

Decoding logic implication of in-memory processing in advanced CMOS technology node

Contact Persons: Prof. David Atienza (david.atienza@epfl.ch)

Dr. Alexandre Levisse (alexandre.levisse@epfl.ch)

Project Description

With the incoming innovations in *Artificial Intelligence* (AI) such as self-driving cars or natural language recognition, *Convolutional Neural Networks* (CNN) have been gaining popularity. Their main asset lays in their natural capability to process data with spatial or temporal interrelationships. However, nowadays CNNs implementations are extremely energy hungry, forbidding their use in embedded systems.

In this context, the ESL laboratory proposed an innovative solution featuring *In-Memory Processing* (IMP) to execute CNNs inside an embedded system. The solution developed in ESL relies on a breakthrough trendy approach named *Bitline Computing* and consists in performing logic operations directly in the cache periphery. It enables fast, low power and low area overhead computation, but opens new questions at all the design levels.

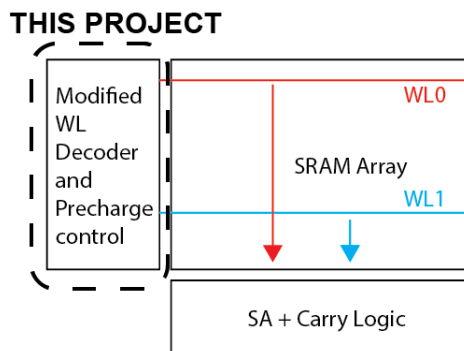


Figure 1: operation of an in-memory processing array with highlighted project objectives.

In this project, we propose to focus on the design of the IMP memory under development in ESL. And particularly on one of the new features of IMP solutions, namely the parallel access of several WordLines (WL) in a SRAM memory array (cf Figure 1).

This project proposes to explore various topologies of WL decoders (such as multiple input static and dynamic programmable decoders or multiple step decoders), benchmark their performances for the proposed memory and then to propose for each, a physical layout exploration. For this purpose, the student will use an advanced CMOS technology PDK (28nm HPM technology) and industrial tools such as Cadence Virtuoso tool for circuit and layout

edition, Synopsys (Calibre, Hspice) and Mentor Graphic (Eldo) tools for electrical simulations and post layout verification.

The project will be carried out at the Embedded Systems Laboratory of EPFL under the supervision of Prof. David Atienza, Dr. Alexandre Levisse and PhD student William Simon.

Project objectives:

1. Understanding of the memory operation and features.
2. Exploration of the standard static and dynamic decoding logic topologies for parallel address decoding
3. Exploration of the scaling-up behavior and benchmark of each proposed solution.
4. If the previous objectives are filled, this work will be part of a submission in a peer-reviewed conference or journal.

Required knowledge and skills:

- Good understanding of MOS transistors behavior
- Advanced knowledge on transistor-level circuit design
- Good analytical skills
- Good background on computer architecture

Appreciated skills:

- Scientific curiosity
- Good communication skills
- Advanced English
- Autonomous work ability

Type of work: 20% theory analysis, 80% hardware design and exploration