# Computational SRAM: Towards Efficient Near-Memory Computing through Tightly Coupled HW/SW Design

***Emanuele Valea***, *Jean-Philippe Noel, Thaddée Bricout, Henri-Pierre Charles, Leo De La Fuente, Bastien Giraud, Maha Kooli, Benjamin Lacour, Manuel Pezzin, Maria Ramirez Corrales*

2nd In-Memory Architectures and Computing Applications Workshop

# Computational SRAM (C-SRAM)

- The proposed **NMC** solution is compatible with any existing memory macros (**off-the-shelf memory compilers, custom design**)

- Coupled with a **Vector Processing Unit (VPU)** executing NMC instructions issued from the host CPU

- A **specific C-SRAM ISA** is supported

- In SoC architectures **data RAM** can be advantageously **replaced by a C-SRAM**

### C-SRAM Instruction Set Architecture

| Category | Mnemonic | Description |
|---|---|---|
| Memory | copy | Copy a line into another |
| | bcast | Broadcast 8/16/32-bit value to the whole Line |
| | hswap | Horizontal 32/64-bit word swap |
| Logical | slli, srli | Shift Left or Right Logical Immediate |
| | (n)and, (n)or, (n)xor | Logical AND, OR & XOR (and negation) |
| Arithmetic | add, sub | Arithmetic 8/16/32-bit Addition & Subtraction |
| | mullo, mulhi | Arithmetic 8-bit integer Multiply |
| | maclo | Arithmetic 8-bit integer Multiply-Accumulate |

### C-SRAM Architecture



### C-SRAM Testchip



- What about programming and **SW compatibility**?

**[1]** M. Kooli *et al.*, "*Towards a Truly Integrated Vector Processing Unit for Memory-bound Applications Based on a Cost-competitive Computational SRAM Design Solution*", ACM JETC, Vol. 18, Issue 2, No. 40, 2022, pp. 1-26.
**[2]** J.-P. Noel *et al.*, "*A 35.6 TOPS/W/mm2 3-stage pipelined computational SRAM with adjustable form factor for highly data-centric applications*", IEEE SSCL, Vol. 3, 2020, pp. 286–289.
**[3]** J.-P. Noel *et al.*, "*Computational SRAM Design Automation using Pushed-Rule Bitcells for Energy-Efficient Vector Processing*", DATE Conference, 2020, pp. 1187–1192.

# HybroGen compilation toolchain (open-source)

- **HybroGen** provides a **compilation flow** from source C code to **parallelized** binary employing specific **NMC instructions**

- **Dynamic compilation** approach in order to **adapt NMC scheduling** to format and size of input data



Dedicated HybroLang language, inspired by C language

[1] M. Kooli *et al.*, "*Towards a Truly Integrated Vector Processing Unit for Memory-bound Applications Based on a Cost-competitive Computational SRAM Design Solution*", ACM JETC, Vol. 18, Issue 2, No. 40, 2022, pp. 1-26.
[4] K. Mambu et al., "*Dedicated Instruction Set for Pattern-based Data Transfers: an Experimental Validation on Systems Containing In-Memory Computing Units*," in IEEE TCAD

# Come to the poster!



- We show some **preliminary results** on several applications:
  - Image processing
  - Artificial Intelligence
  - Cryptography (post-quantum)

|  | Sobel Filter (Edge Detection) | Matrix Multiplication | Frame Differencing | Convolution (Tensor Flow) |
|---|---|---|---|---|
| **Speed-Up** | ~x4.3 | ~x4.9 | ~x7.7 | ~x3.5 |
| **Energy Reduction** | ~x8.3 | ~x13 | ~x10.3 | ~x4 |

[5] Mambu et al. 2022 *Towards Integration of a Dedicated Memory Controller and Its Instruction Set to Improve Performance of Systems Containing Computational SRAM.J. Low Power Electron. Appl.*,12, 18.



**FrodoKEM (PQC)**

- 3% and 15% time reduction compared to the state of the art
- 12% time reduction on key generation function
- 4x speed up on matrix product execution

# Thank you!

Emanuele Valea

Email: emanuele.valea@cea.fr