# Memory devices for neuromorphic computing

## Fabien ALIBART

## IEMN-CNRS, Lille

Trad: Toutes les questions que je me suis posé sur le neuromorphique sans jamais (oser) les poser

Increase of Fault (nanoscale engineering)

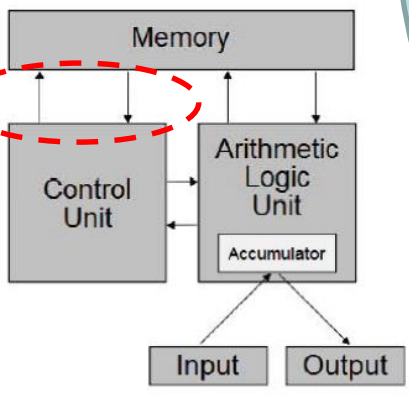Saturation of clock frequency
+
Energy consumption

New needs for computing Recognition, Mining, Synthesis (Intel)

SEMICONDUCTOR TECHNOLOGY CHALLENGES

Von Neumann bottleneck



Shift toward a new paradigm for computation

BIO-INSPIRED COMPUTING to match the brain performances (low power consumption, fault tolerant, performances for RMS)
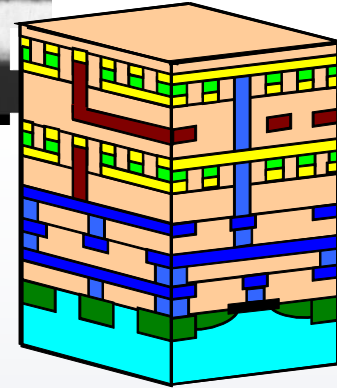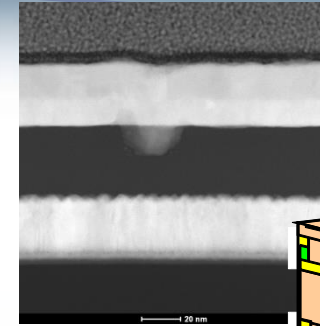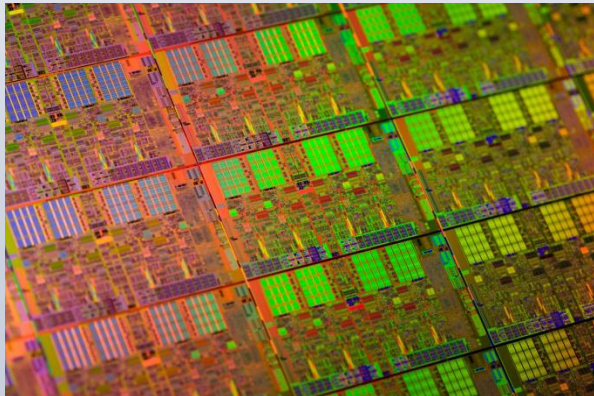
# NNET directions

Supercomputer resources
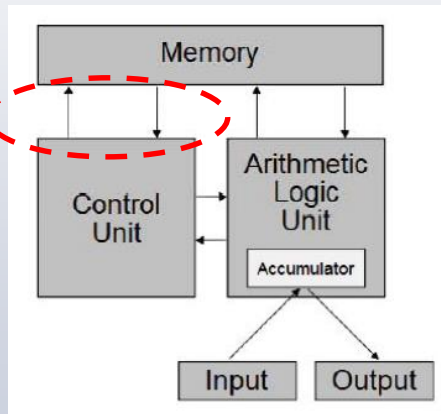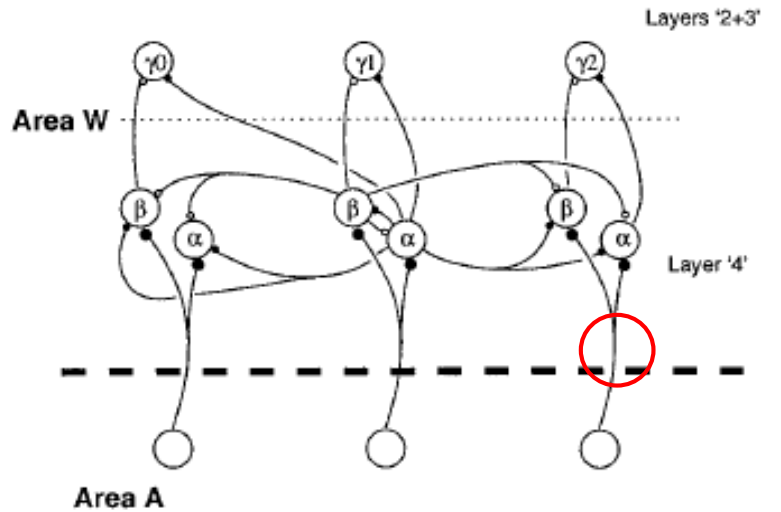Purely digital

$10^{11}$ neurons
$10^{15}$ synapses

Emerging nanotechnologies
New architecture concepts
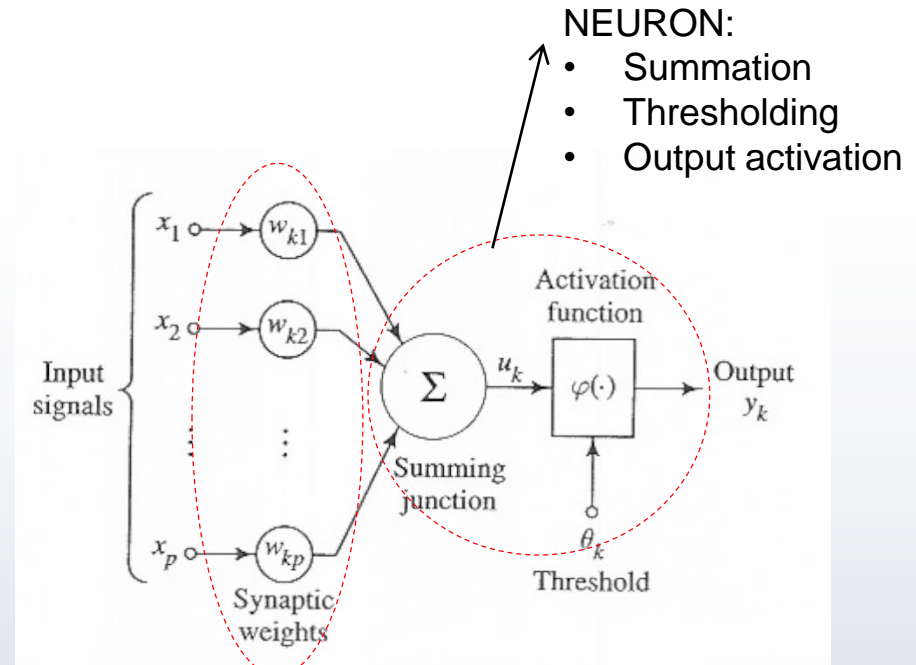and integration strategies

Custom IC, Mix analog/digital
Multichip approach,…
i.e. with conventional technologies

INSTITUT CARNOT IEMN

iemn
UMR CNRS 8520
Recherche
Formation
Transfert

# BNNs vs ANNs

**Biological neural network**





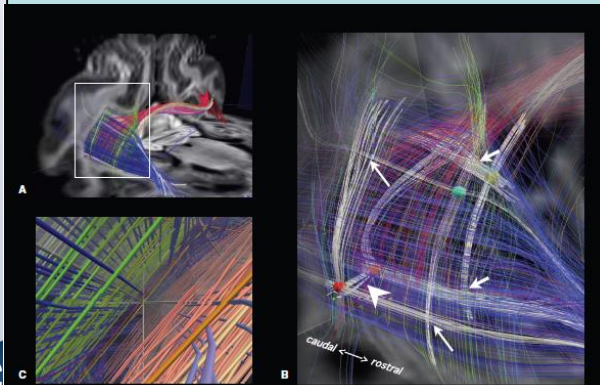**Artificial neural network**



NEURON:
- Summation
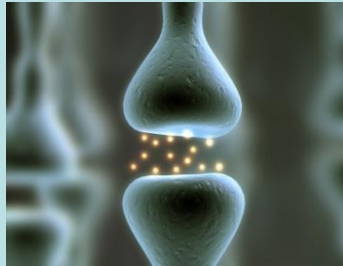- Thresholding
- Output activation
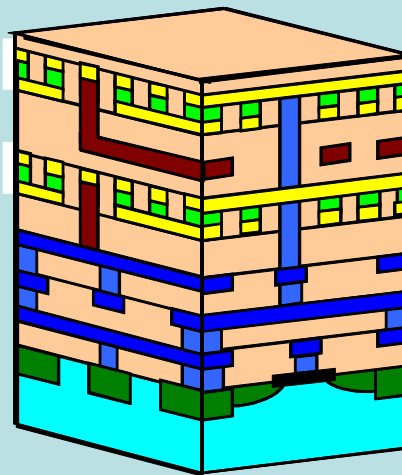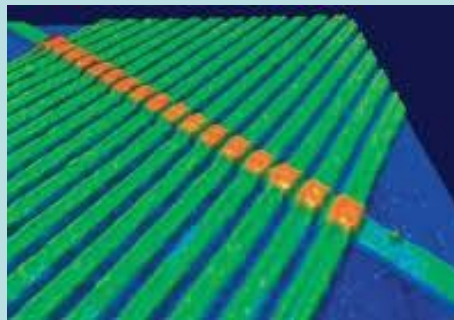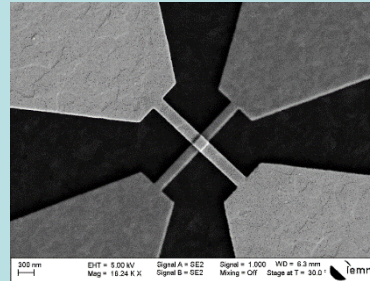
SYNAPSES:
- Input weighting
- Weight adaptation

<u>The memory is in the processing unit</u>
(Direct solution to Von Neumann bottleneck!)
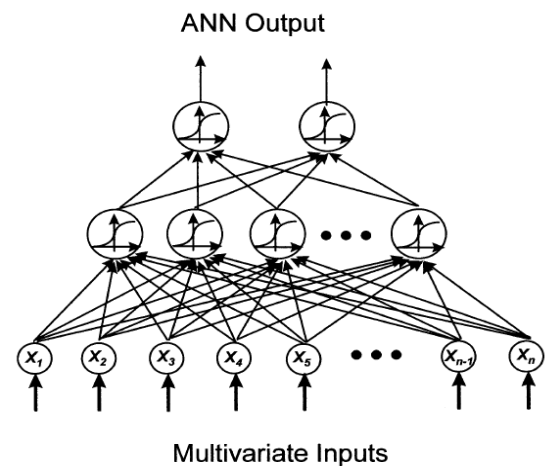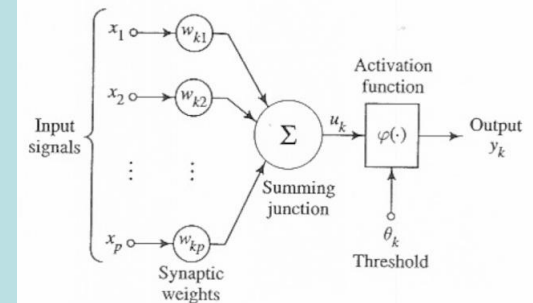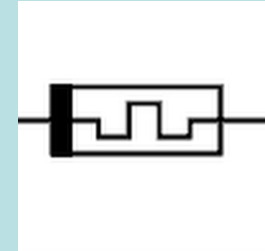
# Neuromorphic in between
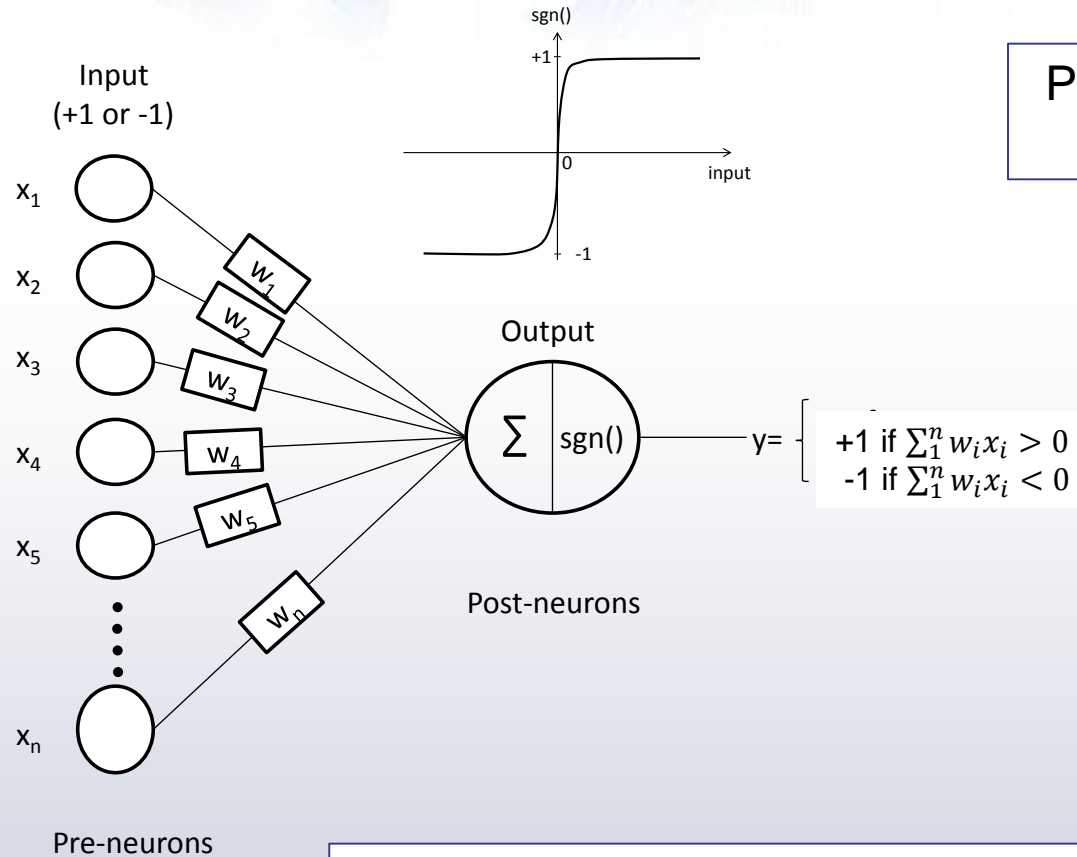
## BNNs

## Neuromorphic Eng.

## ANNs

- Intro to Artificial Neural Network (ANNs)

- Intro to biological Neural Network (BNNs)

- Nanodevices for ANNs

- Nanodevices for BNNs
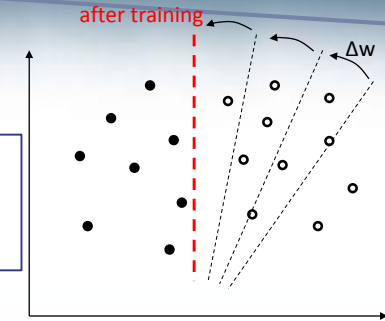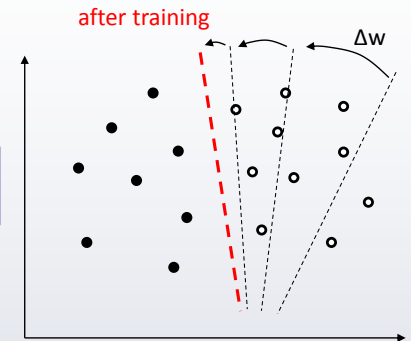
- STDP: somewhere in between

Rosenblatt, 1957

Classification of vectors (datas)

sgn()

+1

0

input

-1

Input
(+1 or -1)

$x_1$

$x_2$

$x_3$

$x_4$

$x_5$

$x_n$

$w_1$

$w_2$

$w_3$

$w_4$

$w_5$

$w_n$

Output

$\Sigma$  sgn()

$y=$
$+1$ if $\sum_1^n w_i x_i > 0$
$-1$ if $\sum_1^n w_i x_i < 0$

Post-neurons

Pre-neurons

Perceptron rule

after training
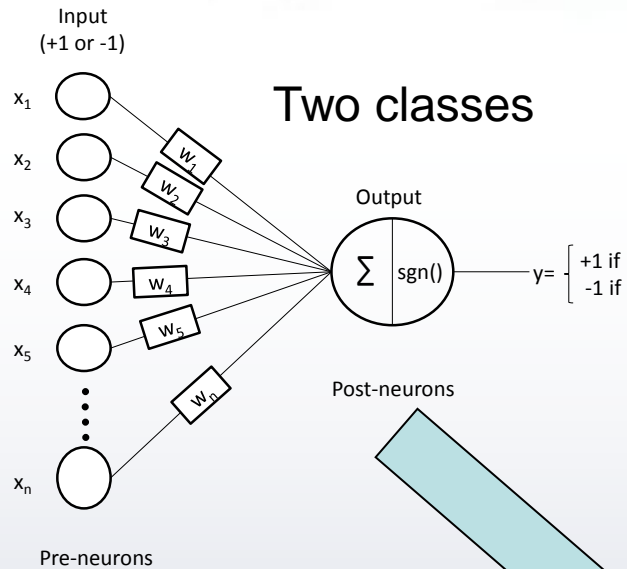
$\Delta w$

Delta rule

after training

$\Delta w$

Backpropagation

Training/learning: find the optimal weight
- No analytical solution
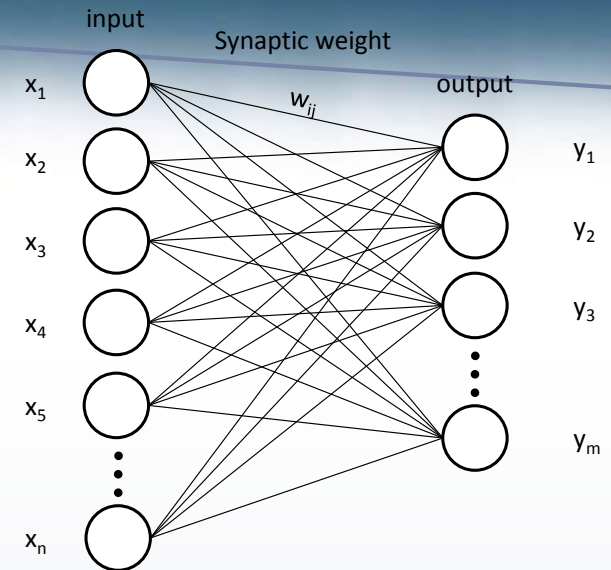- Learning algorithms: iterative correction based on known data

Practically, we can do: Pattern classification, Clustering, Prediction…

Multiple classes

Input
(+1 or -1)

Two classes

$x_1$
$x_2$
$x_3$
$x_4$
$x_5$
$x_n$

$w_1$
$w_2$
$w_3$
$w_4$
$w_5$
$w_n$

Output

$\Sigma$ | sgn()

$y= \begin{cases} +1 \text{ if} \\ -1 \text{ if} \end{cases}$

Post-neurons

Pre-neurons

Non linearly separable ensembles

input

Synaptic weight

output

$x_1$
$x_2$
$x_3$
$x_4$
$x_5$
$x_n$

$w_{ij}$

$y_1$
$y_2$
$y_3$
$y_m$

| Structure | Description of decision regions | Exclusive-OR problem | Classes with meshed regions | General region shapes |
|---|---|---|---|---|
| Single layer | Half plane bounded by hyperplane | | | |
| Two layer | Arbitrary (complexity limited by number of hidden units) | | | |
| Three layer | Arbitrary (complexity limited by number of hidden units) | | | |

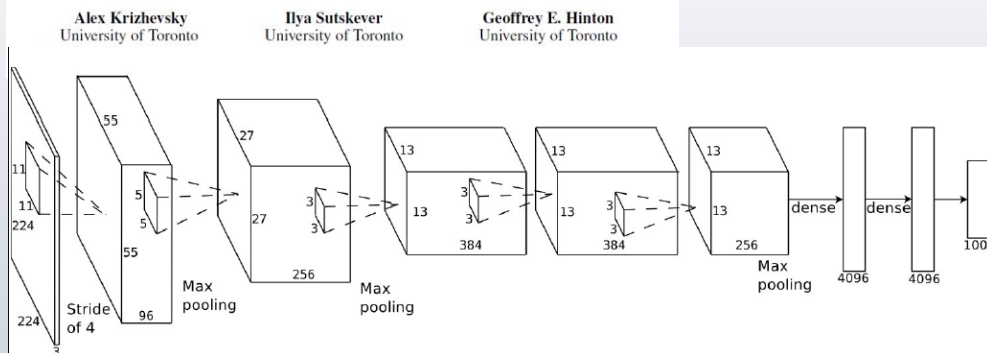**Figure 8. A geometric interpretation of the role of hidden unit in a two-dimensional input space.**

From the MINST Database of Hand–written Digits



- 28x28 pixel array
- 10 classes
- 7840 symapses

**ImageNet Classification with Deep Convolutional Neural Networks**

Alex Krizhevsky
University of Toronto

Ilya Sutskever
University of Toronto

Geoffrey E. Hinton
University of Toronto

- 256x256 pixel array
- 1000 classes



- Trained with stochastic gradient descent on two NVIDIA GPUs for about a week
- 650,000 neurons
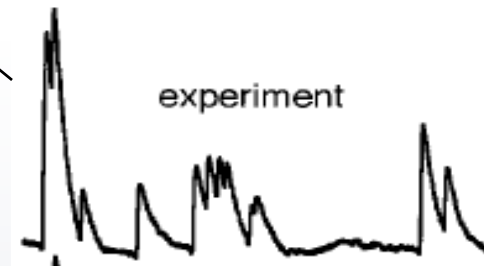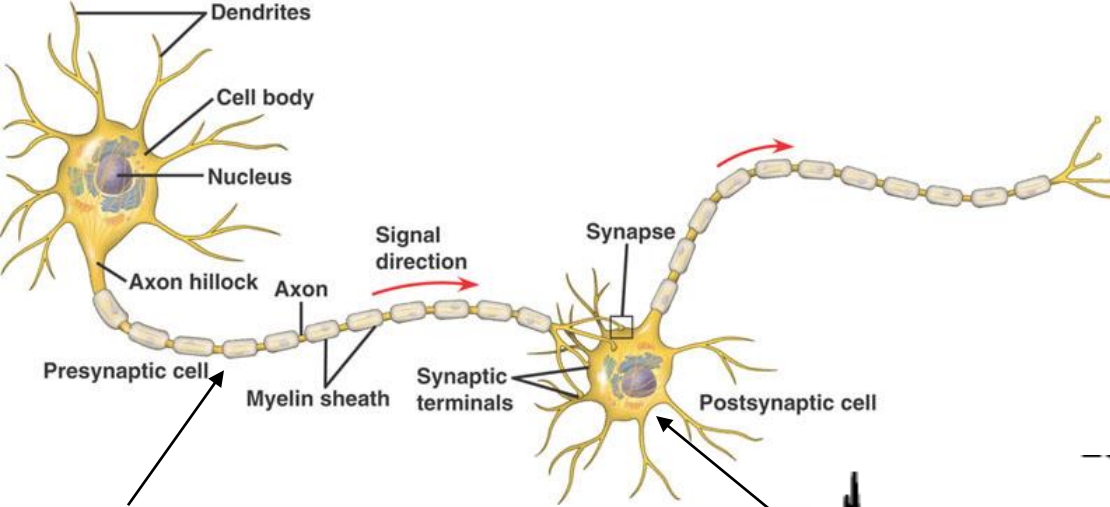- 60,000,000 parameters
- 630,000,000 connections     synapses



| mite | container ship | motor scooter | leopard |
|---|---|---|---|
| mite | container ship | motor scooter | leopard |
| black widow | lifeboat | go-kart | jaguar |
| cockroach | amphibian | moped | cheetah |
| tick | fireboat | bumper car | snow leopard |
| starfish | drilling platform | golfcart | Egyptian cat |

| grille | mushroom | cherry | Madagascar cat |
|---|---|---|---|
| convertible | agaric | dalmatian | squirrel monkey |
| grille | mushroom | grape | spider monkey |
| pickup | jelly fungus | elderberry | titi |
| beach wagon | gill fungus | ffordshire bullterrier | indri |
| fire engine | dead-man's-fingers | currant | howler monkey |

ANNs:

- Practical application will required ulra high density of nanodevices

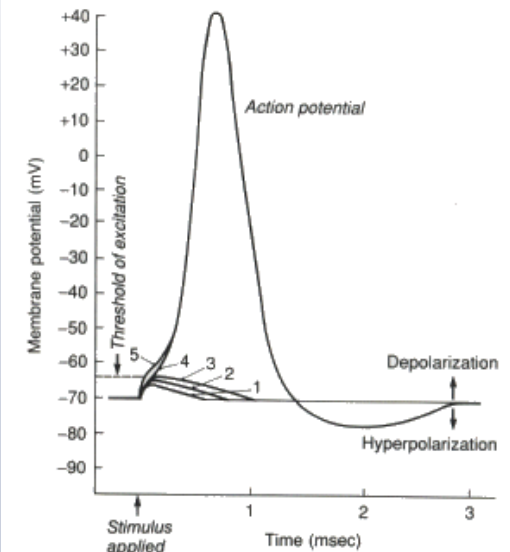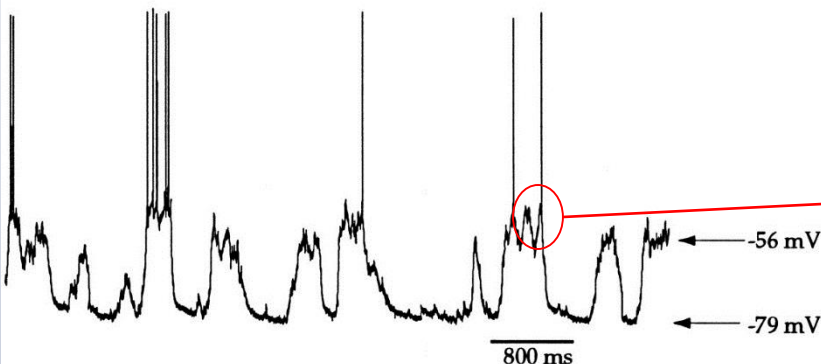- Main challenge: GPU or FPGA are serious challengers

# BNNs: neurons

The neuron membrane can be seen as a transmission line (RC) When the membrane potential reach a threshold, a spike is triggered
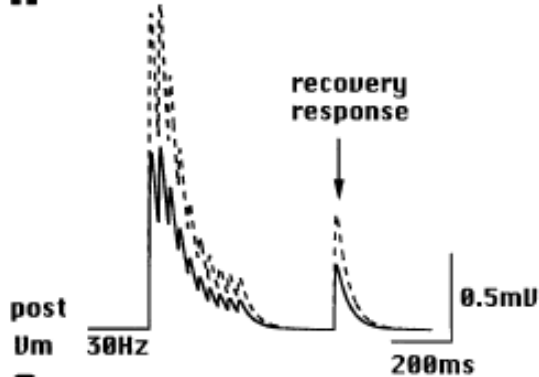
- Currents are ionic
- Low speed of propagation
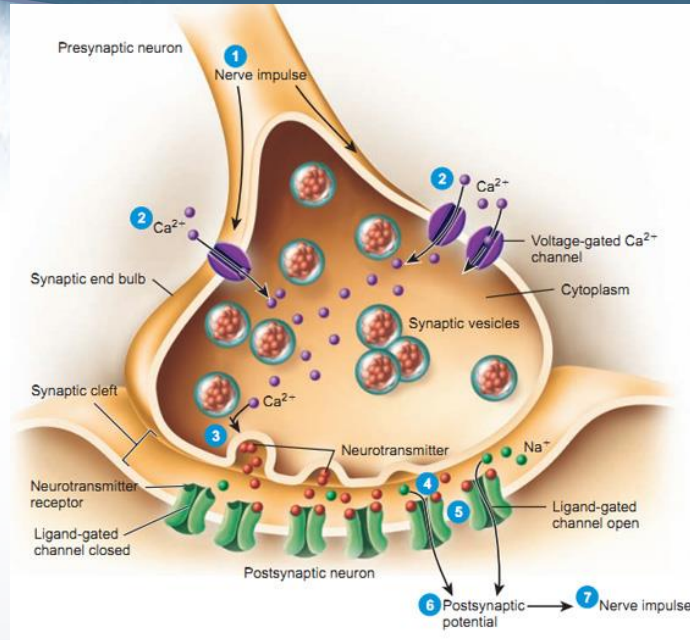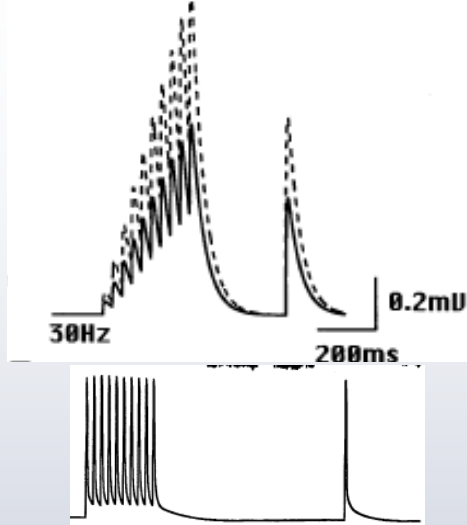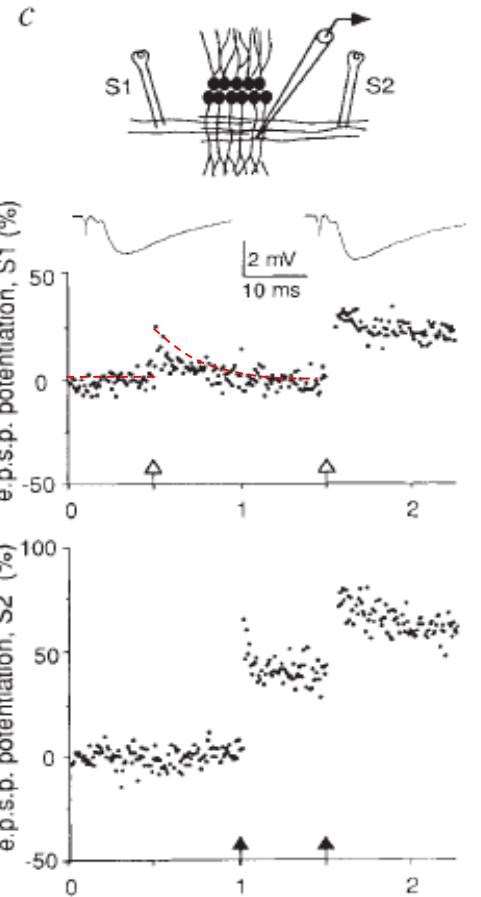- Active devices (soma and Ranvier's node)

**A** Depressing Synapses

recovery response

post Um 30Hz

0.5mV

200ms

**B** Facilitating Synapses

30Hz

0.2mV

200ms

Presynaptic neuron
- Nerve impulse
- Ca²⁺
- Voltage-gated Ca²⁺ channel
- Synaptic end bulb
- Cytoplasm
- Synaptic vesicles
- Synaptic cleft
- Ca²⁺
- Neurotransmitter
- Na⁺
- Neurotransmitter receptor
- Ligand-gated channel closed
- Ligand-gated channel open
- Postsynaptic neuron
- Postsynaptic potential
- Nerve impulse

- The AP release neurotransmitters from the pre-neuron to the post neuron receptors
- Neurotransmitters open ionic channels
- Ionic concentration change the polarization of the post-neuron

- Various time scale of plasticity
- Complex dynamics with many species (ions,NT,..)
- Unidirectionnal

C

S1  S2

2 mV
10 ms

e.p.s.p. potentiation, S1 (%)

e.p.s.p. potentiation, S2 (%)

a

12 h

5 mV

Learning in BNNs:



$a_i$ $\qquad\qquad\qquad$ $a_j$

Who fire together wire together, (Hebbs)

Synaptic learning

$$\frac{dw_{ij}}{dt} \propto a_i \cdot a_j$$

Synaptic adaptation

$$\frac{dw_{ij}}{dt} \propto a_i + a_j$$

Example, the BCM learning rule:

$$\frac{dw_{ij}}{dt} = \varphi(a_j(t)) \cdot a_i(t) - \varepsilon w_{ij}$$

$\varphi(a_j) < 0 \;\; for \;\; a_j < \theta_m \qquad \& \qquad \varphi(a_j) > 0 \;\; for \;\; a_j > \theta_m$



Synaptic Change

$\theta^{Hebb}$

$\theta^{Anti-Hebb}$

$\theta_m$

Stimulation Rate  @ fixe $a_j$



Synaptic Change

LTD $\qquad$ LTP

$\theta^{BCM}$

$\theta_m$ $\qquad\qquad$ $\theta_m$

$\theta_m$

Stimulation Rate  @ fixe $a_j$

Variation of Hebbs rule
- Unsupervised
-

BNNs:

- No charges (i.e. electrons), only ions

- Slow, with rich dynamics

- Still unsolved issues

  – Basics of computing in the brain (coding,…)

  – What do we really need for computing (for practical applictions)

The crossbar structure is the perfect architecture for massively parallel processing

It is compatible with back end process on top of a CMOS substrate

top wire level

similar two-terminal devices at each crosspoint

bottom wire level

via translation layer

crossbar layer

CMOS layer

Strukov, PNAS 2009

Development SoC

Footprint $4F^2$

$10^{12}$ devices/cm$^2$

input
Synaptic weight
$w_{ij}$
output

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$ $x_n$

$y_1$ $y_2$ $y_3$ $y_m$

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$

$y_1$ $y_2$ $y_3$ $y_4$ $y_5$

Crossnet

L= 2n² +1 states

From binary…

… to multilevel…

… to analog

Beck, 2000

Chanthbouala, 2012

(a) Adaptive programming algorithm → $V_{prog}(k)$ → PCM cell → $R(k)$ → Σ ← $R_{TARGET}$, $e(k)$

(b) Switching: LPS → RPS

Papandreou, 2011

**Implemented algorithm**

Start
(inputs: desired state $I_{desired}$, desired accuracy $A_{desired}$; initialize: write voltage to small non-disturbing value $V_{WRITE}$ = 200 mV, voltage step $T_{VSTEP}$ = 10 mV;

Read
(apply $V_{READ}$ = 200 mV and read current $I_{current}$)

Processing
Is state reached within required precision, i.e. $(I_{desired} - I_{current})/I_{desired} < A_{desired}$ ?

Processing
check for overshoot and set the sign of increment, i.e. sign = $I_{current} - I_{desired}$; if $V_{WRITE}$ !=$V_{READ}$ and sign !=oldsign then initialize $V_{WRITE}$ = 200 mV

Write
apply pulse $V_{WRITE}$

Processing
$V_{WRITE} = V_{WRITE}$ + sign * $T_{VSTEP}$ oldsign = sign

Finish

F. Alibart et al. Nanotechnology, 23 075201, 2012

# ANNs: implementations









- 50x50 pasive crossbar array
- Binary operation (multilevel, more or less)
- CBRAM technology

Strukov, 2015



12x12 Xbar array (TiO2 memristive devices)
- Online training (variation of delta rule)
- Three classes, 60 synapses

# ANNs: implementations





$$\Delta W_{ij} = \eta * X_i * \delta_j$$

"learning rate"

- Demo of multilayer perceptron (with backprop)
- 165000 PCM analog synapses

Crossbar, IEDM, 2014



- Record density 4Mb
- Binary only
- 100nm half pitch
- Fully functionnal!

- Still a huge gap betwen requirements and what is available (from the memory perspectives)

- Do we need learning (i.e. smart memory) or pure memory (i.e. storage + analog)

- Co-integration still not demonstrated

# NANO FOR BNNs

STP



$N_n$ neurotransmiter (NT) activated at spike $n$

Recovery of NT with a time $\tau_d$

$V_{IN}$

$I_n = f(N_n)$

$V_{OUT}$

charge of $N_n$ holes in the NP at spike $n$

Pentacene thin film

$I_n = f(N_n)$

Discharge of the holes with a time $\tau_d$

Gold nanoparticles

(Alibart, Adv. Funct. Mat, 2011

- Still very emerging

- No clear idea of applications but a very complex (and exciting!) field

- Can we map BNNs (ionic systems) with electronic devices (instead of ionic)
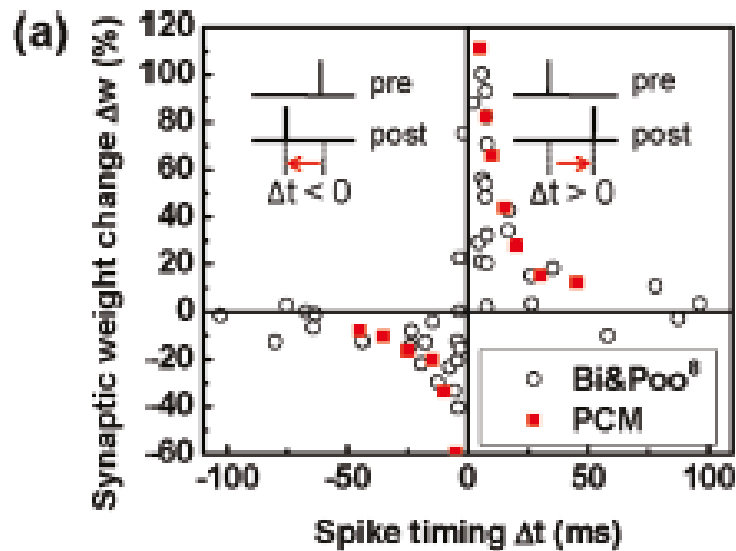
# NANO IN BETWEEN: STDP
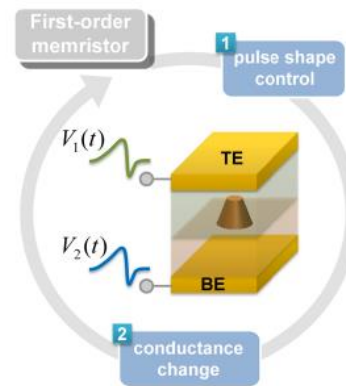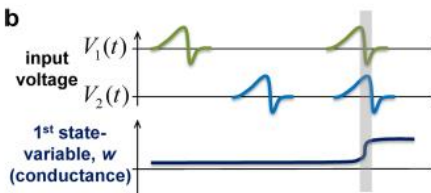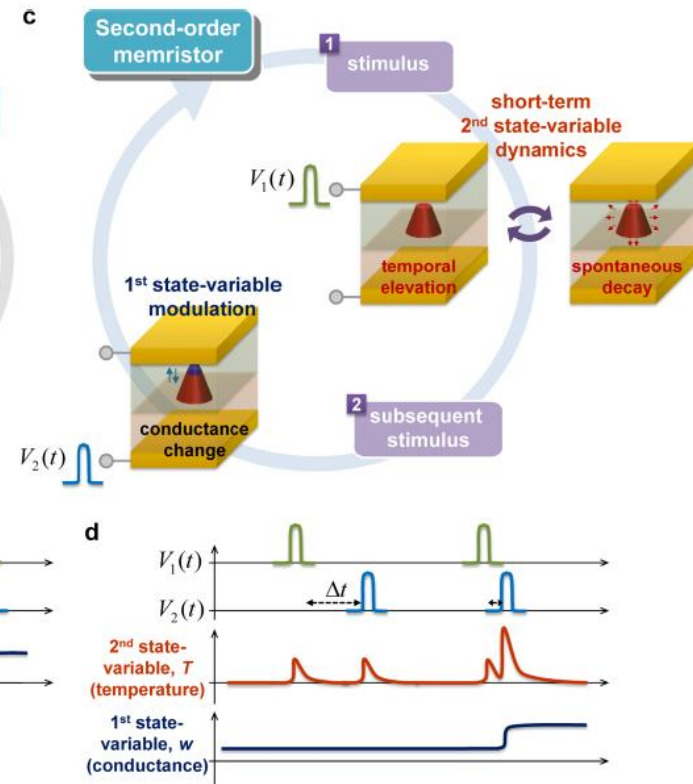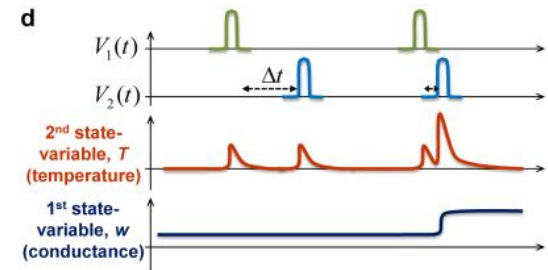
- No hardware demo at the network level

- One layer OK, what about multi-layers?

- STDP but what else?

- Neuromorphic in between BNNs and ANNs. Maybe an issue for visibility (what is our community?)

- Neuromorphic as a bridge between ANNs and BNNs?

- Doing more than identifying neuromorphic in nano, it is time to built it!

# Thank you!

(One final tip: The hippocampus is not an animal)