# A Convolutional Neural Network Accelerator Architecture with Fine-Granular Mixed Precision Configurability

**Xian Zhou** , Li Zhang , Chuliang Guo , Xunzhao Yin , Cheng Zhuo

College of Information Science and Electronic Engineering, Zhejiang University, Hangzhou, China
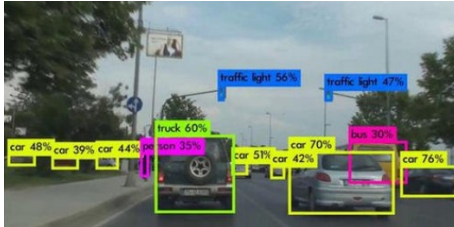
浙江大学
Zhejiang University

# Outline

- Motivation

- Architecture and Dataflow

- Optimization and Tradeoff

- Experimental Results

# CNN in Mobile Applications

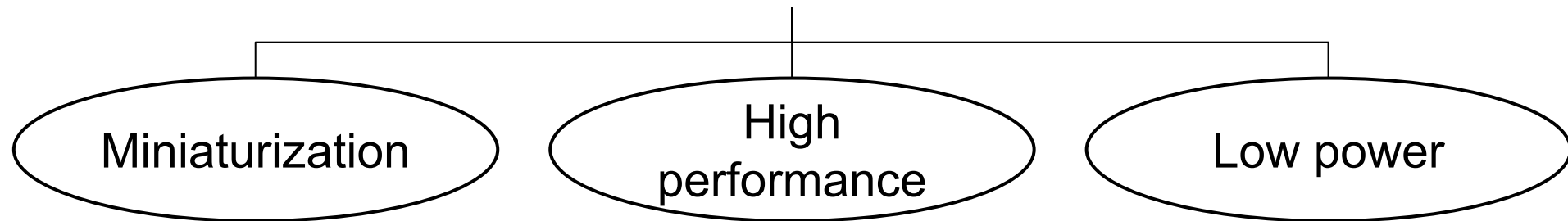CNN has been widely deployed in various deep learning domains

Traffic

Industry

Monitor



Characteristics

Miniaturization

High performance

Low power

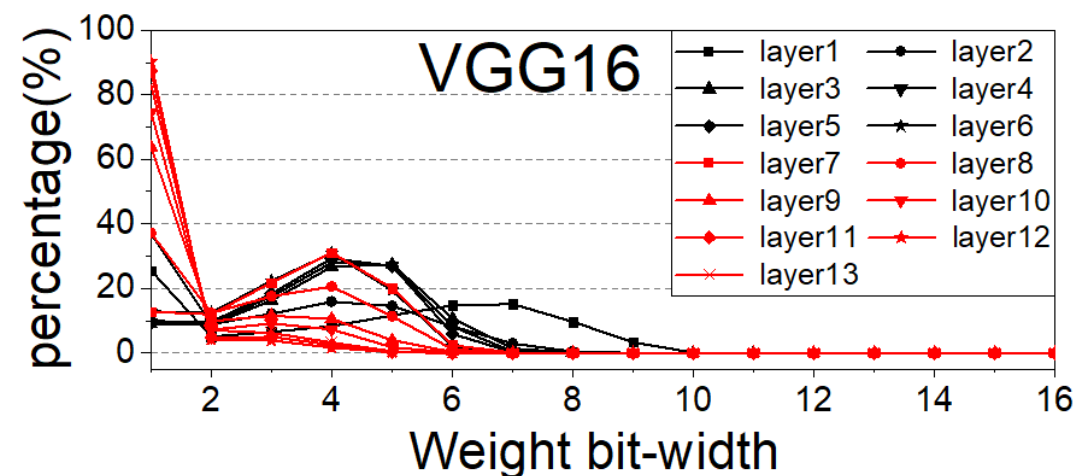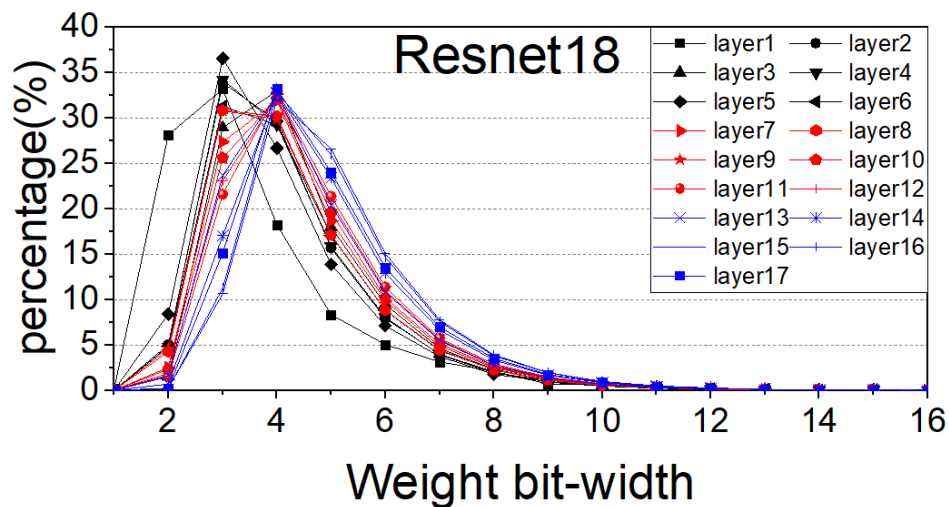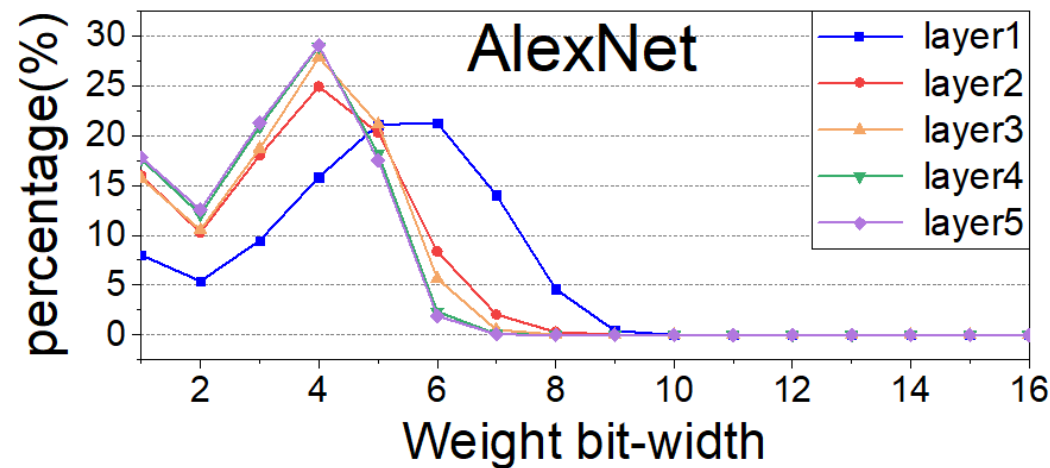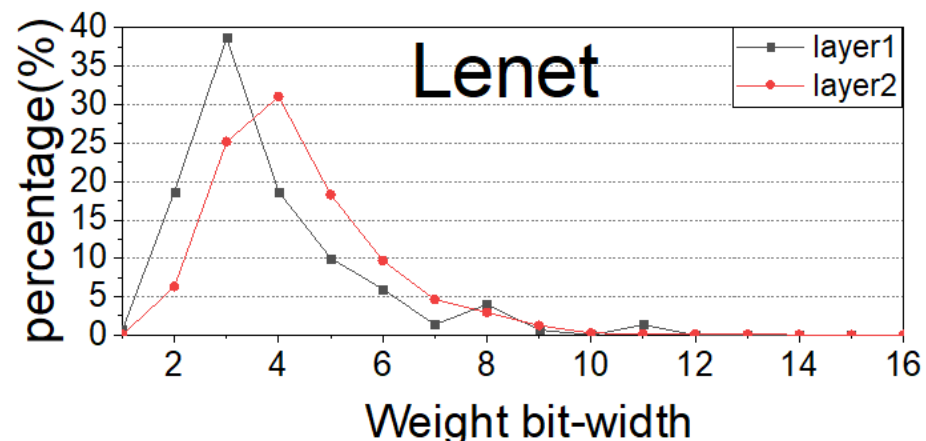CNN hardware accelerator become a popular solution

# Difficulties in CNN Deployment

CNN algorithm complexity continues growing and significant arithmetic and storage consumption

| Model | Top-1 | Top-5 | Ops (Bn) | GPU (ms) | CPU (s) | Weights (MB) |
|---|---|---|---|---|---|---|
| Densenet 201 | 77 | 93.7 | 10.85 | 32.6 | 1.38 | 66 |
| Darknet19 | 72.9 | 91.2 | 7.29 | 6.2 | 0.87 | 80 |
| Darknet53 | 77.2 | 93.8 | 18.57 | 13.7 | 2.11 | 159 |
| Resnet 18 | 70.7 | 89.9 | 4.69 | 4.6 | 0.57 | 44 |
| Resnet 34 | 72.4 | 91.1 | 9.52 | 7.1 | 1.11 | 83 |
| Resnet 50 | 75.8 | 92.9 | 9.74 | 11.4 | 1.13 | 87 |
| Resnet 101 | 77.1 | 93.7 | 19.7 | 20 | 2.23 | 160 |
| Resnet 152 | 77.6 | 93.8 | 29.39 | 28.6 | 3.31 | 220 |
| ResNeXt 50 | 77.8 | 94.2 | 10.11 | 24.2 | 1.2 | 220 |
| AlexNet | 57 | 80.3 | 2.27 | 3.1 | 0.29 | 238 |
| VGG-16 | 70.5 | 90 | 30.94 | 9.4 | 4.36 | 528 |

High consumption and limited memory size in hardware deployment

# Weight Distribution



Most weights with a low effective bit-width and few with a high effective bit-width

# Weight Storage Pattern in A Single Precision System

Format of 16-bit fixed-point data in Memory

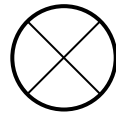| | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| data1: **20** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| data2: **-32** | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| data3: **8** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| data4: **3** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 |
| data5: **0** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| data6: **12** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 |
| data7: **0** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

noneffective bit

effective bit

Too many noneffective bit makes a significant  waste of storage resources!

# Operation in A Single Precision System

$A$(16bit) $\quad$ $A_{15}\ A_{14}\ A_{13}\ A_{12}\ A_{11}\ A_{10}\ A_9\ A_8\ A_7\ A_6\ A_5\ A_4\ A_3\ A_2\ A_1\ A_0$

$*$ $\quad$ $\bigotimes$ ——a 16*16 bit multiplier in single precision system

$B$(16bit) $\quad$ $B_{15}\ B_{14}\ B_{13}\ B_{12}\ B_{11}\ B_{10}\ B_9\ B_8\ B_7\ B_6\ B_5\ B_4\ B_3\ B_2\ B_1\ B_0$

$A$(16bit) $\quad$ $A_{15}\ A_{14}\ A_{13}\ A_{12}\ A_{11}\ A_{10}\ A_9\ A_8\ A_7\ A_6\ A_5\ A_4\ A_3\ A_2\ A_1\ A_0$

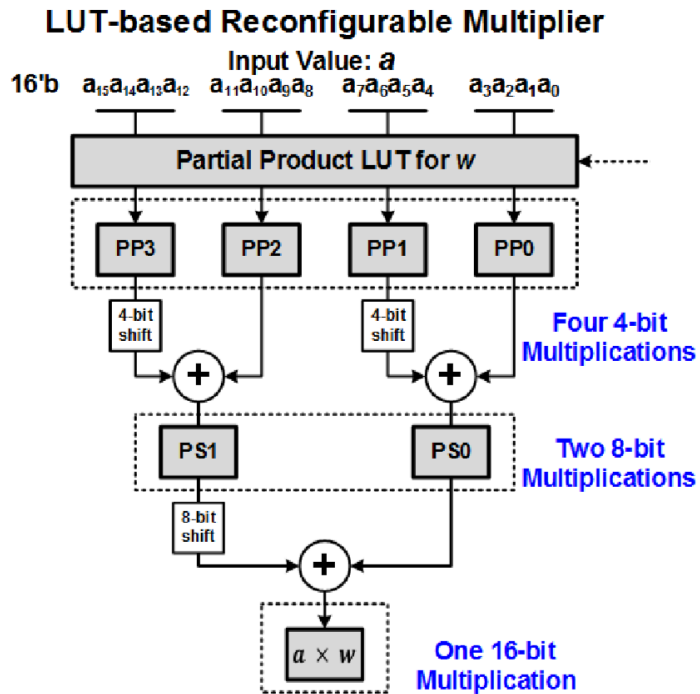$*$ $\quad$ $\bigotimes$ ——a 16*4 bit multiplier in best bit-width

11 $\quad$ 0 0 0 0 0 0 0 0 0 0 0 0 1 0 1 1

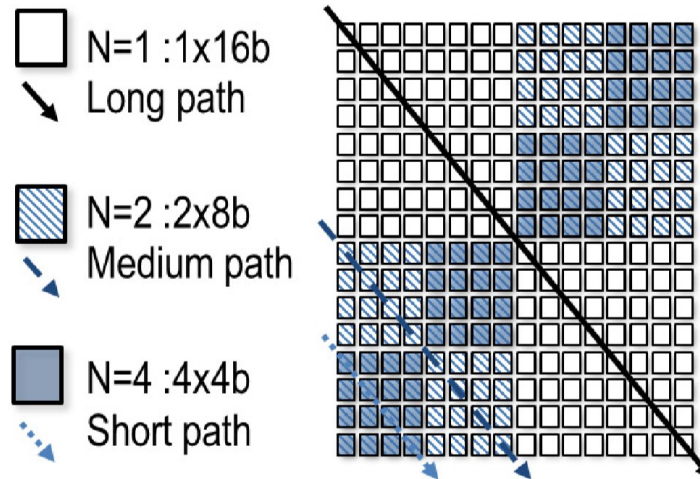All data is processed using the multiplier with the largest bit-width

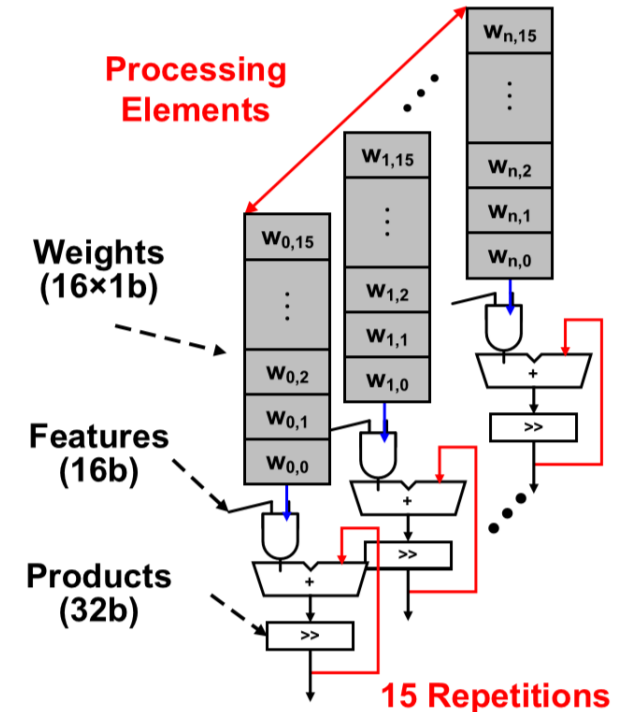A significant  waste of operation resources!

# Related Work

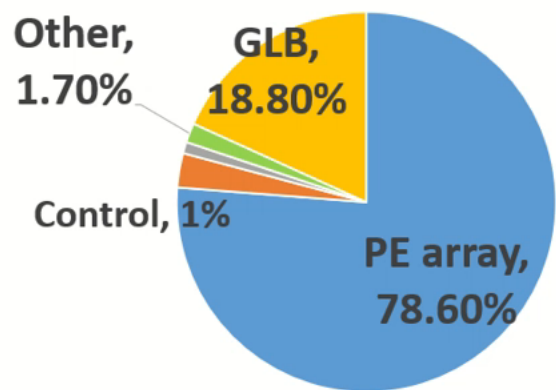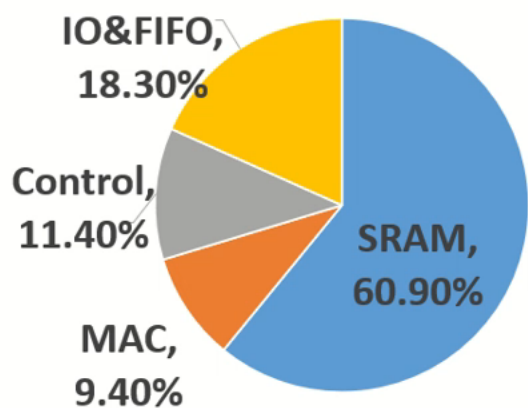DNPU(ISSCC2017)    ENVISION(ISSCC2017)    UNPU(JSSC2019)



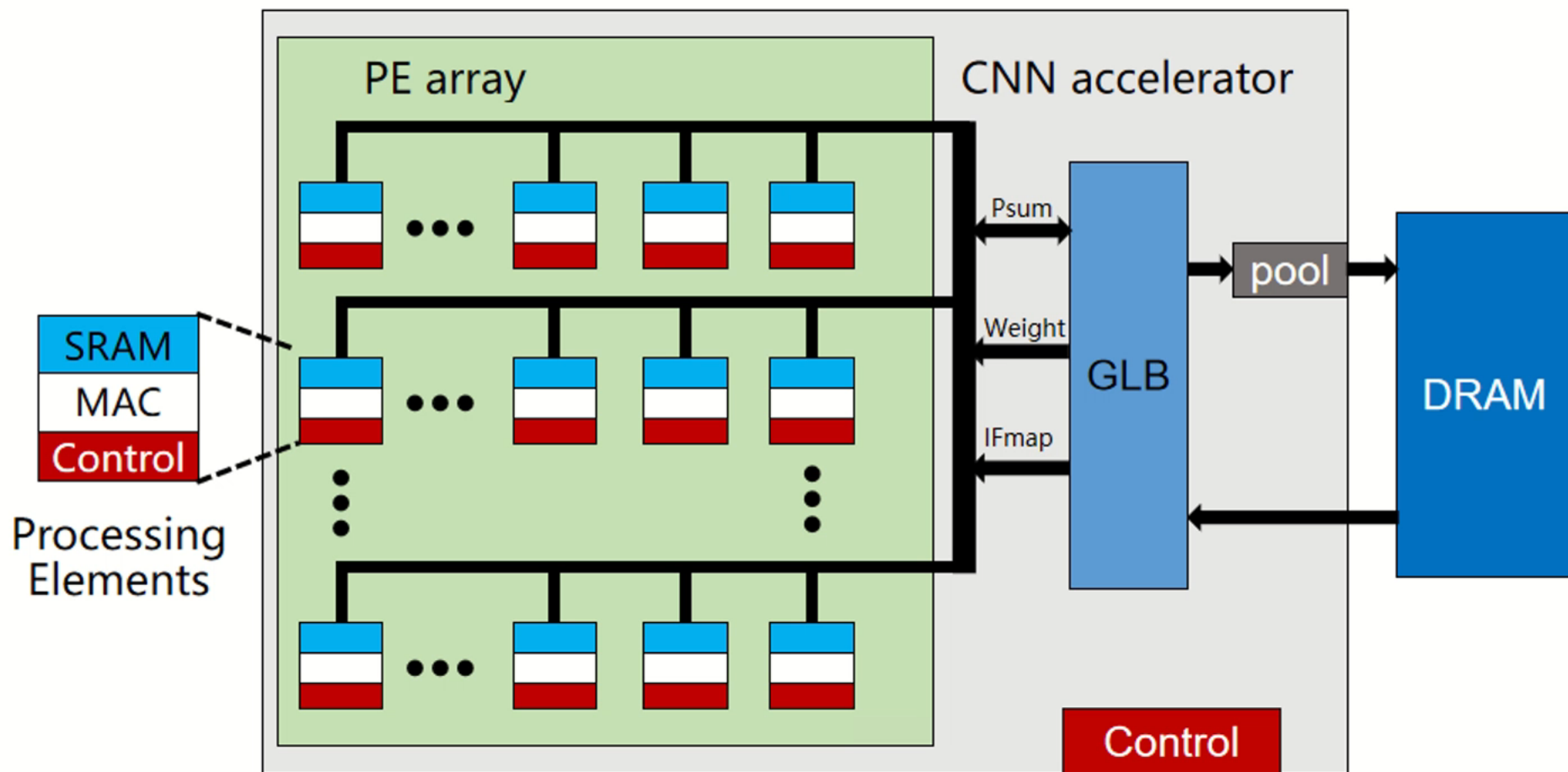The prior work do **not** support mixed precision computing within the same layer

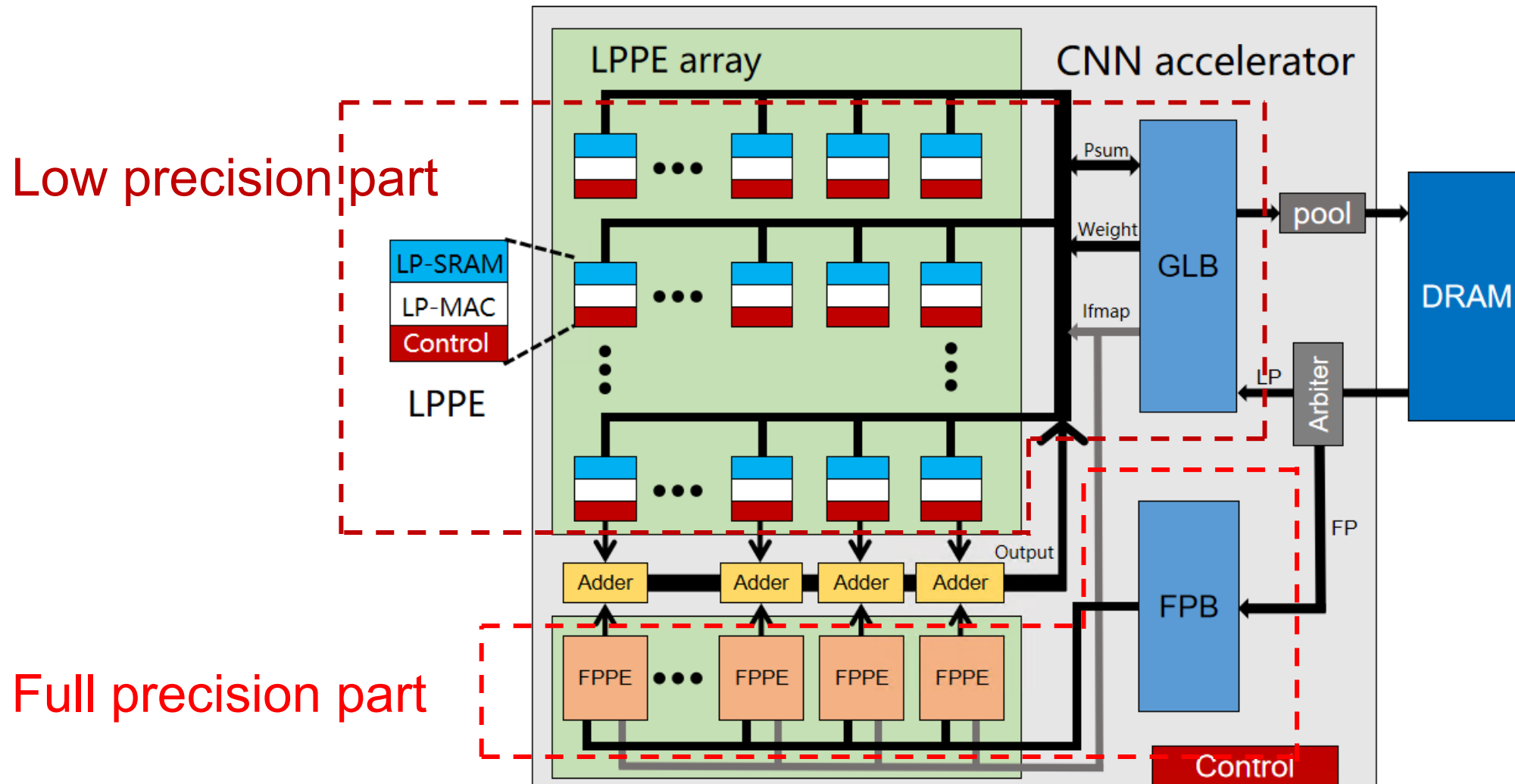# Commonly-Used CNN Accelerators Architecture
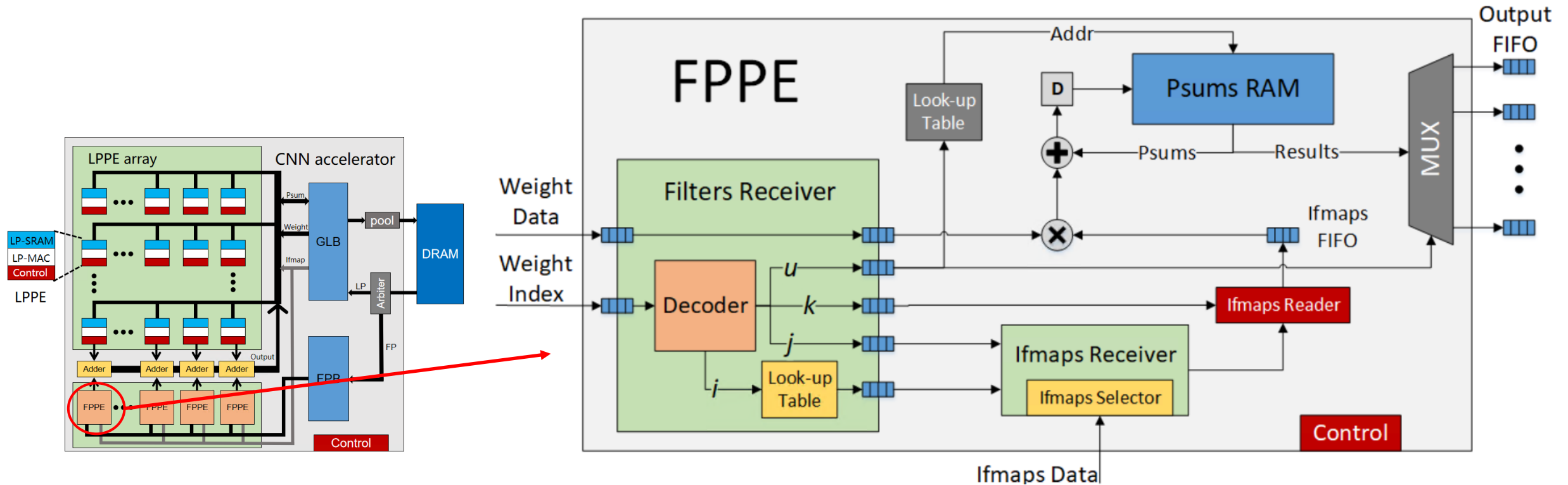


Accelerator*

PE*

*Data from Y. H. Chen, etc, "Eyeriss: An energy efficient reconfigurable accelerator for deep convolutional neural networks," IEEE J. Solid-State Circuits

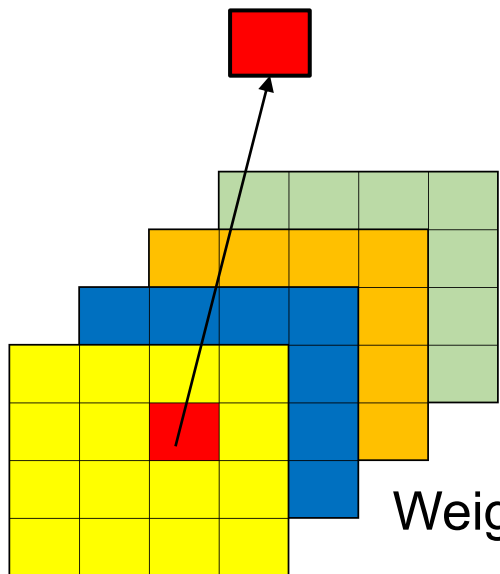# Mixed Precision CNN Accelerators Architecture

# Structure of The Proposed FPPE

# Decoder in FPPE

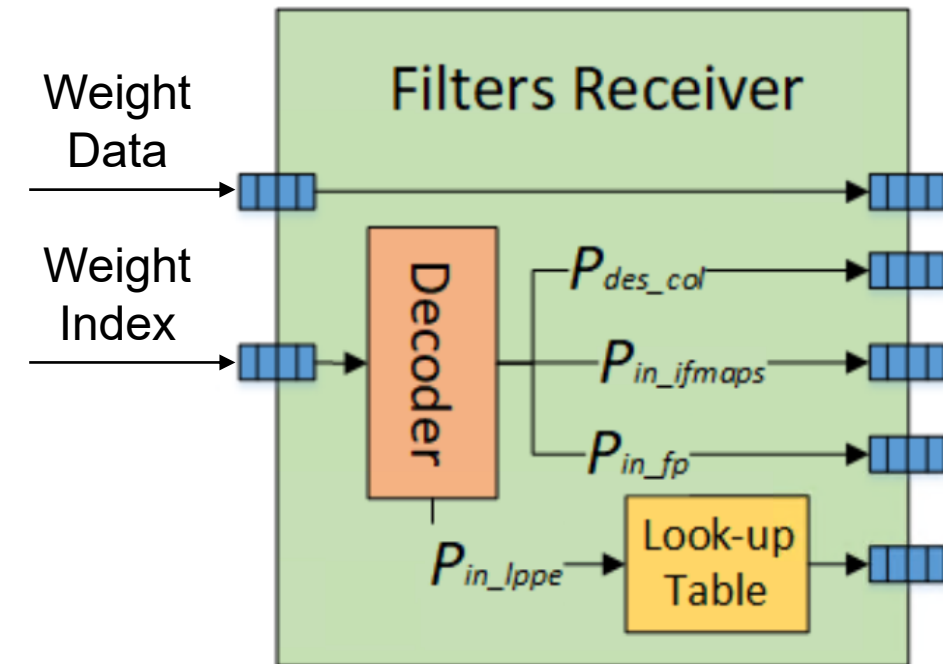Decoder translates weight position into FPPE parameter by PE and CNN shape



Weight position
Channel:1
Row:2
Col:3

Weights plane

**decoder**

FPPE parameter

$P_{des\_col}$

$P_{in\_ifmap}$

$P_{in\_fp}$

$P_{in\_lppe}$

# Arbiter Logic

Assume that a 16-bit signed fixed-point weight is fetched as $\{W_{15}, W_{14} \ldots W_1, W_0\}$

For a given bit-width threshold $W_i$ :

    low precision weights : $\{W_{14}, W_{13} \ldots W_{i-1}\} = (16 - i)\{W_{15}\}$

    full precision weights : $\{W_{14}, W_{13} \ldots W_{i-1}\} \neq (16 - i)\{W_{15}\}$
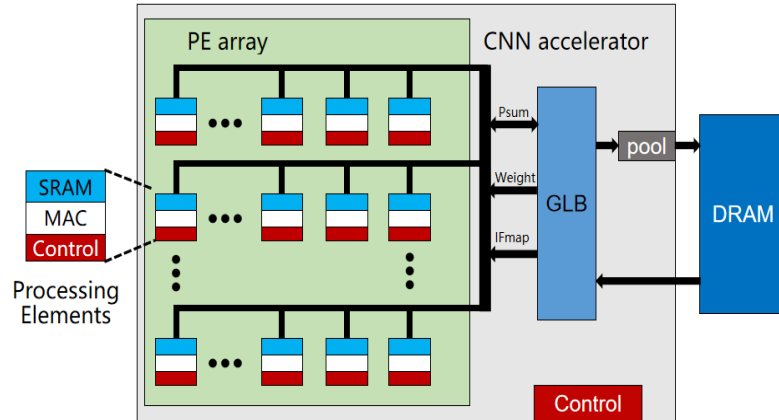


Sign bit

| bit-width threshold($W_i$) | "-68" is low precision |
|---|---|
| 10 | **Yes** |
| 8 | Yes |
| 5 | No |

# Comparison of Weight Storage Size

Weight storage for **single** precision architecture



Weight storage for **mixed** precision architecture



VS

$$S_{PE} = N \times W \times k$$
$$S_{GLB} = M \times W$$

Full precision part

$$S'_{FPPE} = N \times k \times (1 - p(W_i)) \times W$$
$$S'_{FPB} = M \times (1 - p(W_i)) \times W$$

Low precision part

$$S'_{LPPE} = N \times k \times W_i$$
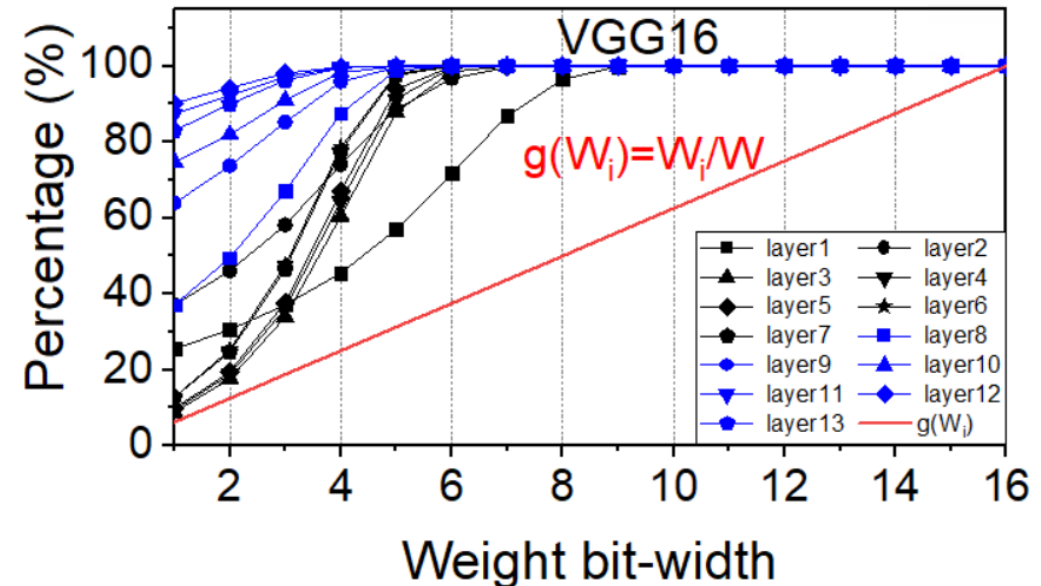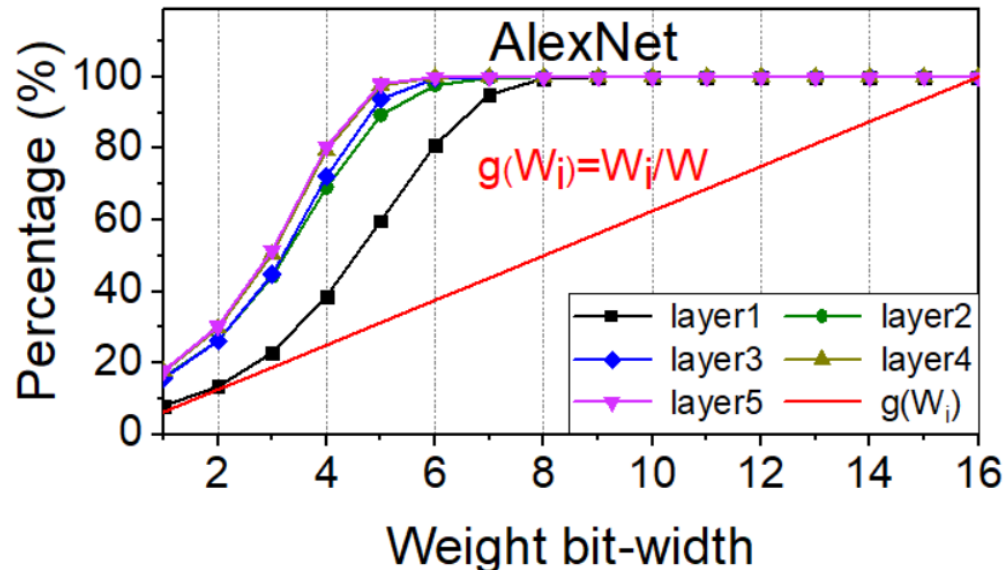$$S'_{GLB} = M \times W_i$$

# Bit-Width Threshold Selection

Relative storage saving for PE (or GLB) can be approximately calculated by above EQs, and simplified to:
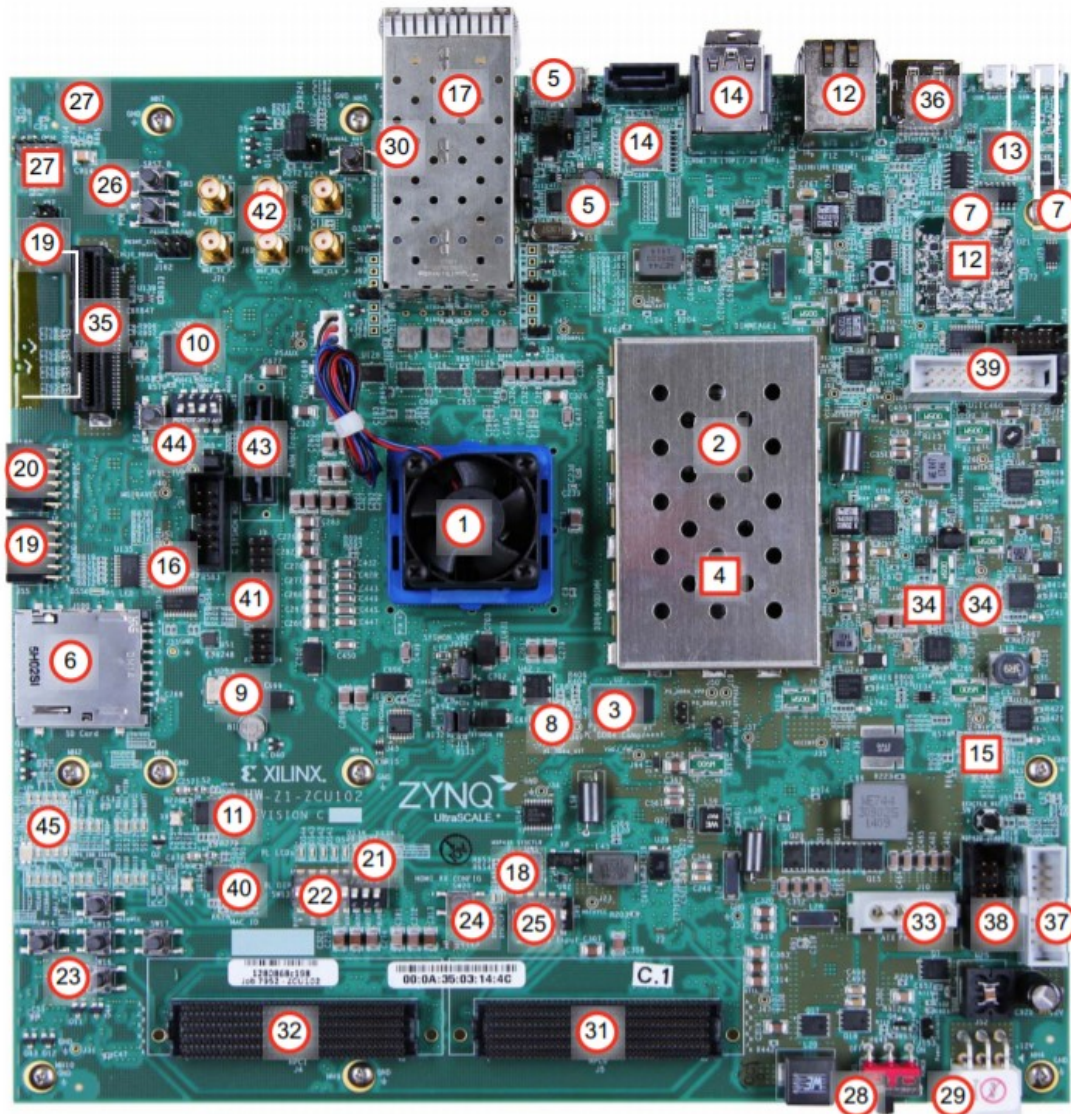
Area saving rate of PE:

$$\alpha \approx \frac{S_{PE}-S'_{PE}-S'_{FP\,PE}}{S_{PE}} = p(W_i) - \frac{W_i}{W}$$

Area saving rate of GLB:

$$\beta \approx \frac{S_{GLB}-S'_{GLB}-S'_{FPB}}{S_{GLB}} = p(W_i) - \frac{W_i}{W} = \alpha$$

# Experimental Setup



-Low power
-Data reuse
-Reconfigurable

| Baseline Accelerator | |
|---|---|
| Board | ZCU102 |
| Dataflow | Row Stationary |
| Arithmetic Precision | 16-bit fixed point |
| Clock Rate | 200Mhz |
| GLB storage size | 224KB |
| #PE in a row | 14 |
| #PE in a col | 12 |
| PE storage size | 608B |
| #LUT used for MAC | 280 |
| #LUT used for PE | 1313 |

# Experimental Result——Area Reduction



Baseline/PE x168 — Prop./LPPE x168 — Prop./FPPE x14

| | Baseline/PE | Prop./LPPE | Prop./FPPE |
|---|---|---|---|
| **Number** | 168 | 168 | 14 |
| **Storage(Byte)** | 608 | 384 | 1003 |
| **#LUTs(MAC)** | 280 | 128 | 280 |
| **#LUTs(PE)** | 1313 | 777 | 1606 |
| **Total** **Storage(Byte)** | 102,144(1x) | 78,554(0.77x) | |
| **#LUTs** | 220,584(1x) | 153,020(0.69x) | |

# Experimental Result

| | | Storage (KB) | | | # LUT | | Norm. power |
|---|---|---|---|---|---|---|---|
| | | Weight | Ifmap | Psum | MAC | System | |
| Alex Net | Baseline | 116.8 | 69.3 | 137.7 | 47k | 204k | 1x |
| | Proposed | 62.0 | 69.5 | 142.9 | 25k | 192k | 0.88x |
| | Saving | 46.9% | -0.3% | -3.8% | 45.9% | 6.2% | 12.1% |
| VGG 16 | Proposed | 58.9 | 69.5 | 142.5 | 23k | 183k | 0.88x |
| | Saving | 49.6% | -0.3% | -3.5% | 49.8% | 10.5% | 12.1% |



(a) AlexNet     (b) VGG16

- Weight storage area reduced by nearly 50%
- Number of LUTs for calculation is reduced by nearly 50%
- Number of LUTs in the system is reduced by 6% to 10%
- Save about 12% power than baseline
- Actual total storage saving in AlexNet and VGG16 is almost 17.8% and 16.8%

# Conclusion

- Proposed architecture of PE to simultaneously store and calculate with different employ two separate groups precisions

- Implement the proposed CNN accelerator using an FPGA platform

- Weight storage area in PE and GLB reduced by nearly 50%

- Total storage area reduced by almost 17.8%

- Critical path delay is reduced by almost 28%

- Dynamic power saving by 12.1% without timing penalty

# Q & A