

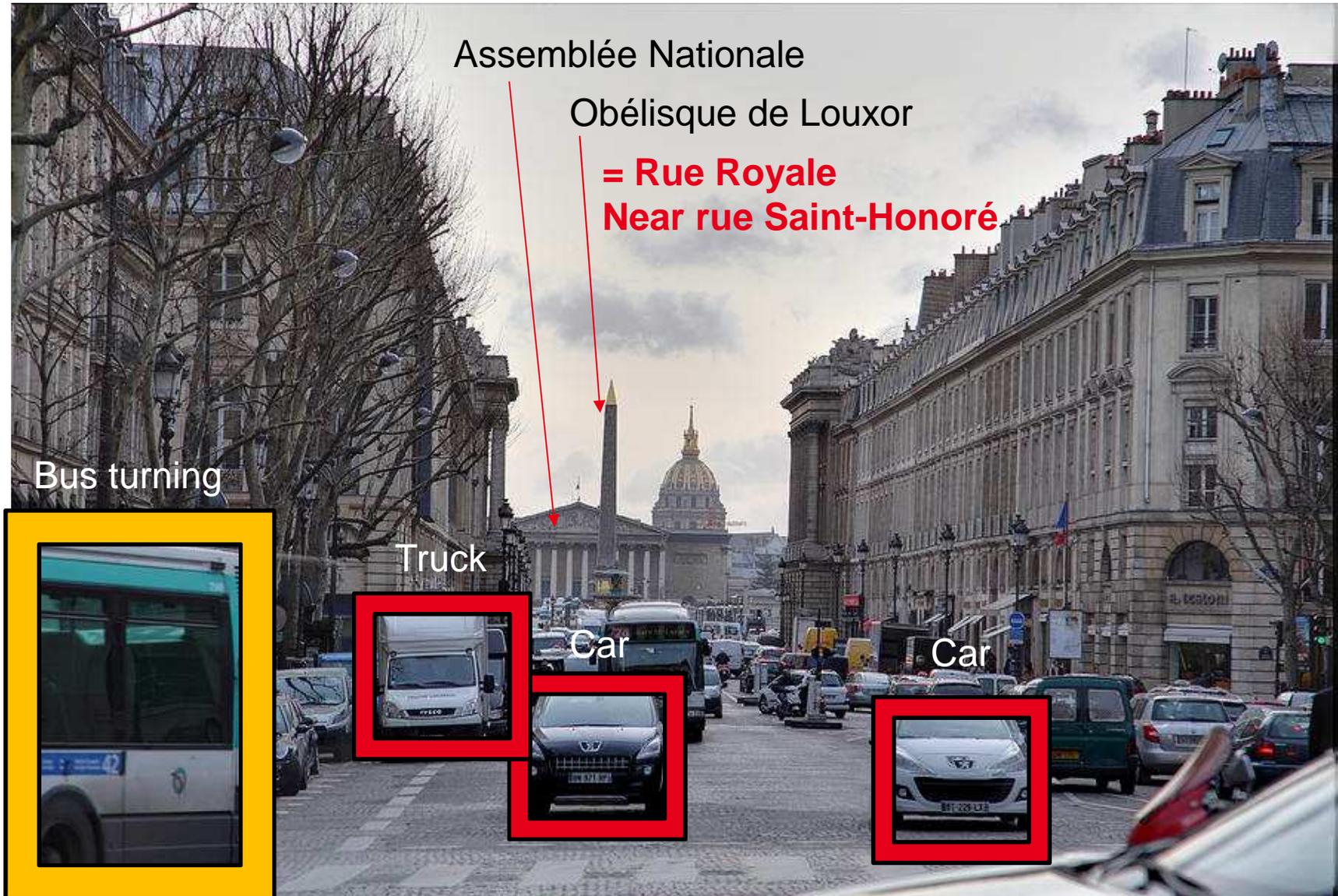


# BRAIN-INSPIRED COMPUTING FOR ADVANCED IMAGE AND PATTERN RECOGNITION

Leti Devices Workshop | Marc DURANTON | December 4, 2016



# IMAGE RECOGNITION: KEY FOR FUTURE APPLICATIONS



# ImageNet: Classification

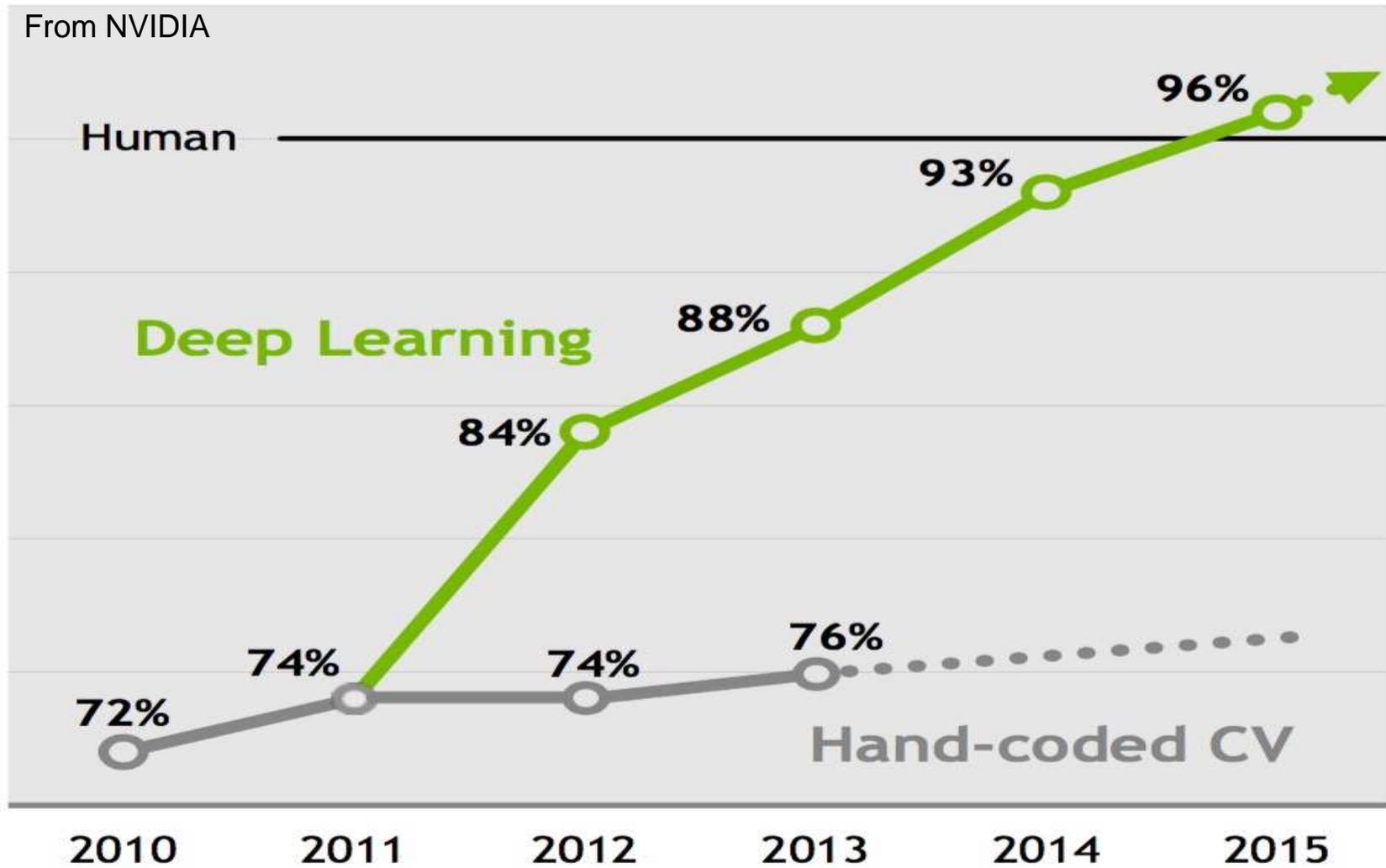
Y LeCun

- Give the name of the dominant object in the image
- Top-5 error rates: if correct class is not in top 5, count as error
  - ▶ Black: ConvNet, Purple: no ConvNet

2012 Teams	%error	2013 Teams	%error	2014 Teams	%error
Supervision (Toronto)	15.3	Clarifai (NYU spinoff)	11.7	GoogLeNet	6.6
ISI (Tokyo)	26.1	NUS (singapore)	12.9	VGG (Oxford)	7.3
VGG (Oxford)	26.9	Zeiler-Fergus (NYU)	13.5	MSRA	8.0
XRCE/INRIA	27.0	A. Howard	13.5	A. Howard	8.1
UvA (Amsterdam)	29.6	OverFeat (NYU)	14.1	DeeperVision	9.5
INRIA/LEAR	33.4	UvA (Amsterdam)	14.2	NUS-BST	9.7
		Adobe	15.2	TTIC-ECP	10.2
		VGG (Oxford)	15.2	XYZ	11.2
		VGG (Oxford)	23.0	UvA	12.1



# COMPETITION ON IMAGENET: SINCE 2012, CONVOLUTIONAL NEURAL NETWORKS (CNN) ARE LEADING!





# Deep Learning is Everywhere (ConvNets are Everywhere)

Y LeCun

## ■ Lots of applications at Facebook, Google, Microsoft, Baidu, Twitter, IBM...

- ▶ Image recognition for photo collection search
- ▶ Image/Video Content filtering: spam, nudity, violence.
- ▶ Search, Newsfeed ranking

## ■ People upload 800 million photos on Facebook every day

- ▶ (2 billion photos per day if we count Instagram, Messenger and Whatsapp)

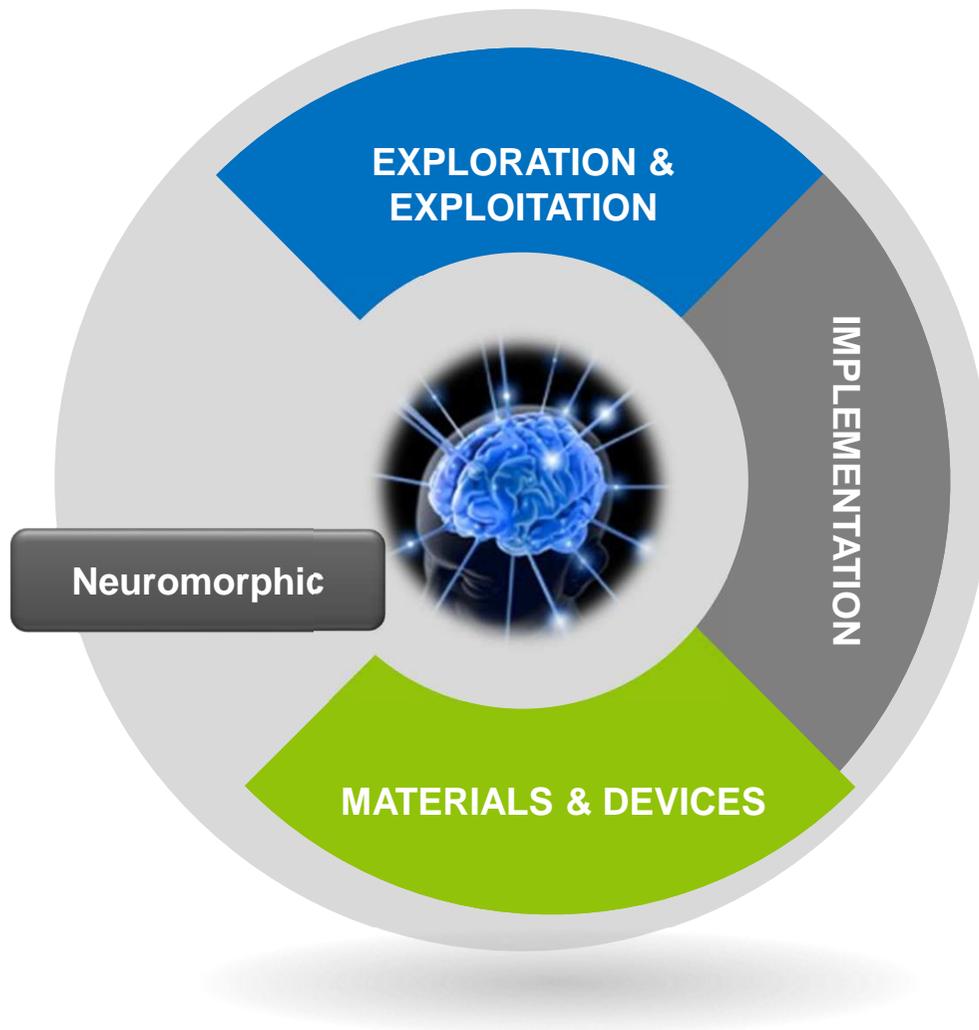
## ■ Each photo on Facebook goes through two ConvNets within 2 seconds

- ▶ One for image recognition/tagging
- ▶ One for face recognition (not activated in Europe).

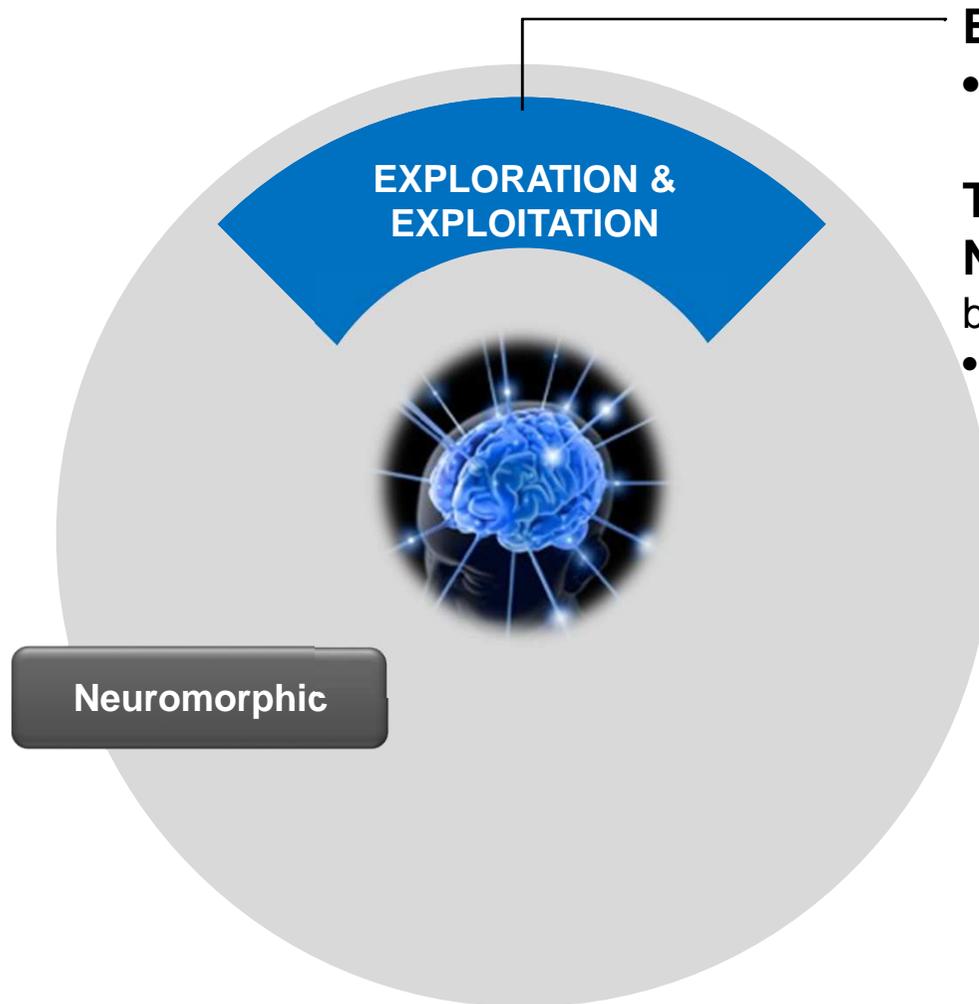
## ■ Soon ConvNets will really be everywhere:

- ▶ self-driving cars, medical imaging, augmented reality, mobile devices, smart cameras, robots, toys.....

# DEEP LEARNING AND NEUROMORPHIC SYSTEMS AT LETI AND LIST



# DEEP LEARNING AND NEUROMORPHIC SYSTEMS AT LETI AND LIST



## Exploitation of Deep Neural Networks

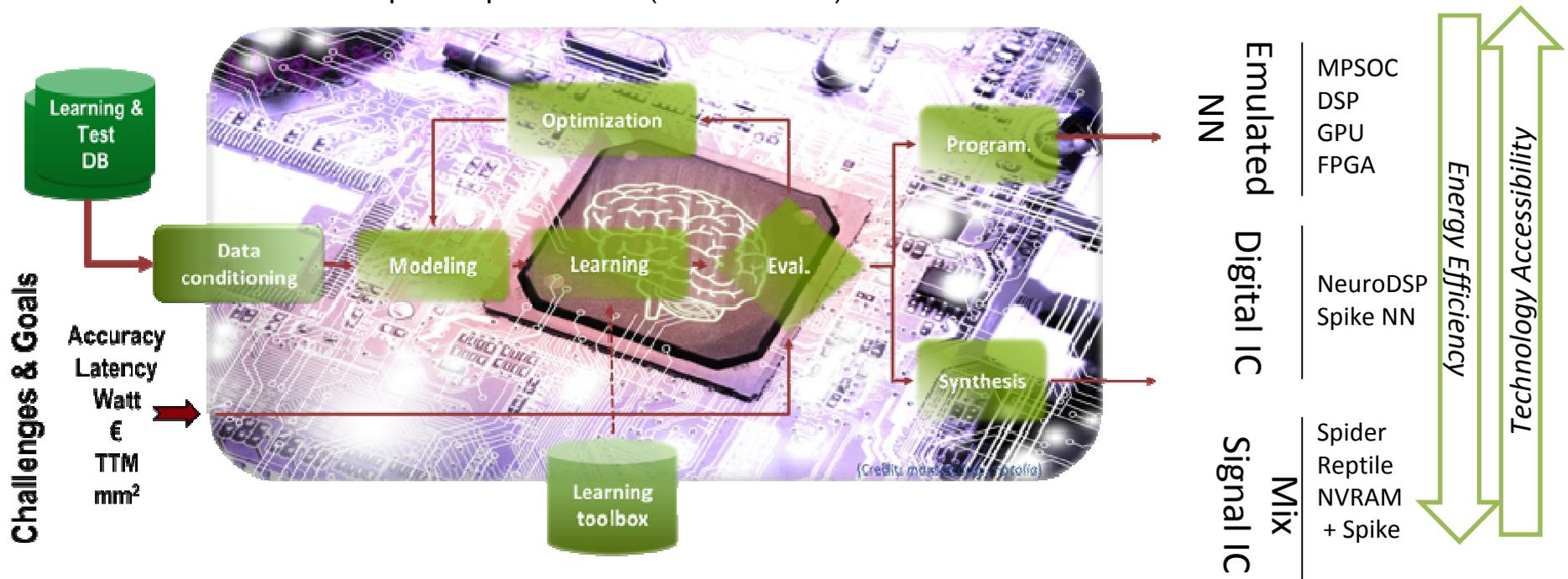
- Image recognition, annotation and indexing

**Tools for fast and accurate Neural Network (NN) exploration & Architecture benchmarking: *N2D2***

- Neural Network exploration (including with spike coding and new materials)

## N2D2: PLATFORM FOR DEVELOPING DEEP NEURAL NETWORK APPLICATIONS

- **N2D2** is a **platform to design and generate deep neural network (DNN)** and to select the computing platform which fit best application needs
- Fast benchmarking of Components Off the Shelf and exports to dedicated ASIC:
  - Parallel processors (OpenCL, OpenMP)
  - GPU (OpenCL, Cuda, CuDNN)
  - FPGA (RTL, HLS)
  - Leti & List specific processors (like **P-Neuro**)



# FAST AND ACCURATE NN EXPLORATION

## Automated architecture mapping and benchmarking tool flow

### 1) Deep network builder

```

; Environment
[env]
SizeX=8
SizeY=8
ConfigSection=env.config
Type=Pool
PoolWidth=2
PoolHeight=2
NbChannels=32
Stride=2

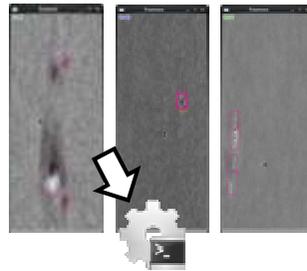
; First layer (convolution)
[conv1]
Input=env
Type=Conv
KernelWidth=3
KernelHeight=3
NbChannels=32
Stride=1

; Second layer (pooling)
[pool1]
Input=conv1
Type=Pool
PoolWidth=2
PoolHeight=2
NbChannels=32
Stride=2

; Third layer (fully connected)
[fc1]
Input=conv2
Type=Fc
NbChannels=1000

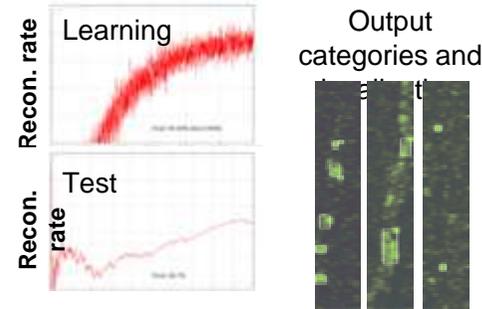
```

### 2) Learning a database



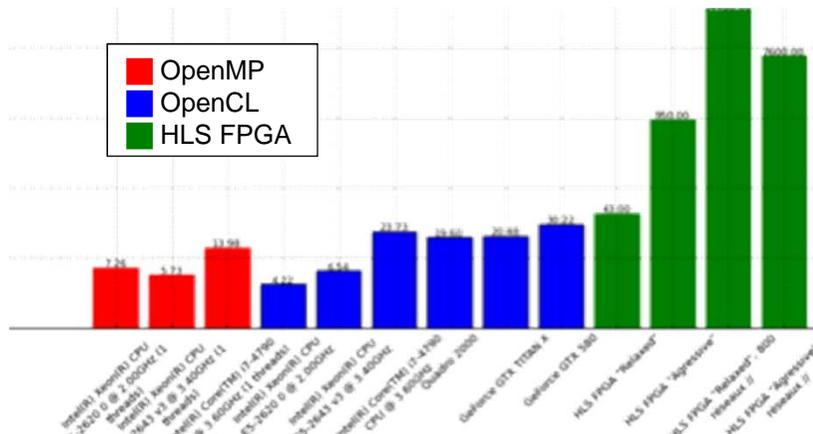
N2D2 software framework

### 3) Analysis of network performances



↓ Inference phase

### 4) CPU, GPU and FPGA-based real-time implementation



→ Wide targets range, perfs and power metrics

```

signal mrv_val1435_fu_08752_p0 : STD_LOGIC_VECTOR (7 downto 0);
signal mrv_val1436_fu_08753_p0 : STD_LOGIC_VECTOR (7 downto 0);
signal mrv_val1437_fu_08754_p0 : STD_LOGIC_VECTOR (7 downto 0);
signal mrv_val1438_fu_08755_p0 : STD_LOGIC_VECTOR (7 downto 0);
signal ap_M0_Fork : STD_LOGIC_VECTOR (7 downto 0);
begin
-- The current state lag_C5_fork of the state machine, --
ap_C5_fork_assign_proc : process(ap_clk)
begin
if ap_clk'event and ap_clk = '1' then
if lag_rst = '1' then
ap_C5_fork = ap_S7_S13_Fork_0;
else
ap_C5_fork = ap_M0_Fork;
end if;
end if;
end process;

-- ap_done_reg assign process --
ap_done_reg_assign_proc : process(ap_clk)
begin
if ap_clk'event and ap_clk = '1' then
if lag_rst = '1' then
ap_done_reg = ap_const_logic_0;
else
if (lag_const_logic_1 = ap_continue) then
ap_done_reg = ap_const_logic_1;
else
ap_done_reg = ap_const_logic_0;
end if;
end if;
end process;

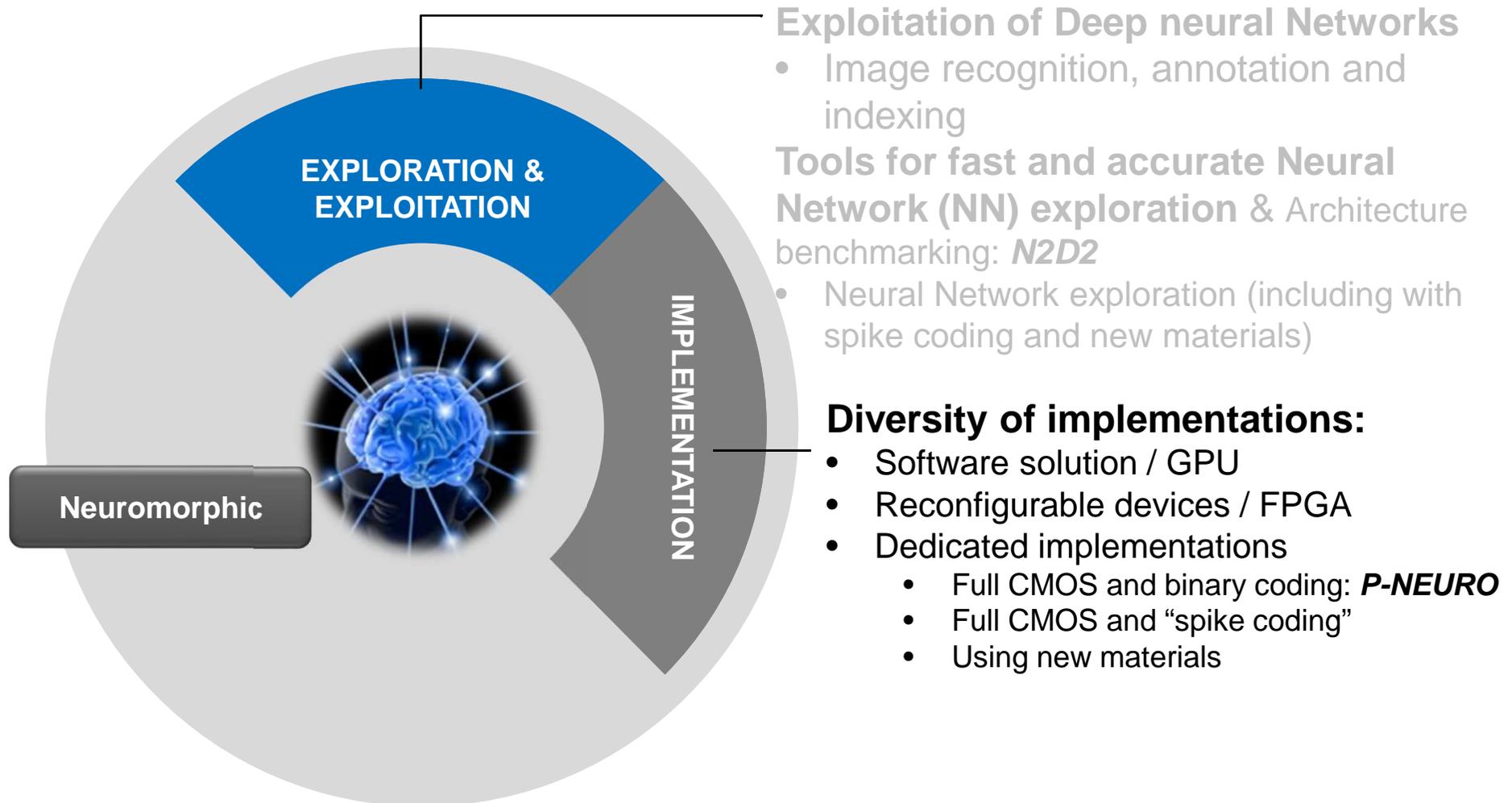
-- ap_reg_posten_p0p_110 assign process --
ap_reg_posten_p0p_110_assign_proc : process(ap_clk)

```





# DEEP LEARNING AND NEUROMORPHIC SYSTEMS AT LETI AND LIST

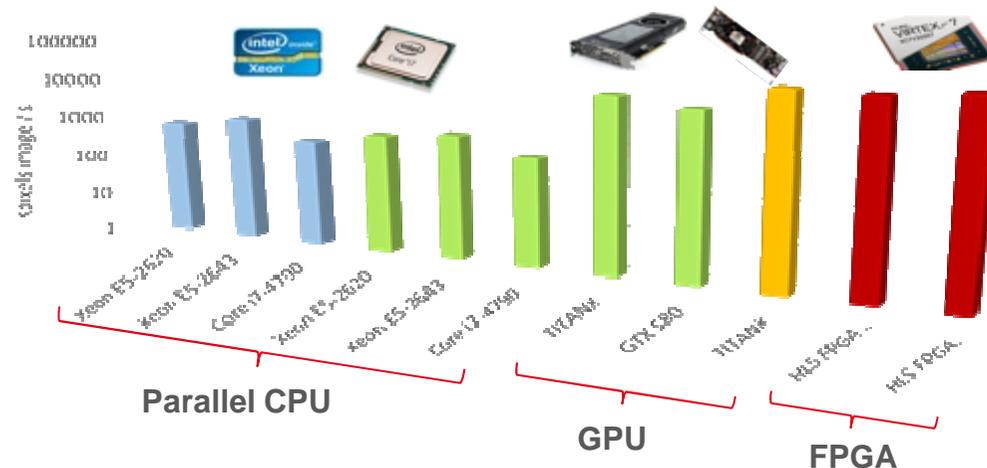


# N2D2 and P-Neuro: complete solution for Deep Learning in smart nodes

## Fast benchmarking of Components Off The Shelf:

- Parallel processors
- GPU
- FPGA (HLS)

- OpenMP
- OpenCL
- CUDA
- HLS FPGA



## Performance of *P-Neuro* neural network processing unit

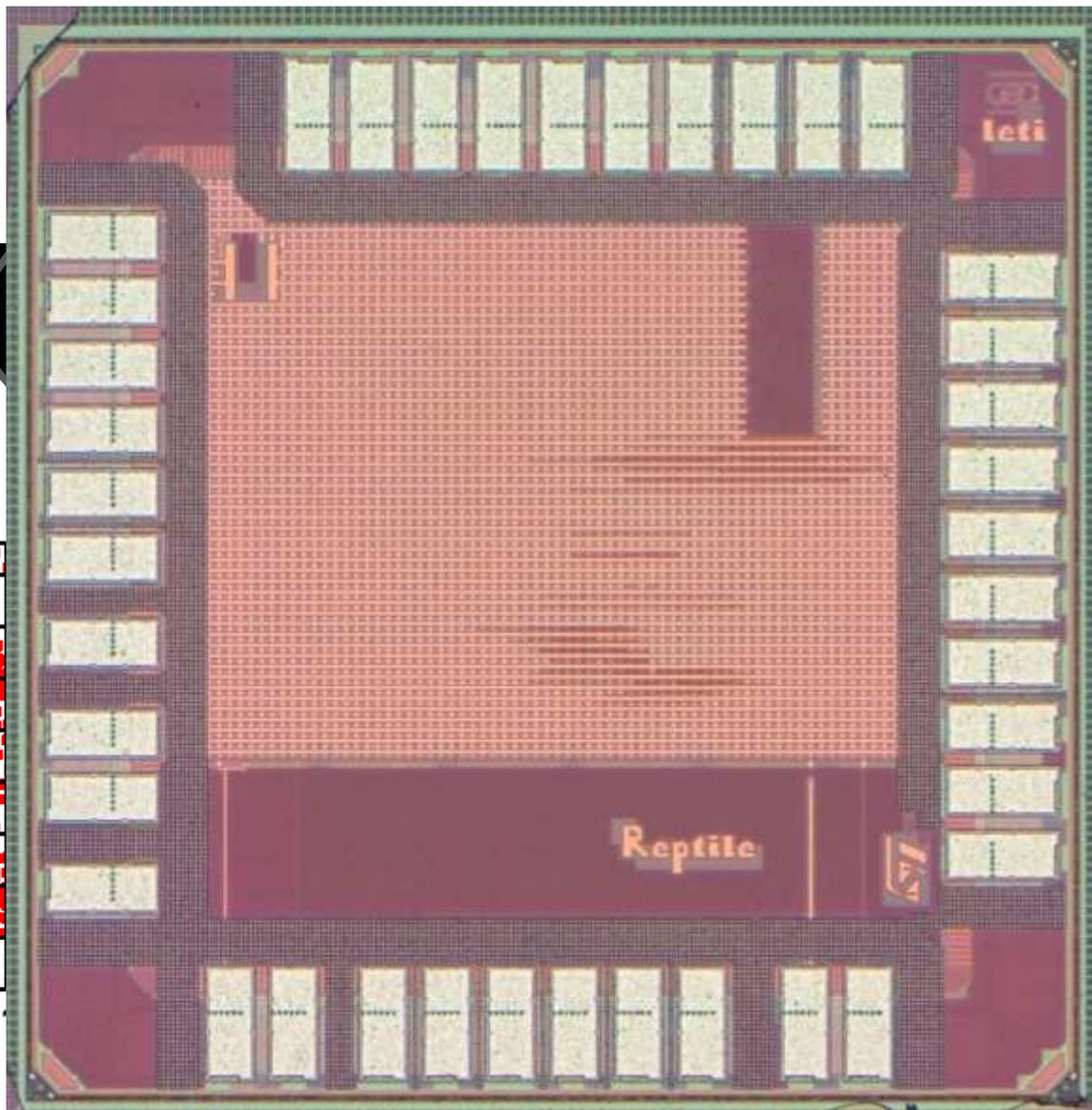
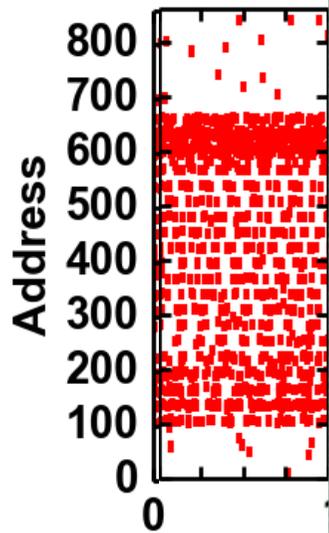
- Example on Faces extraction,
  - Database of 18000 images
- Comparison of 5 different architectures
- Focus on energy efficiency
- Expected performance of *P-Neuro*:
  - FDSOI 28nm, 1GHz
  - 1.8 TOPs/W, <0.5 mm<sup>2</sup> (4 cores)
  - Fully scalable from 1 to 1024 cores
  - Ready for integration in smart nodes

Target	Frequency	Energy efficiency
Quad ARM A7	900 MHz	380 images/W
Quad ARM A15	2000 MHz	350 images/W
Tegra K1	850 MHz	600 images/W
Intel I7	3400 MHz	160 images/W
<b>P-Neuro (FPGA)</b>	<b>100 MHz</b>	<b>2 000 images/W</b>
<b>P-Neuro (ASIC)</b>	<b>500 MHz</b>	<b>125 000 images/W</b>



# SPIKE-BASED CODING

29x29 pixels  
841 addresses



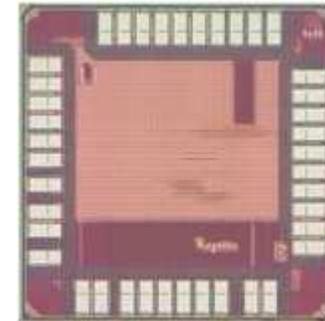
 **Correct Output**



# THE PROMISES OF SPIKE-CODING NN

- Reduced computing complexity and natural temporal and spatial parallelism
- Simple and efficient performance tunability capabilities
- Spiking NN best exploit NVMs such as RRAM, for massively parallel synaptic memory

	Formal neurons	Spiking neurons
Base operation	- Multiply-Accumulate (MAC)	+ Accumulate only
Activation function	- Non-linear function	+ Simple threshold
Parallelism	- Spatial multiplexing	+ Spatial and temporal multiplexing

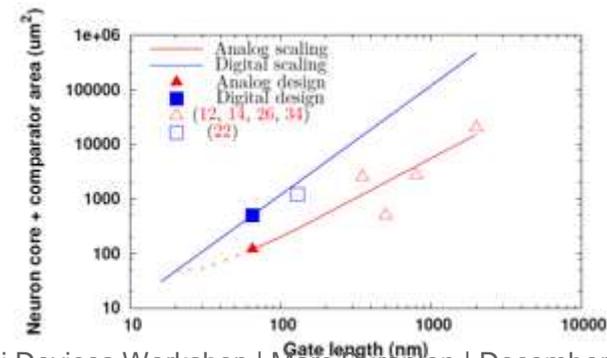
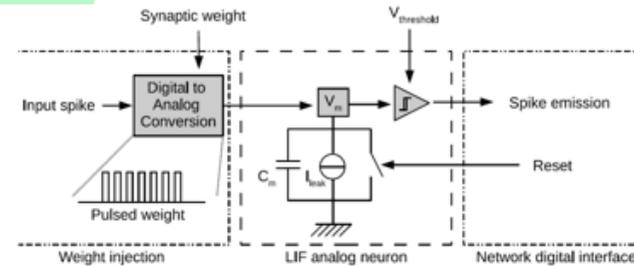


## Two test chips implemented in 65nm

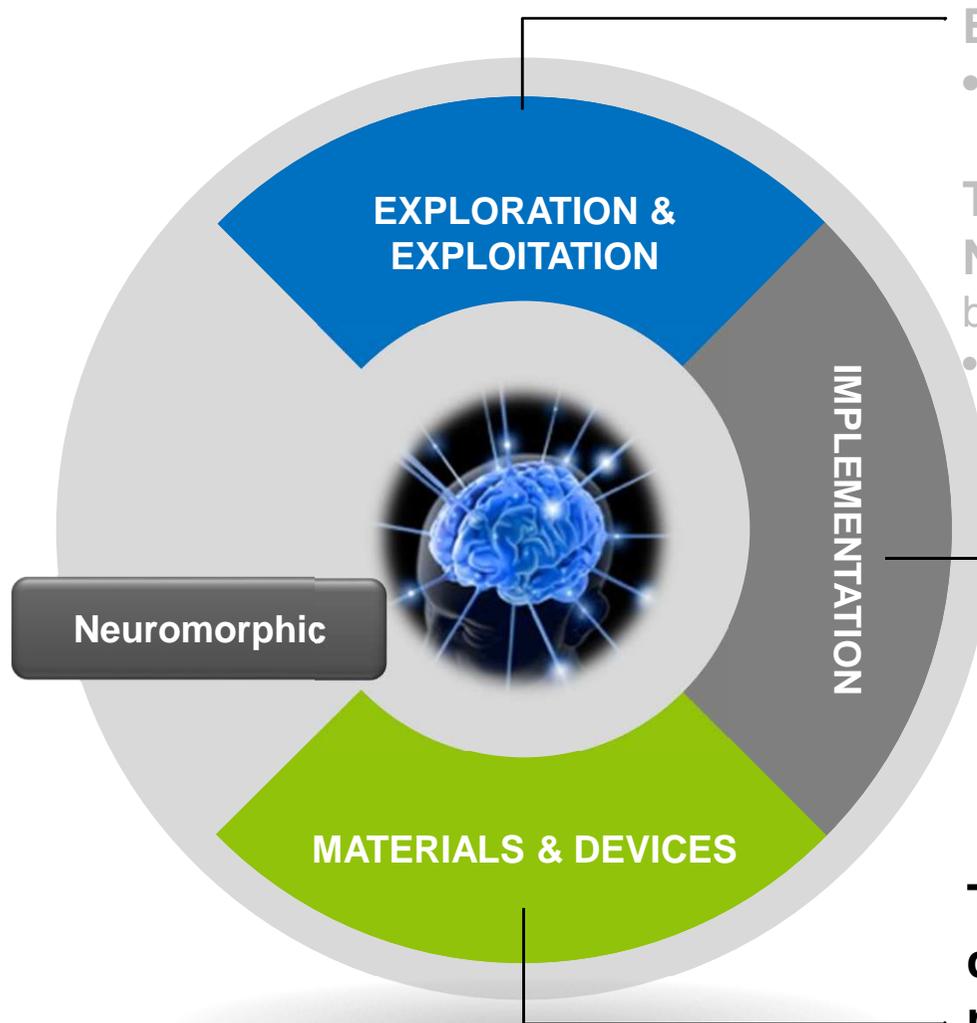
- Reptile: 3 tiles of 12 neurons
- Spider: 25 tiles of 12 neurons

## Advanced technology nodes

- Comparison of Analog and Digital neurons
- Gain of Analog neuron (less area) reduces → Curves cross at 22nm node



# DEEP LEARNING AND NEUROMORPHIC SYSTEMS AT LETI AND LIST



## Exploitation of Deep neural Networks

- Image recognition, annotation and indexing

## Tools for fast and accurate Neural Network (NN) exploration & Architecture benchmarking: *N2D2*

- Neural Network exploration (including with spike coding and new materials)

## Diversity of implementations:

- Software solution / GPU
- Reconfigurable devices / FPGA
- Dedicated implementations
  - Full CMOS and binary coding: *P-NEURO*
  - Full CMOS and “spike coding”
  - Using new materials

## Take full advantage of advanced devices to break the density and power issues:

- 3D integration, CoolCube™.
- RRAM, PCM and new devices,

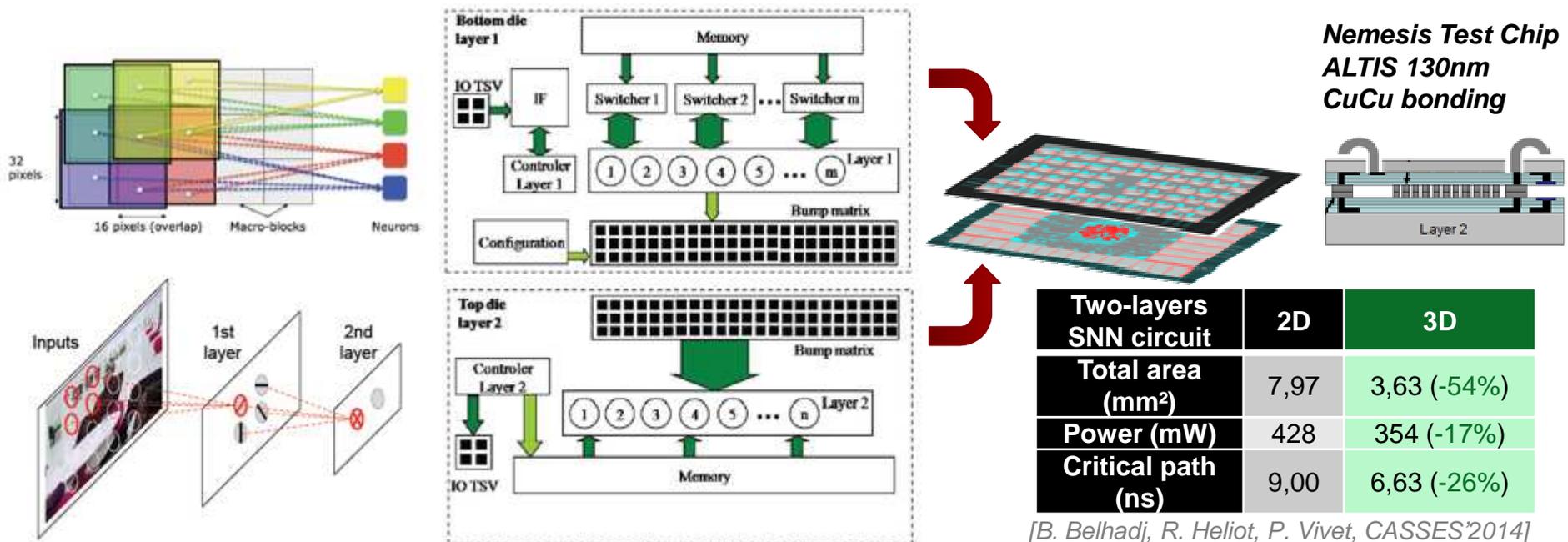
# 3D SPIKING NEURAL NETWORK

## Neural Networks

- Naturally 3D for 2D inputs, layers optimally distributed in stacked dies
- Vertical connections between layers: minimizes interconnect length, avoid routing congestion

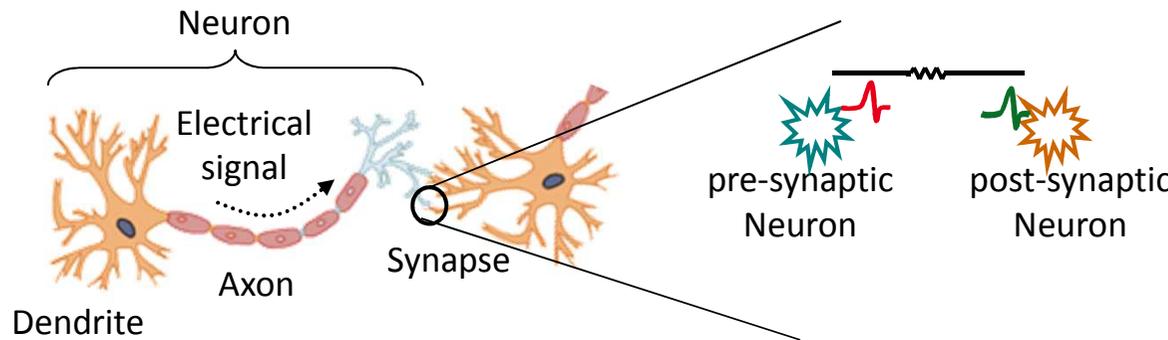
## NEMESIS 3D two-layers SNN test chip

- 1<sup>st</sup> layer: 48 macro-block neurons, 1024 synapses per neuron (49 152 total)
- 2<sup>nd</sup> layer: 50 fully connected neurons, 2 400 synapses

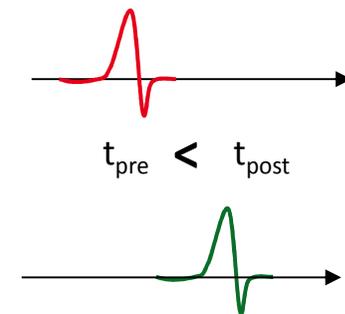
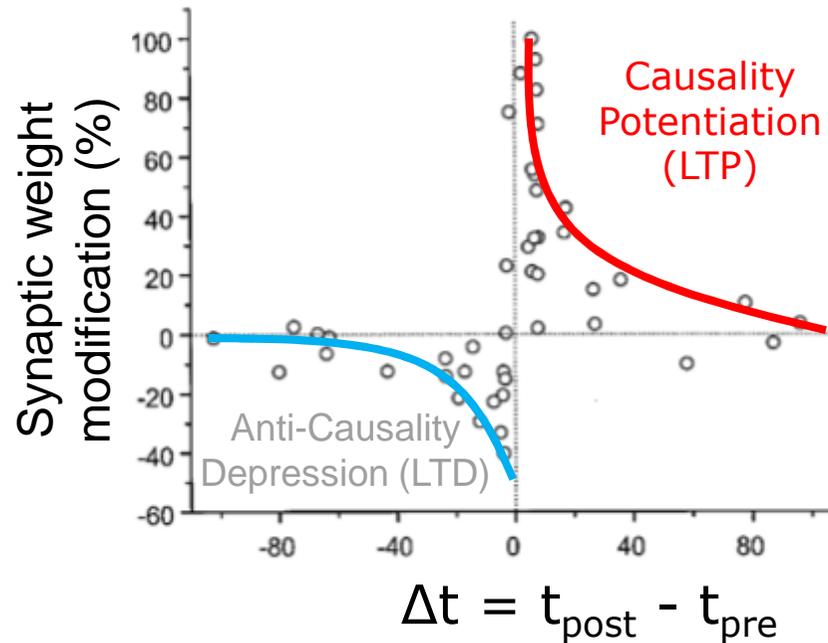
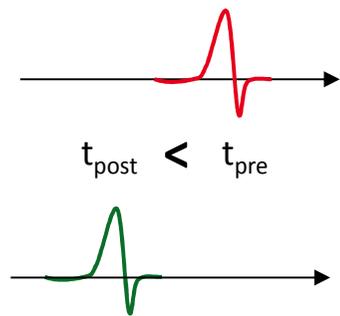


**→ 3D offers 2x better total area and 25% better power efficiency vs 2D**

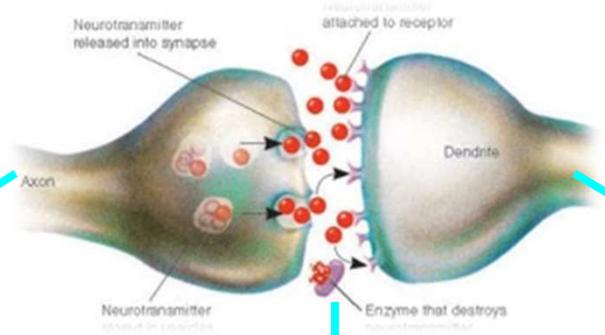
# LEARNING FROM NEUROSCIENCE: A STDP (SPIKE TIMING DEPENDENT PLASTICITY) PRIMER



STDP = correlation detector  
 → Possible learning model of the brain?



# NEW ELEMENT: RRAM AS SYNAPSES



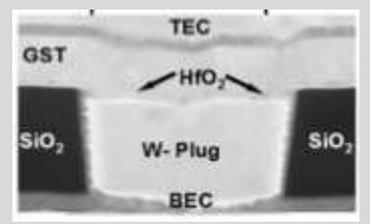
**Thermal effect**

**Electrochemical effect**

**Electronic effect  
oxygen vacancies**

## PCM

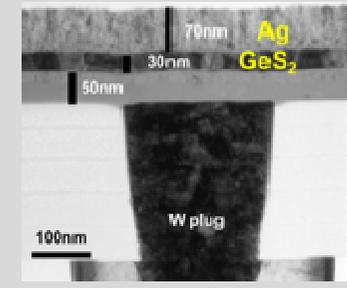
**GST  
GeTe  
GST + HfO<sub>2</sub>**



*M.Suri, et. al, IEDM 2011*  
*M.Suri, et. al, IMW 2012 , JAP 2012*  
*O.Bichler et al. IEEE TED 2012*  
*M.Suri et al., EPCOS 2013*  
*D.Garbin et al., IEEE Nano 2013*

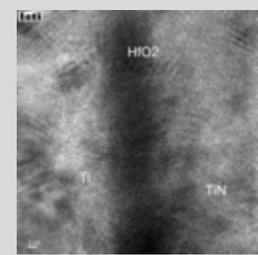
## CBRAM

**Ag / GeS<sub>2</sub>**



## OxRAM

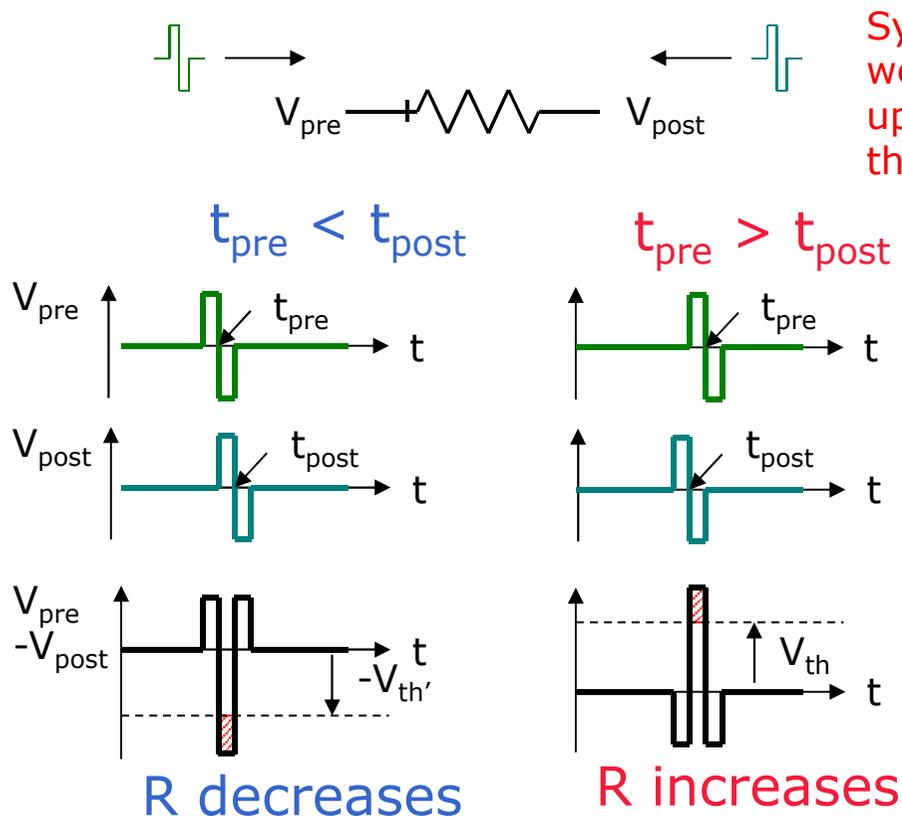
**TiN/HfO<sub>2</sub>/Ti/TiN**



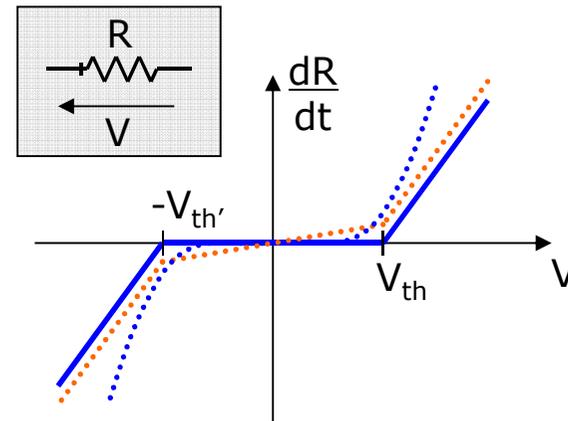
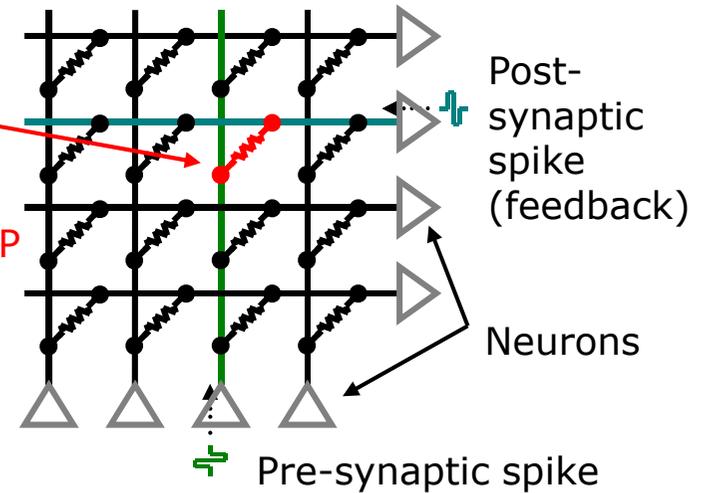
*D.Garbin et al. IEDM 2014*  
*D.Garbin et al., IEEE TED 2015*

# PRINCIPLE CROSSBARS OF MEMRISTORS

## First Proposed by Snider(1)

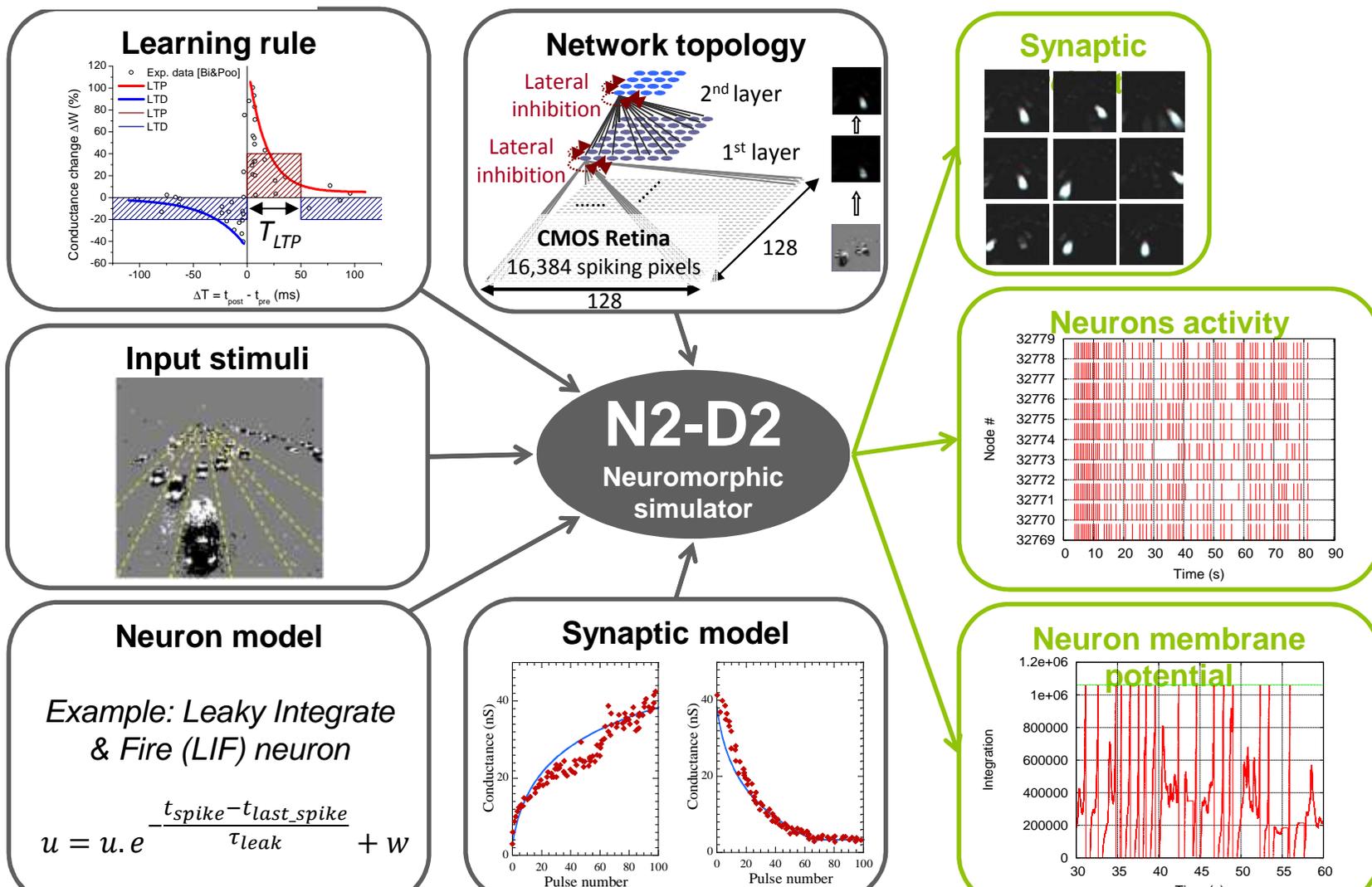


Synaptic weight update through STDP



1. G. Snider, *Nanoscale Architectures*, 2008
2. B. Linares-Barranco et al, *Nature Precedings*, 2009

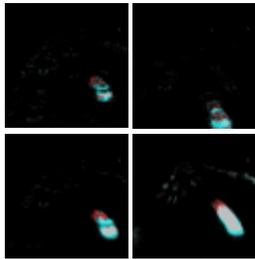
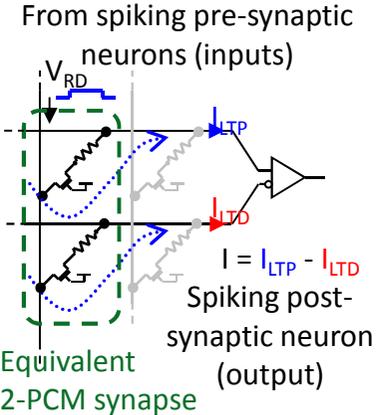
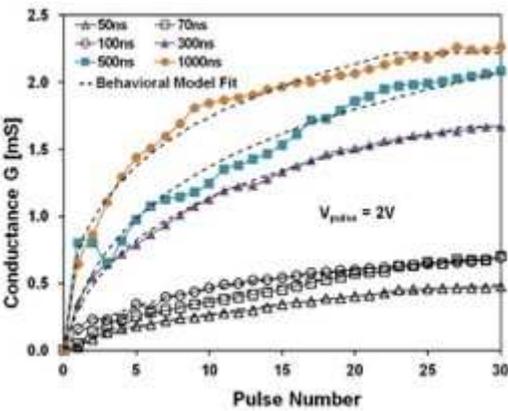
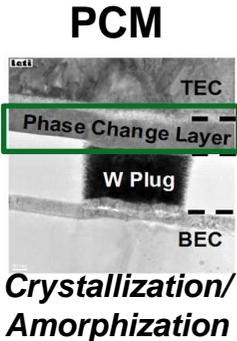
# BIO-INSPIRED MODELS EXPLORATION



→ Complete tool flow for bio-inspired synapses, neurons and learning rules network simulations

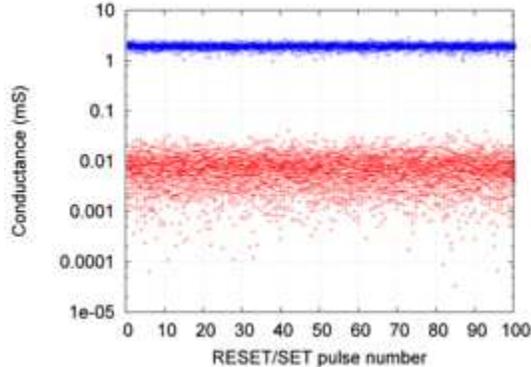
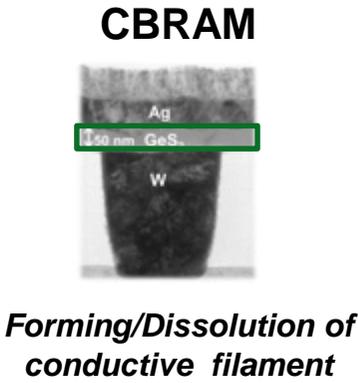
# NVM SYNAPSES IMPLEMENTATIONS

- 2-PCM synapses for unsupervised cars trajectories extraction



[O. Bichler et al., Electron Devices, IEEE Transactions on, 2012]

- CBRAM binary synapses for unsupervised MNIST handwritten digits classification with stochastic learning

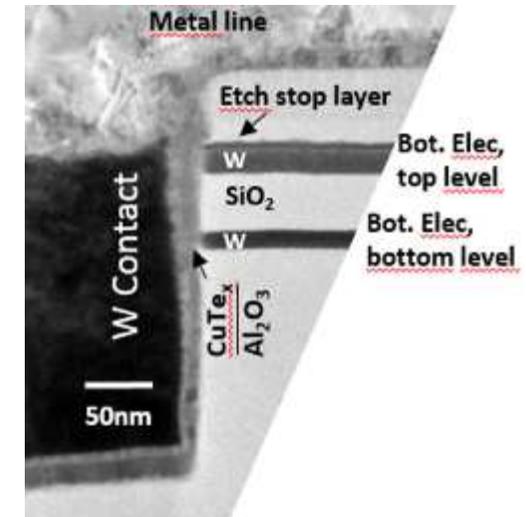


[M. Suri et al., IEDM, 2012]

## EXAMPLE OF ON-GOING INVESTIGATIONS: VRRAM FOR NEUROMORPHIC APPLICATIONS

- **Investigation of VRRAM based on CBRAM stack**

- 2 levels (proof of concept)
- 16 levels (goal)
- 1 select transistor per level (proof of concept)
- Integrated selector (goal)
- CBRAM most suitable R for neuromorphic
- OxRAM also analysed



- **Design: support development for VRRAM**

- **High Density**: Estimate the maximum size of a VRRAM-based array supposing to have an integrated selector [E. Cha, ISCAS 2014]
- **Neuromorphic**: propose a circuit dimensioning for the neuromorphic approach presented at IEDM 2015 (1TnR pillar ~ Synapse, NO Selector)

## ■ Objective:

- Fabricate a chip implementing a neuromorphic architecture that supports state-of-the-art machine learning algorithms and spike-based learning mechanisms.

## ■ Features:

- 28nm FDSOI technology with RRAM synapses
- Ultra low power scalable and reconfigurable architecture
- 50x lower dissipation than digital equivalent
- TFT based scalable multichip architecture platform
- A technology to implement on-chip learning, using native adaptive characteristics of electronic synaptic elements

## A NEW EU COLLABORATIVE PROJECT: NEURAM3

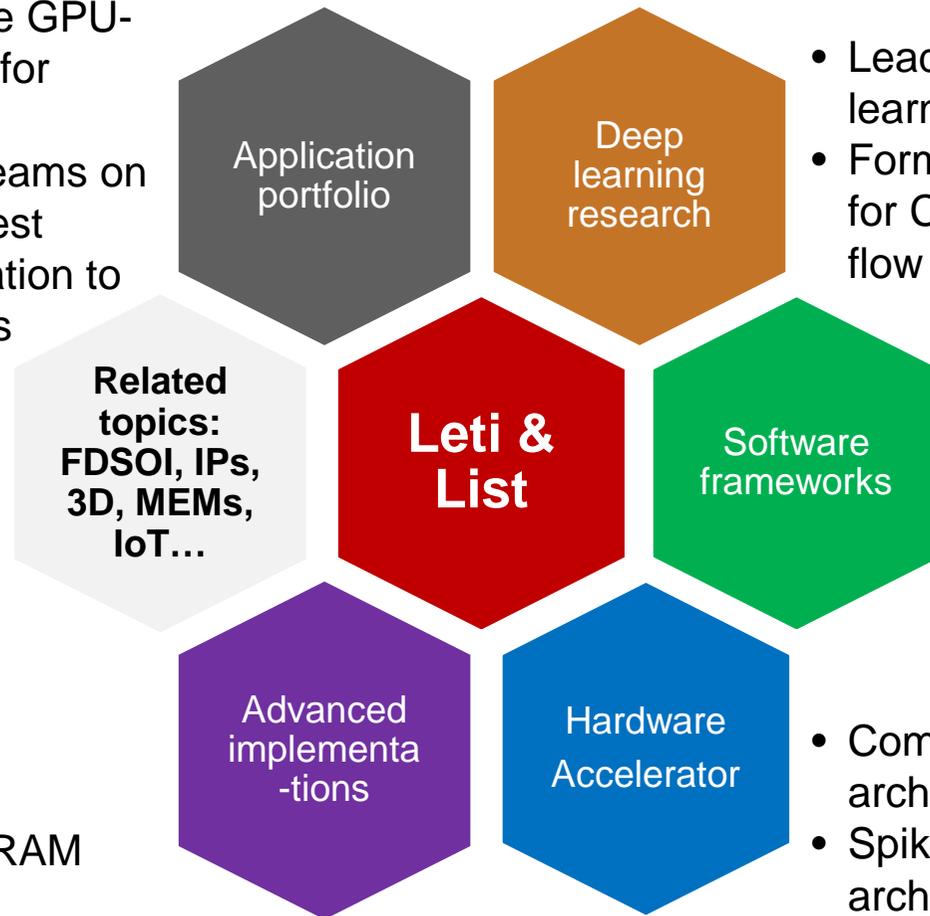


Participant no.	Organization name	Short name	Country
1 (Coordinator)	Commissariat a l'energie atomique et aux energies alternatives	CEA	France
2	Interuniversitair Micro-Electronica Centrum IMEC VZW	IMEC	Belgium
3	Stichting IMEC Nederland	IMEC-NL	Netherlands
4	IBM Research GmbH	IBM	Switzerland
5	University of Zurich, Institute of Neuroinformatics	UZH	Switzerland
6	Agencia Estatal Consejo Superior de Investigaciones Cientificas, Instituto de Microelectronica de Sevilla	CSIC	Spain
7	Consiglio Nazionale delle Ricerche	CNR	Italy
8	Jacobs University Bremen	JAC	Germany
9	ST-Microelectronics S.A.	STM	France

# LETI AND LIST ASSETS IN DEEP LEARNING

## Summary of key points

- Large-scale database GPU-accelerated learning for CNN
- Among the leading teams on ImageClef2015 contest
- From scratch exploration to industrial applications



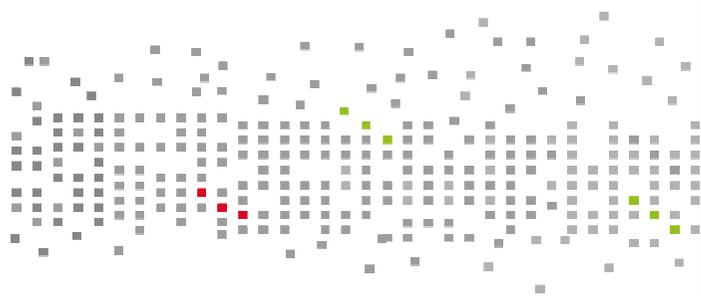
- Lead in bio-inspired STDP learning (IEDM'11,12,14)
- Formalized spike-coding for CNN, complete tool flow for co-simulation

- Complete framework with C, OpenCL, CUDA and HLS exports
- Complete tool flow for spike-coding DSP

- 2-PCMs synapse (patented) scheme (IEDM'15)
- Lead in SNN with RRAM devices (IEDM'14)

- Competitive reconfigurable architecture with P-Neuro
- Spike-coding DSP architecture
- Increased efficiency with 3D

***Thank you  
for your  
attention***



**Leti, technology research institute**  
Commissariat à l'énergie atomique et aux énergies alternatives  
Minatec Campus | 17 rue des Martyrs | 38054 Grenoble Cedex | France  
[www.leti.fr](http://www.leti.fr)

