

# Cambricon-C: Efficient 4-bit Matrix Unit via Primitivization

Yi Chen

USTC & ICT, CAS



# Cambricon-C: Efficient 4-bit Matrix Unit via Primitivization

Yi Chen	Yongwei Zhao	Yifan Hao	Yuanbo Wen	Yuntao Dai
<i>USTC SKLP, ICT, CAS</i>	<i>SKLP, ICT, CAS</i>	<i>SKLP, ICT, CAS</i>	<i>SKLP, ICT, CAS</i>	<i>USTC</i>
chen_yi@mail.ustc.edu.cn	zhaoyongwei@ict.ac.cn	haoyifan@ict.ac.cn	wenyuanbo@ict.ac.cn	daiyuntao@mail.ustc.edu.cn

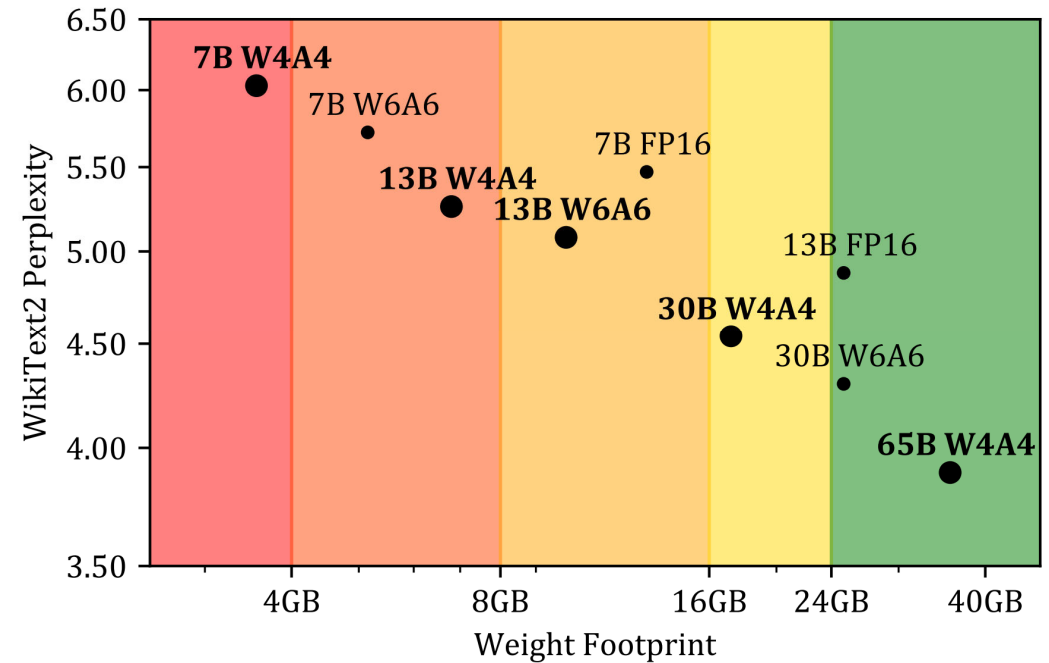
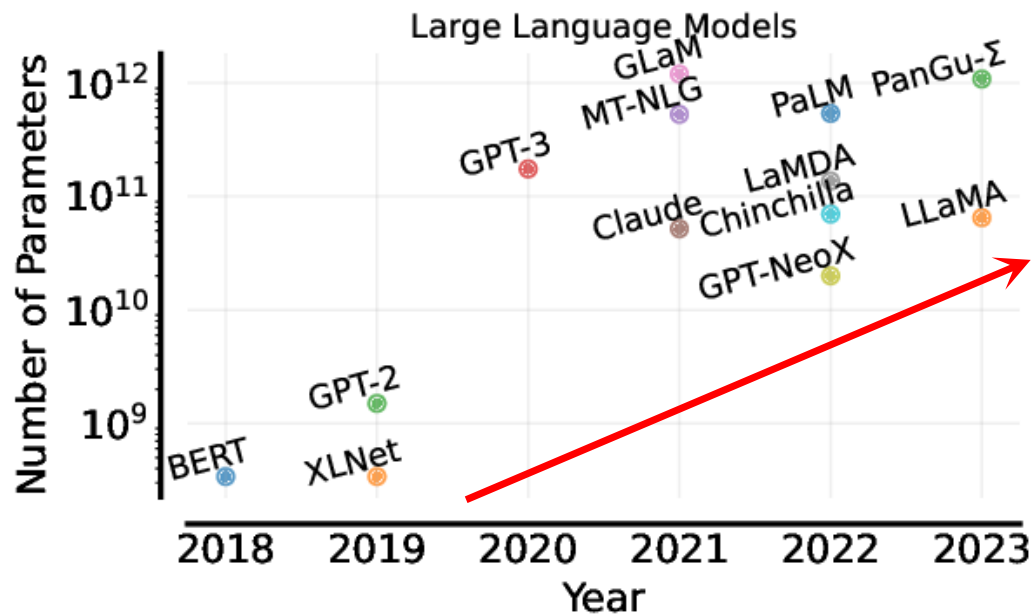
Xiaqing Li	Yang Liu	Rui Zhang	Mo Zou	Xinkai Song
<i>SKLP, ICT, CAS</i>	<i>SKLP, ICT, CAS UCAS</i>	<i>SKLP, ICT, CAS</i>	<i>SKLP, ICT, CAS</i>	<i>SKLP, ICT, CAS</i>
lixiaqing@ict.ac.cn	liuyang22z1@ict.ac.cn	zhangrui@ict.ac.cn	zoumo@ict.ac.cn	songxinkai@ict.ac.cn

Xing Hu	Zidong Du	Huaping Chen	Qi Guo*	Tianshi Chen
<i>SKLP, ICT, CAS SHIC</i>	<i>SKLP, ICT, CAS SHIC</i>	<i>USTC</i>	<i>SKLP, ICT, CAS</i>	<i>Cambricon Technologies</i>
huxing@ict.ac.cn	duzidong@ict.ac.cn	hpchen@ustc.edu.cn	guoqi@ict.ac.cn	tchen@cambricon.com

# Outlines

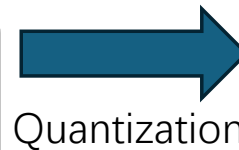
- Background & Motivation
- Algorithm
  - Primitivized Matrix Multiplication
  - Quarter Square Multiplication
- Architecture
  - Overview
  - PE design
  - Converter
- Evaluation and Results
- Conclusion

# Background: Emerging 4-bit Large Model



## Model Size Increases

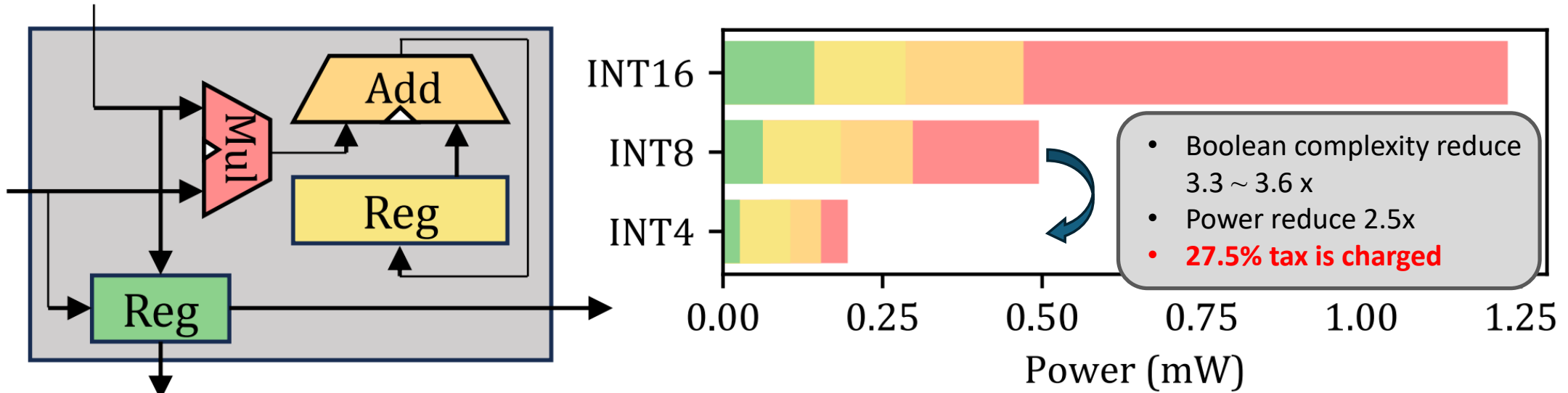
- Model performance is related to model scale
- Model size imposes significant demands on the computational resources



## Model Deployment

- 4-bit variants get best performance under the same memory budget.
- GMMML, vLLM, TensorRT-LLM, Open-VINO

# Background: Hardware Matrix Unit



## 1. Repeating Computation

- INT4 multiplication presents < 256 unique value.
- Matrix dimensions are often far > 256

## 2. Accumulator dominating

- $n$  dimension INT4 matrix needs  $8 + \log_2 n$  bits Accumulator

**27.5% Tax !**

# Algorithm: Primitivized Matrix Multiplication

## Key1: Avoid duplicated multiplications by pre-computing a Look-up Table

Suppose an INT4 Multiplication:

$$y = \sum_{i=0}^{N-1} w_i x_i, \quad N \gg 256$$

Define pre-computed  $\mathbf{m}$  includes all  $w_i x_i$  possible values

$$\mathbf{m} = \begin{bmatrix} 0 \times 0, & 0 \times 1, & 0 \times 2, & \dots, & 0 \times -1, \\ 1 \times 0, & 1 \times 1, & 1 \times 2, & \dots, & 1 \times -1, \\ 2 \times 0, & 2 \times 1, & 2 \times 2, & \dots, & 2 \times -1, \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ -1 \times 0, & -1 \times 1, & -1 \times 2, & \dots, & -1 \times -1 \end{bmatrix}$$



Formula can also be presented as:

$$y = \sum_{i=0}^{N-1} \mathbf{p}_i \mathbf{m}^T,$$

Where  $\mathbf{p}_i$  is the one-hot vector

$$\mathbf{p}_i = [p_0, p_1, \dots, p_{255}], \quad p_j = \begin{cases} 1, & \text{if } j = x_i \circ w_i \\ 0, & \text{otherwise} \end{cases}$$

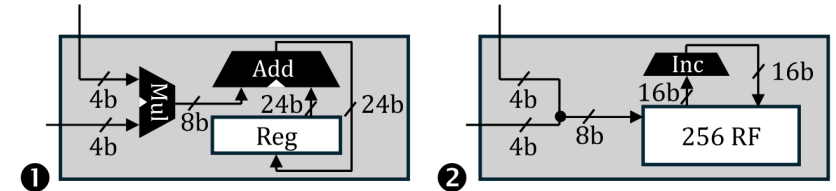
## Key2: Reduce the strength of additions

With the distributive law of multiplication:

$$y = \sum_{i=0}^{N-1} \mathbf{p}_i \mathbf{m}^T = (\sum_{i=0}^{N-1} \mathbf{p}_i) \mathbf{m}^T$$

Major computation part: summation of **one-hot** vector

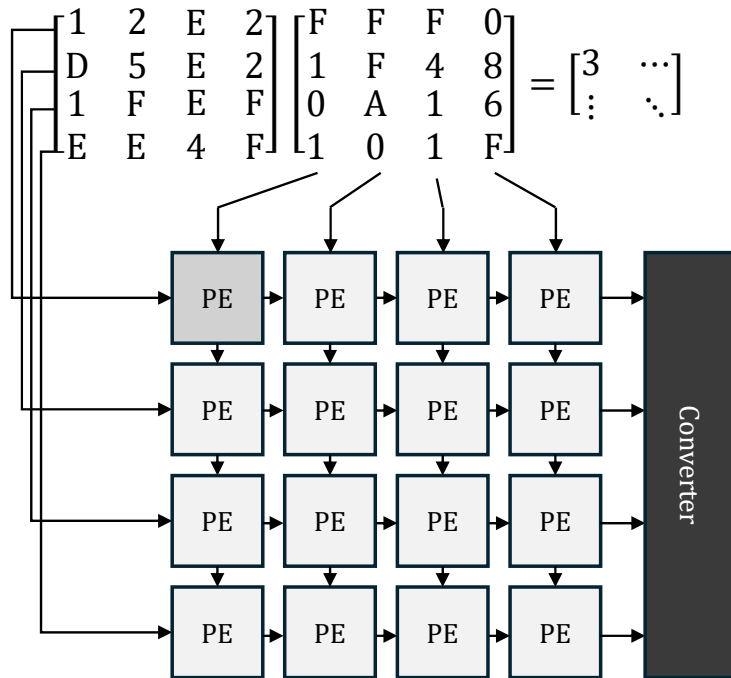
⇒ 256 number increment + final once multiplication



	Tr. Count	Energy	Conventional MM	Primitivized MM
4b Mul	241 T	0.227 pJ	× 11,008	× 256*
24b Add	679 T	0.584 pJ	× 11,007	× 255
16b Inc	122 T	0.028 pJ	/	× 11,008

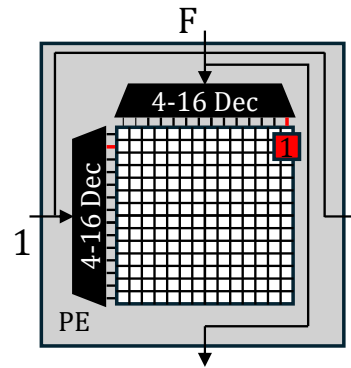
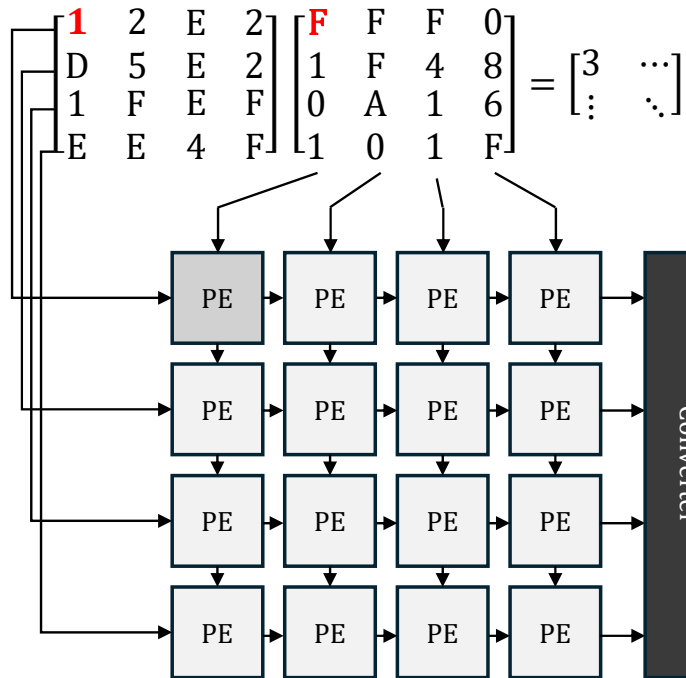
(\* required in the last step multiplying  $\mathbf{m}^T$ , 8b×16b to be precise.)

# Algorithm: Primitivized Matrix Multiplication



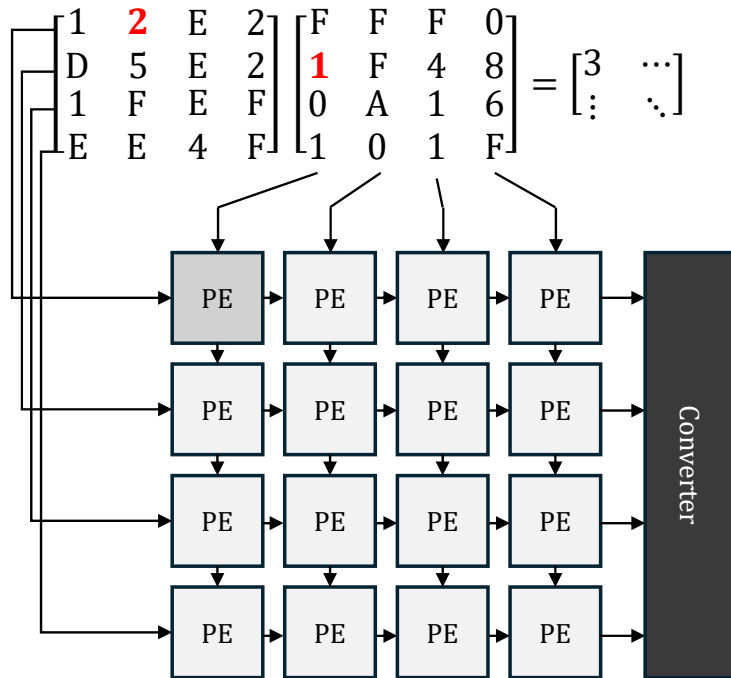
# Algorithm: Primitivized Matrix Multiplication

Compute  $\sum_{i=0}^{N-1} p_i$

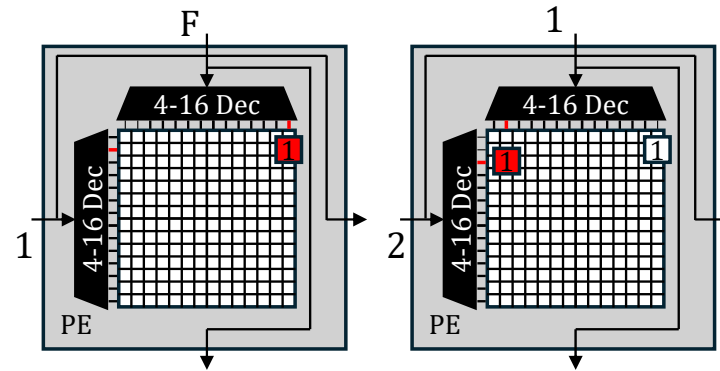




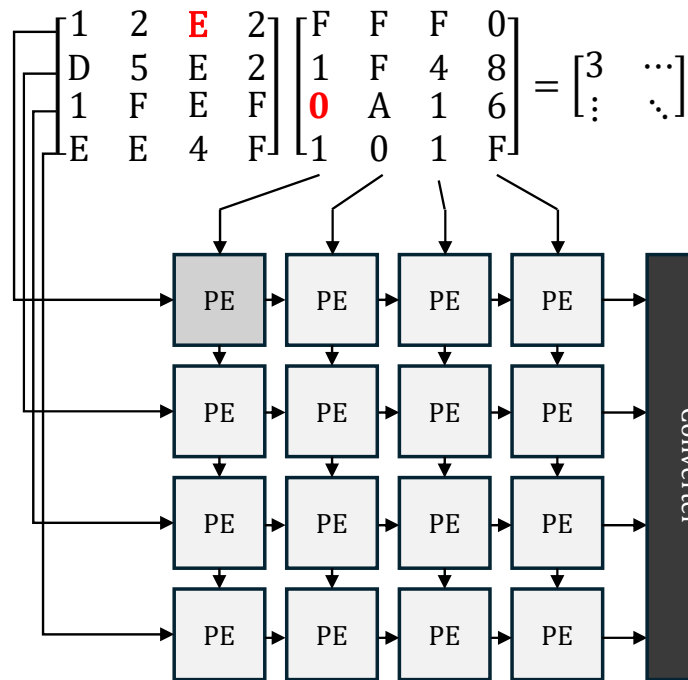
# Algorithm: Primitivized Matrix Multiplication



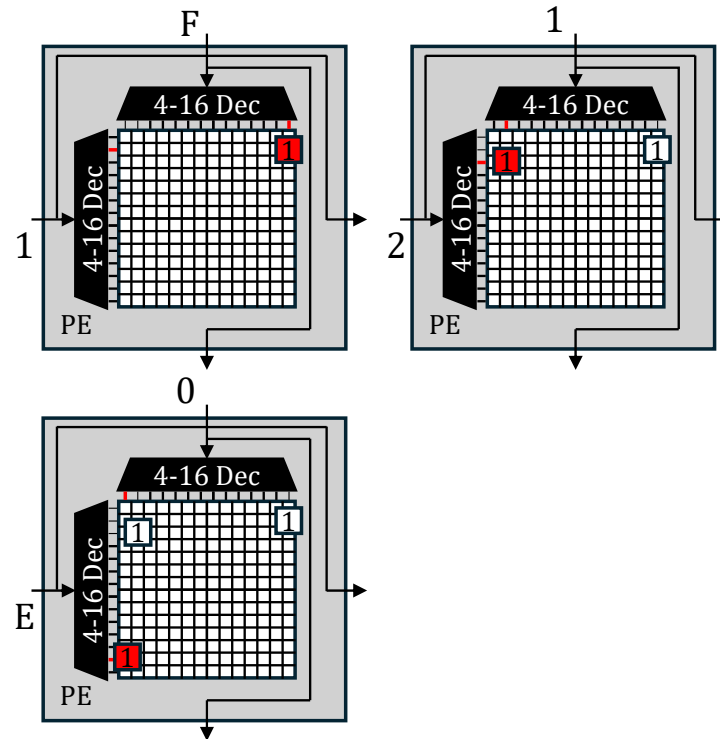
Compute  $\sum_{i=0}^{N-1} p_i$



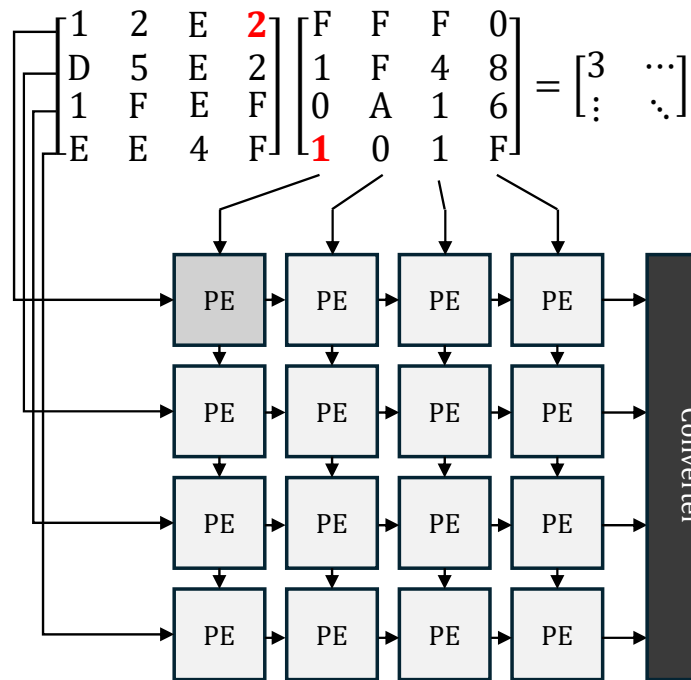
# Algorithm: Primitivized Matrix Multiplication



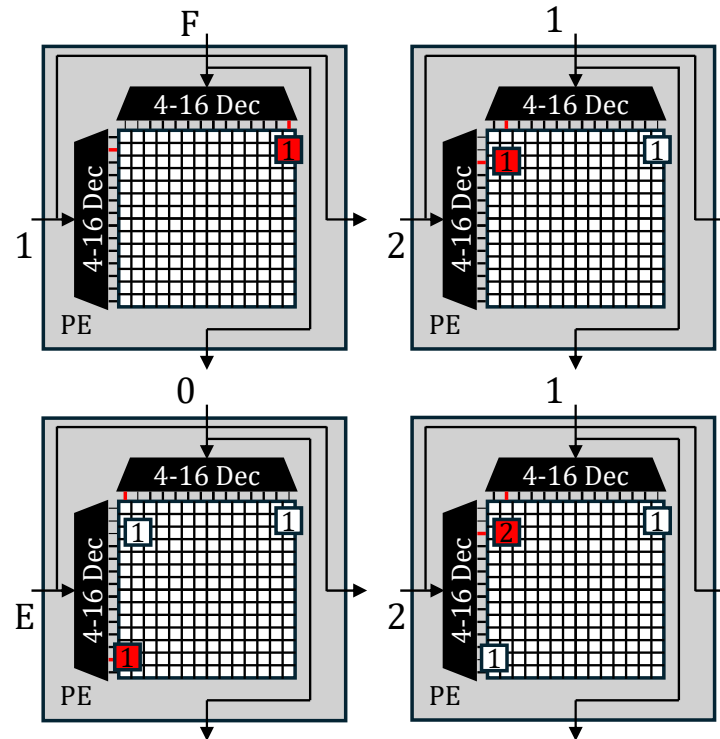
Compute  $\sum_{i=0}^{N-1} p_i$



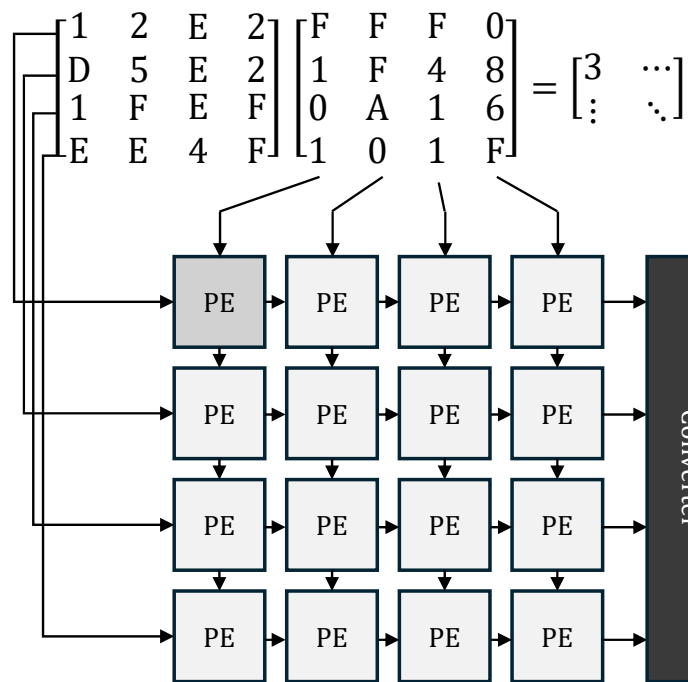
# Algorithm: Primitivized Matrix Multiplication



Compute  $\sum_{i=0}^{N-1} p_i$

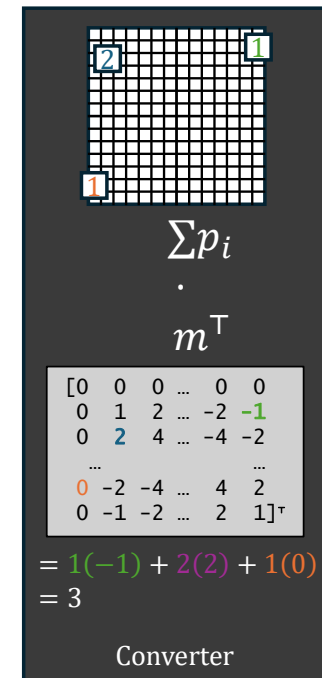
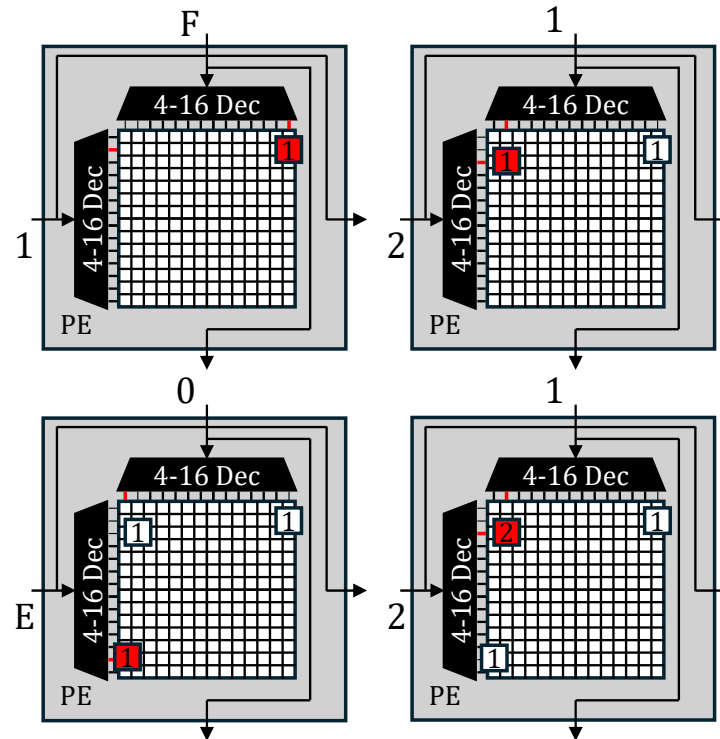


# Algorithm: Primitivized Matrix Multiplication



Compute  $\sum_{i=0}^{N-1} p_i$

Compute  $(\sum_{i=0}^{N-1} p_i) m^T$



# Algorithm: Quarter Square Multiplication

## Key: Reduce the number of counters

It is intuitively obvious that

$$\frac{(x+y)^2}{4} - \frac{(x-y)^2}{4} = xy$$

As  $x$  and  $y$  are integers.

$\Rightarrow x+y$  and  $x-y$  are both odd or even

$\Rightarrow$  The remainders will be cancelled.

$$\text{i.e. } (x+y)^2 \equiv (x-y)^2 \pmod{4}$$

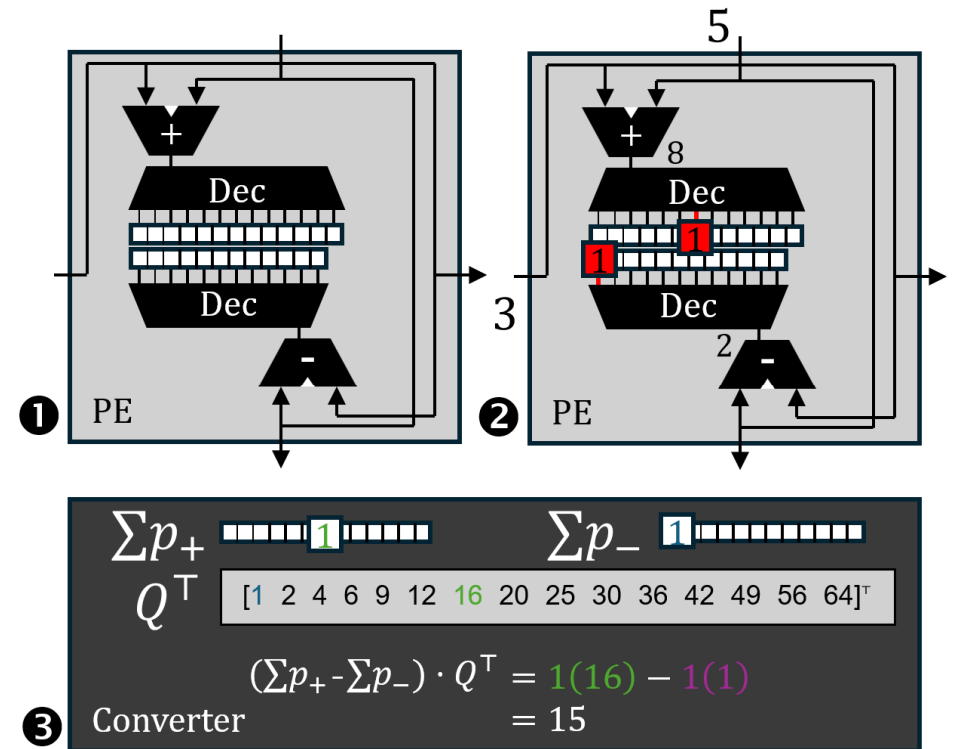
$$\Rightarrow \left\lfloor \frac{(x+y)^2}{4} \right\rfloor - \left\lfloor \frac{(x-y)^2}{4} \right\rfloor = xy, \quad x, y \in \mathbb{Z}$$

The pre-compute look-up table becomes  $Q(n) = \lfloor n^2/4 \rfloor$

As  $Q(0) = Q(1) = 0$  and  $Q(n) \equiv Q(-n)$

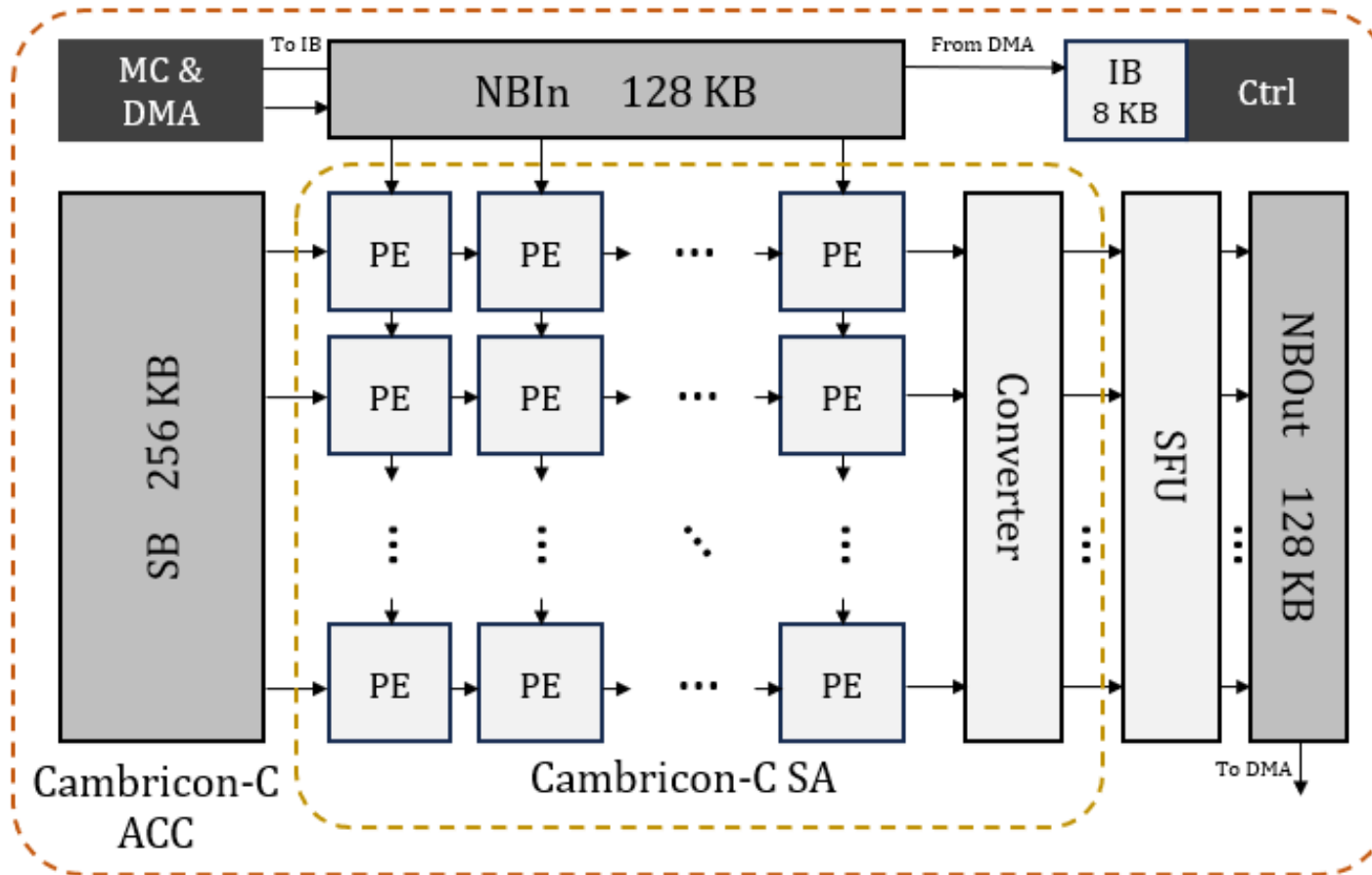
$x+y \in [-16, 14]$ ,  $x-y \in [-15, 15]$ ,

$$\Rightarrow Q = [1, 2, 4, 6, 9, 12, 16, 20, 25, 30, 36, 42, 49, 56, 64]$$



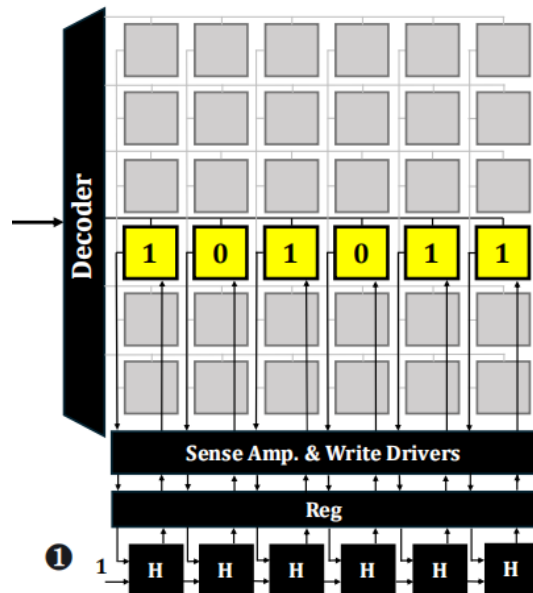


# Architecture: Overview



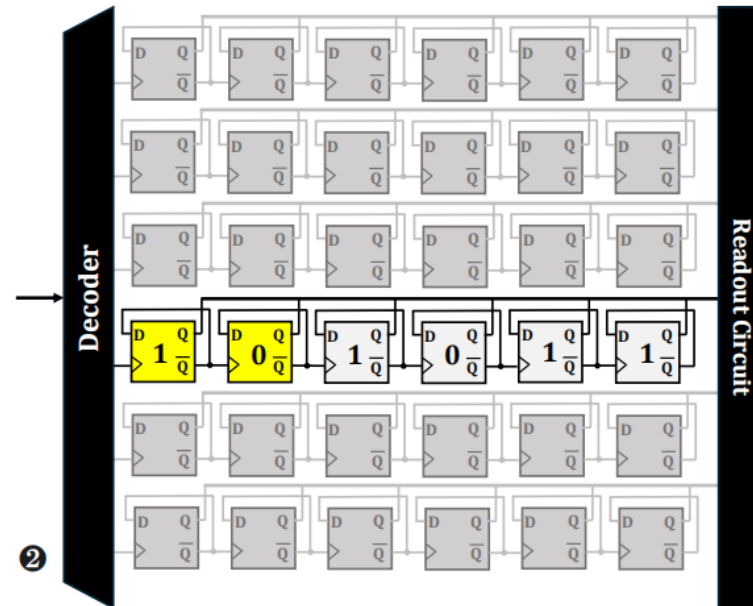
- Cambricon-C ACC
  - Cambricon-C SA
    - 32x32 PE array
    - Converter
  - 128KB NBIn
  - 128KB NBOut
  - 256KB SB
  - Controller (Ctrl)
  - SFU
  - 8KB Instruction buffer (IB)
  - Direct memory access module (DMA)

# Architecture: Counters



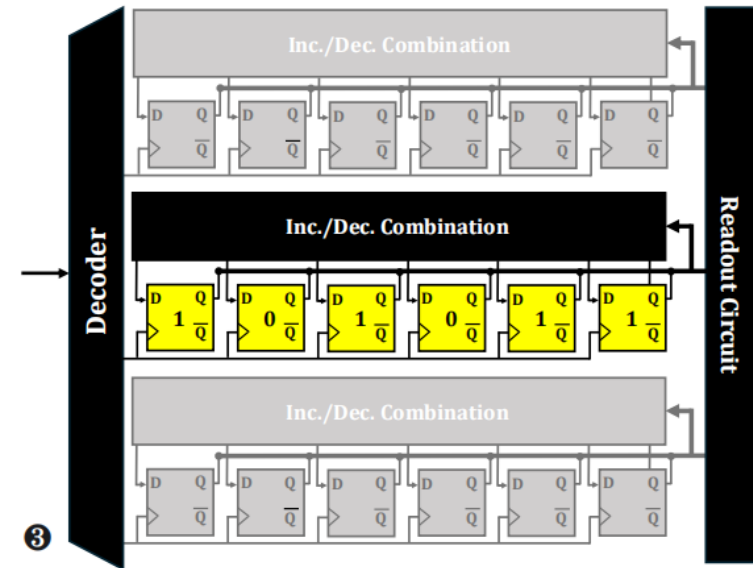
## SRAM-based Counters

- + Lower leakage power
- + Lower area
- Inefficient incrementing
- Hard to customize



## Ripple Counters

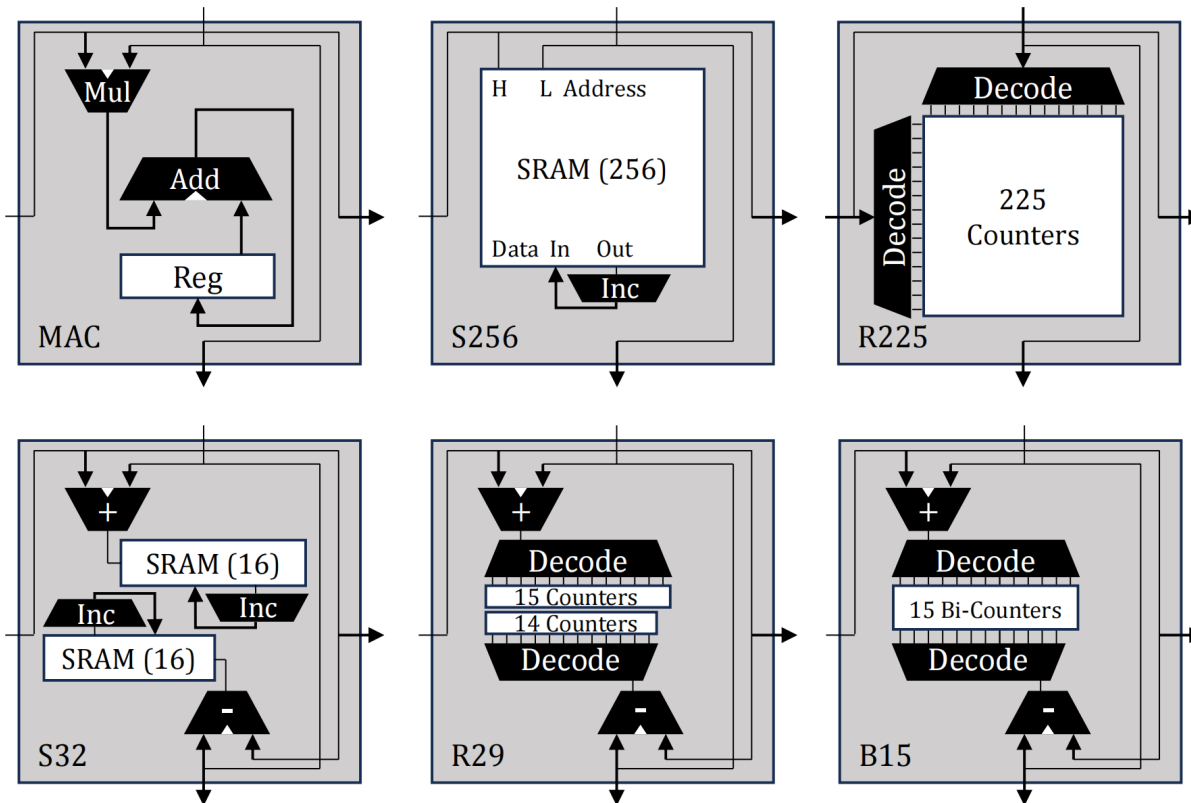
- + Efficient incrementing
- + High speed
- + Easy to customize
- Large area and leakage power
- Unstable while incrementing



## Bi-directional Counters

- + Combine up/down counters
- + Easy to customize
- High leakage power
- Less efficient.

# Architecture: PE Designs

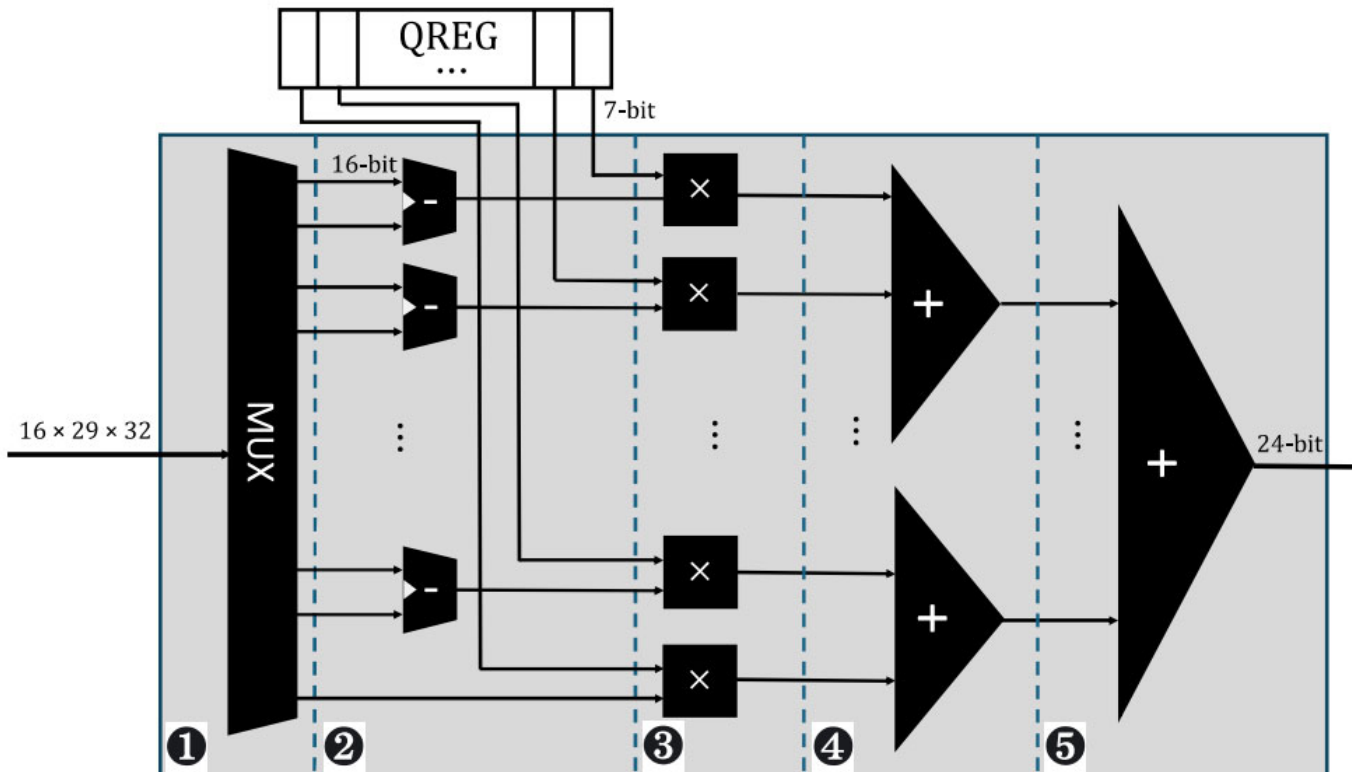


PE Design	Area ( $\mu\text{m}^2$ )	Energy (pJ)	Frequency (MHz)
MAC	<b>1422</b>	2.508	357
S256	21480	3.367	<b>400</b>
R225	57113	2.525	<b>400</b>
S32	13045	2.439	<b>400</b>
R29	8914	<b>1.274</b>	<b>400</b>
B15	9821	2.734	300

- Sacrifice area to gain **lower energy**
- **QSM** significantly increases the PE performance.
- **R29** outperform other schemes
- **S32** and **B15** has larger area and power. As it introduces more combination circuit.



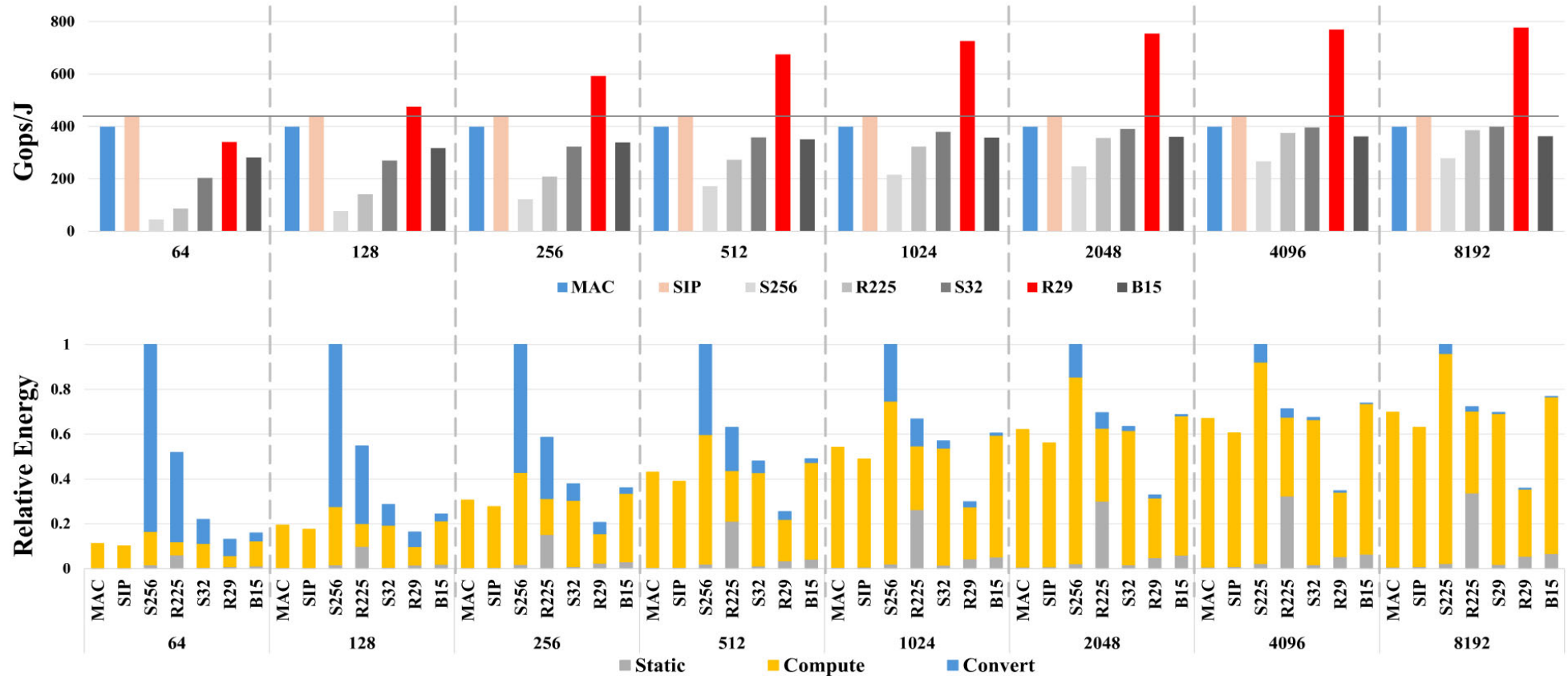
# Architecture: Converter



- Converter for R29-PE array
  - Counter Readout
  - Substraction of counter values
  - Multiplication by Pre-loading Value
  - Final summation

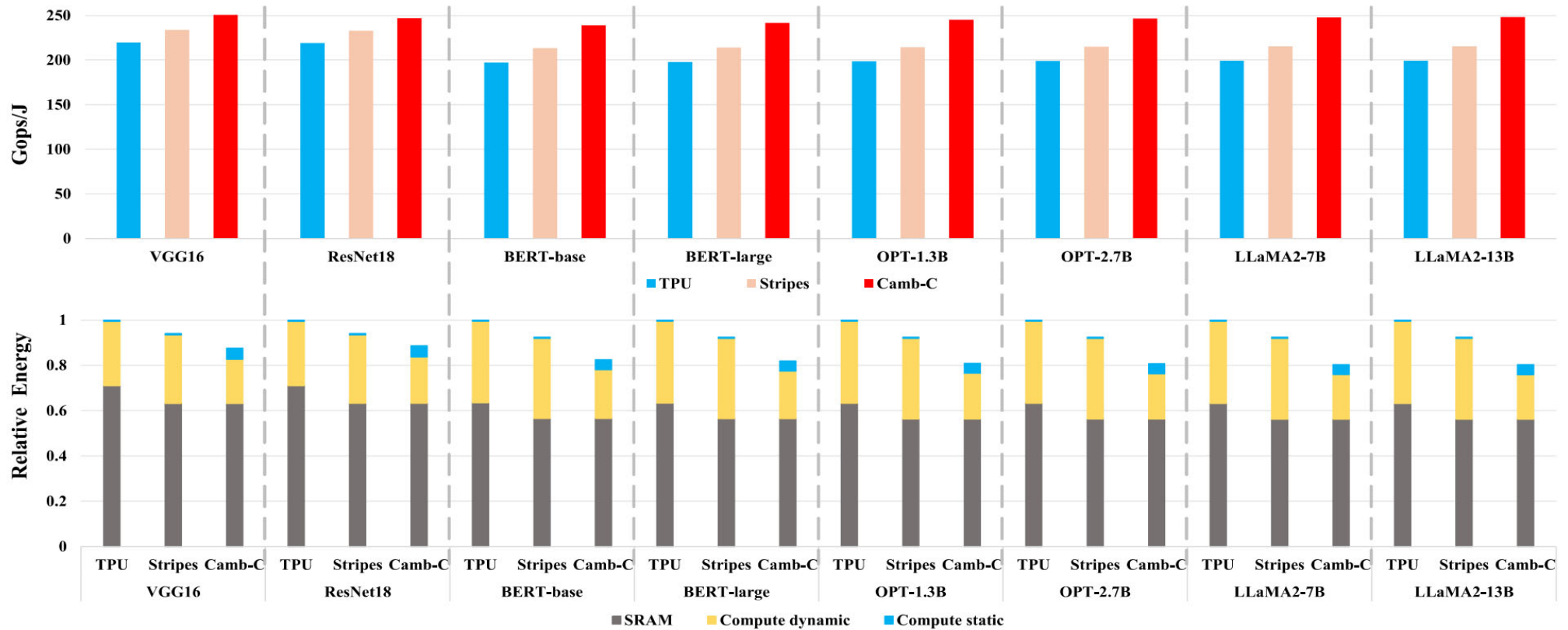


# Evaluation: Matrix dimension



- R29 outperform other PEs when matrix dimensions reach 128.
- The proportion of conversion energy decreases with the dimension increasing.
- The maximum benefit of R29 reaches 1.95x than MAC and 1.76x than SIP

# Evaluation: Performance on DNN



- Cambricon-C consistently outperforms TPU and Stripes in various DNN models.
- Cambricon-C gains 1.15x and 1.25x improvement compared with Stripes and TPU.

# Evaluation: Comparison with Other Methods

Design	Energy (pJ)	Frequency (MHz)
MAC [33]	2.508	357
SIP [23]	2.265	400
PIP [1]	2.549	400
FuseKNA [40]	6.450	400
CARAT [34]	6.638	400
R29	<b>1.273</b>	400

- PIP: Skip zero-bit in bit-serial compute.
  - Extra decoder and shifter
- FuseKNA: compute with Look-up table
  - LUT-based AGU involves multiple registers
- CARAT: Avoid repeating value computation
  - Suitable for float point data, but not INT.
  - longer array size to shed through put loss.

The circuit design for INT4 Computation is quite simple and neat, making it not easy to achieve net gains.

# Conclusion

- Algorithm:
  - Primitivized matrix multiplication
  - Quarter square multiplication
- Architecture
  - 5 different PE design
  - Converter
  - Cambricon-C accelerator
- Evaluation and Results
  - R29 gets 1.95x energy saving than MAC
  - Cambricon-C ACC get 1.25x energy saving than TPU



MICRO 2024

---

# Thanks!

## Q & A

Feel free to contact me on  
Whova app or e-mail: [chen\\_yi@mail.ustc.edu.cn](mailto:chen_yi@mail.ustc.edu.cn)