

# CR2: A Compressed ReRAM-based DNN Accelerator by Combining Computation and Read operation



香港中文大學(深圳)  
The Chinese University of Hong Kong, Shenzhen

Shi-Hao Hong and Yeh-Ching Chung  
The Chinese University of Hong Kong, Shenzhen

## Background

- ReRAM-based PIM: A non-von neumann architecture performing Matrix-Vector Multiplication in memory cell arrays.

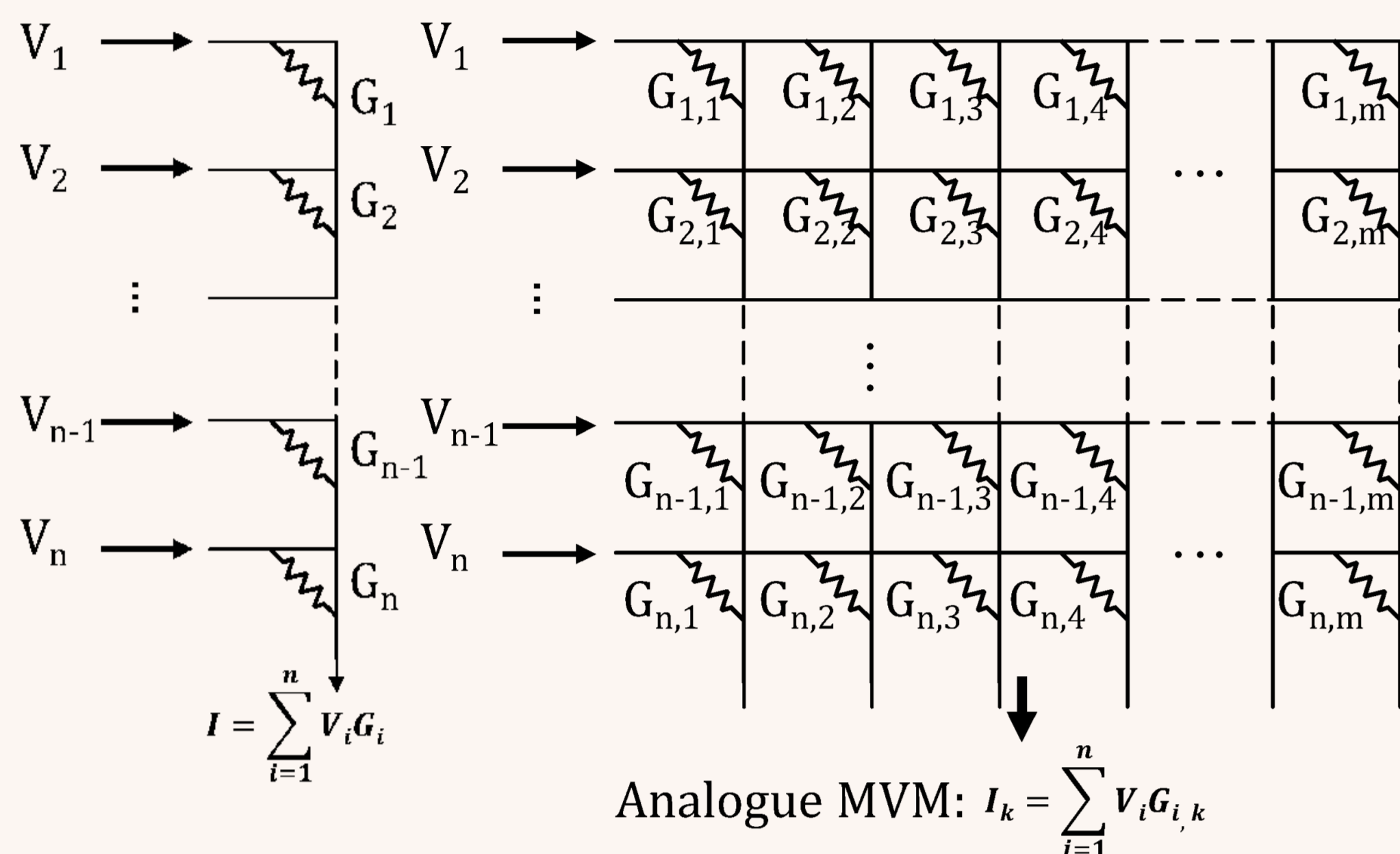
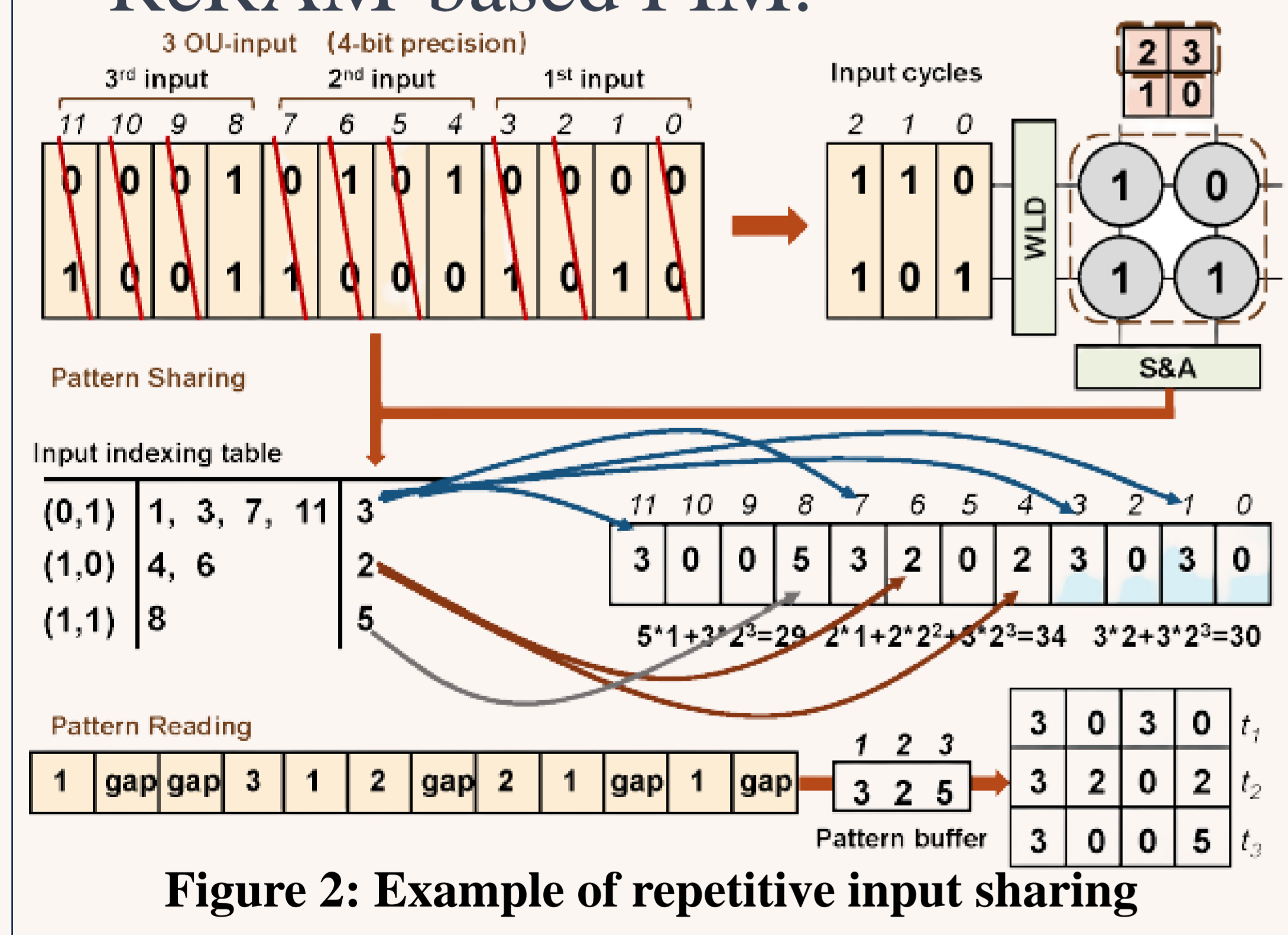


Figure 1: ReRAM-based MVM computation

- Operation Unit (OU): Due to hardware limits, it is necessary to split an MVM into multiple operation units (OUs).
- Sparsity pruning and repetition sharing are common methods to compress crossbar arrays in ReRAM-based PIM.

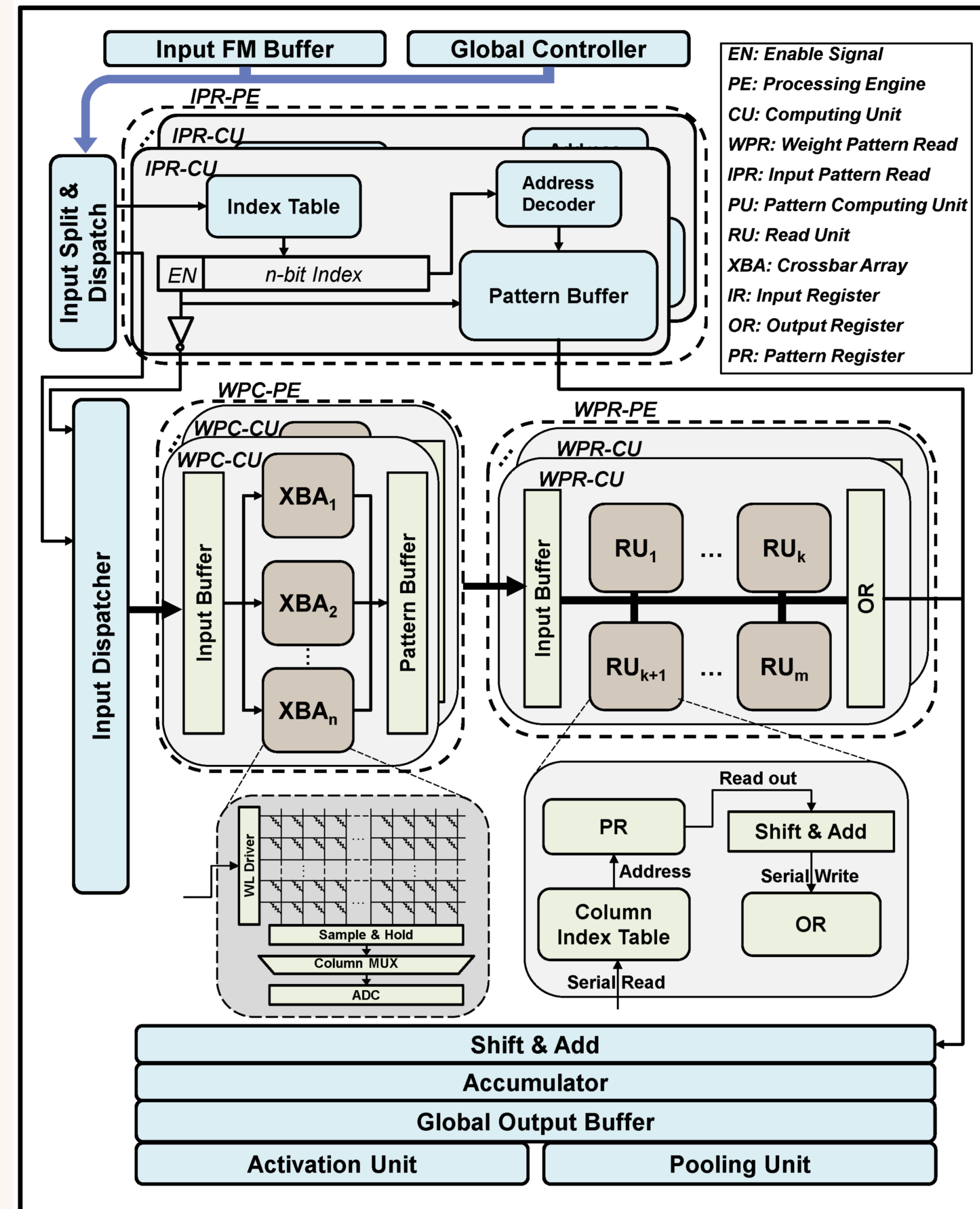


## Motivation

Current studies have not explored the concentrated distribution of inputs. Preliminary experiments reveal that certain input patterns are frequently repeated while changing datasets. It allows us to replace repetitive computations with reading buffer.

## Design of CR2

- Architecture: CR2 consists mainly of weight pattern compute (WPC), weight pattern reuse (WPR) and input pattern reuse (IPR).



Size of Learning data (SL), Size of Testing data (ST) are parameters in learning allocation strategy

we propose a buffer allocation strategy tailored to repeating input patterns at each layer

$$\max \sum_j P_j [B_j] \text{ s.t. } B_j \geq 0, \sum_j B_j W_j \leq B$$

Recurrence searching:  $P_{max}(i, w) = \max(P_{max}(i-1, w - kwi) + Pi[k]), k \in [0, \min(xi, \lfloor \frac{w}{w_i} \rfloor)]$

## Evaluation

Benchmark: Imagenet task  
Models: AlexNet, VGG16 ResNet50 and GoogLeNet, where weights are quantized to 8-bit.

