

## TOWARDS NEXT-GENERATION COMPUTING





### TECHNOLOGICAL SOLUTIONS FOR IMPROVING THE FIGURES OF MERIT OF EDGE AI APPLICATIONS

Alexandre VALENTIAN

## leti

## **TRENDS IN AI COMPUTING**



### 

## **TRENDS IN AI COMPUTING**





## The High Cost of Data Movement

Fetching operands costs more than computing on them



Bill Dally, "To ExaScale and Beyond", 2010



## **TRENDS IN EDGE COMPUTING**

Increased computing efficiency

#### Weight quantization

#### Reduced bit accuracy

- Smaller memory footprint
- Lighter operations

#### Variable bit precision

Handling higher bit accuracy when needed

• For higher inference precision

#### Sparsity

Skip MAC operations

• When weight or intermediate result is 0

#### Increased storage efficiency

#### **Near memory computing**

Avoid external memory accesses

#### Weights

- Embedded Non-Volatile Memory Intermediate results
  - SRAM or Embedded DRAM

#### **In-Memory computing**

SRAM or Embedded NVM

Digital or analog



#### **Von Neumann architecture**



#### **In-Memory Computing (IMC) architecture**



Technological solutions for improving the figures of merit of Edge AI applications | imec - Leti workshop@IEDM | December 8, 2019



Multi valued RRAM technology compatible with advanced logic

- Cell density of 40F<sup>2</sup> is available at 28nm
- And compatible to sub-20nm

Roadmap from single bit to multi-bit\* and beyond

• From 1 bit to 4 bit and beyond

\* [T. Wu, ISSCC 2019]

Roadmap for increasing embedded cell density

 From 40F<sup>2</sup> down to 4F<sup>2</sup> thanks to new selector technology

Examples with technology available today

- ResNet50 (74 MB of weights) → 15 mm<sup>2</sup> of memory
- YoloV3 (101MB of weights)  $\rightarrow$  20 mm<sup>2</sup> of memory



28nm RRAM integration



Selector and RRAM integration [IEDM 2019]

### 

## **TRENDS IN AI COMPUTING**





**Back-propagation algorithm** 

- Necessitates to keep all intermediate results (activations)
- With a batch size of more than one
  - To not cycle too much the non-volatile memories

This requires a tremendous amount of activation memory

- Example YoloV3
  - A batch of 20 images requires 800MB of memory





**Advantages** 

Increasing computing & memory capacity

Trends

- « Chipletization » : Generic computing templates, Heterogeneous technologies



# Leti DISTRIBUTED MEMORY-CENTRIC EDGE AI COMPUTING ARCHITECTURE

**Memory-Centric architecture** 

- No more global buffers
- No more power-hungry caches
- Fully distributed memory and control
- Energy efficient use of memory using 3D technology

### Edge Al architecture, using

- Generic PE engines
- Vertically and horizontally connected computing clusters
- In-Memory Computing tiles (IMC)
- Dense NVM for storage
- DRAM for online learning





**Trends in Edge AI applications** 

- Inference first
- Then lifelong local learning

Main challenge is to reduce data movement

This can be solved thanks to a combination of architecture and technology

- Combination of In-Memory Computing
- Non-volatile memory for synaptic weights
- 3D technology for heterogeneous integration

LETI is working to advance those technologies

## THANK YOU FOR YOUR ATTENTION

