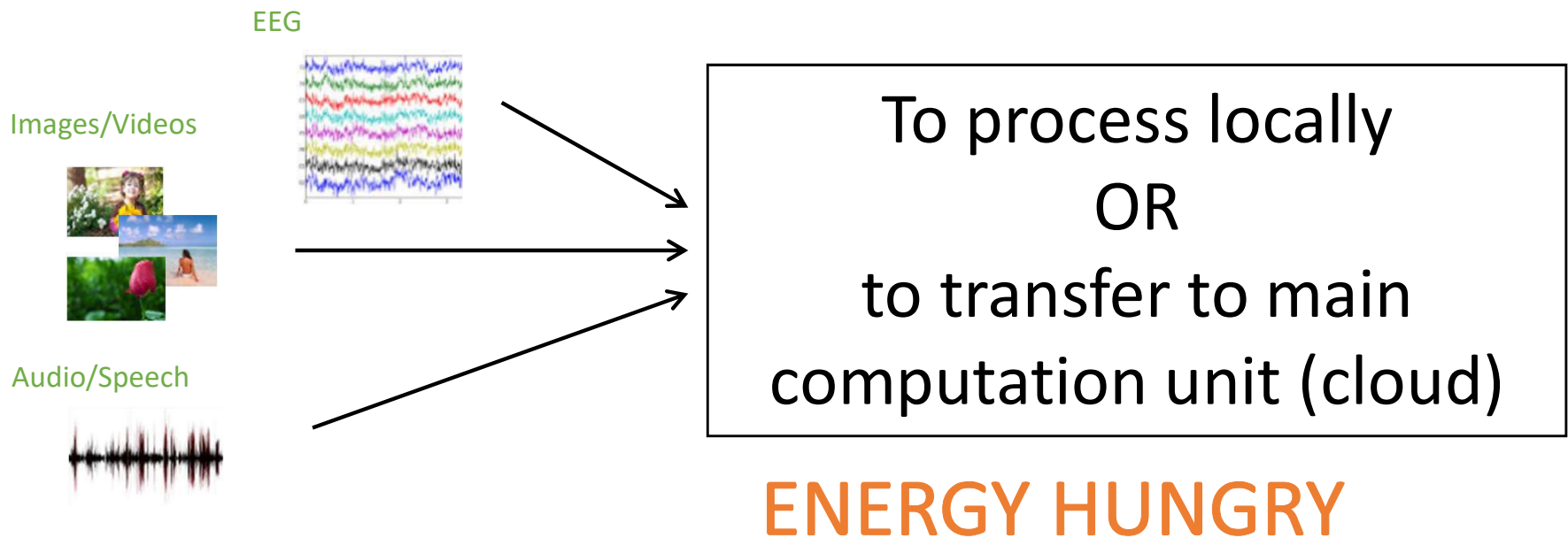# LEOPAR: Low-Energy On-chip Pre-processing for Activity Recognition

Benoit LARRAS, Kévin HÉRISSÉ, Antoine FRAPPÉ, Bruno STEFANELLI, Andreas KAISER

14/05/2019 – Colloque du GdR BioComp

# Project context

- **Massive amounts of data**
- **Always-on sensing**

EEG

Images/Videos



Audio/Speech

To process locally
OR
to transfer to main
computation unit (cloud)

ENERGY HUNGRY

Small, cheap, no battery replacement
→ **Towards Near-Sensor Computing**

iemn
Institut d'Electronique, de Microélectronique
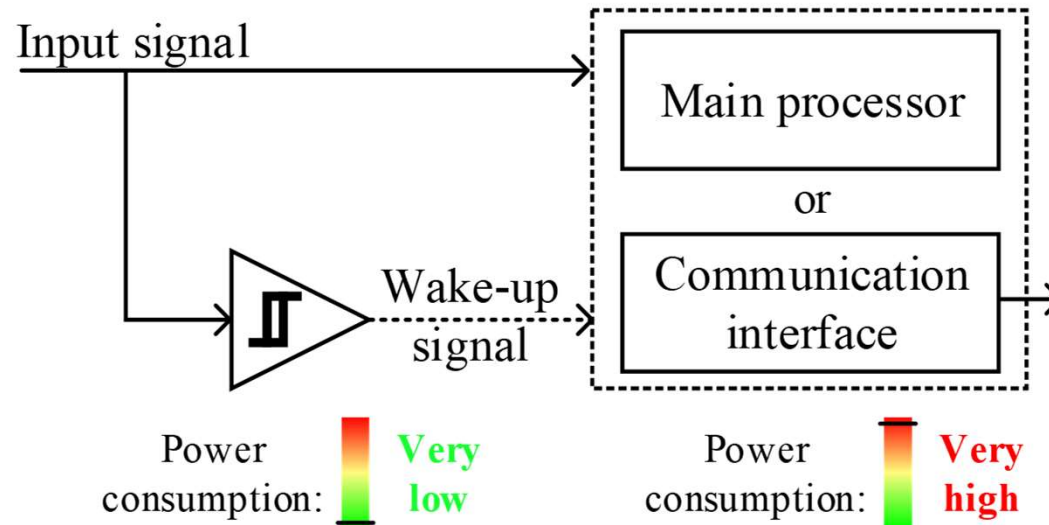et de Nanotechnologie
UMR CNRS 8520

# Application fields

- Audio processing
  - Voice Activity Detection in noisy context
  - Vowels, words, language recognition
  - Specific feature extraction

- Human-body signal classifications
  - ECG, EEG, etc…

- Vibration and movement recognition

- Image processing
  - Motion-triggered cameras
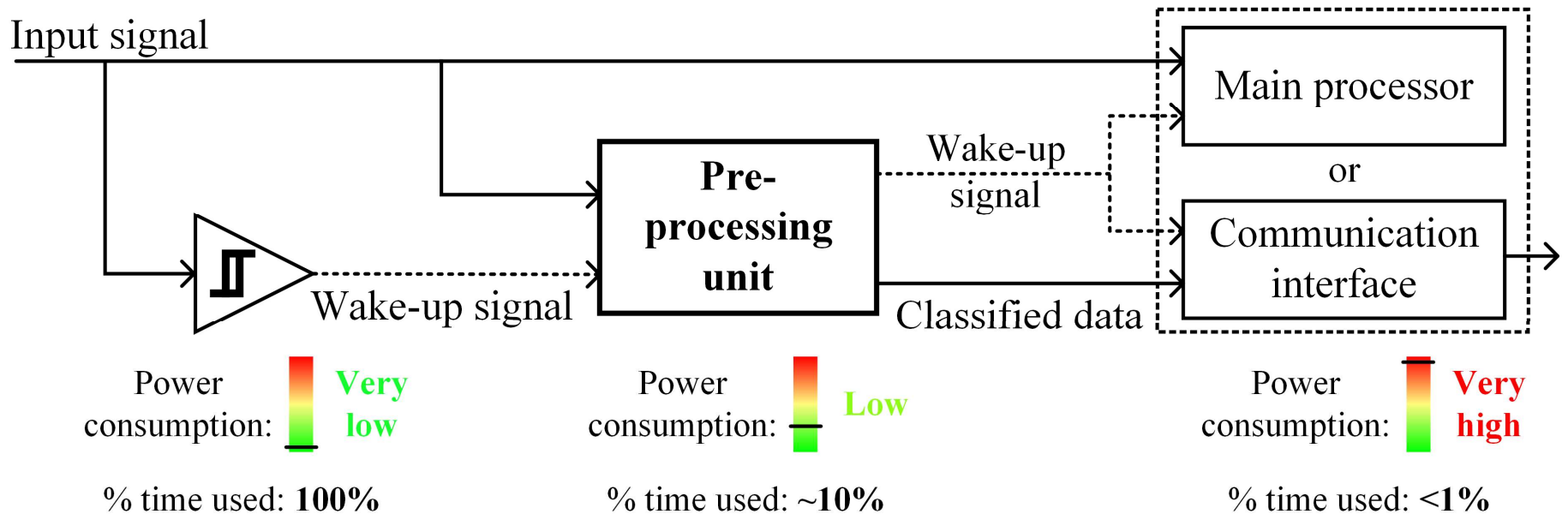  - Face detection / Owner-activated devices

- Automotive

iemn
Institut d'Electronique, de Microélectronique
et de Nanotechnologie
UMR CNRS 8520

3

# Project objectives

**Standard scheme:**



Non relevant data is processed if it exceeds the threshold…

# Project objectives



Input signal → Pre-processing unit → Main processor or Communication interface

Wake-up signal

Wake-up signal

Classified data

Power consumption: **Very low**  •  % time used: **100%**

Power consumption: **Low**  •  % time used: **~10%**

Power consumption: **Very high**  •  % time used: **<1%**

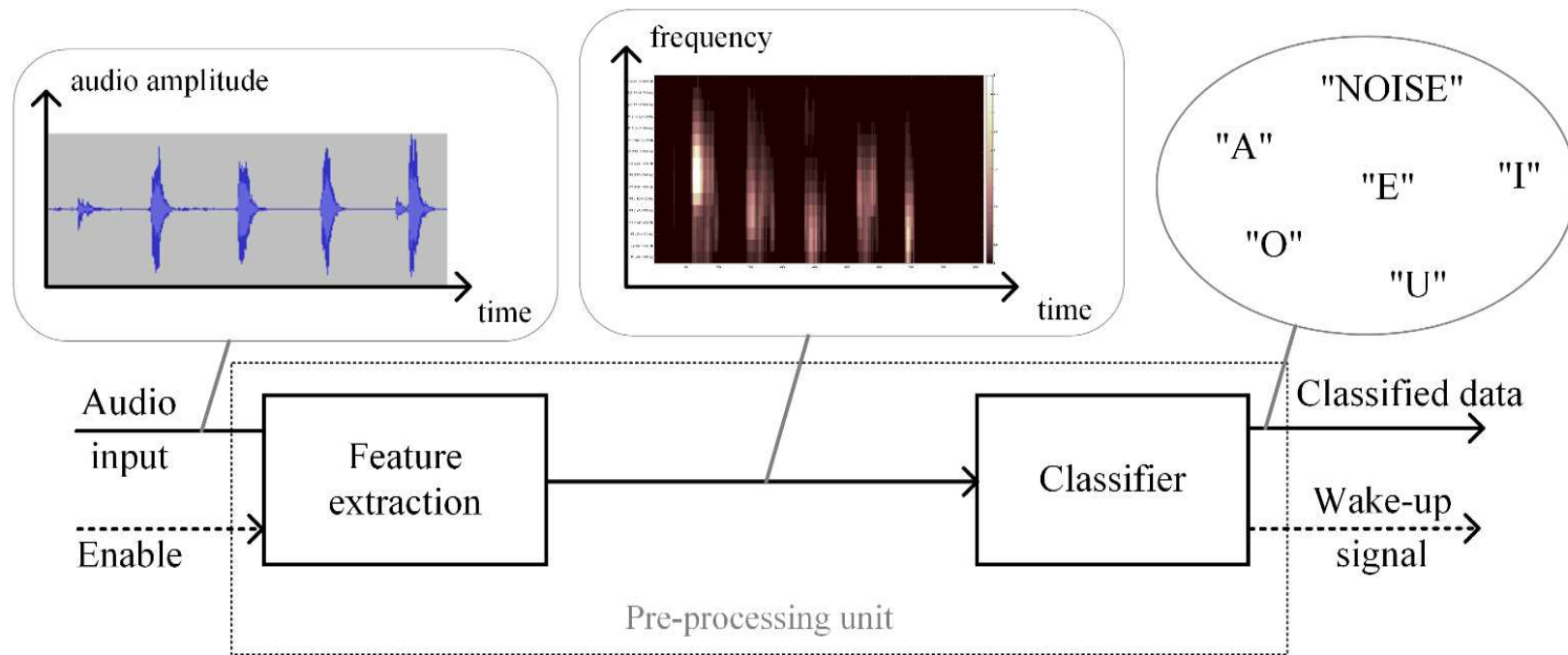**Near-sensor Computing**: process **relevant** data **as close as the sensor** as possible

- Aggregation of a **lower amount of data**
- Need of energy-hungry processing during a **lower amount of time**

# Envisionned demonstration



- Focus on **audio applications**: voice activity detection, vowels recognition, keyword detection.

- **On-chip event-driven** feature extraction

- **Small-scale neuro-inspired** classification unit

# Feature extraction
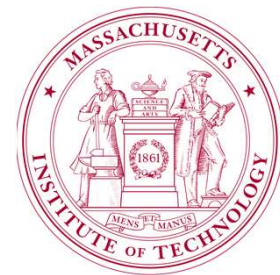
Objective: **extract energy in different frequency bands**

- Analog filter bank [Badami, JSSC 2016]

  - **Low energy**

  - **Non configurable filters**

  - **High silicon area**

- Digital FFT [Price, JSSC 2017]

  - **Configurability**

  - **Audio fidelity**

  - **Latency**

  - **High complexity**

  - **High energy**

# Feature extraction

- **Digital filter bank**

    - **Configurability**

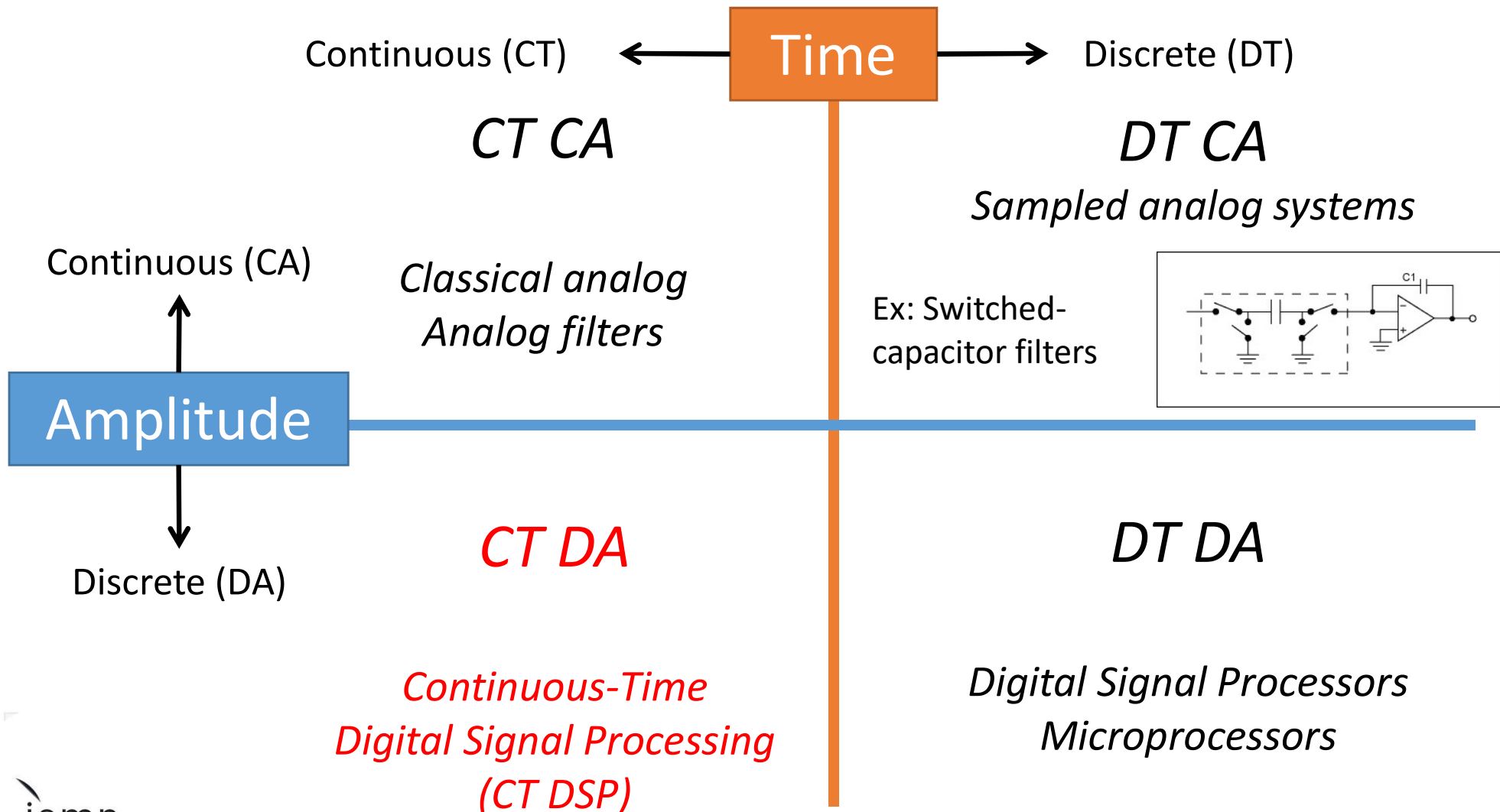    - **Low latency**

    - **Implementation capability**

Requires **preliminary always-on** A-to-D conversion and signal processing of the complete spectrum
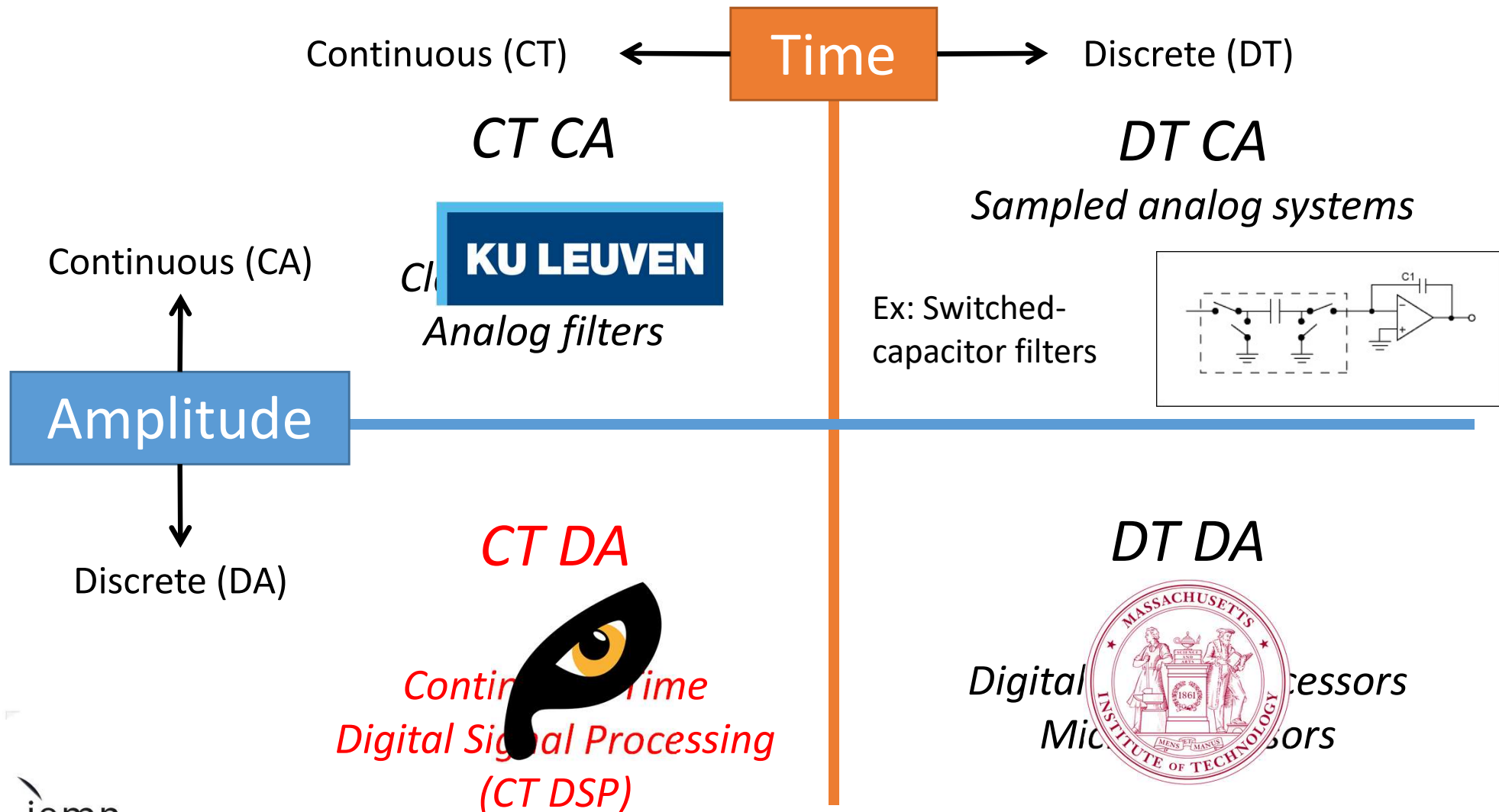
**Event-driven / Clockless ?**

→ Advantages of both analog and digital implementations

iemn
Institut d'Electronique, de Microélectronique et de Nanotechnologie
UMR CNRS 8520

# Opportunity: Continuous-Time Digital Signal Processing (CTDSP)

Continuous (CT) ← **Time** → Discrete (DT)

*CT CA*

*DT CA*

*Sampled analog systems*

Continuous (CA)

*Classical analog*
*Analog filters*

Ex: Switched-capacitor filters

**Amplitude**

Discrete (DA)

*CT DA*

*DT DA*

*Continuous-Time*
*Digital Signal Processing*
*(CT DSP)*

*Digital Signal Processors*
*Microprocessors*

iemn
Institut d'Electronique, de Microélectronique
et de Nanotechnologie
UMR CNRS 8520

# Opportunity: Continuous-Time Digital Signal Processing (CTDSP)

Continuous (CT) ← **Time** → Discrete (DT)

*CT CA*

*DT CA*

*Sampled analog systems*

Continuous (CA)

**KU LEUVEN**

*Cl...*

*Analog filters*

Ex: Switched-capacitor filters

**Amplitude**

Discrete (DA)

*CT DA*

*DT DA*

*Conti... ...Time Digital Sig...al Processing (CT DSP)*

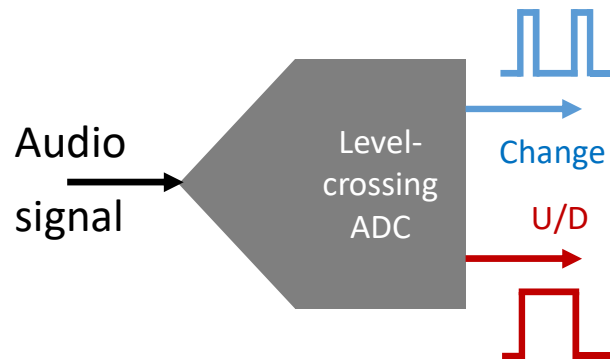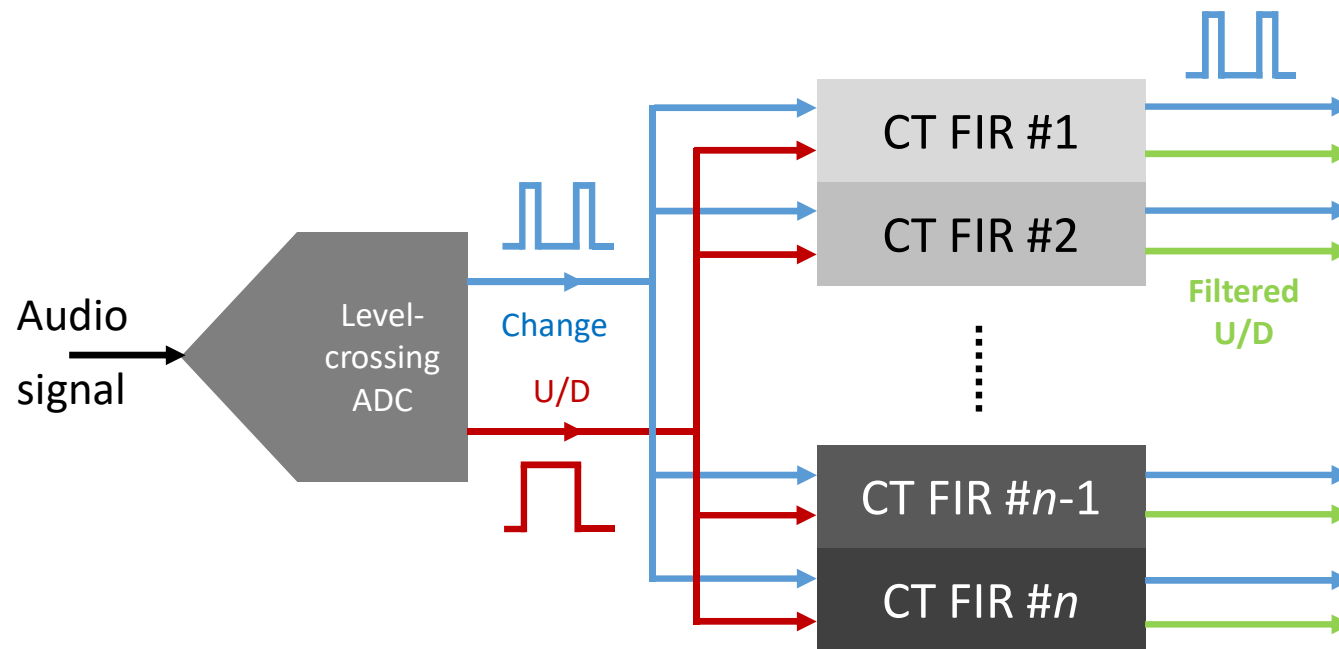*Digital... ...cessors Mic...ssors*

# Continuous-Time (CT) advantages

- **Event-driven system**
  - No clock
  - Event-driven power consumption

- **CMOS Digital System**
  - Configurability
  - Scalability
  - High integration level

# Continuous-Time (CT) advantages

- **Event-driven system**
  - No clock
  - Event-driven power consumption

- **CMOS Digital System**
  - Configurability
  - Scalability
  - High integration level

Audio signal → Level-crossing ADC → Change, U/D
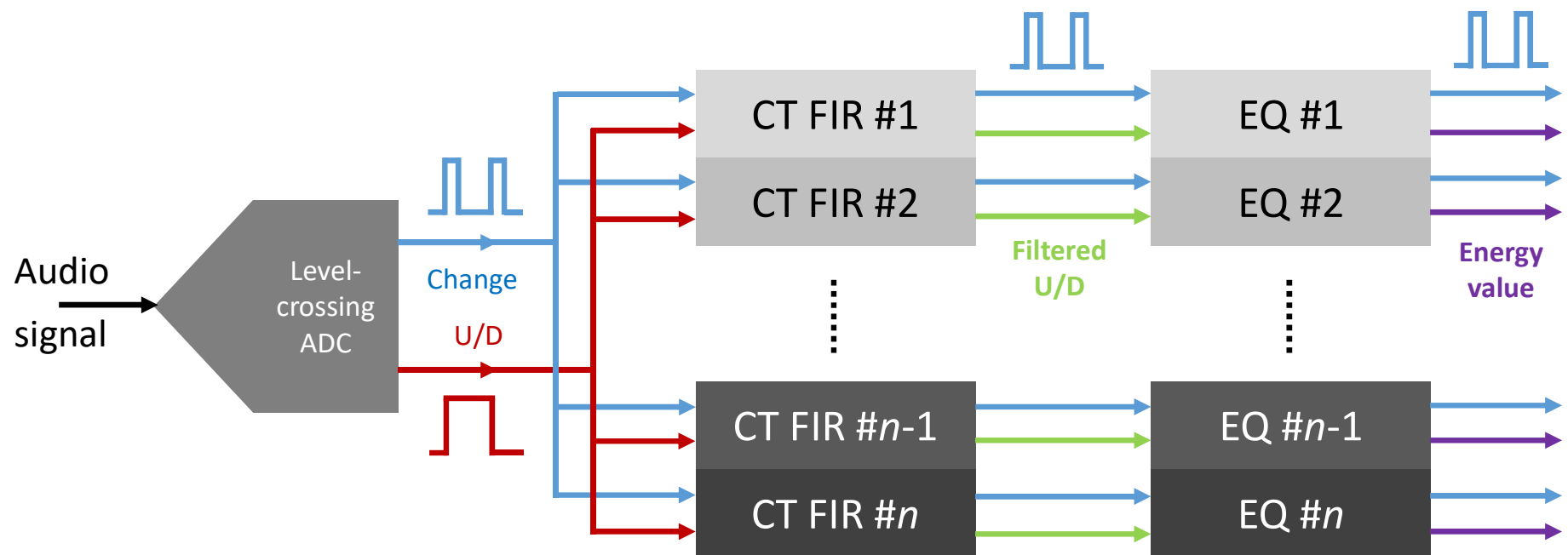
# Continuous-Time (CT) advantages

- **Event-driven system**

  - No clock

  - Event-driven power consumption

- **CMOS Digital System**

  - Configurability

  - Scalability

  - High integration level



[Kurchuk, JSSC 2012]

# Continuous-Time (CT) advantages

- **Event-driven system**
  - No clock
  - Event-driven power consumption

- **CMOS Digital System**
  - Configurability
  - Scalability
  - High integration level

Audio signal → Level-crossing ADC → Change / U/D → CT FIR #1, CT FIR #2, ... CT FIR #$n$-1, CT FIR #$n$ → Filtered U/D → EQ #1, EQ #2, ... EQ #$n$-1, EQ #$n$ → Energy value

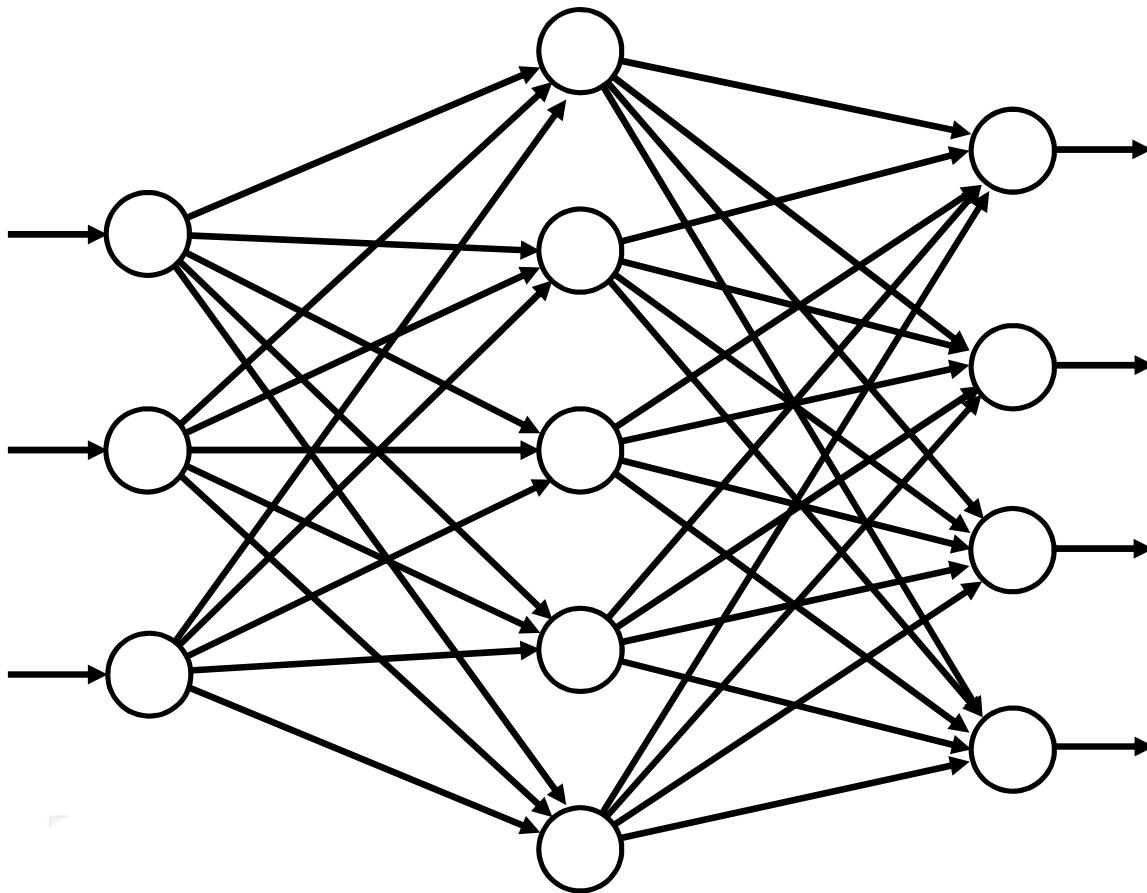EQ: Energy Quantifier

# Classification

- Detection of a **small number of specific patterns**: voice activity, vowels, specific sounds, etc.

- **Limited amount of features** → limited amount of computing units (neurons)

- Embedded environment: **energy and complexity requirements**

→ Towards a **binarized**, **small-scale** classifier with **determined data storage**

iemn
Institut d'Electronique, de Microélectronique
et de Nanotechnologie
UMR CNRS 8520

# Opportunity: Small-scale classifiers

- **Only necessary functions implemented**
  - Online inference only, towards binary synaptic weights
  - Activation function: e.g. local Winner-Takes-All

- **Asynchronous behavior →** Event-driven compatible

- **Short reaction time →** Real-time compatible

- **Envisioned classifier models:**
  - LSTM
  - Spiking neural networks
  - Clique-based networks

iemn
Institut d'Electronique, de Microélectronique
et de Nanotechnologie
UMR CNRS 8520

# Neural networks models

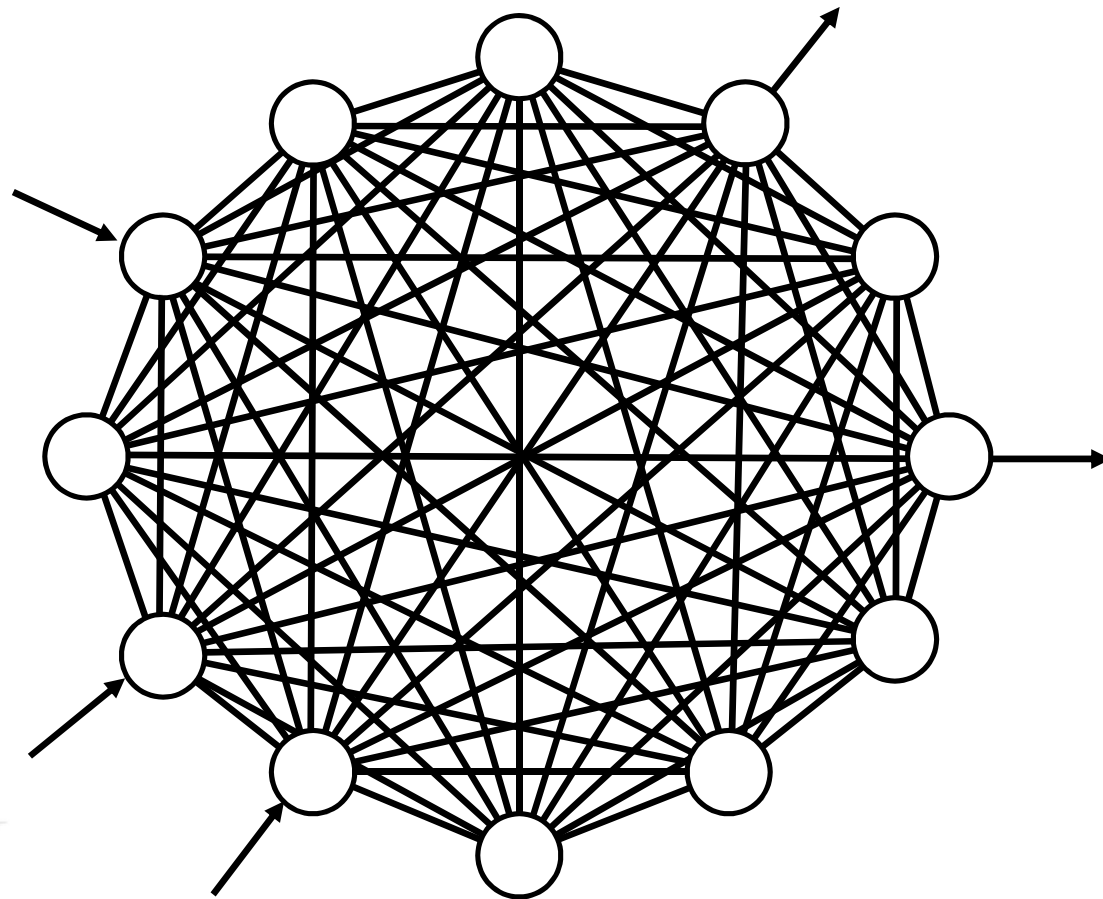Several organizations for the neurons:



**Feedforward neural networks**

- Full connectivity from a layer to the next one

- Unidirectional links

# Neural networks models

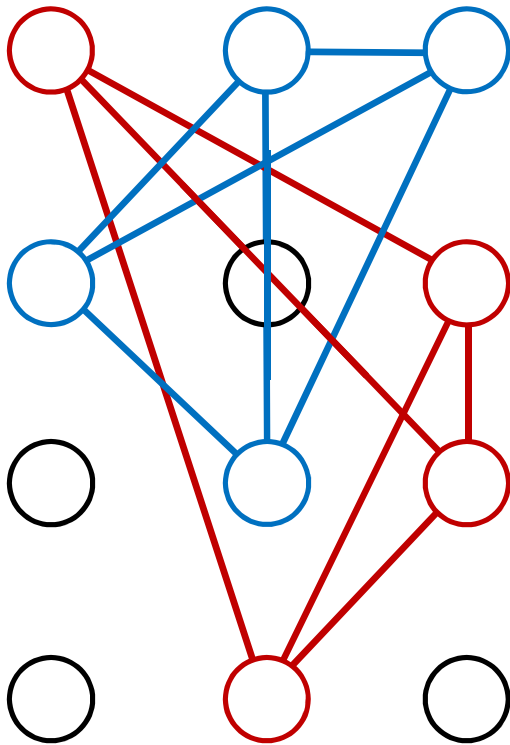Several organizations for the neurons:



**Recurrent neural networks (Hopfield)**

- Full connectivity between the neurons

- Bidirectional links

iemn
Institut d'Electronique, de Microélectronique
et de Nanotechnologie
UMR CNRS 8520

# Neural networks models

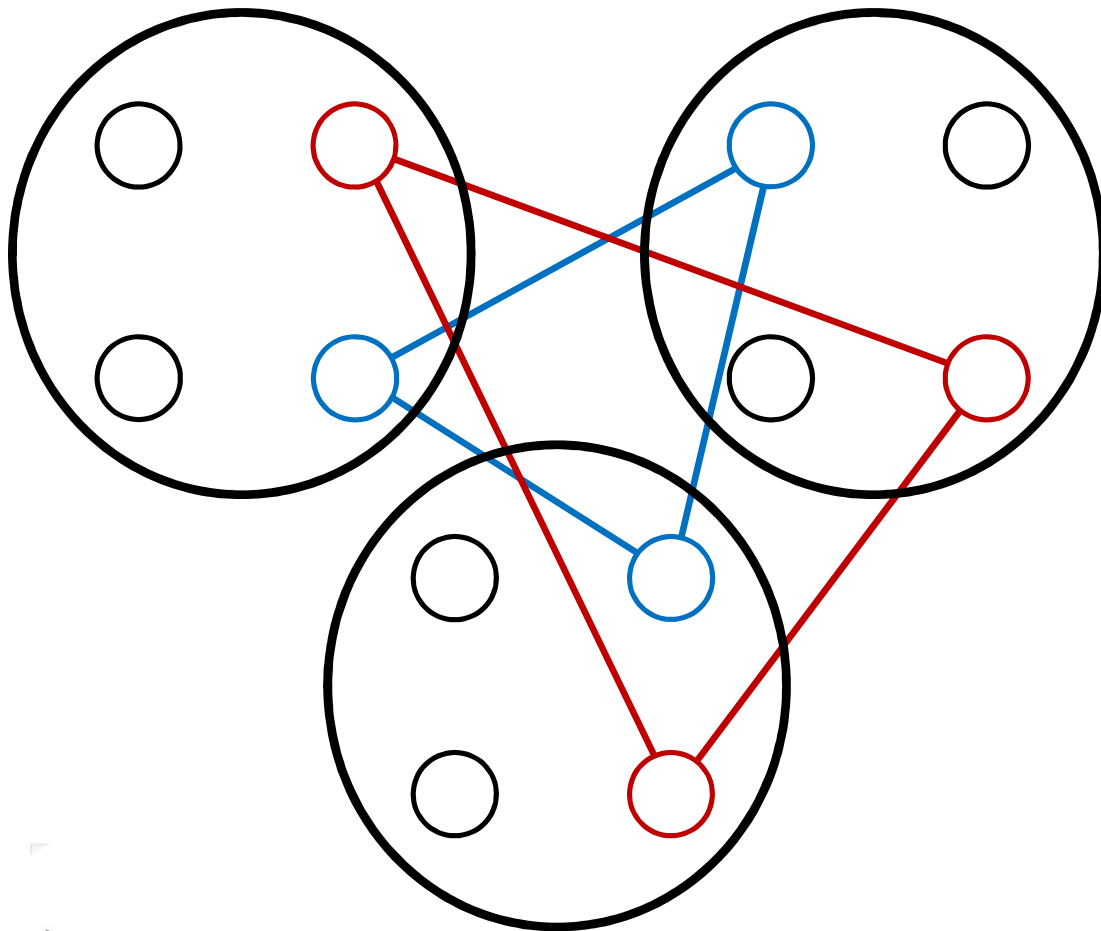Several organizations for the neurons:

**Clique-based neural networks**

- Connections between neurons only through cliques

- Bidirectional links

# Neural networks models
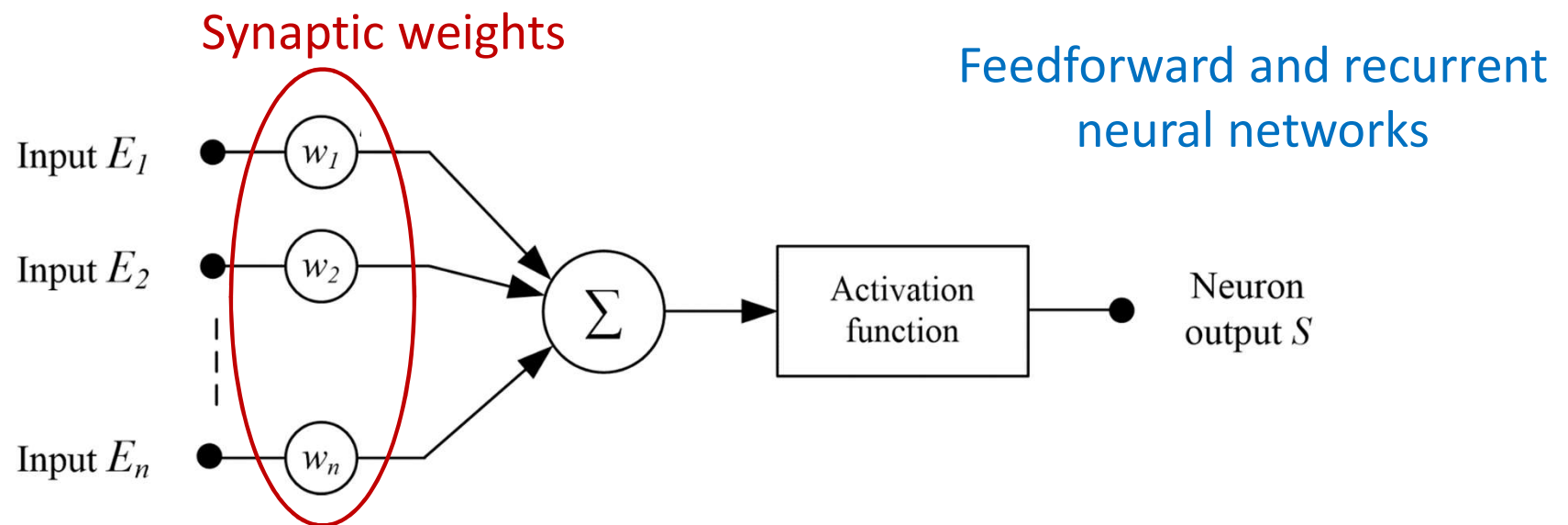
Several organizations for the neurons:



**Clustered clique-based networks**

- Division in clusters

- Connections between neurons from different clusters
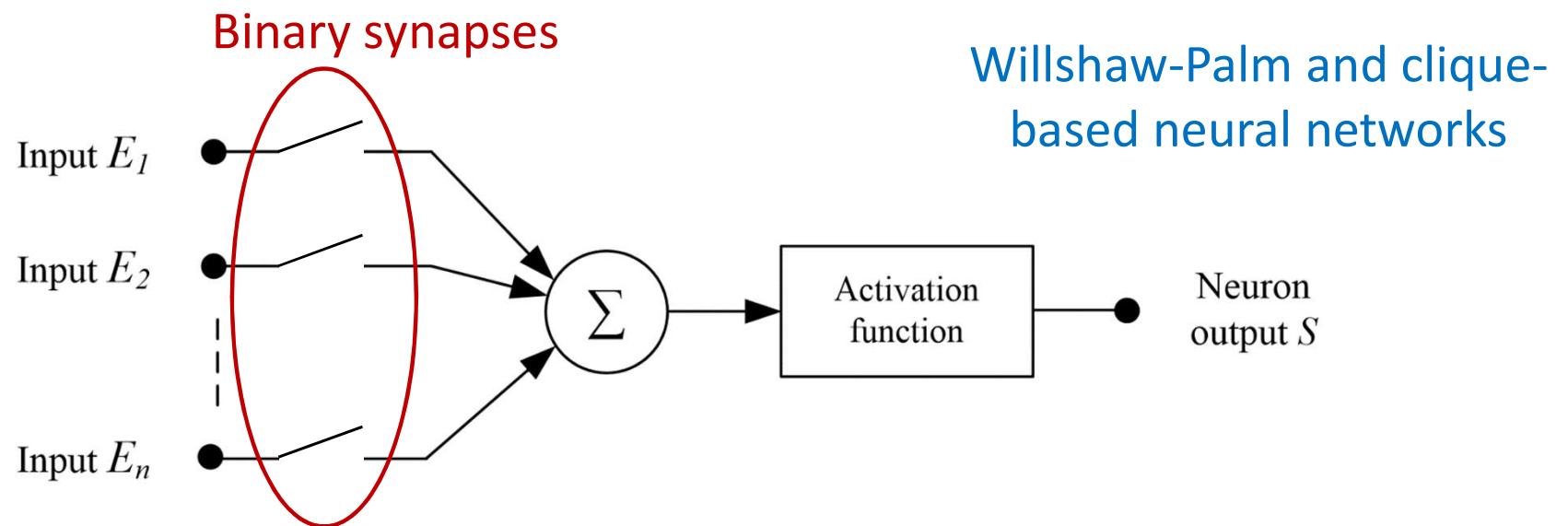
[Gripon and Berrou, TNNLS 2011]

# Inside a neuron

Structure of a neuron:

Synaptic weights
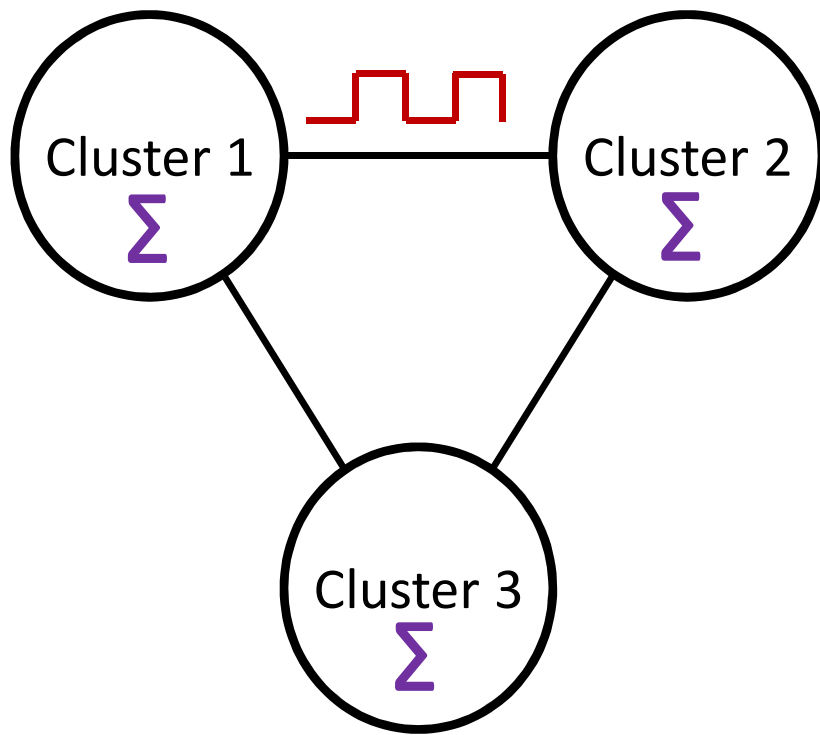
Feedforward and recurrent neural networks

Input $E_1$ — $w_1$

Input $E_2$ — $w_2$

Input $E_n$ — $w_n$

$\Sigma$ → Activation function → Neuron output $S$

iemn
Institut d'Electronique, de Microélectronique
et de Nanotechnologie
UMR CNRS 8520

# Inside a neuron

Structure of a neuron:

Binary synapses

Willshaw-Palm and clique-based neural networks



Less complex activation function: WTA rule

→ comparison + activation

iemn
Institut d'Electronique, de Microélectronique
et de Nanotechnologie
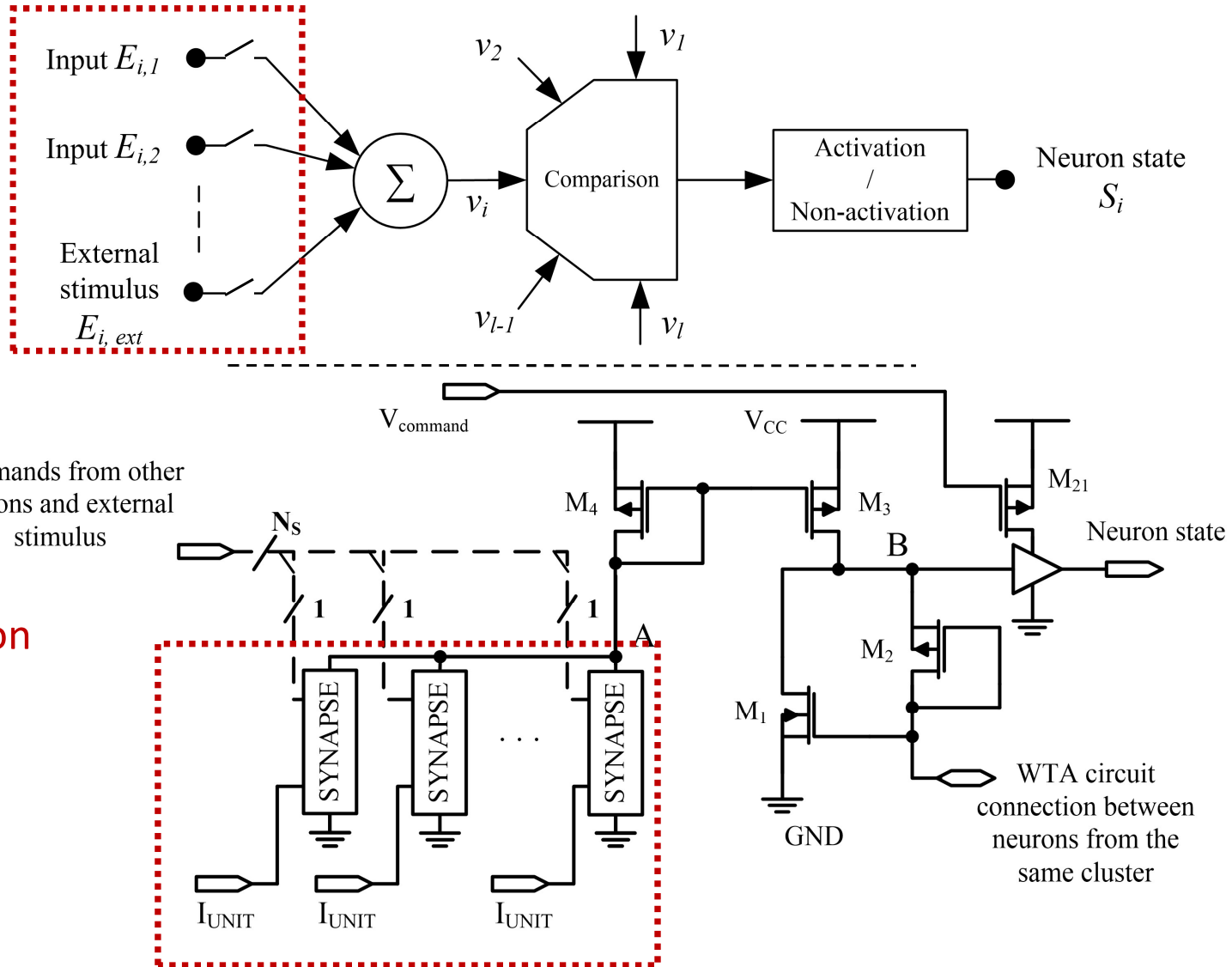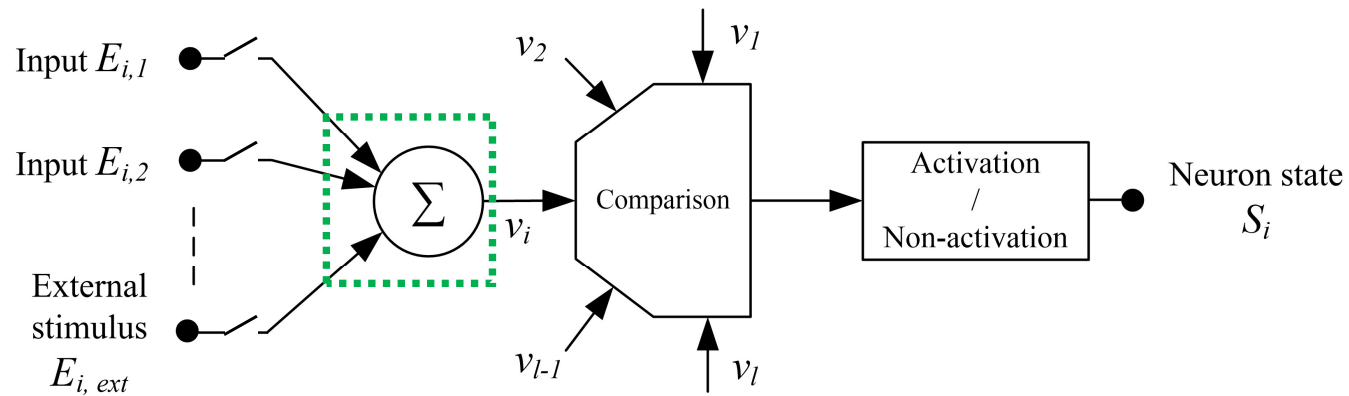UMR CNRS 8520

# Implementation choices



- **Binary information** exchanged by the neurons

  → **Communication: digital signals**

- **Simple analog circuits** adapted to the functions in a neuron

  → **Computations: analog signals**

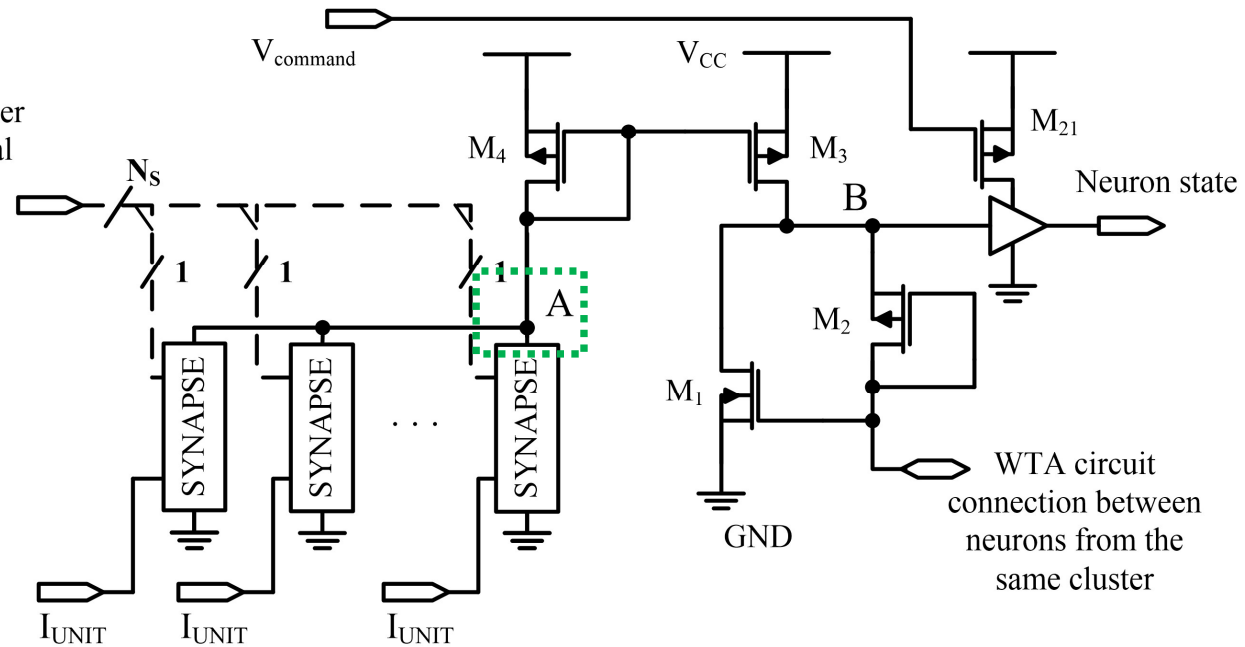  → **Mixed-signal asynchronous implementation**

# What about circuit ?



Input $E_{i,1}$

Input $E_{i,2}$
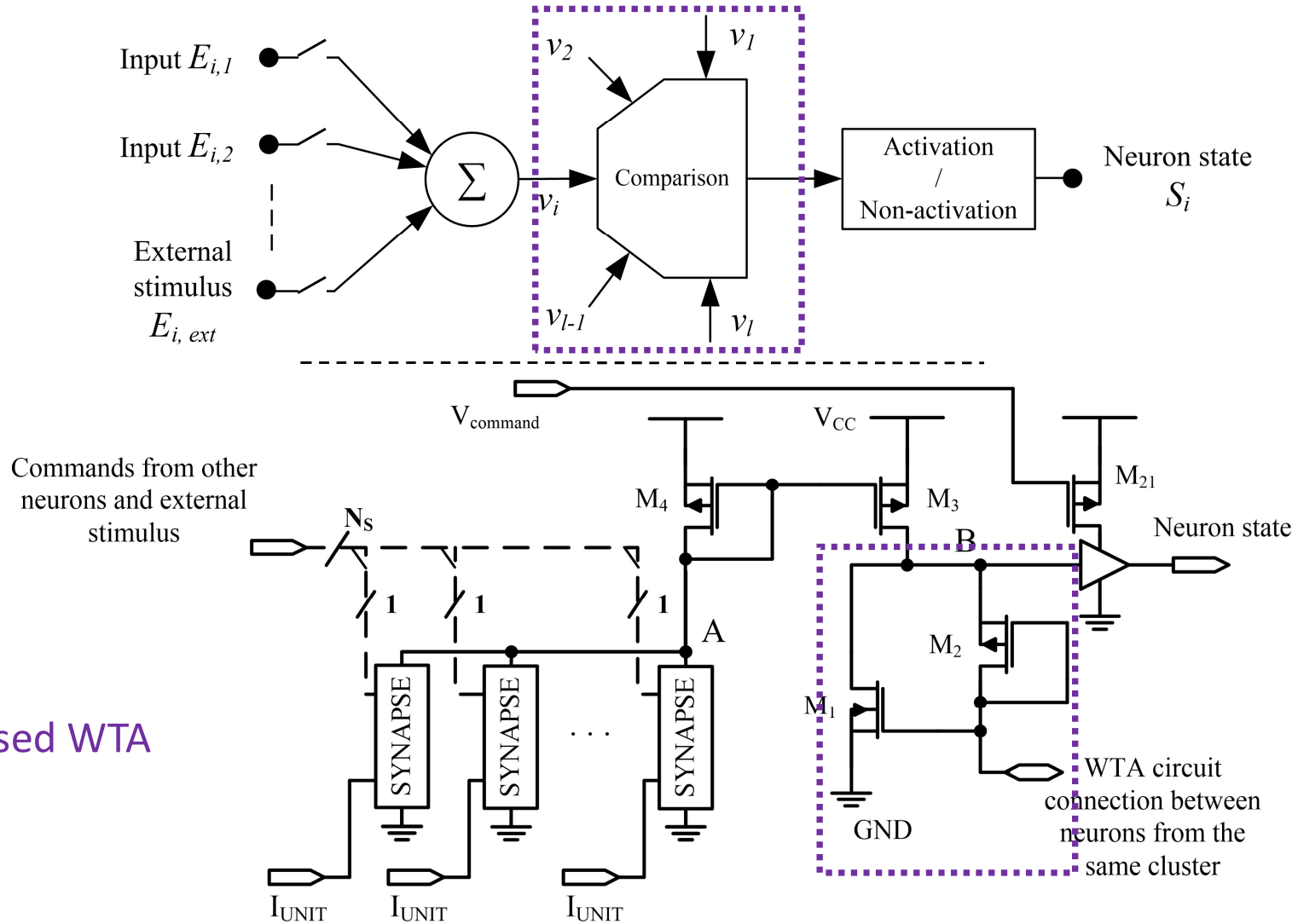
External stimulus $E_{i,ext}$

$v_2$     $v_1$

$\Sigma$   $v_i$   Comparison

Activation / Non-activation

Neuron state $S_i$

$v_{l-1}$   $v_l$

$V_{command}$

$V_{CC}$

Commands from other neurons and external stimulus

$N_S$

V-to-I conversion

1   1   1

$M_4$   $M_3$   $M_{21}$

B

Neuron state

A

$M_2$

$M_1$

SYNAPSE   SYNAPSE   . . .   SYNAPSE

$I_{UNIT}$   $I_{UNIT}$   $I_{UNIT}$

GND

WTA circuit connection between neurons from the same cluster

iemn
Institut d'Electronique, de Microélectronique et de Nanotechnologie
UMR CNRS 8520
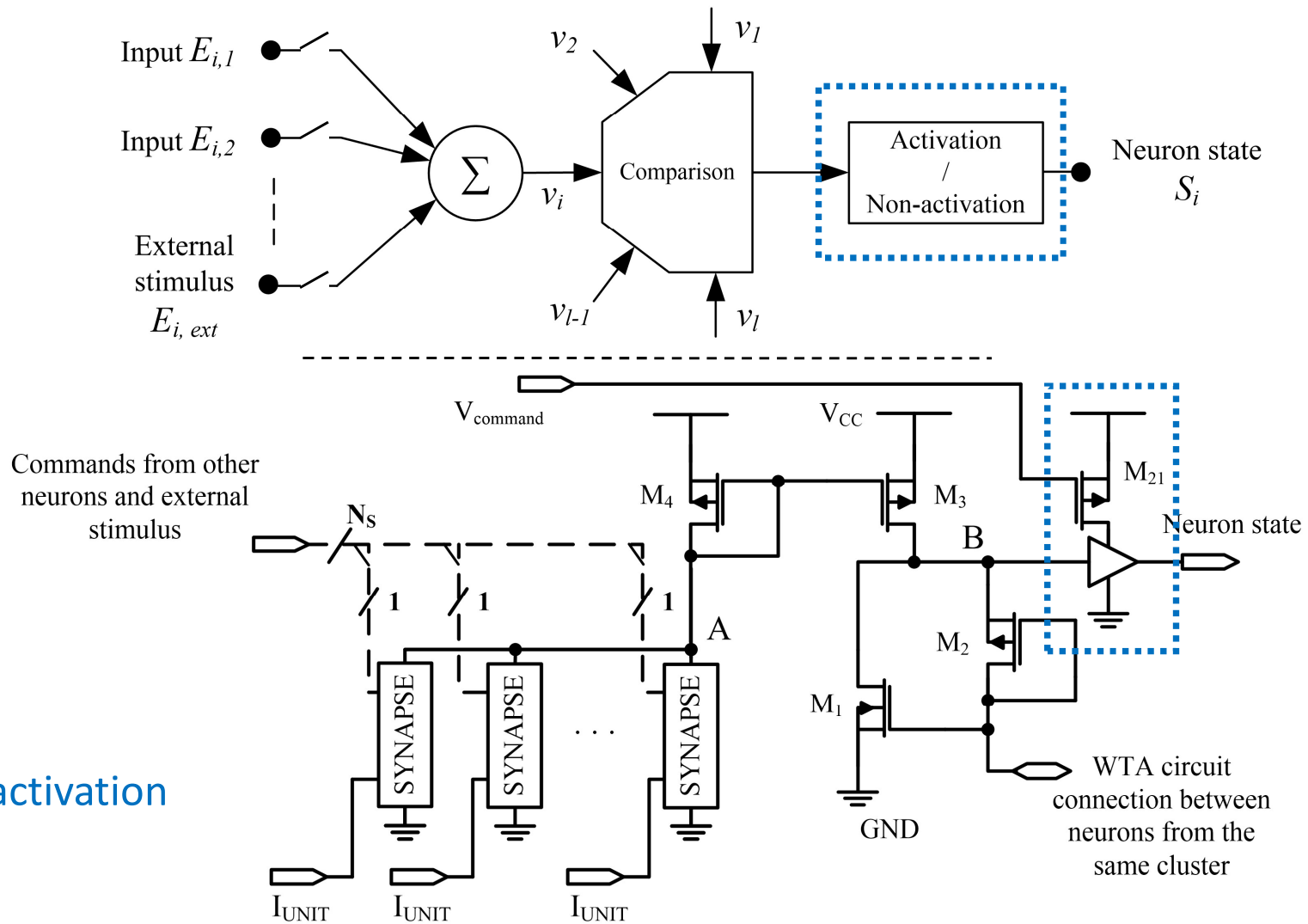
24

# What about circuit ?



**Current addition**

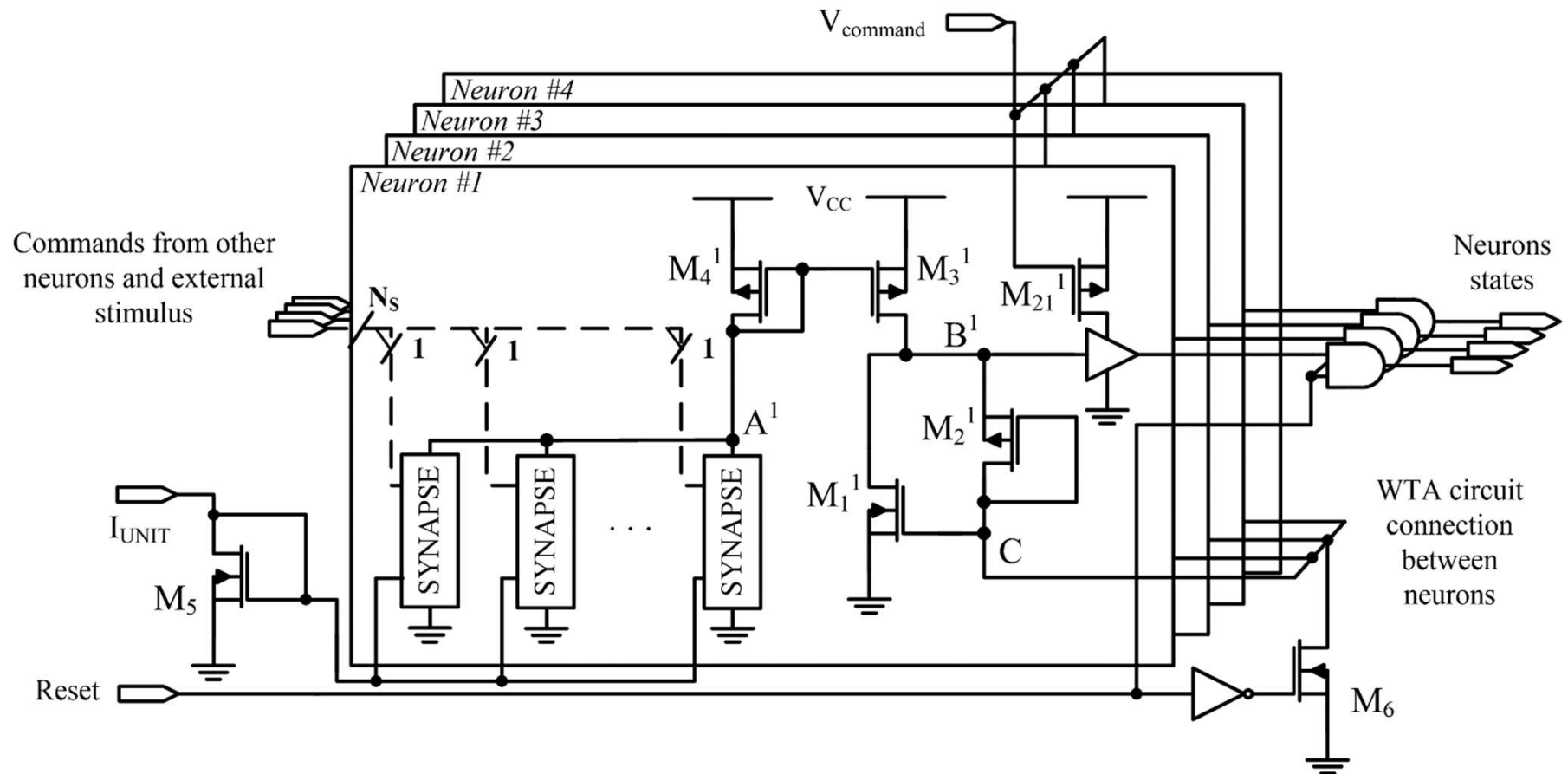# What about circuit ?



Current-based WTA

# What about circuit ?



Threshold activation

# What about circuit ?

**Schematic of a cluster of 4 neurons:**
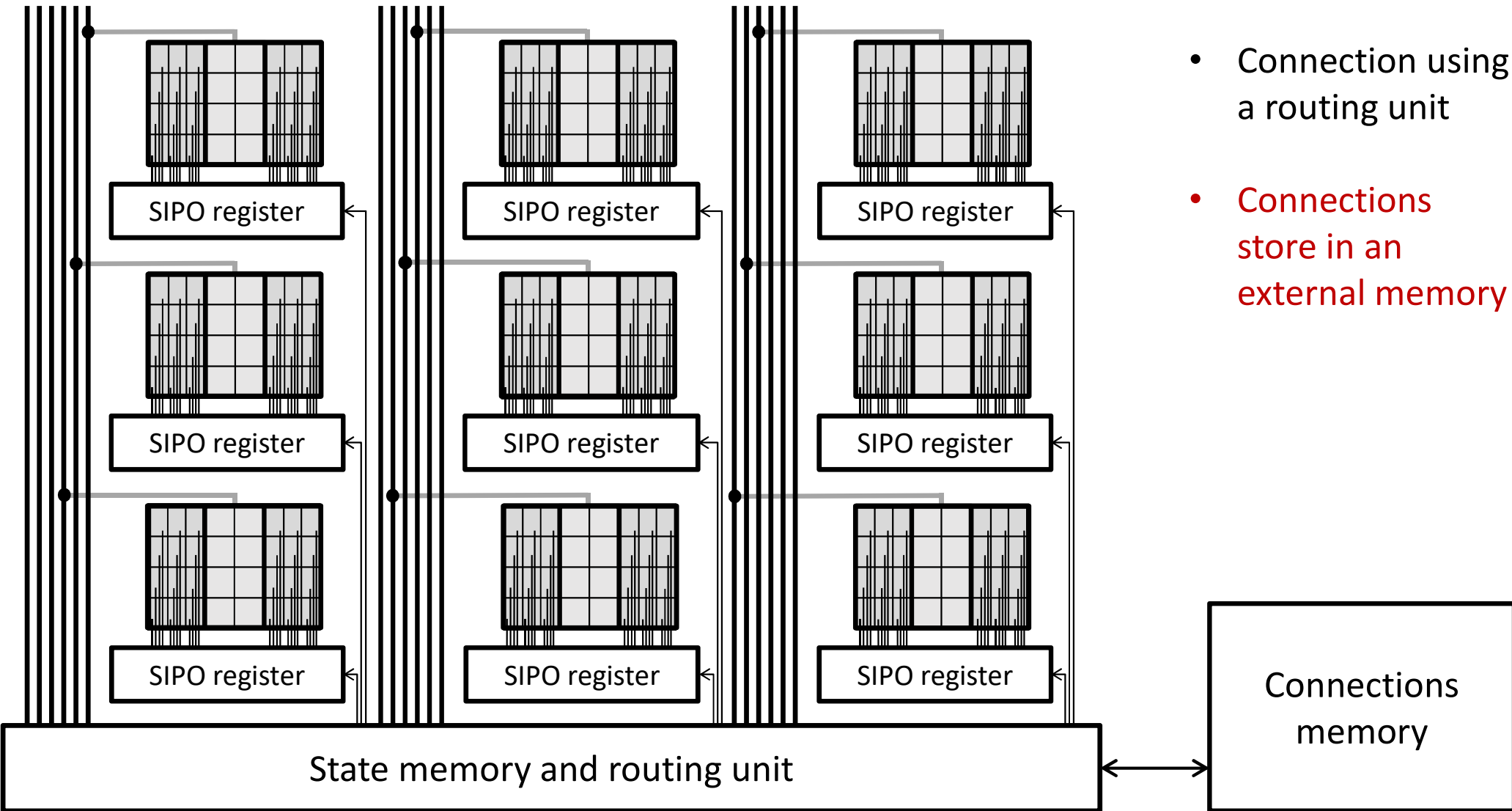
# Network topologies



Neurons          Synapses

- Cluster matrix

- Hardwired connections between neurons

- Fastest response

- No flexibility

iemn
Institut d'Electronique, de Microélectronique et de Nanotechnologie
UMR CNRS 8520

# Network topologies



- Connection using a routing unit

- Connections store in an external memory

SIPO register

SIPO register

SIPO register

SIPO register

SIPO register

SIPO register

SIPO register

SIPO register

SIPO register

State memory and routing unit

Connections memory

iemn
Institut d'Electronique, de Microélectronique et de Nanotechnologie
UMR CNRS 8520

30

# Network topologies

- Iteration of the process on one cluster

- Flexibility: topology changes with the number of iterations
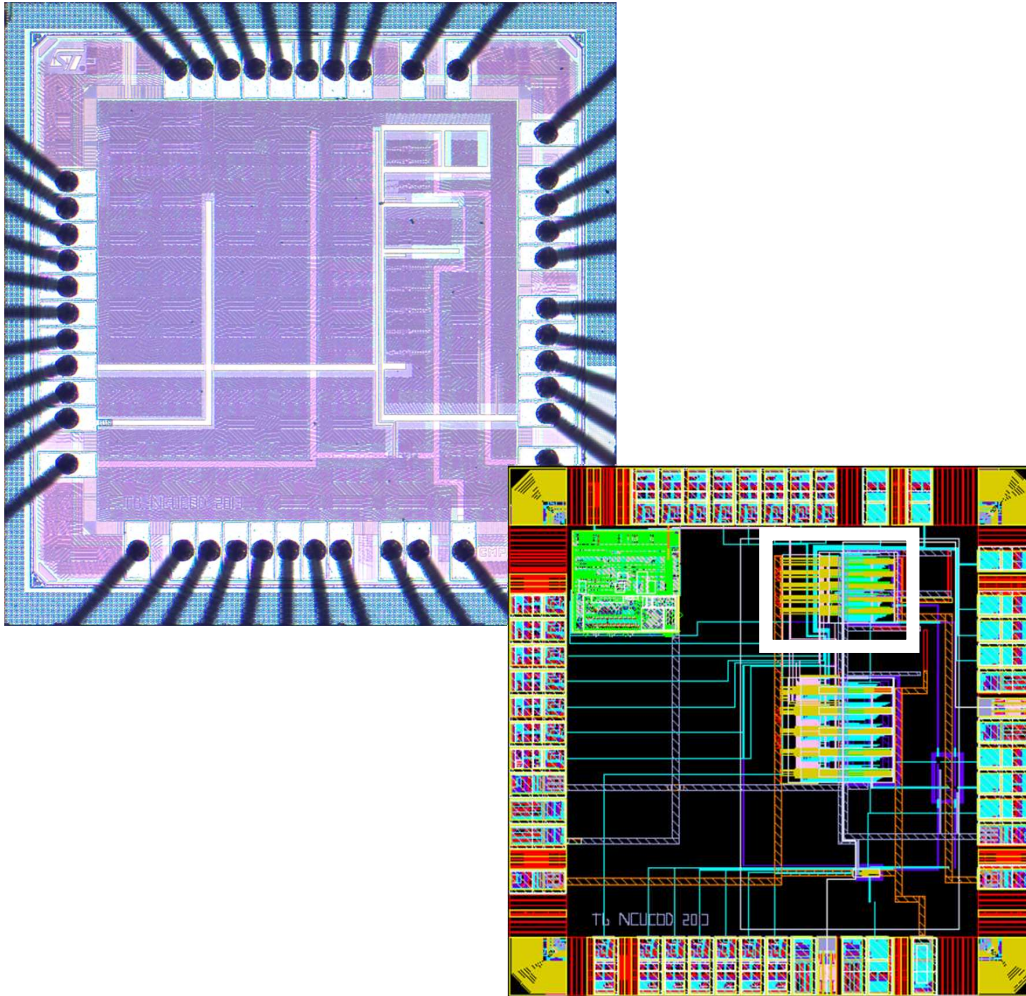
- Latency

SIPO register

State memory and routing unit

Connections memory

iemn
Institut d'Electronique, de Microélectronique
et de Nanotechnologie
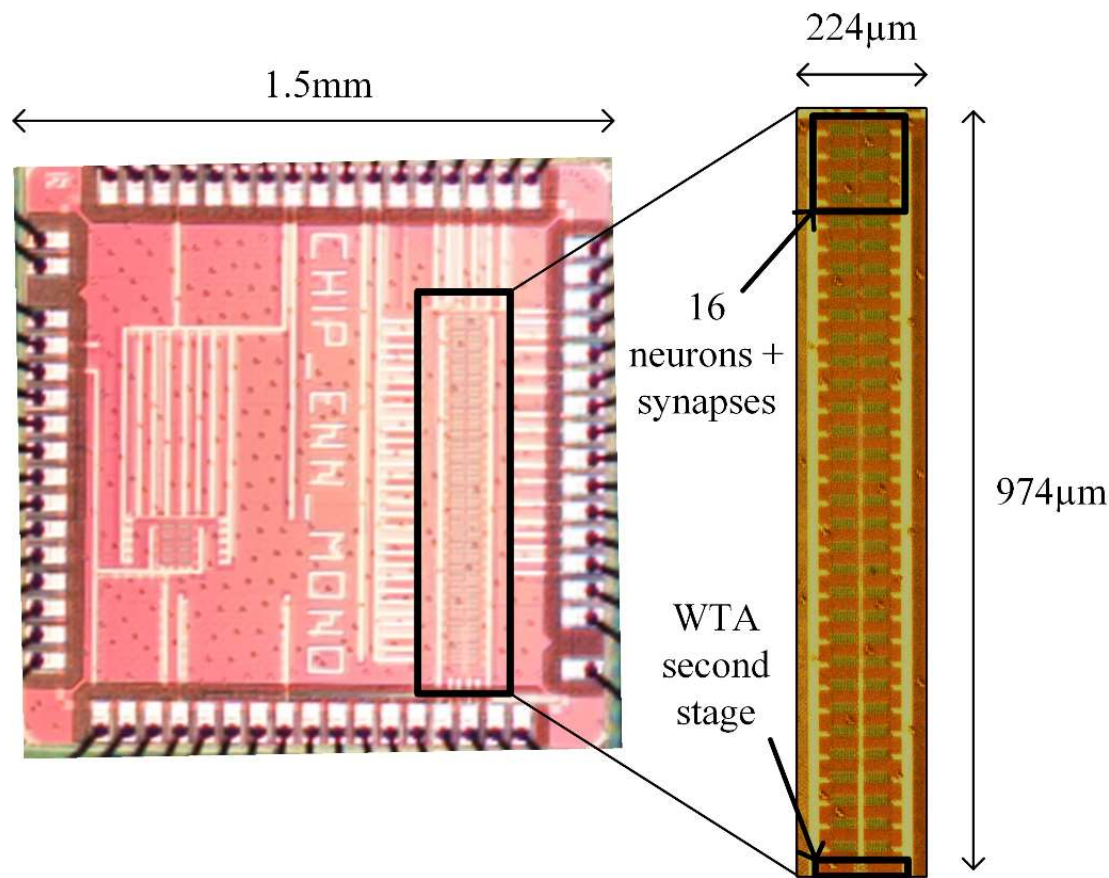UMR CNRS 8520

# Hardware realizations (1/2)



- 5 cluster of 6 neurons => 30 neurons

- Hardware connections => asynchronous

- Control signals generated by an FPGA

| Technology node | 65-nm CMOS |
|---|---|
| Silicon area occupation | 0,019 mm² |
| Supply voltage | 1 V |
| Synaptic current | 300 nA |
| Static current | 5,4 µA |
| Network response time | 58 ns |
| Energy consumption per synaptic event per neuron | 48 fJ |

[Larras, TCAS-I 2016]

iemn
Institut d'Electronique, de Microélectronique et de Nanotechnologie
UMR CNRS 8520

32

# Hardware realizations (2/2)



1.5mm

224µm

16 neurons + synapses

974µm

WTA second stage

CHIP_ENN_MONO

[Larras, TCAS-I 2019]
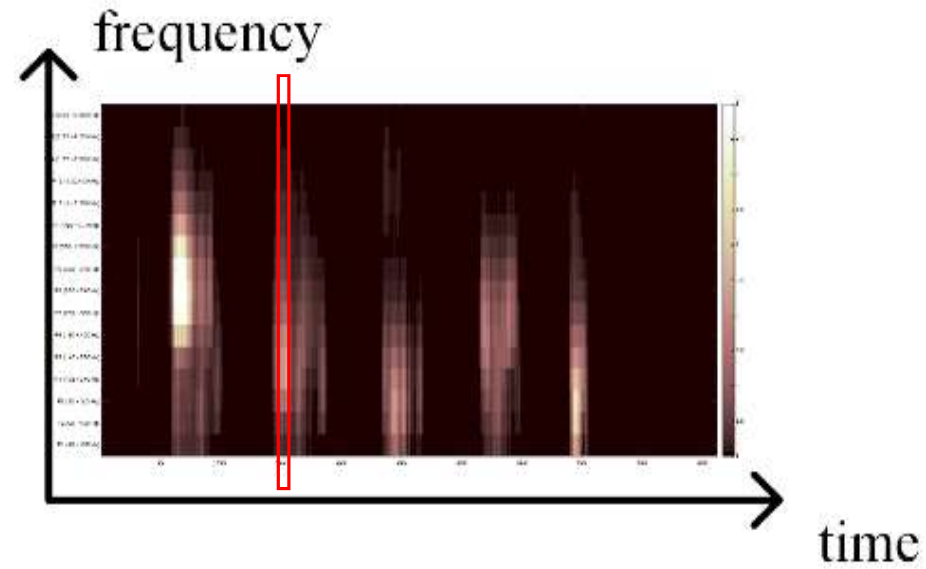
- One cluster of 128 neurons
- Time multiplexing
- Maximum of 3968 emulated neurons
- Driven by an FPGA

| Technology node | 65-nm CMOS |
|---|---|
| Silicon area occupation | 0,21 mm² |
| Supply voltage | 1 V |
| Synaptic current | 300 nA |
| Static current | 23,4 µA |
| Cluster response time | 83 ns |
| Energy consumption per synaptic event per neuron | 115 fJ |

iemn
Institut d'Electronique, de Microélectronique et de Nanotechnologie
UMR CNRS 8520

# Envisionned scheme

# Envisionned scheme

frequency



| 75 Hz | 324 Hz | 1,5 kHz | 3,2 kHz | 4,8 kHz |

- One feature = one cluster

# Envisionned scheme



frequency

75 Hz    324 Hz    1,5 kHz    3,2 kHz    4,8 kHz

- One feature = one cluster
- One neuron per quantization level

iemn
Institut d'Electronique, de Microélectronique
et de Nanotechnologie
UMR CNRS 8520

# Envisionned scheme



frequency

75 Hz     324 Hz     1,5 kHz     3,2 kHz     4,8 kHz

« [ə] »

- One feature = one cluster
- One neuron per quantization level
- **Instantaneous detection of speech formants (cliques)**

# Further opportunities

- **Asynchronous formant extraction**
  - Applications: voice activity detection, phonemes detection

- **Data reduction**
  - From 2-D data to 1-D data
  - Use with LSTM stage to extract keywords

- **Circuit integrability ?**

- **Compatibility with real time ?**

iemn
Institut d'Electronique, de Microélectronique
et de Nanotechnologie
UMR CNRS 8520

# Challenges

- **Feature extraction unit**
  - Event-driven processing with no clocks is difficult to handle and design (concepts, tools)
  - Timing is critical…

- **Classification unit**
  - Generic topology vs. diversity of applications
  - Bridging the gap from theory to efficient hardware

- **Latency and energy consumption!**

- **Integration in advanced CMOS technology**

iemn
Institut d'Electronique, de Microélectronique et de Nanotechnologie
UMR CNRS 8520

# Conclusion

- ANR LEOPAR project targeting a breakthrough in the audio processing domain, in terms of energy efficiency

- Circuit implementation leveraging analog and digital domains

- Targeted hardware demonstration: hardware prototype and integrated circuit in 28-nm FDSOI CMOS

iemn
Institut d'Electronique, de Microélectronique
et de Nanotechnologie
UMR CNRS 8520

# Thank you !

Any questions ? Feel free to ask or send an e-mail to
[benoit.larras@yncrea.fr](mailto:benoit.larras@yncrea.fr)

iemn
Institut d'Electronique, de Microélectronique
et de Nanotechnologie

UMR CNRS 8520