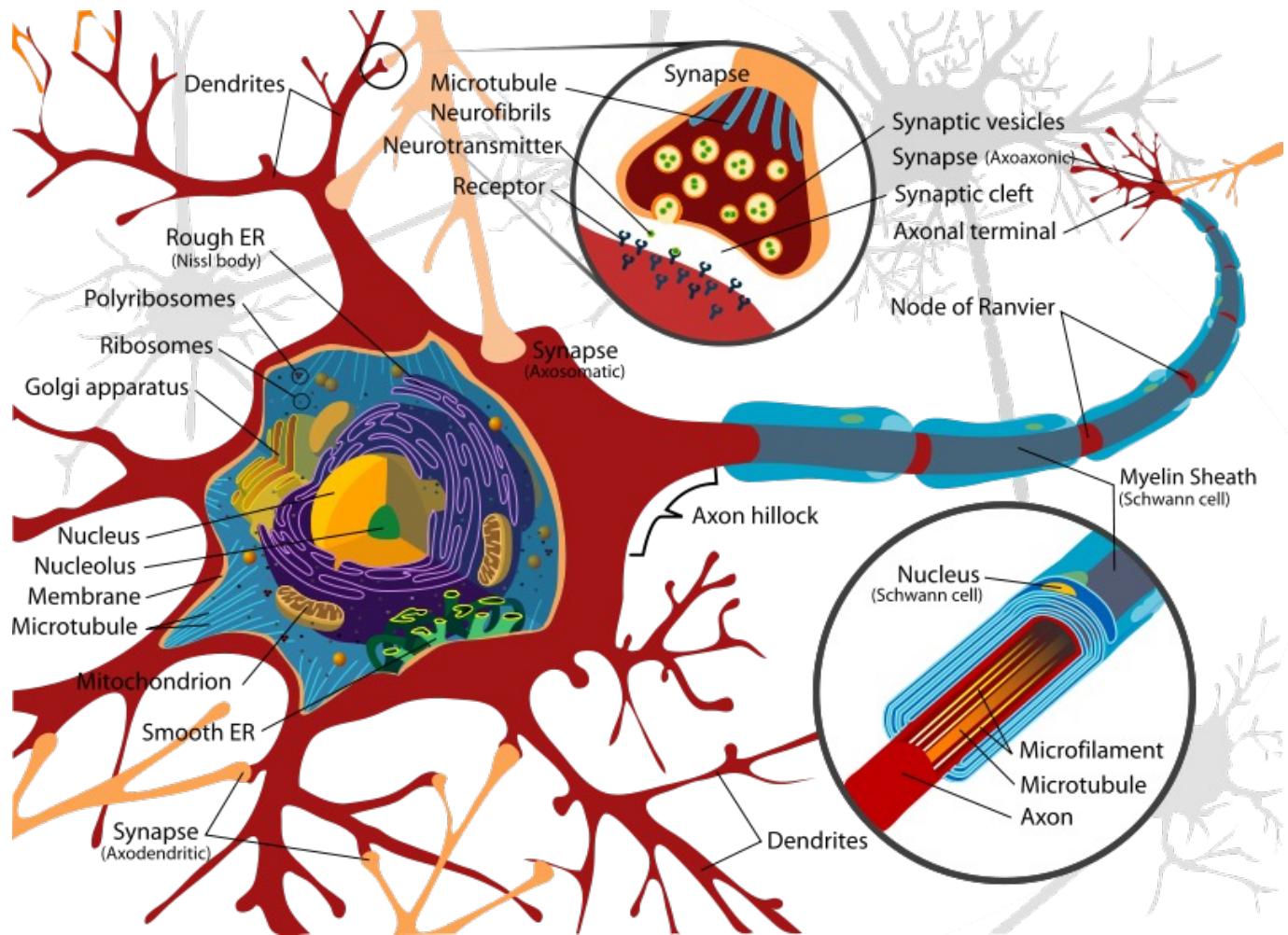


BIOLOGICAL NEURONS

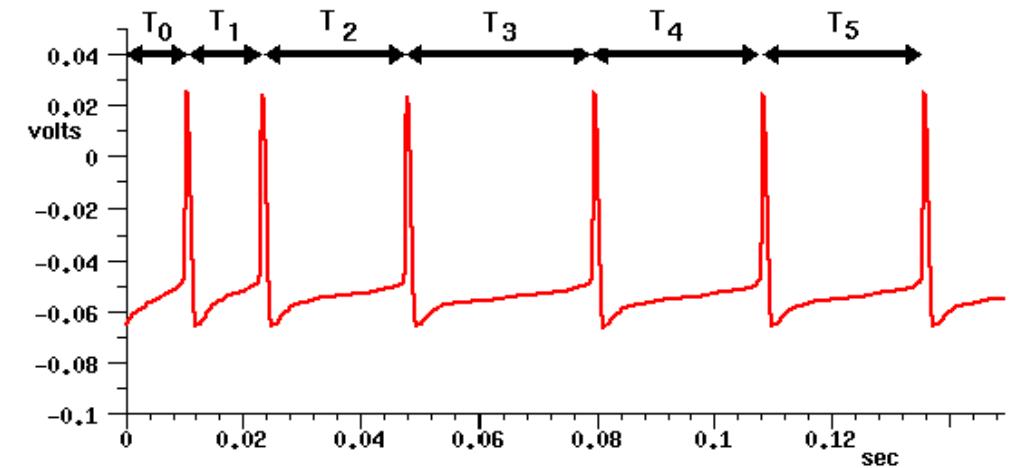
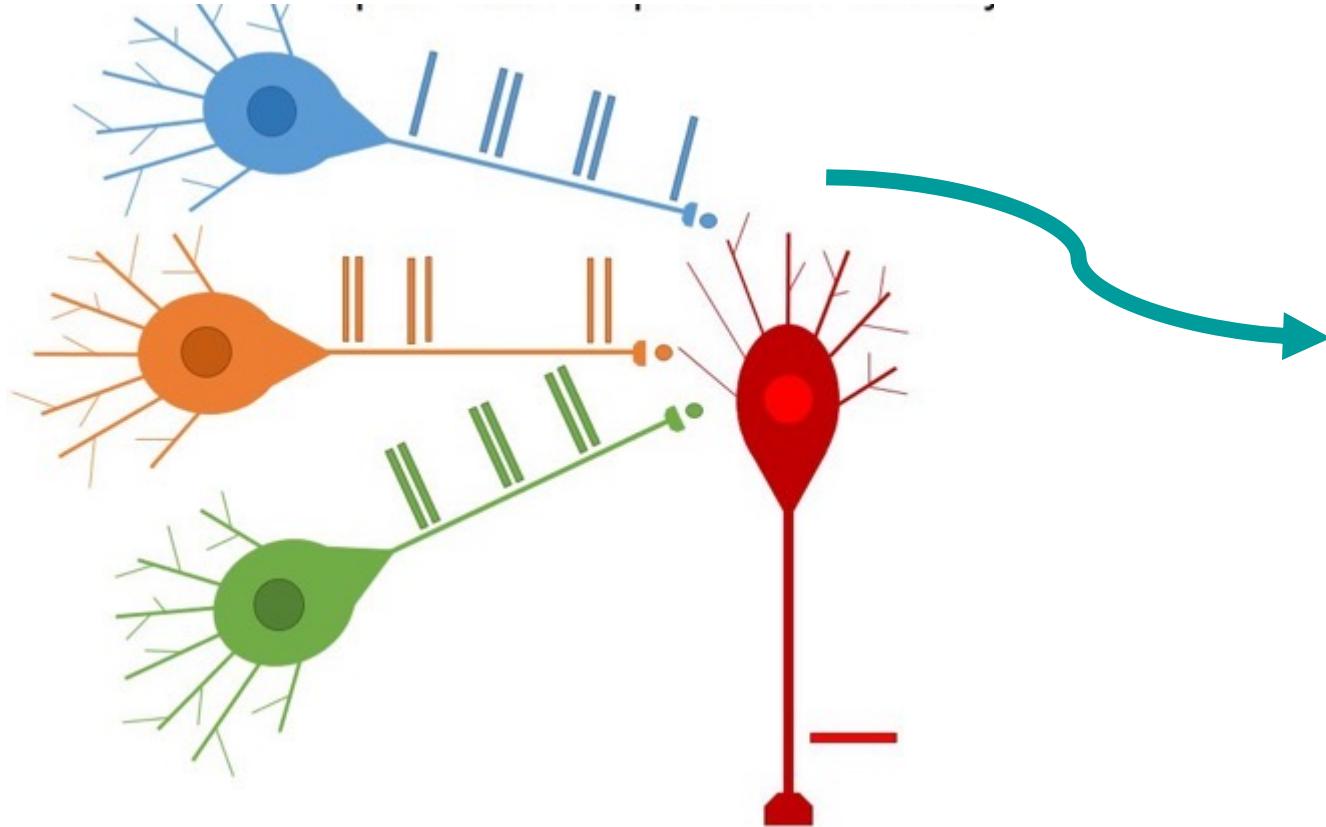


10um

Megias et al. Neuroscience, 2001



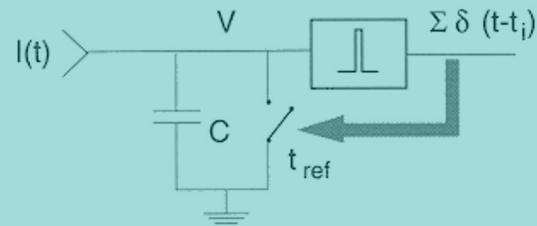
NEURONS



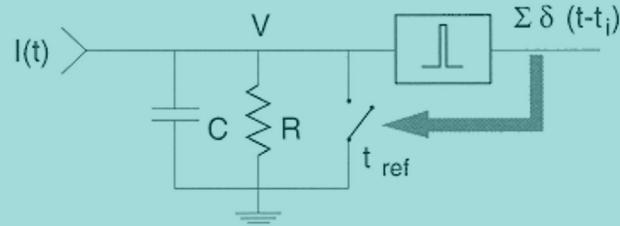
MODELS



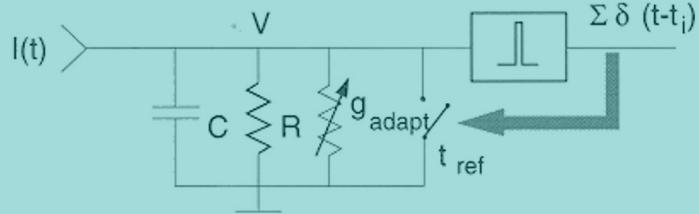
Perfect Integrate-and-Fire Unit



Leaky Integrate-and-Fire Unit

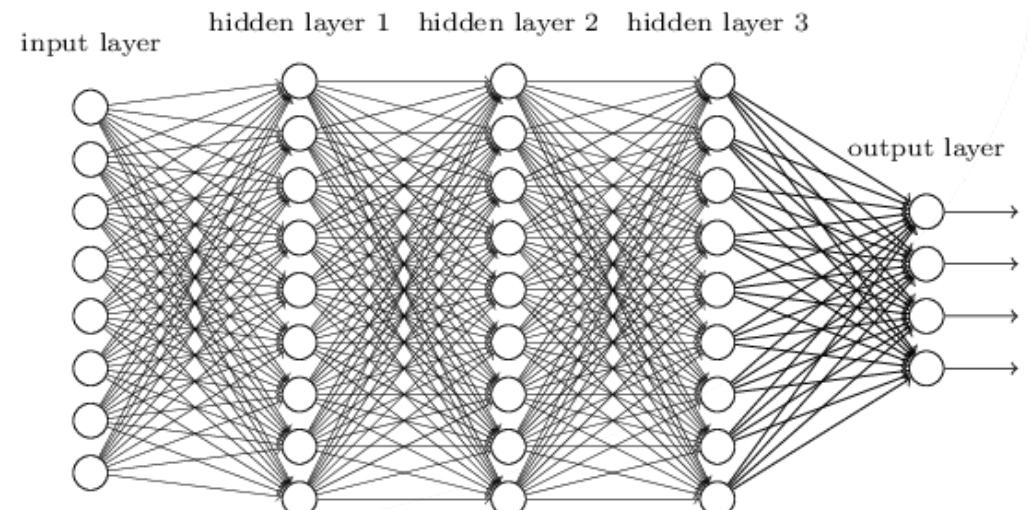
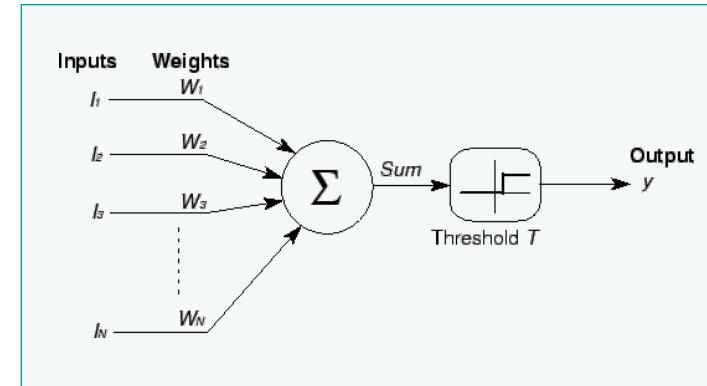


Adapting Integrate-and-Fire Unit

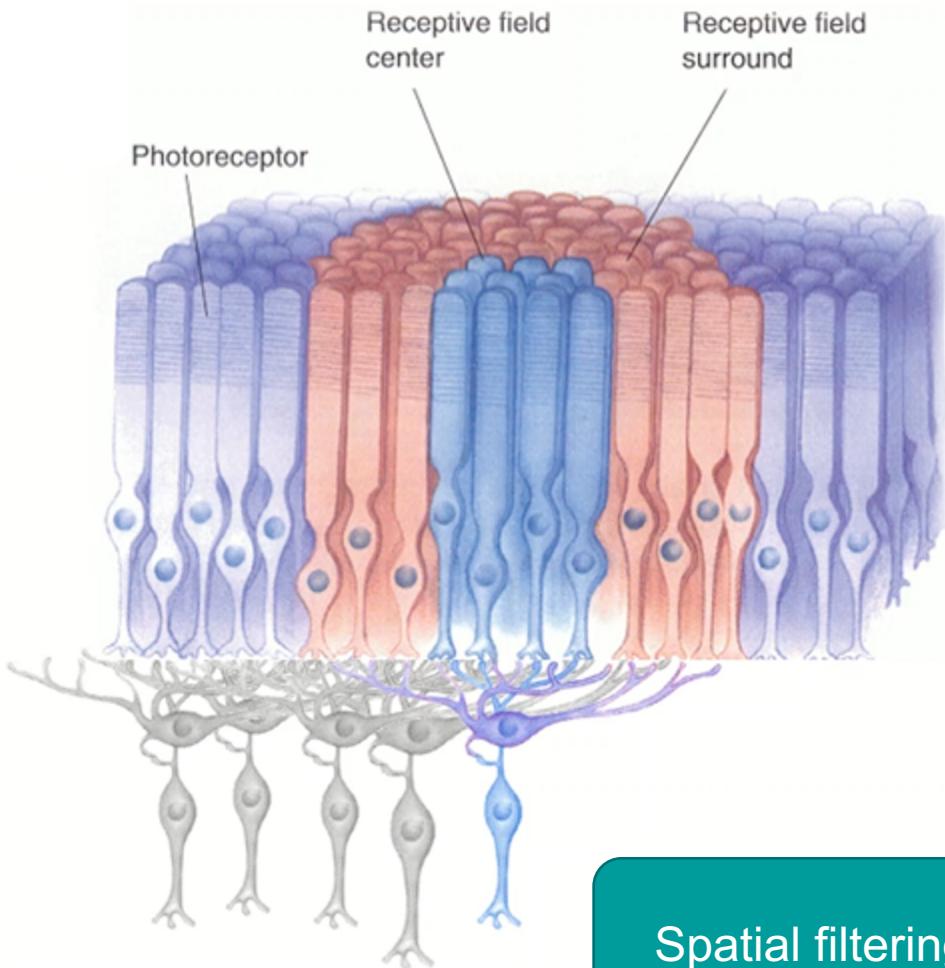


Spiking NN

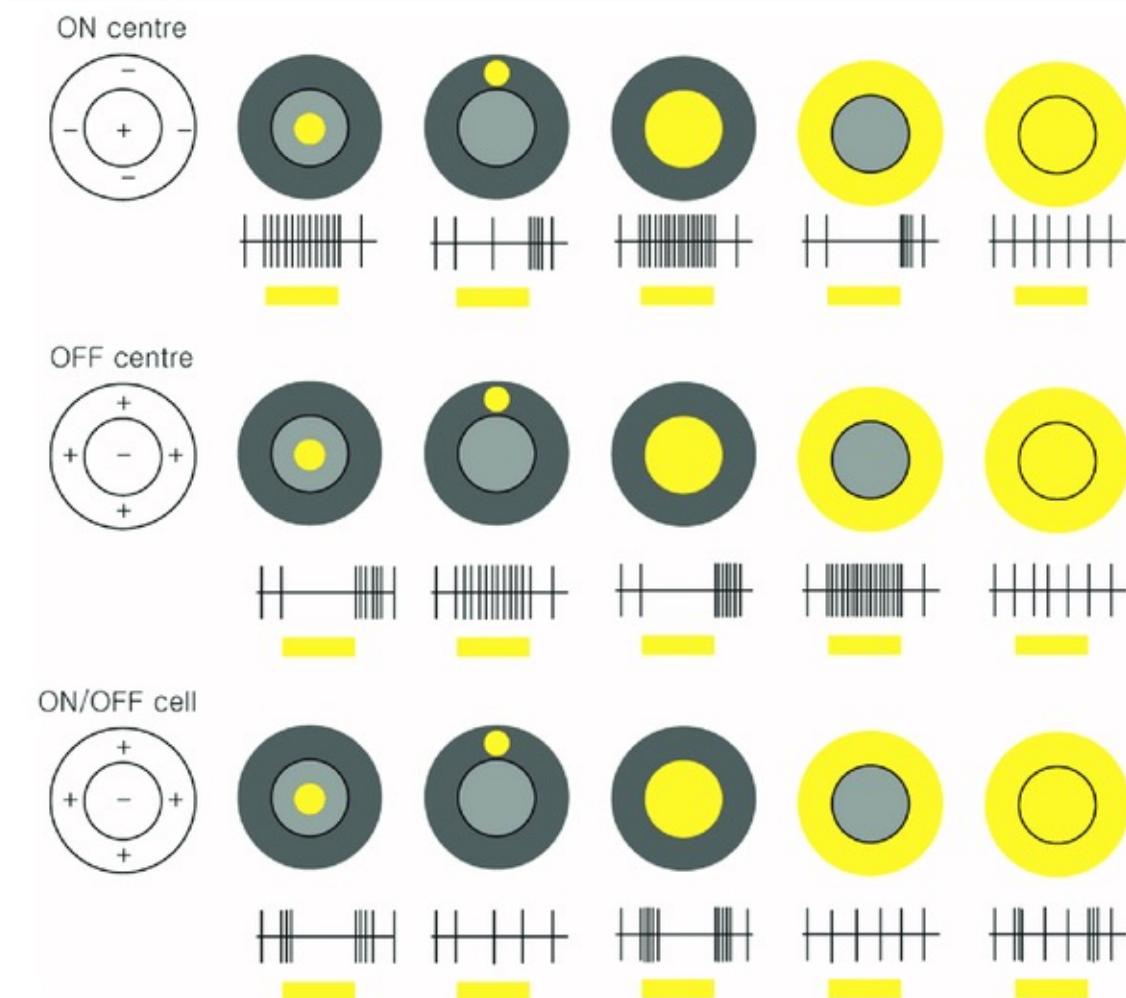
Artificial NN



RETINAL IMAGE PROCESSING



Spatial filtering

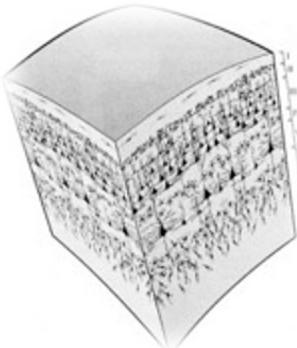
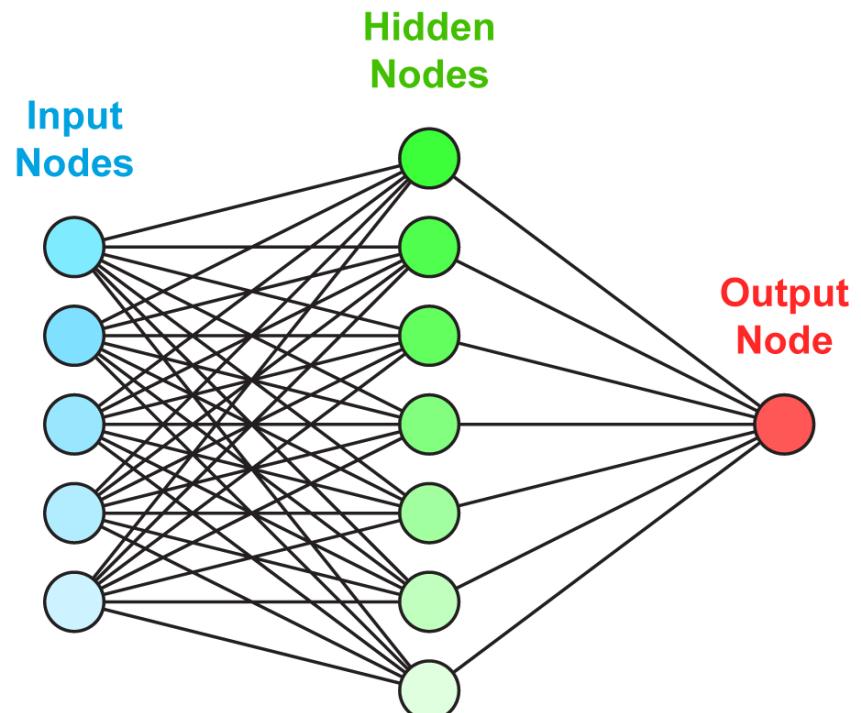




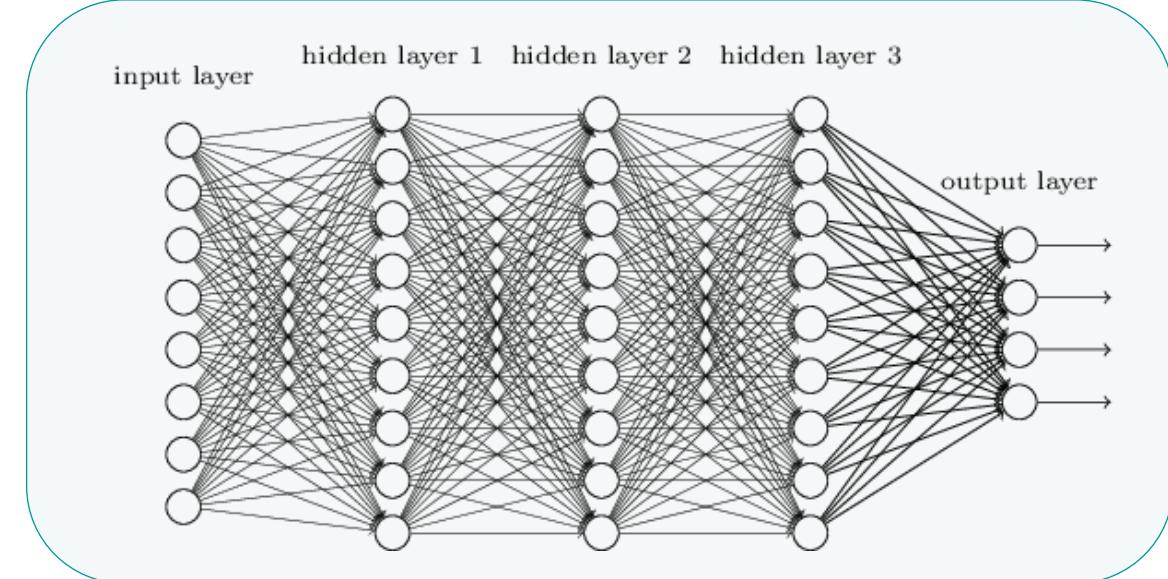
ARTIFICIAL NN

COMPUTING IN MACHINES WITH
BIOLOGICAL INSPIRED ARCHITECTURES

FULLY CONNECTED NN



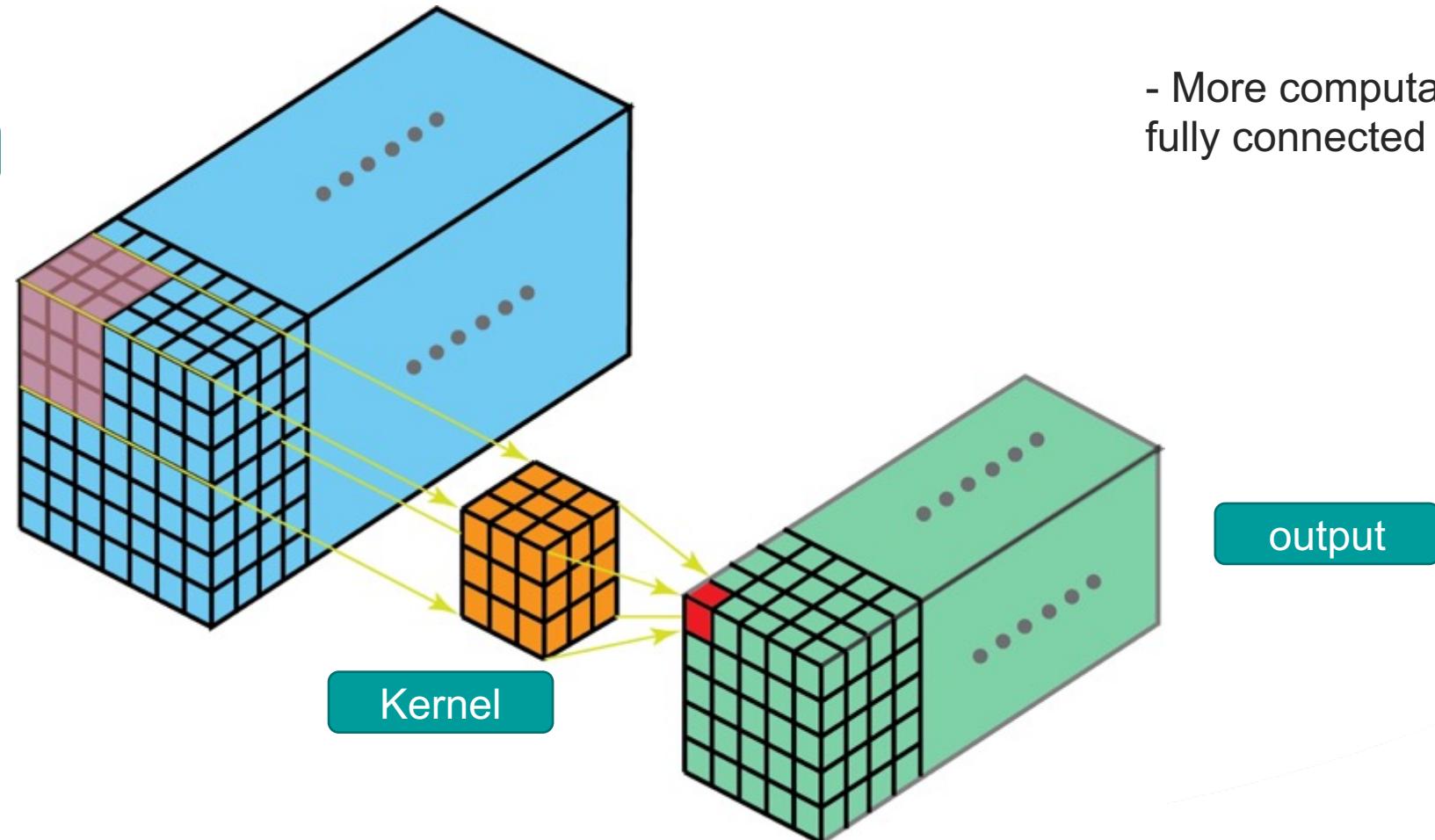
- All neurons from one layer connect to each neuron of the next layer



CONVOLUTIONAL NN

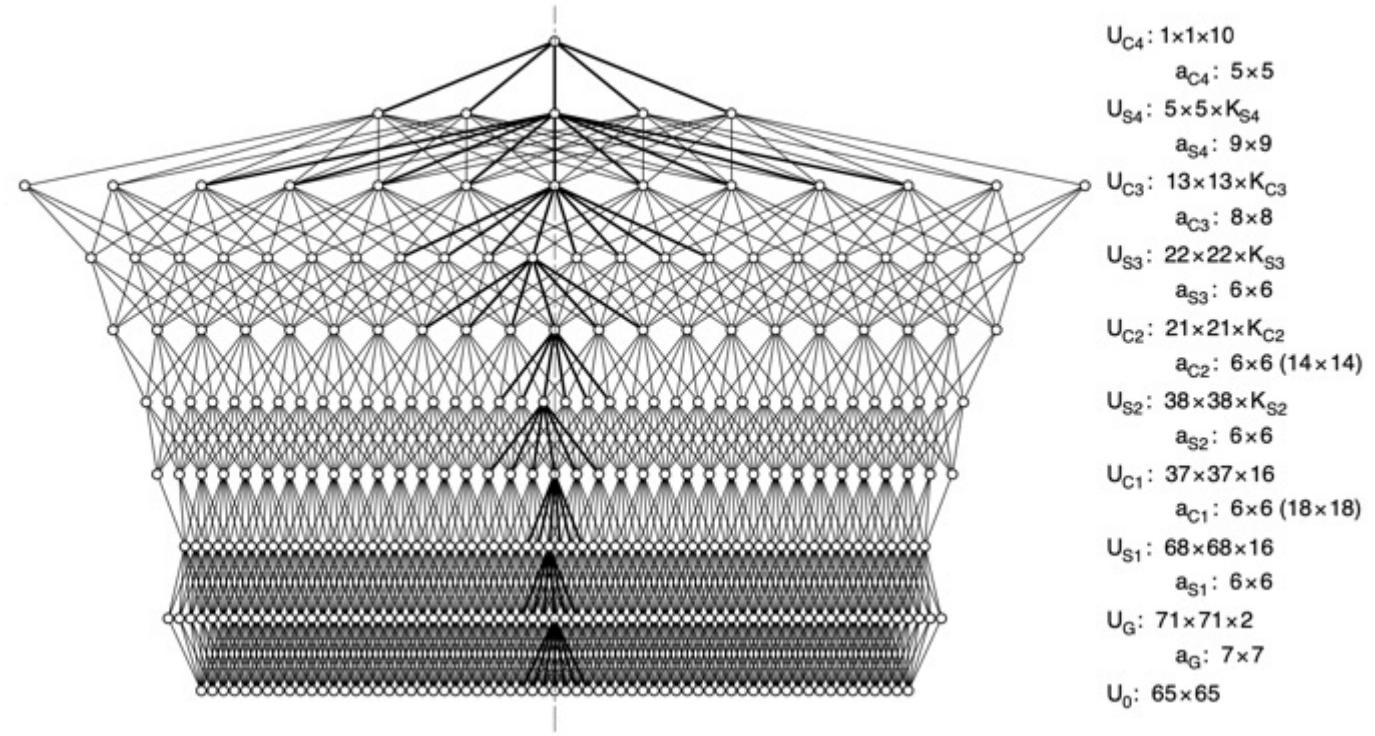
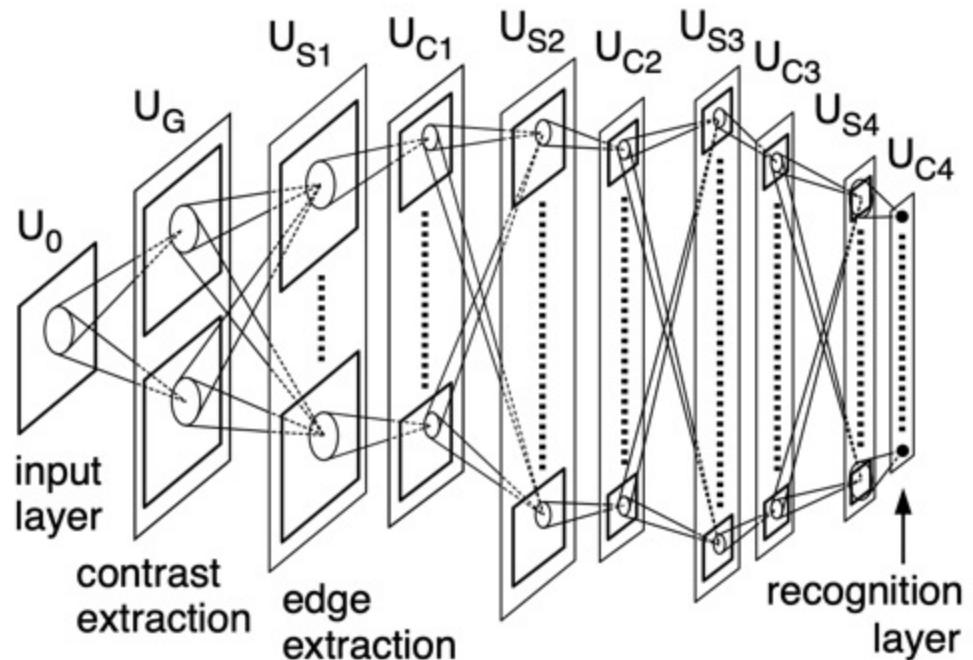


- ≡ A Convolutional layer performs a spatial (3D in general) convolution (spatial multiplication) with a kernel (filter)
- ≡ Resembles a more elaborate ON/OFF center cell



FROM SMALL TO HUGE

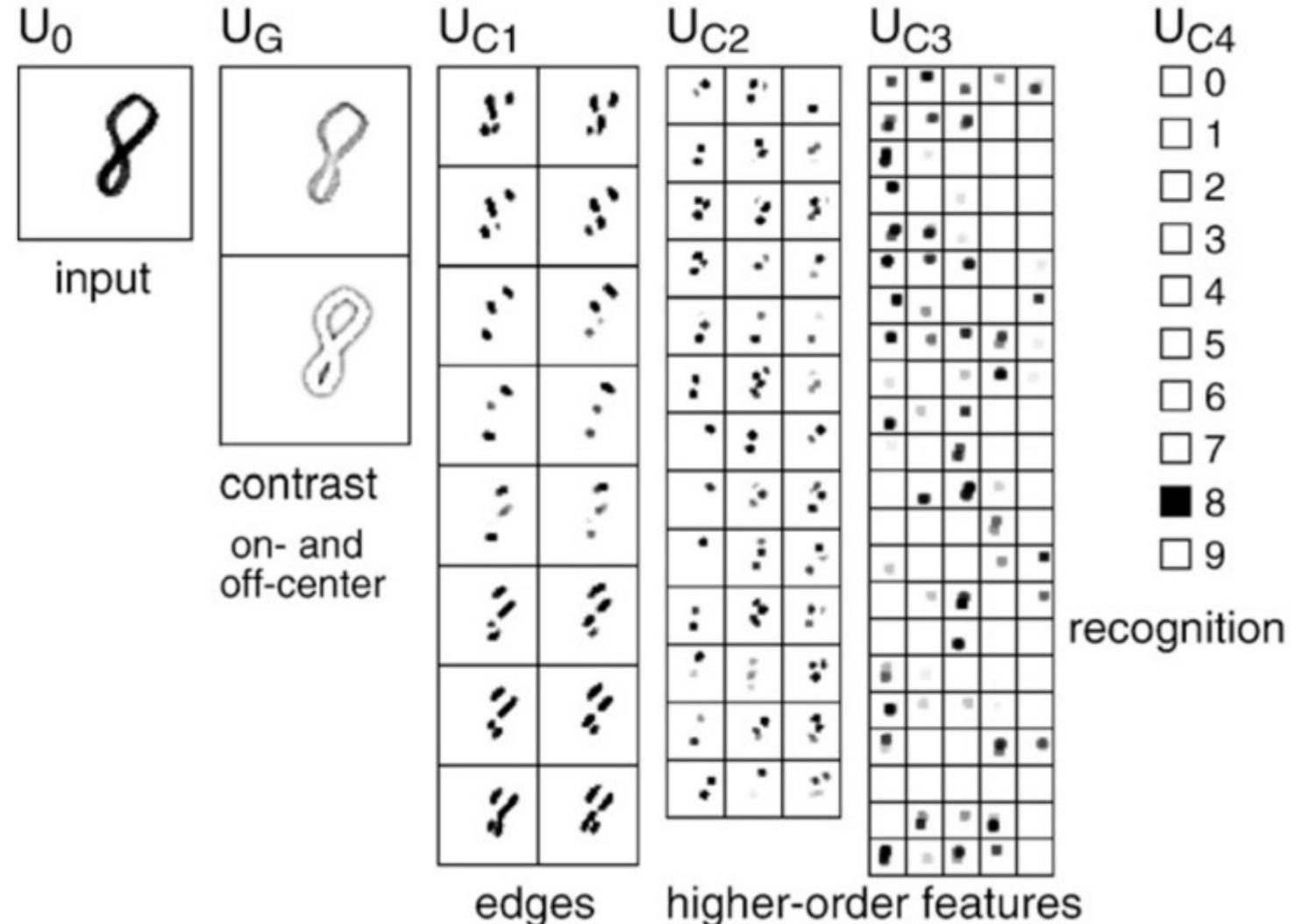
- ☰ 1979: The first “[convolutional neural networks](#)” were used by Kunihiko Fukushima (Neocognitron). Fukushima designed neural networks with multiple pooling and convolutional layers.



FROM SMALL TO HUGE

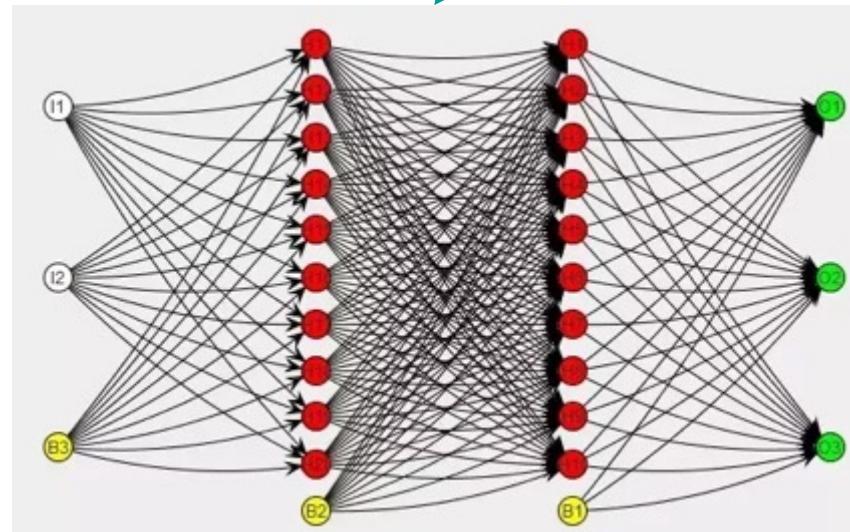


SAL
SILICON AUSTRIA LABS



FROM SMALL TO HUGE

- 1985: Rumelhart, Williams, and Hinton demonstrated back propagation in a neural network could provide "interesting" distribution representations



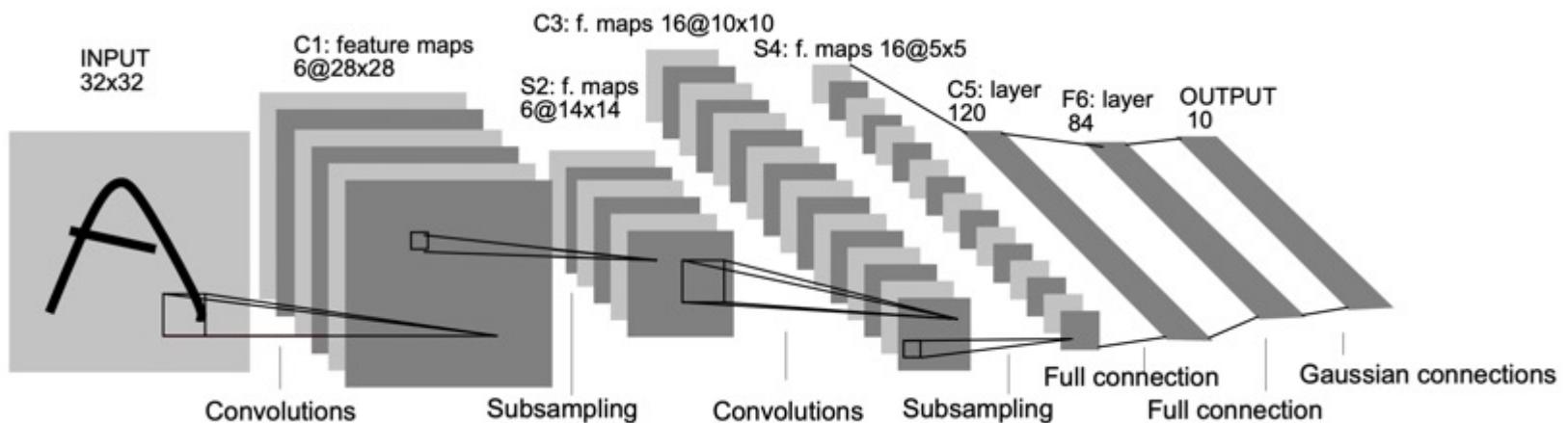
	Out	Truth	Error
Messi	2	10	-8
Mbappe	8	0	8
Dibu	4	0	4

Backpropagation propagates the error back, using the gradients to adjust each parameter

FROM SMALL TO HUGE



- In 1989, Yann LeCun provided the first practical demonstration of backpropagation. He combined convolutional neural networks with back propagation onto read “handwritten” digits.



3	6	8	1	7	9	6	6	9	1
6	7	5	7	8	6	3	4	8	5
2	1	7	9	7	1	2	8	4	5
4	8	1	9	0	1	8	8	9	4
7	6	1	8	6	4	1	5	6	0
7	5	9	2	6	5	8	1	9	7
1	2	2	2	2	3	4	4	8	0
0	2	3	8	0	7	3	8	5	7
0	1	4	6	4	6	0	2	4	3
7	1	2	8	7	6	9	8	6	1

- LeNet 5 would have taken several weeks of training back then
- LeNet 1 (2600 param.) ~ 100.000 multiply/add operations

- MNIST 60.000 element dataset



FROM SMALL TO HUGE

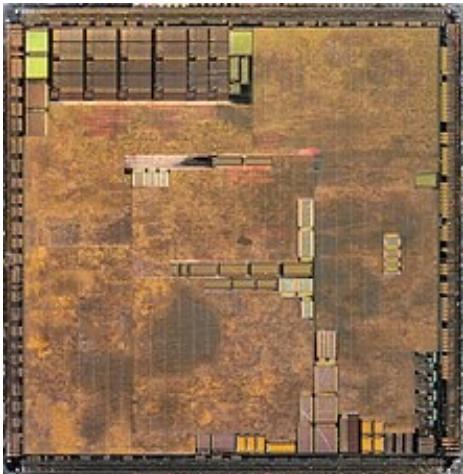


SAL
SILICON AUSTRIA LABS

- 1999: computers started becoming faster and GPU (graphics processing units) available.



1999 Nvidia GeForce 256



Overview	
Manufacturer	NVIDIA
Original Series	GeForce 256
Release Date	October 11th, 1999
Launch Price	\$179 USD
PCB Code	180-P0003-0100
Board Model	NVIDIA P3
Graphics Processing Unit	
GPU Model	NV10
Architecture	NV10
Fabrication Process	220 nm
Die Size	111 mm ²
Transistors Count	17M
Transistors Density	153.2K TRAN/mm ²
Pixel Pipelines	4
TMUs	4
ROPs	4
Clocks	
Base Clock	120 MHz
Boost Clock	TBC MHz
Memory Clock	143 MHz
Effective Memory Clock	143 Mbps
Memory Configuration	
Memory Size	32 MB
Memory Type	SDR
Memory Bus Width	64-bit
Memory Bandwidth	1.1 GB/s

FROM SMALL TO HUGE



SAL
SILICON AUSTRIA LABS

- 2009: Fei-Fei Li from Stanford launched [ImageNet](#), a free database of more than 14 million labeled images.

IMAGENET
Home Download Challenges About
Not logged in. Login | Signup

14,197,122 images, 21641 synsets indexed

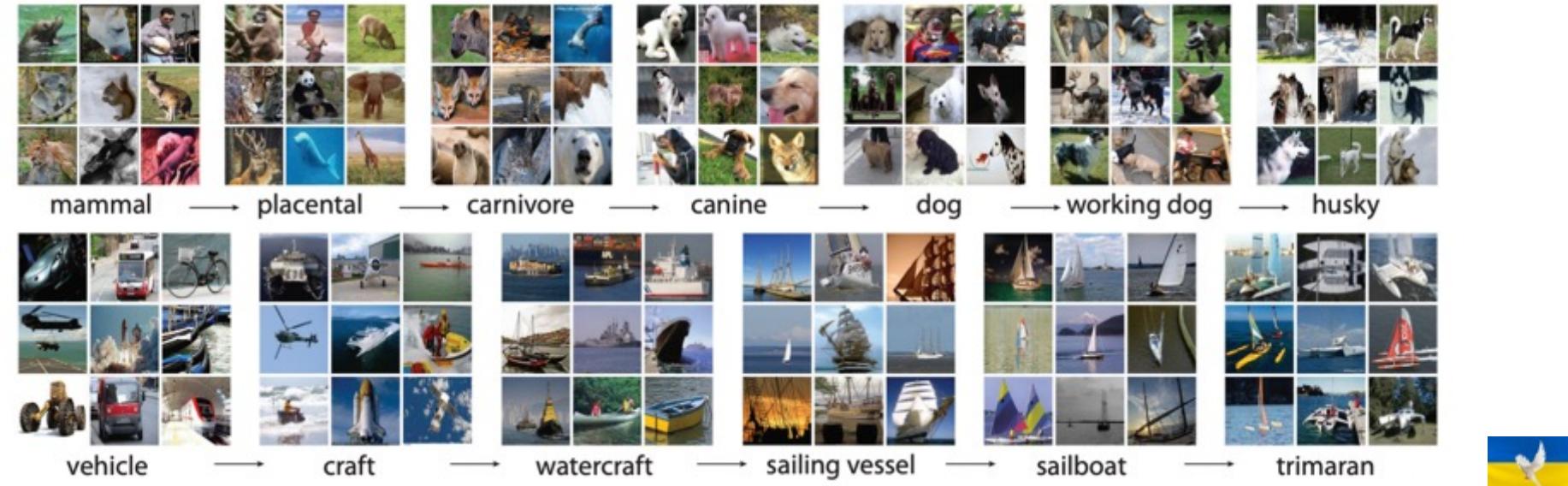
ImageNet is an image database organized according to the [WordNet](#) hierarchy (currently only the nouns), in which each node of the hierarchy is depicted by hundreds and thousands of images. The project has been instrumental in advancing computer vision and deep learning research. The data is available for free to researchers for non-commercial use.

Mar 11 2021. ImageNet website update.

© 2020 Stanford Vision Lab, Stanford University, Princeton University [imagenet.help.desk@gmail.com](#) Copyright infringement

ImageNet: A Large-Scale Hierarchical Image Database

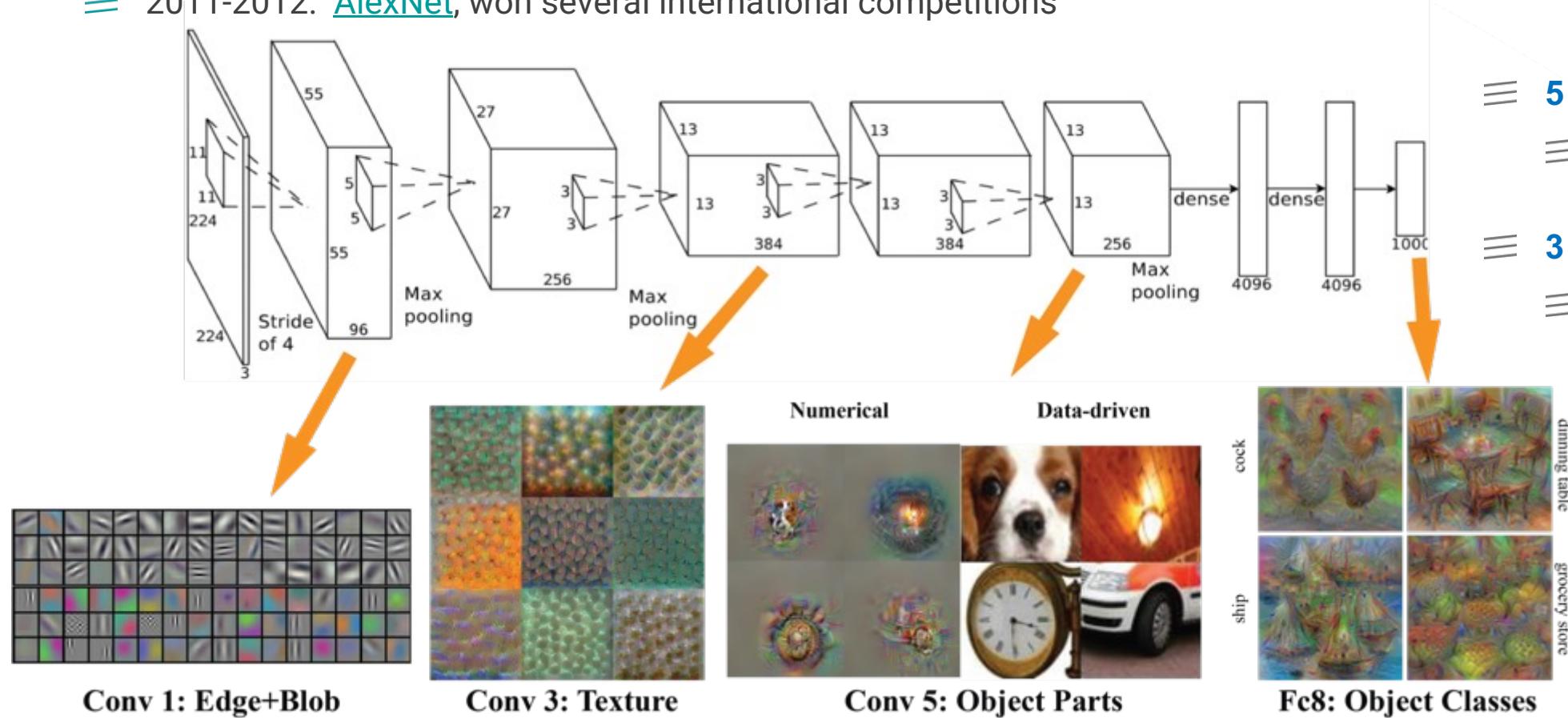
Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li and Li Fei-Fei
Dept. of Computer Science, Princeton University, USA
{jiadeng, wdong, rsocher, jial, li, feifeili}@cs.princeton.edu



3.2 millions

FROM SMALL TO HUGE

- 2011-2012: [AlexNet](#), won several international competitions

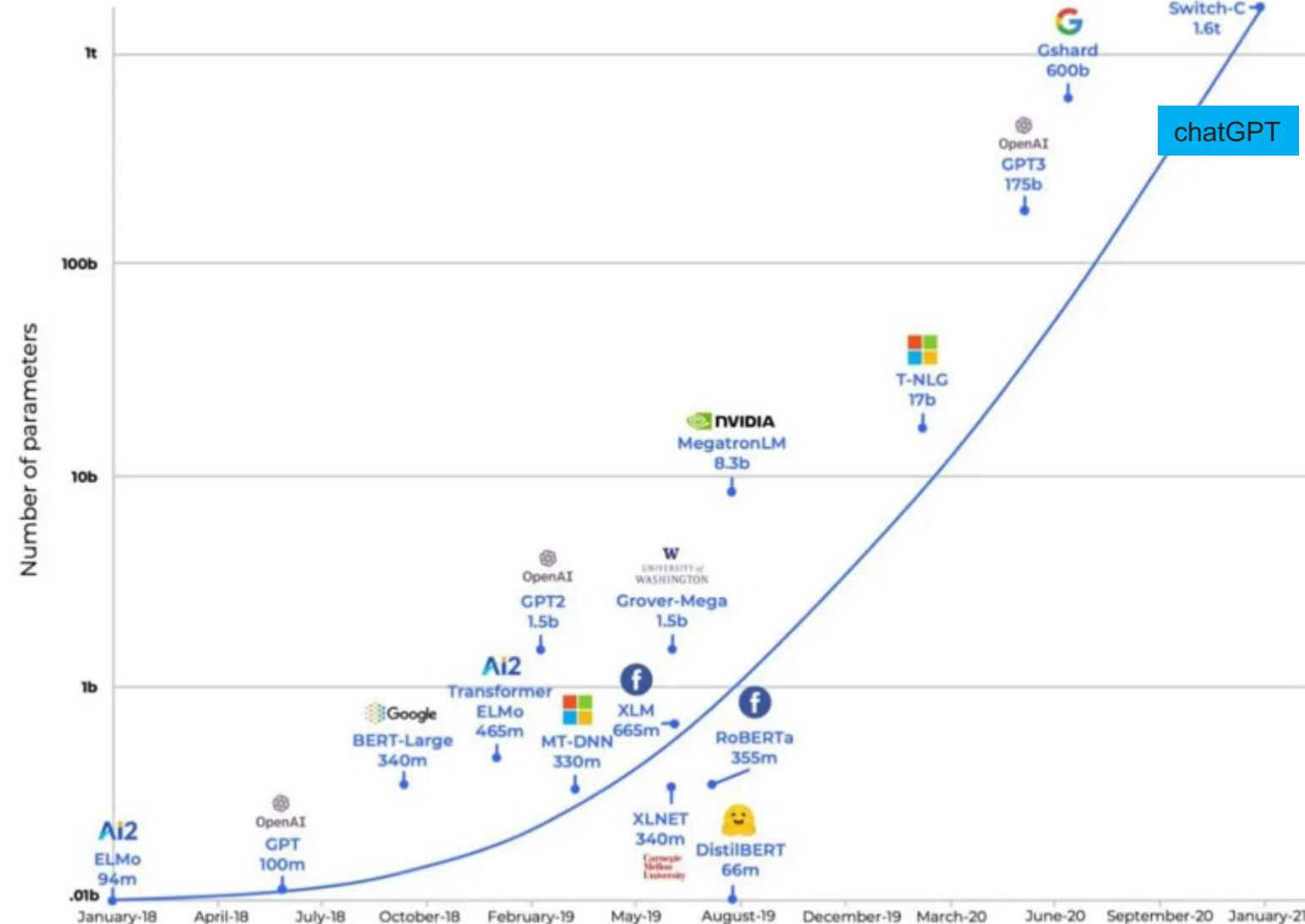


650K neuronas, 60M parámetros, 630M conexiones



FROM SMALL TO HUGE

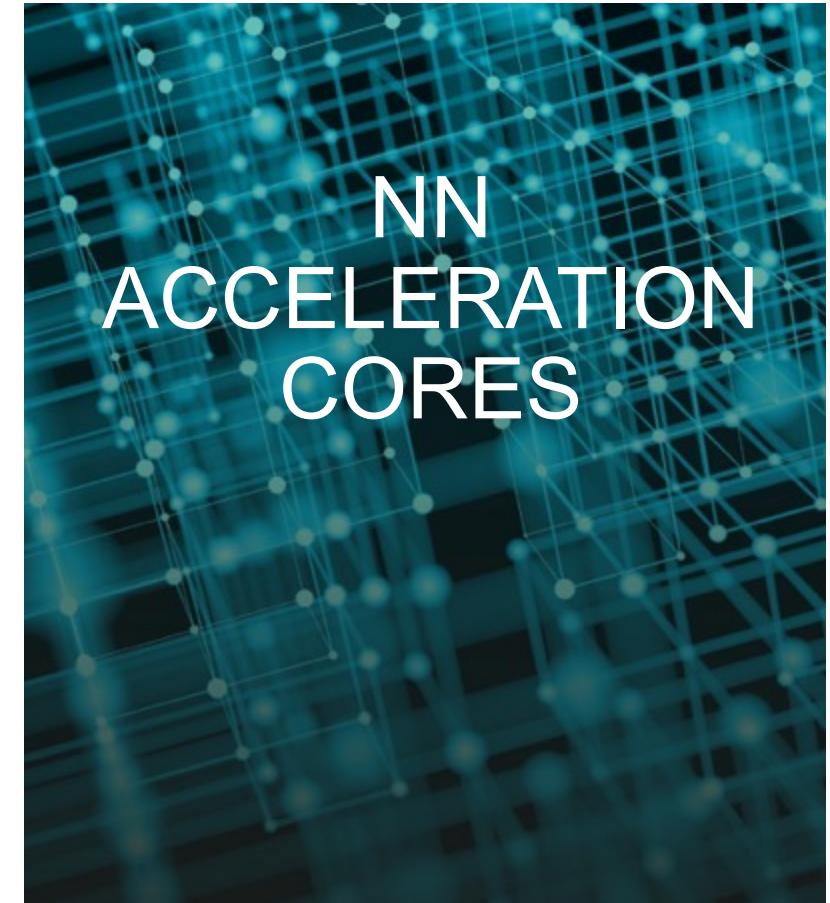
- ≡ 2018 till today
- ≡ Natural language processing engines
- ≡ Models with hundred billion parameters



CORES



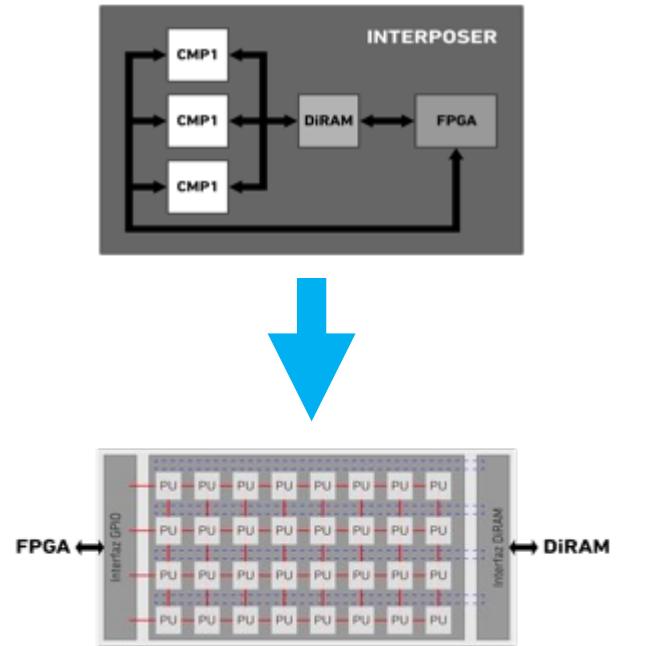
- ≡ A collection of computational cores with improved energy efficient



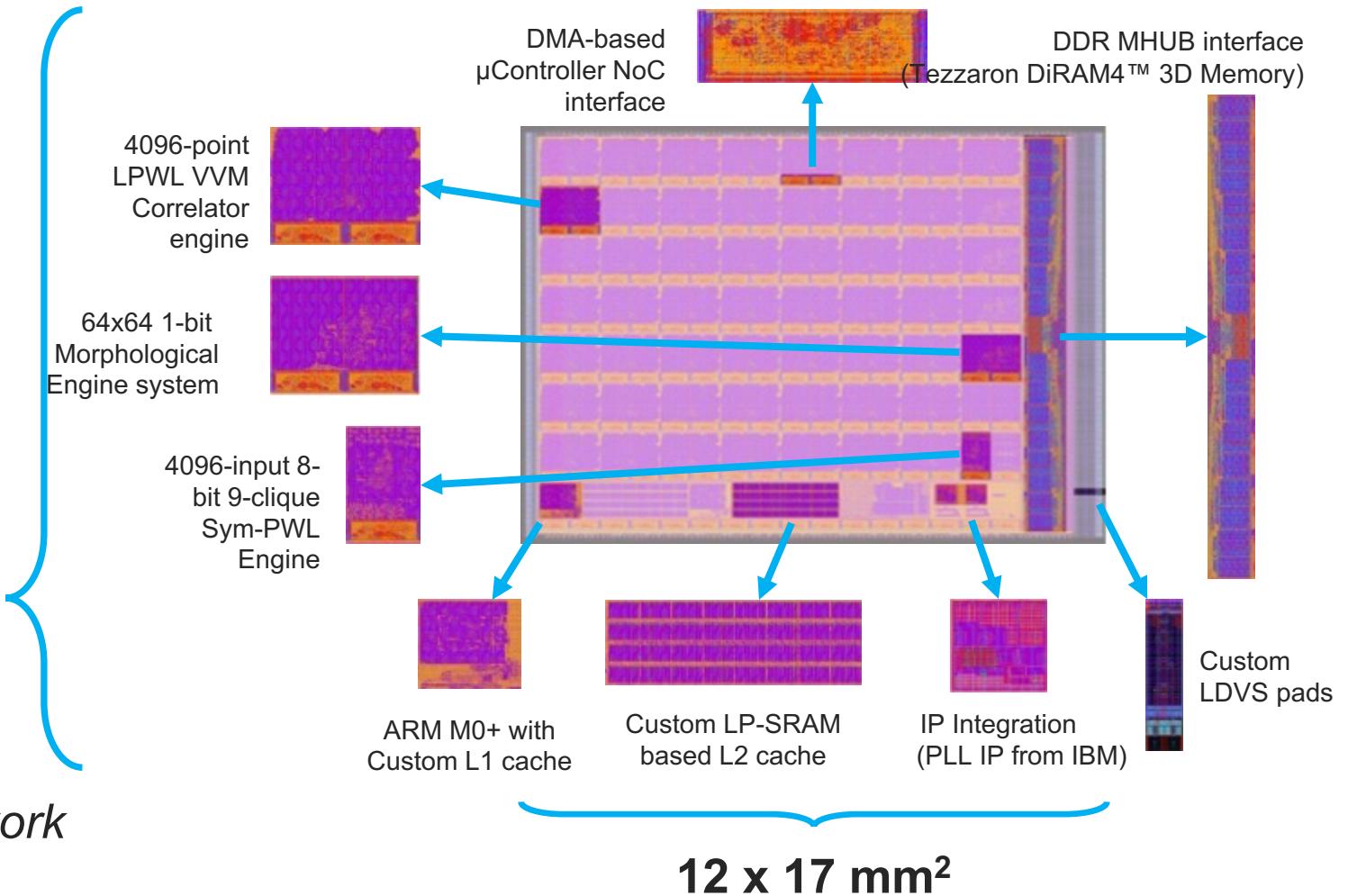


HETEROGENEOUS SALAMI CHIP

SAL
SILICON AUSTRIA LABS



L1-NoC → *token ring*
L2-NoC → *mesh network*



END

