# ReCom: An Efficient Resistive Accelerator for Compressed Deep Neural Networks
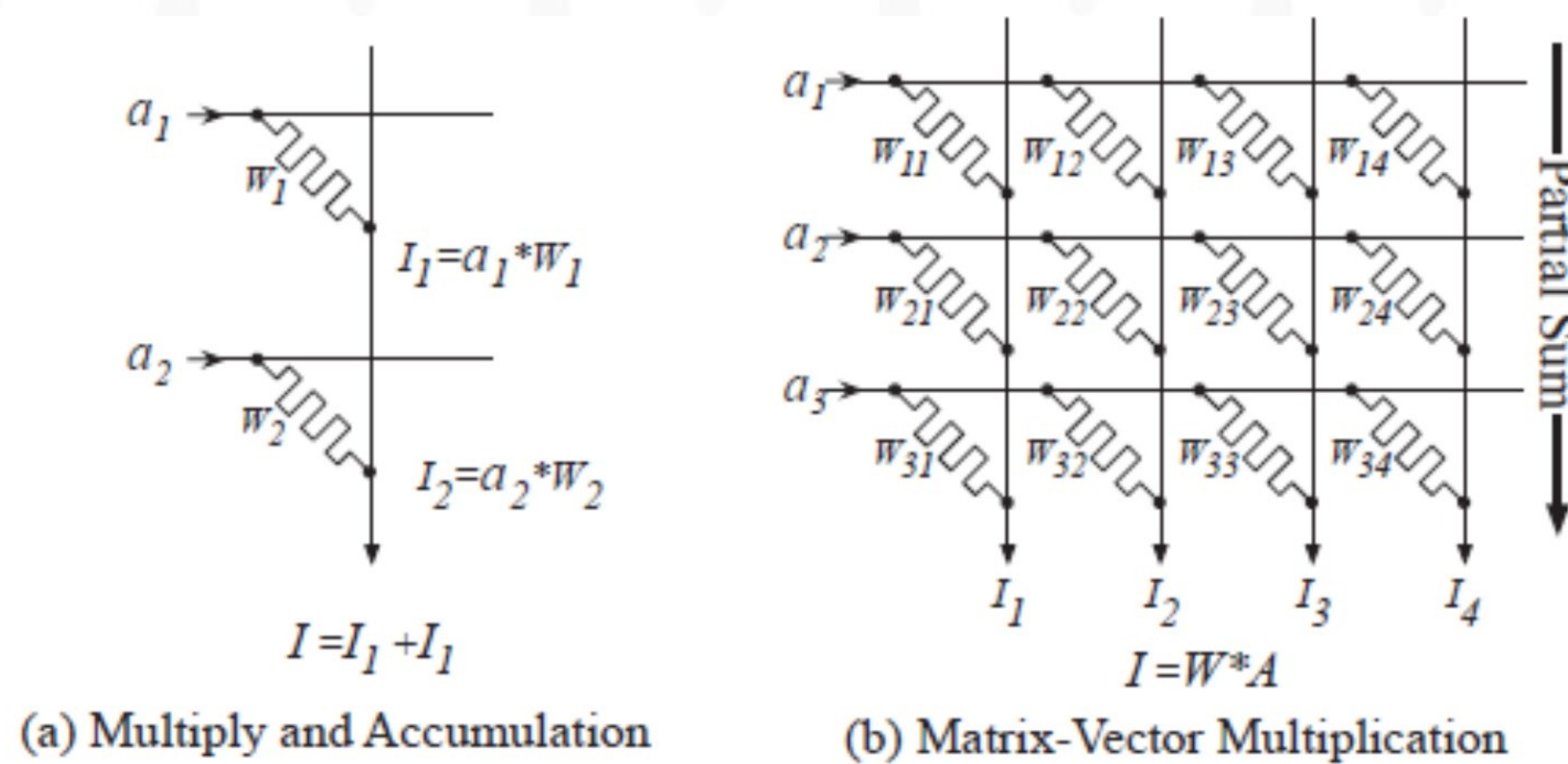
Houxiang Ji[1], Lianghao Song[2], Li Jiang[1], Hai(Halen) Li[2], Yiran Chen[2]
1. Department of CSE, Shanghai Jiao Tong University
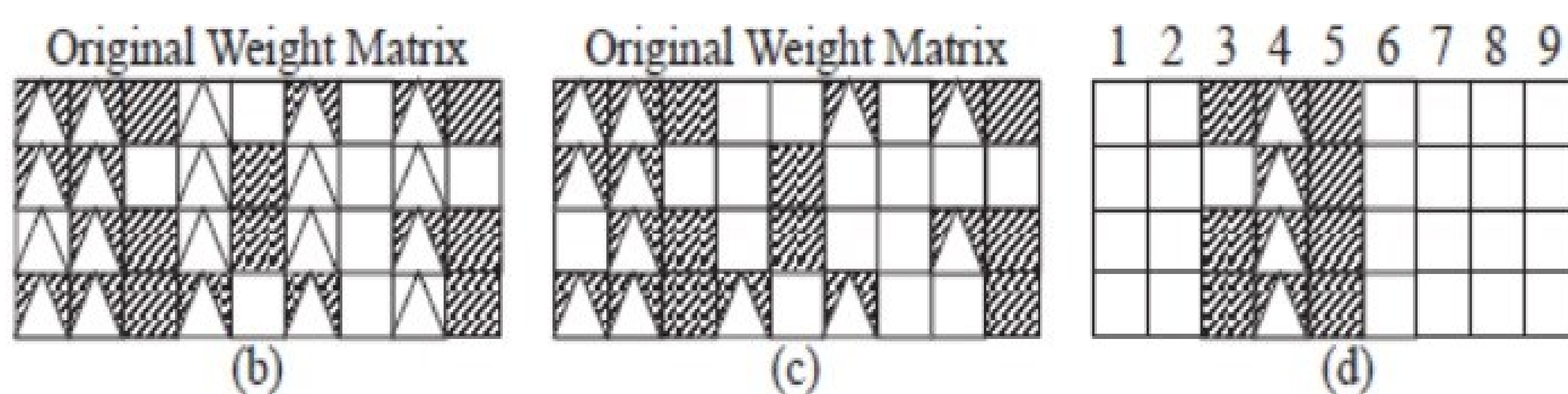2. Department of ECE, Duke University
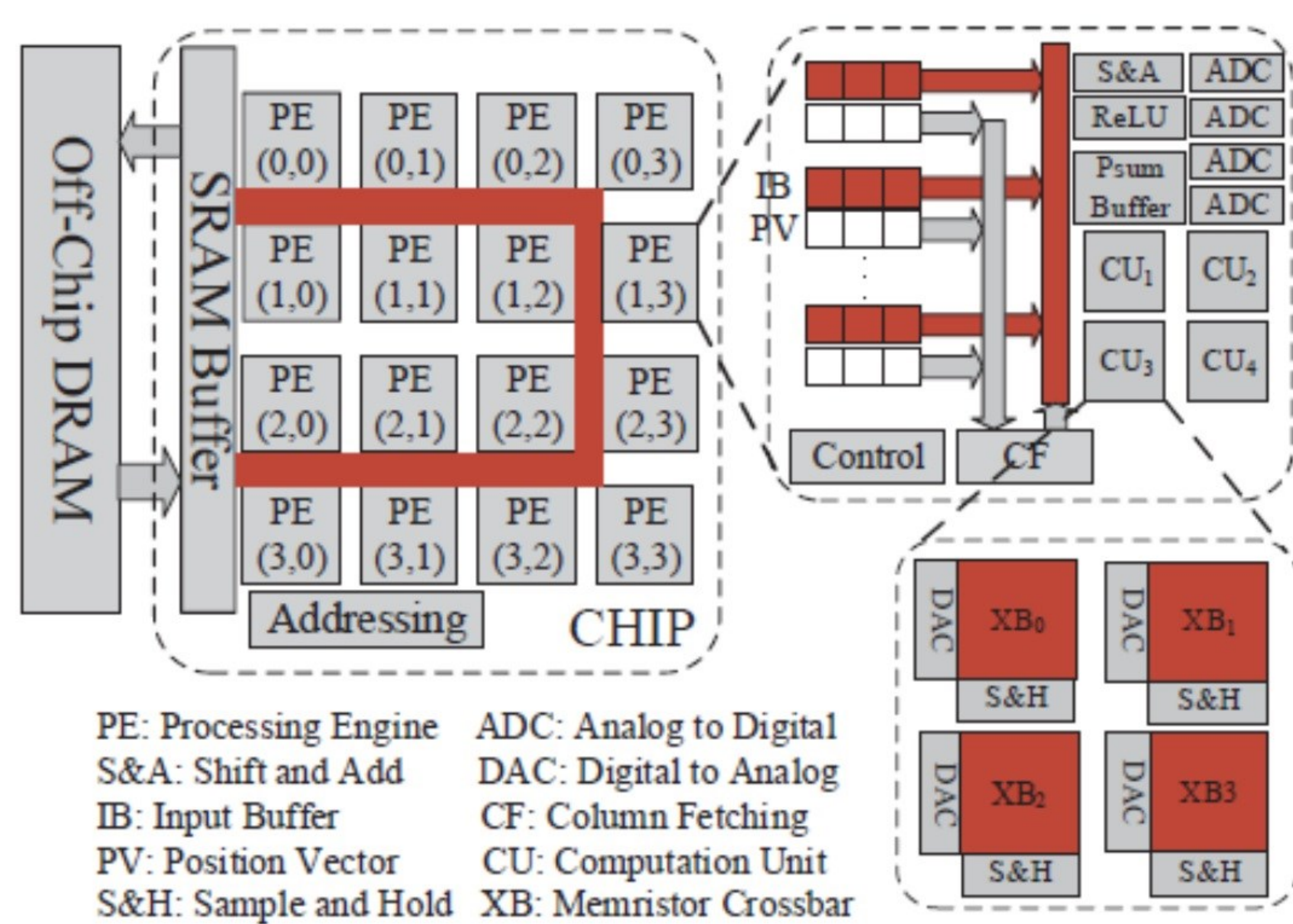
## I. Introduction & Motivation

1. ReRAM-based Accelerator for Deep Learning Algorithms



(a) Multiply and Accumulation    (b) Matrix-Vector Multiplication

## 2. Sparse Neural Network Acceleration



## II. ReCom Architecture



PE: Processing Engine    ADC: Analog to Digital
S&A: Shift and Add    DAC: Digital to Analog
IB: Input Buffer    CF: Column Fetching
PV: Position Vector    CU: Computation Unit
S&H: Sample and Hold    XB: Memristor Crossbar

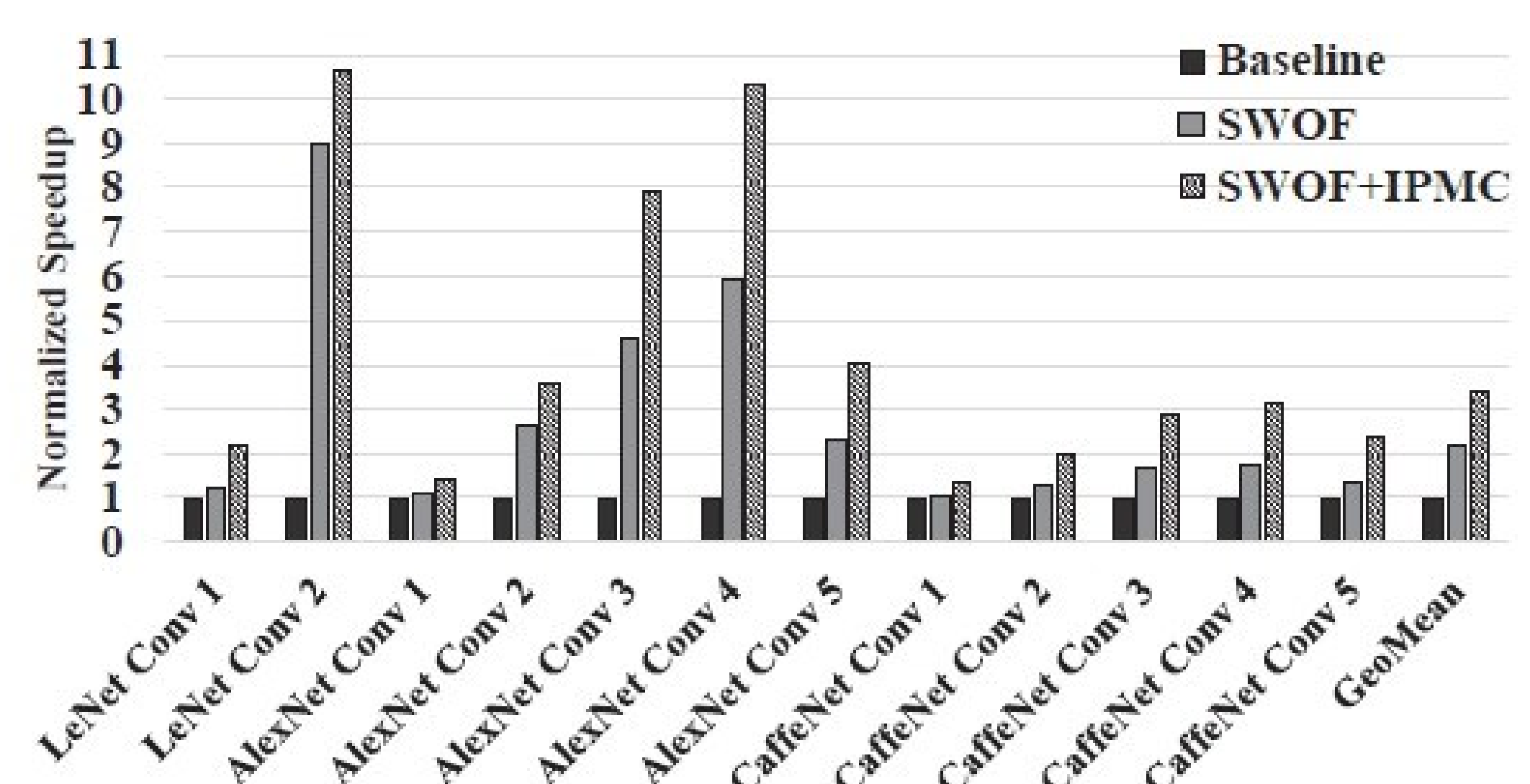1. Structural Compression on Weight Matrix

2. Structurally-compressed Weight Oriented Fetching (SWOF)

3. In-layer Pipeline for Memory and Computation (IPMC)

## III. Evaluation

3.37x speedup (up to 10.66x)
2.41x energy saving (up to 9.43x)



## IV. Conclusion

RECOM, the *first* accelerator to support the sparse DNN processing in ReRAM. Deep neural networks are compressed structurally by specific regularization on each layer with little or no loss in accuracy. Our experiments show notable improvement on processing speed and energy-saving.