



# LOCAL: Low-Complex Mapping Algorithm for Spatial DNN Accelerators

Midia Reshadi, David Gregg  
School of Computer Science and Statistics,  
Lero, Trinity College Dublin

# Outline

---

- Introduction
- Motivation
- Problem formulation
- LOCAL mapping algorithm
- Simulation results

# Introduction

---

- DNN applications
  - Self-driving cars
  - Recommendation systems
  - Language translation



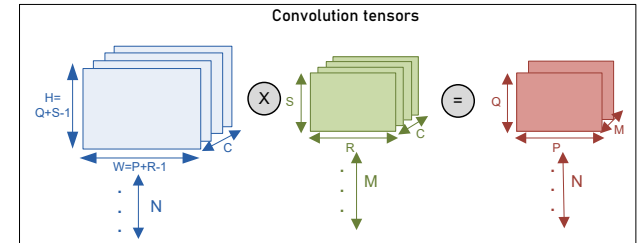
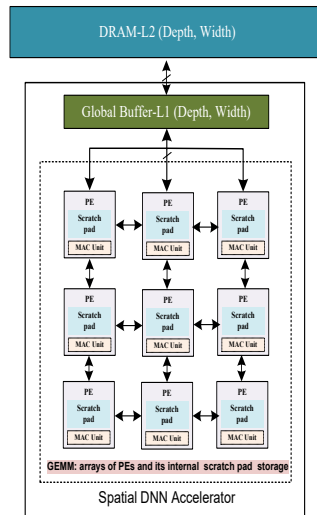
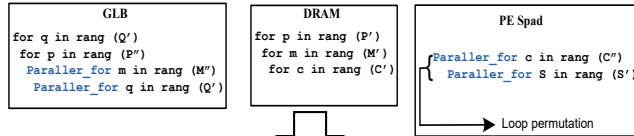
# Mapping algorithm

Loop representation of convolution tensors

```

1 for r in rang (R):
2   for s in rang (S):
3     for p in rang (P):
4       for q in rang (Q):
5         for c in rang (C):
6           for m in rang (M):
7             for n in rang (N):
8               Output[p][q][k][n] += Weight[r][s][k][c] * Input[p+r][q+s][c][n];

```

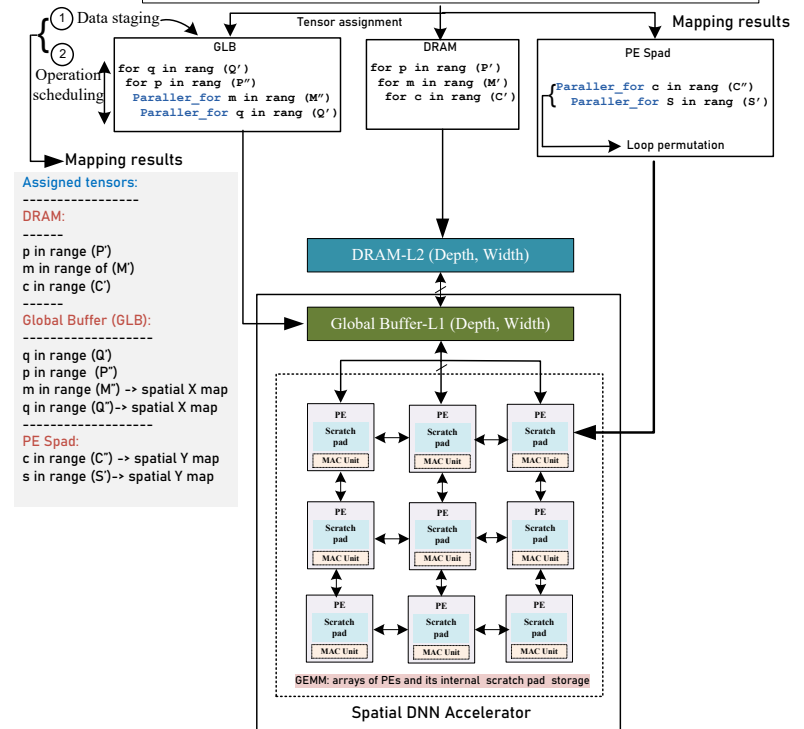


Loop representation of convolution tensors

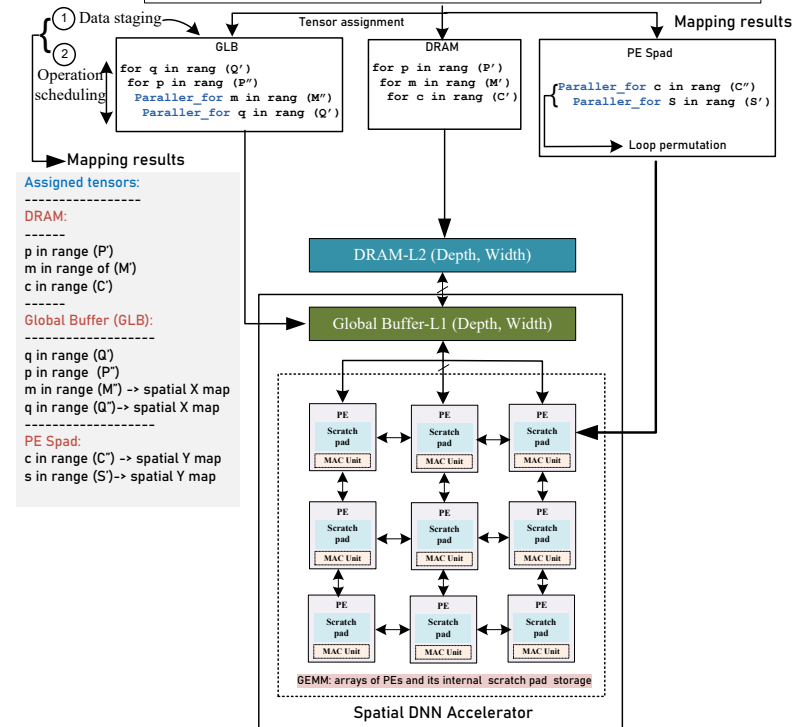
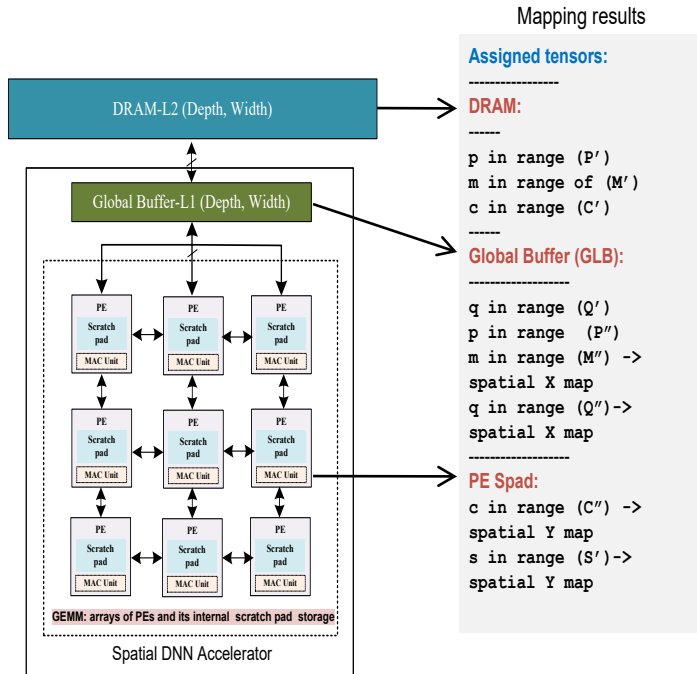
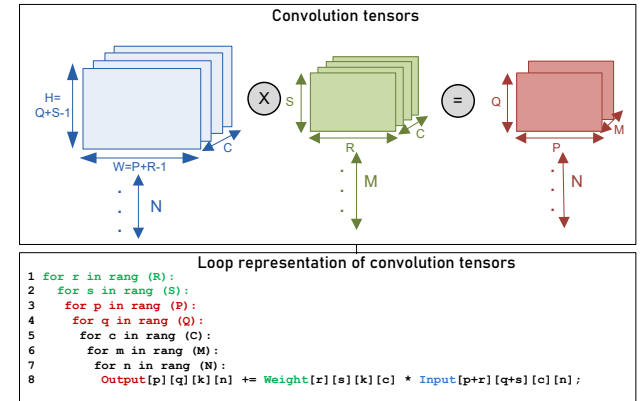
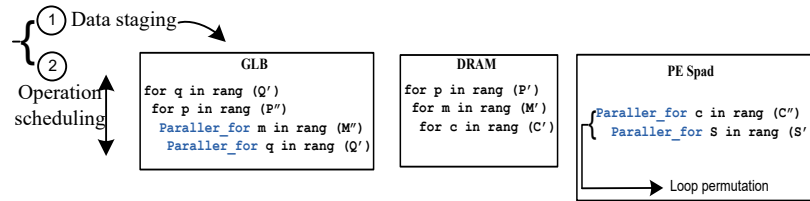
```

1 for r in rang (R):
2   for s in rang (S):
3     for p in rang (P):
4       for q in rang (Q):
5         for c in rang (C):
6           for m in rang (M):
7             for n in rang (N):
8               Output[p][q][k][n] += Weight[r][s][k][c] * Input[p+r][q+s][c][n];

```



# Mapping algorithm

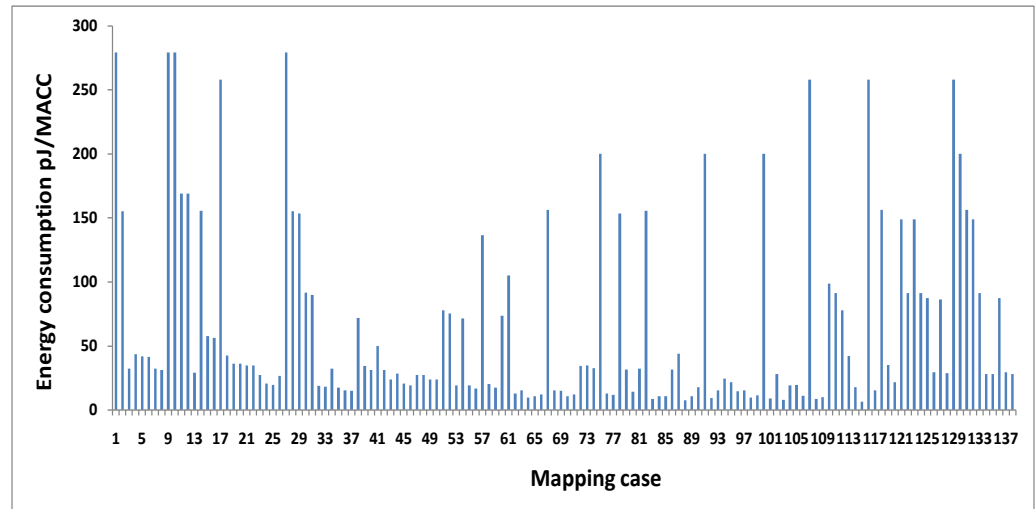


# Motivation

---

- Random mapping of Layer 5 of VGG16 to Eyeriss
  - 3,000 random mapping cases
  - 77% difference between the *Random\_max* and the *Random\_med*

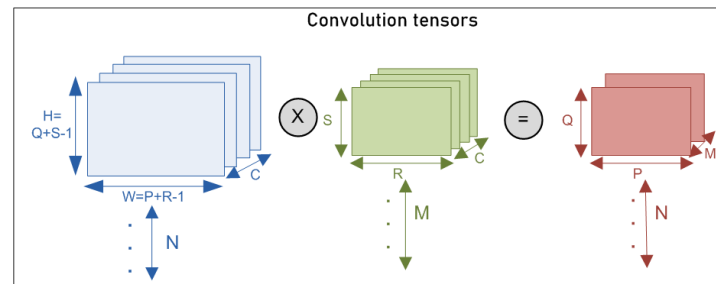
| SPA architecture             |             |
|------------------------------|-------------|
| SPA                          | Eyeriss     |
| On-chip storage levels       | 2           |
| DRAM(width)                  | 64          |
| L <sub>1</sub> (depth,width) | (16384, 64) |
| L <sub>0</sub> (depth,width) | (16,16)     |
| PE array                     | (12,14)     |



# Preliminaries

---

- Convolutions



$CT = \{Input, Weight, Output\}$

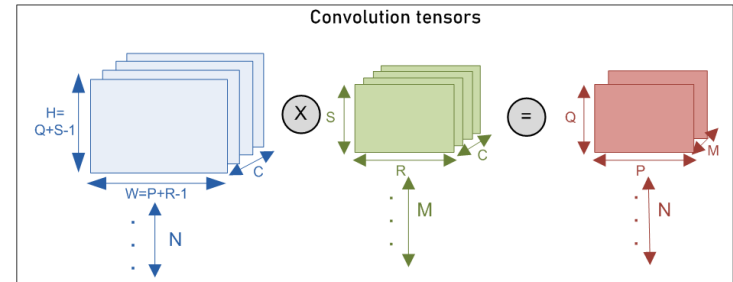
$CT \in R^{dimension}$

$dimension = \{N, M, C, W, H, R, S, P, Q\}$

# Preliminaries

---

- Convolutions



$$\text{dimension} = \{N, M, C, W, H, R, S, P, Q\}$$

$$W \in R^{MCRS}$$

$$I \in R^{NCHW}$$

$$O \in R^{NMPQ}$$

$$ct_i \in CT$$

$$CT = \{ct_1, ct_2, ct_3\}$$

$$CT = \{R^{MCRS}, R^{NCHW}, R^{NMPQ}\}$$

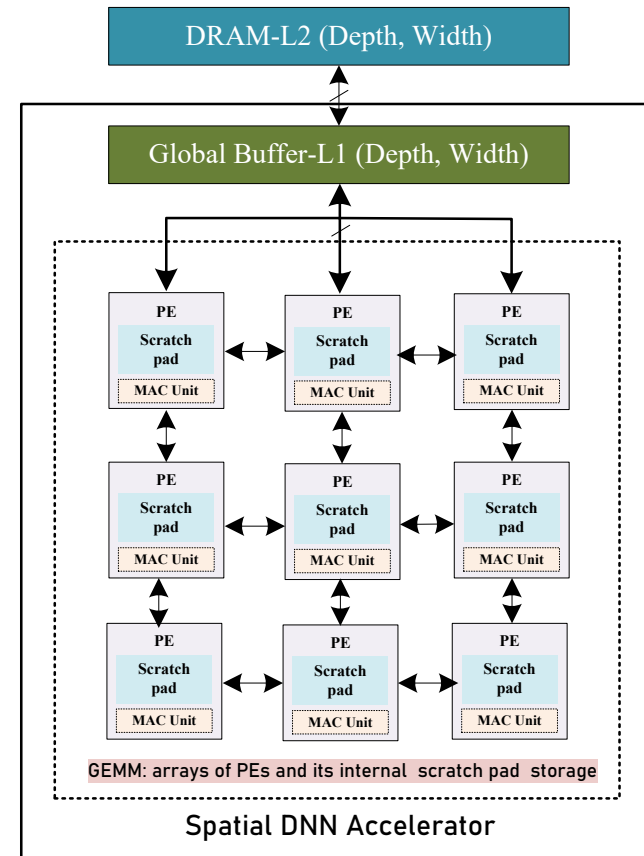


# Preliminaries

- Spatial DNN accelerator

$$\text{SPA} = \{\text{Storage}[i, j, k], \text{PE}[m, n]\}$$

$$\text{PE}_{m,n} \in R^{x,y}$$



# Preliminaries

$$S_{i,j,k} \Rightarrow S \in R^{i,j,k}$$

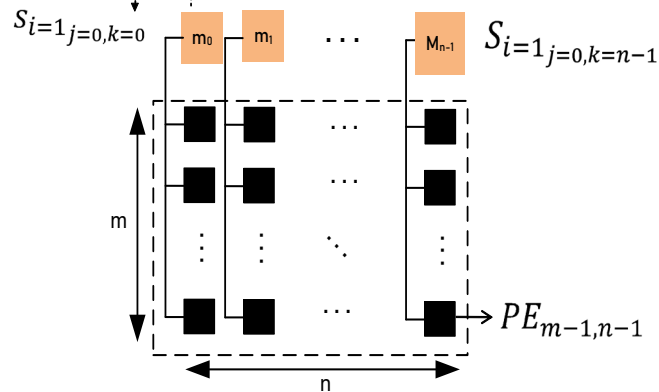
Eyeriss style

Level 3 (L3) e.g. DRAM

$$S_{i=3,j=0,k=0}$$

Level 2 (L2) e.g. GLB

$$S_{i=2,j=0,k=0}$$



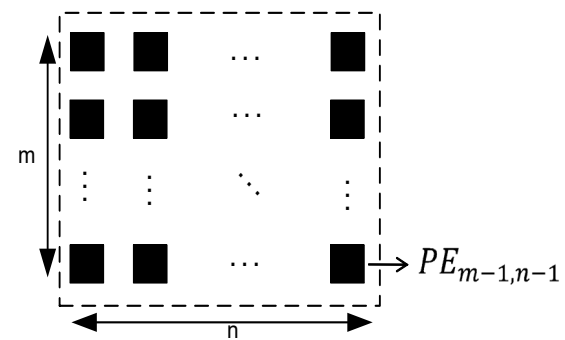
NVDLA style

Level 2 (L2) e.g. DRAM

$$S_{i=2,j=0,k=0}$$

Level 1 (L1) e.g. SRAM

$$S_{i=1,j=0,k=0}$$



# Mapping algorithm

---

- Assignment  $ct_i \in CT$  *assign* to  $s_{i,j,k} \in S$

- Bounding  $ct_i [0, rang) \in CT$  *assign* to  $s_{i,j,k} \in S$   
for  $Q$  in  $[0,5)$  assigned to  $s_{i=1,j=0,k=0}$



**EXAMPLE**

- Scheduling

$ct_i, ct_j, ct_k \in CT$  *assign* to  $L_i$

- Parallelization

Order  $\left\{ \begin{array}{l} ct_i [0, rang_i) \\ ct_j [0, rang_j) \\ ct_k [0, rang_k) \end{array} \right.$

Spatial computing  $\rightarrow ct_i [0, rang) \in CT$  on  $PE_{i \text{ to } j(x|y)} \in PEs$

Parallel\_for  $S$  in  $[0,7)$  on  $PE_{0-7}$  Spatial X dimension

# Mapping Problem Formulation

---

- ***Given:***
  - Convolution tensors  $CT$
  - Spatial DNN accelerator ( $SPA$ )
- ***Find***
  - A mapping function  $map( )$

$\min\{Energy\}$

$\max\{PE\text{ utilization}\}$

$$utilization\ of\ PEs = \frac{number\ of\ active\ PEs}{number\ of\ PEs}$$

# Local Mapping Algorithm

**Algorithm1:** Tensor to memory assignments  $ct_i$  to  $s_{i(j,k)}$

Inputs: All  $ct_i \in CT$  and all  $s_{i(j,k)} \in S$

Size of  $(s_{i(j,k)})$

PE dimension ( $PE_{m,n}$ )

Outputs:  $ct_i \rightarrow s_{i(j,k)}$  (Assignment & Permutation)

$ct_i \rightarrow PE_{i \rightarrow j(x|y)}$  (Parallelization)

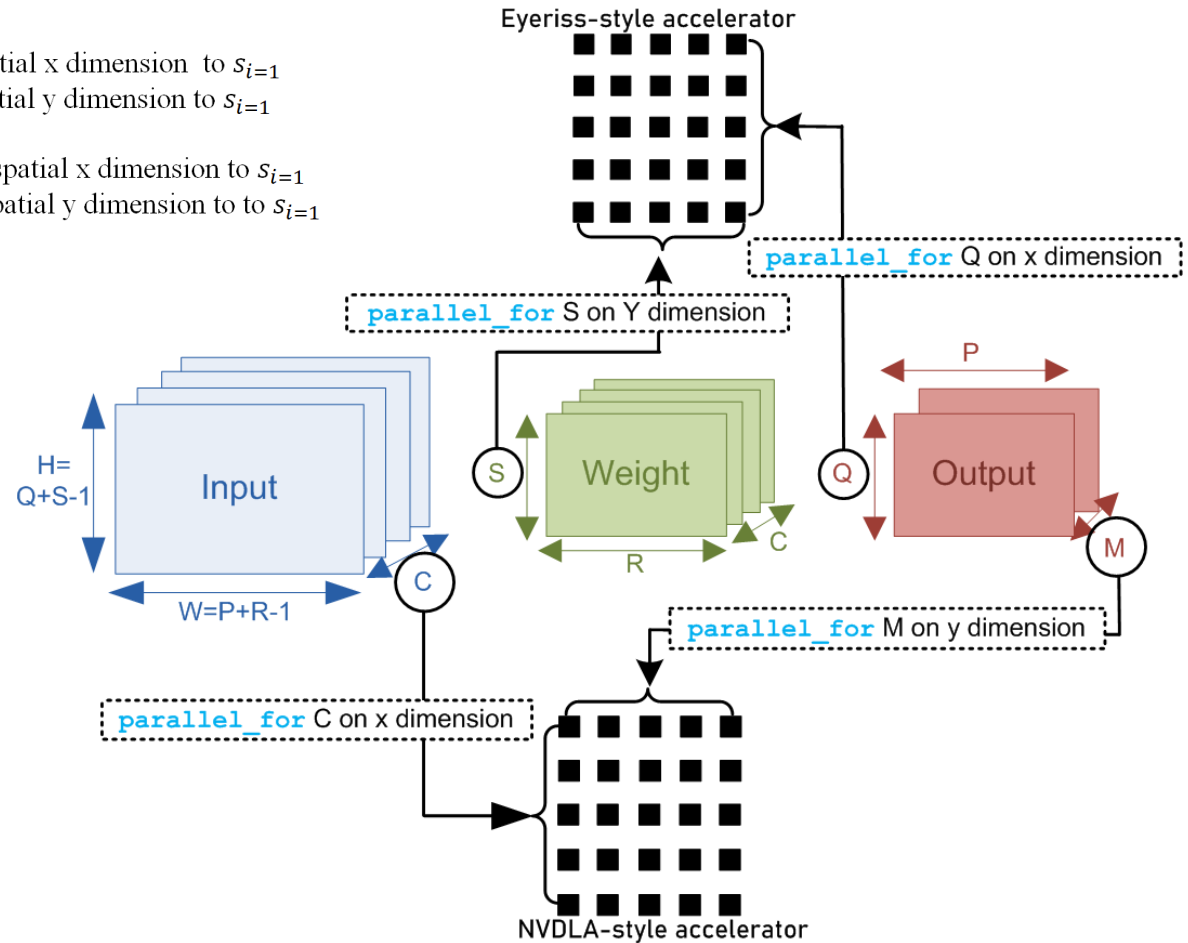
```
1: Parallelization:
2:    $\forall s_{i=1,j,k}$ :
3:     if ( $k=0$ ) //NVDLA style
4:       assign: parallel_for C in  $Rang(m)$  spatial x dimension to  $s_{i=1}$ 
5:       assign: parallel_for M in  $Rang(n)$  spatial y dimension to  $s_{i=1}$ 
6:     else //Eyeriss Style
7:       assign: parallel_for in Q in  $Rang(m)$  spatial x dimension to  $s_{i=1}$ 
8:       assign: parallel_for in S in  $Rang(n)$  spatial y dimension to  $s_{i=1}$ 
9:     end_if
10: Assignment:
11:    $\forall$  all Unassigned tensors:  $U_{ct_i}$ 
12:   for all remaining  $U_{ct_i}$ :
13:     for ( $i=0 \rightarrow i \leq n$ )
14:       Assign  $U_{ct_i}$  to  $s_i$  // Assigning with priority from  $s_{i=0}$  to  $s_{i=n}$ 
15:     end_for
16:   end_for
17: Scheduling:
18:   for all Assigned tensors:  $A_{ct_i}$ 
19:     sort high to low range  $ct_i$ 
20:      $\forall s_i$  for ( $i=0 \rightarrow i \leq \text{max memory level}$ )
21:       do permutation to allocate higher range tensor to lower  $s_i$ 
22:     end_for
```

# More details

```

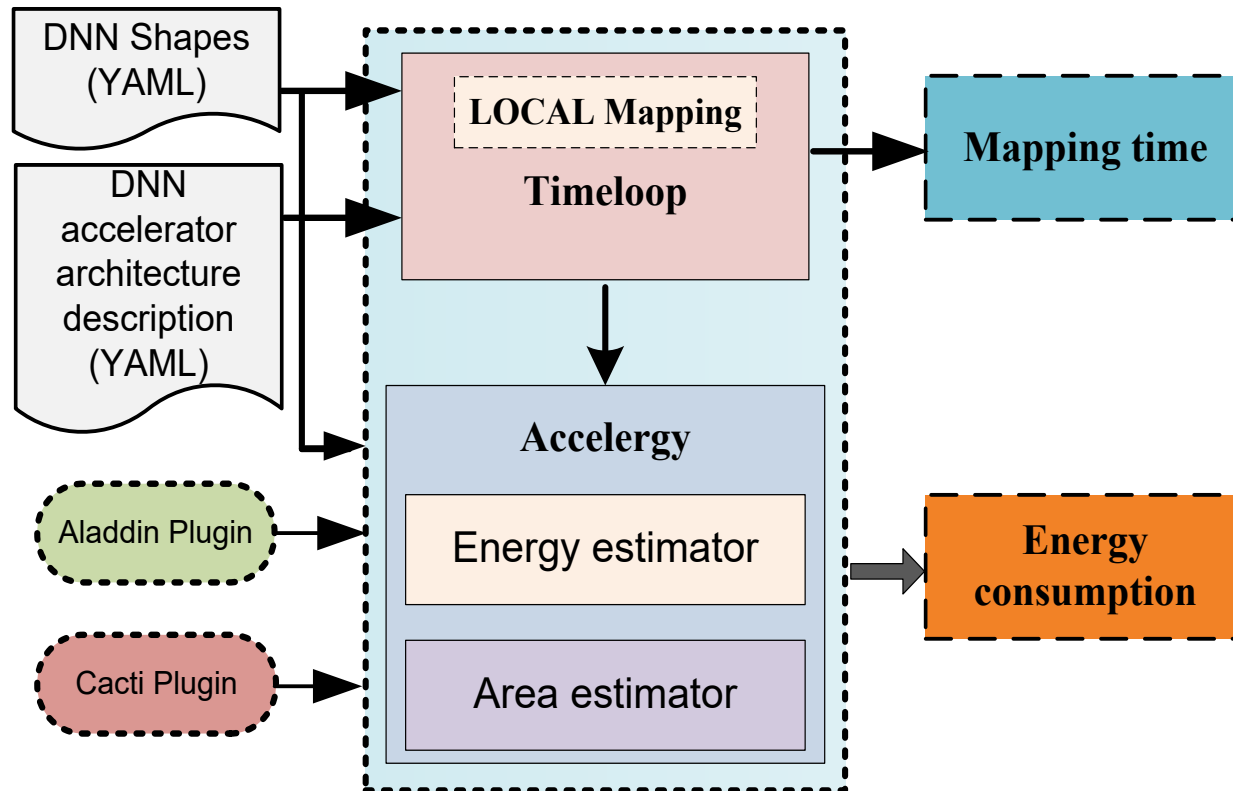
1: Parallelization:
2:  $\forall s_{i=1,j,k}$ :
3:   if (k==0) //NVDLA style
4:     assign: parallel_for C in Rang(m) spatial x dimension to  $s_{i=1}$ 
5:     assign: parallel_for M in Rang(n) spatial y dimension to  $s_{i=1}$ 
6:   else //Eyeriss Style
7:     assign: parallel_for in Q in Rang(m) spatial x dimension to  $s_{i=1}$ 
8:     assign: parallel_for in S in Rang(n) spatial y dimension to  $s_{i=1}$ 
9:   end_if

```



# Evaluation

---



# Simulation Workload

---

**Table 2. Workload categories**

| <b>Category</b>            | <b>Workload</b>                           | <b>Number of MAC operation</b> |
|----------------------------|---|--------------------------------|
| <i>High C value</i>        | 22 <sup>nd</sup> conv layer of Resnet50   | 51380224                       |
|                            | 23 <sup>nd</sup> conv layer of SqueezeNet | 5537792                        |
|                            | 9 <sup>nd</sup> conv layer of VGG16       | 1849688064                     |
| <i>High M value</i>        | 25 <sup>nd</sup> conv layer of SqueezeNet | 24920064                       |
|                            | 24 <sup>nd</sup> conv layer of ResNet50   | 51380224                       |
|                            | 8 <sup>nd</sup> conv layer of VGG16       | 924844032                      |
| <i>High P and Q values</i> | 1 <sup>st</sup> conv layer of SqueezeNet  | 708083712                      |
|                            | 1 <sup>st</sup> conv layer of ResNet50    | 472055808                      |
|                            | 1 <sup>st</sup> conv layer of VGG16       | 86704128                       |

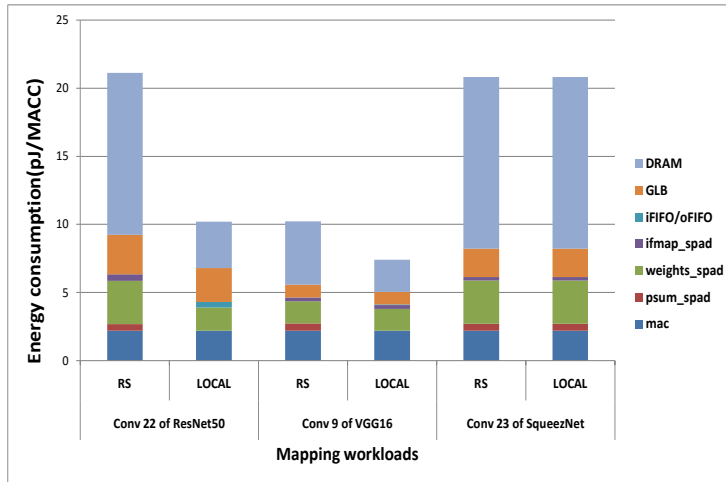


# Mapping time

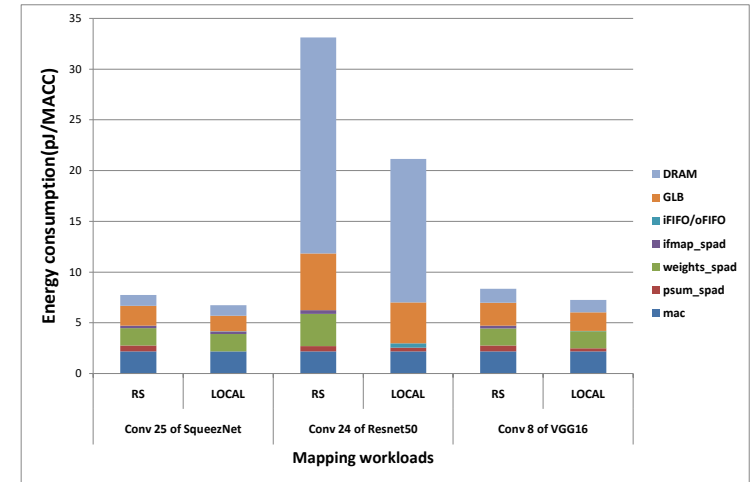
---

| <i>Workload</i>            | <i>Convolution</i>  | <i>Mapping mechanism based on Eyeriss</i> | <i>Mapping time (sec)</i> | <i>Mapping mechanism based on ShiDianNao</i> | <i>Mapping time (sec)</i> | <i>Mapping mechanism based on NVDLA</i> | <i>Mapping time (sec)</i> |
|----------------------------|---------------------|---|---------------------------|--|---------------------------|---|---------------------------|
| <b>High C value</b>        | Resnet50: Conv 22   | RS  | 87                        | OS   | <b>576</b>                | WS                                      | 127                       |
|                            |                     | LOCAL                                     | 16.2                      | LOCAL  | <b>15</b>                 | LOCAL                                   | 6                         |
|                            | VGG16: Conv 9       | RS  | 170                       | OS   | 137                       | WS                                      | 68                        |
|                            |                     | LOCAL                                     | 10                        | LOCAL  | 15                        | LOCAL                                   | 9                         |
|                            | SqueezeNet: Conv 23 | RS  | 17                        | OS   | 125                       | WS                                      | 21                        |
|                            |                     | LOCAL                                     | 16                        | LOCAL  | 67                        | LOCAL                                   | 18                        |
| <b>High M value</b>        | SqueezeNet: Conv 25 | RS  | <b>230</b>                | OS   | 126                       | WS                                      | 996                       |
|                            |                     | LOCAL                                     | <b>6.6</b>                | LOCAL  | 16                        | LOCAL                                   | 31                        |
|                            | Resnet50: Conv 24   | RS  | 74                        | OS   | 116                       | WS                                      | 42                        |
|                            |                     | LOCAL                                     | 22                        | LOCAL  | 28                        | LOCAL                                   | 12                        |
|                            | VGG16: Conv 8       | RS  | 351                       | OS   | 98                        | WS                                      | 411                       |
|                            |                     | LOCAL                                     | 12                        | LOCAL  | 32                        | LOCAL                                   | 24                        |
| <b>High P and Q values</b> | SqueezeNet: Conv1   | RS  | 60                        | OS   | 20                        | WS                                      | <b>2238</b>               |
|                            |                     | LOCAL                                     | 5.1                       | LOCAL  | 7                         | LOCAL                                   | <b>45</b>                 |
|                            | Resnet50: Conv 1    | RS  | 90                        | OS   | 60                        | WS                                      | 140                       |
|                            |                     | LOCAL                                     | 6                         | LOCAL  | 13                        | LOCAL                                   | 23                        |
|                            | VGG16: Conv 1       | RS  | 81                        | OS   | 24                        | WS                                      | 113                       |
|                            |                     | LOCAL                                     | 6.6                       | LOCAL  | 6                         | LOCAL                                   | 17                        |

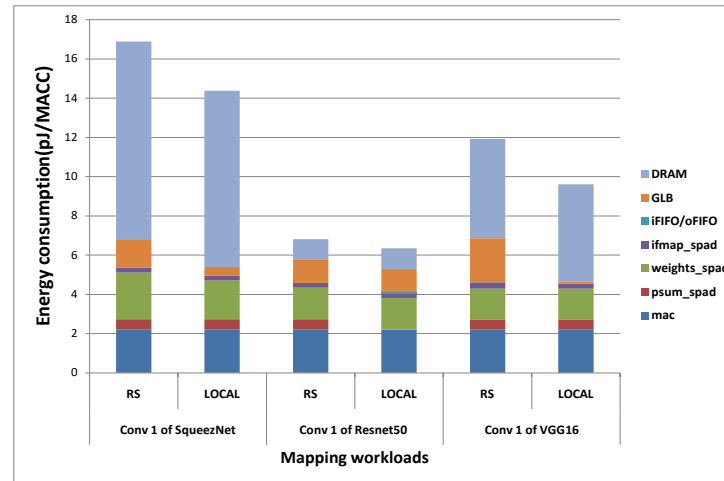
# Simulation Results



**Energy consumption of row stationary and LOCAL mapping in Eyeriss with High C value workload**

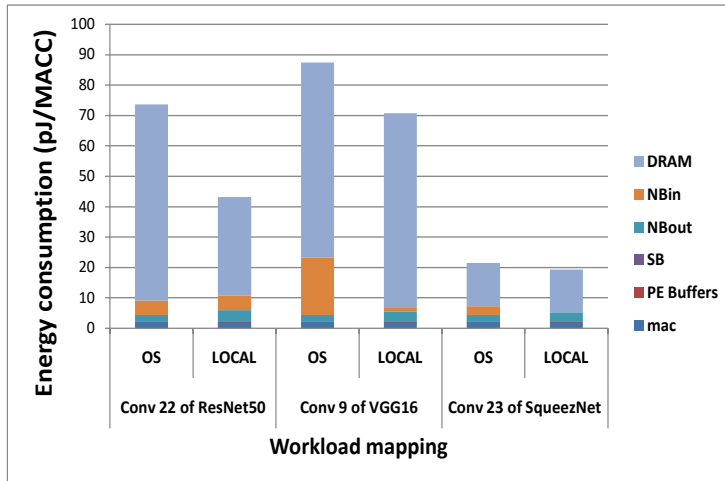


**Energy consumption of row stationary and LOCAL mapping in Eyeriss with High M value workload**

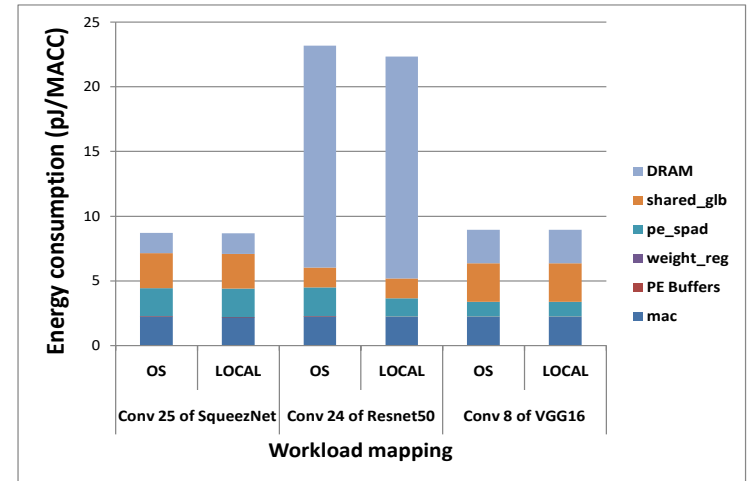


**Energy consumption of row stationary and LOCAL mapping in Eyeriss with High P and Q values workload**

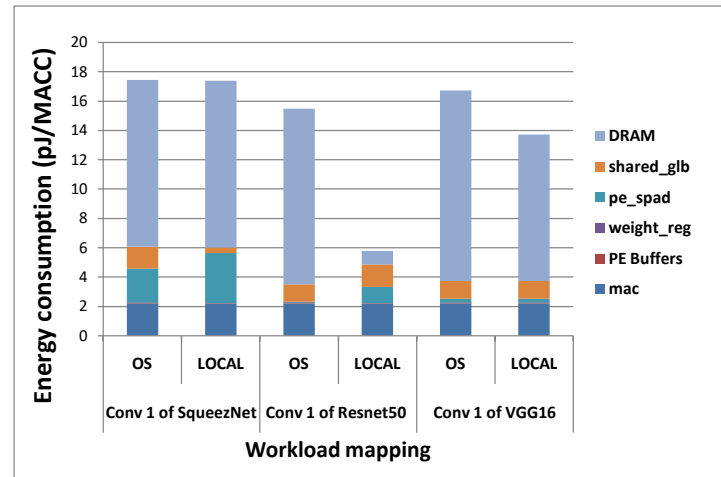
# Simulation Results



Energy consumption of output stationary and LOCAL mapping in Shi-diannao with High C value workload

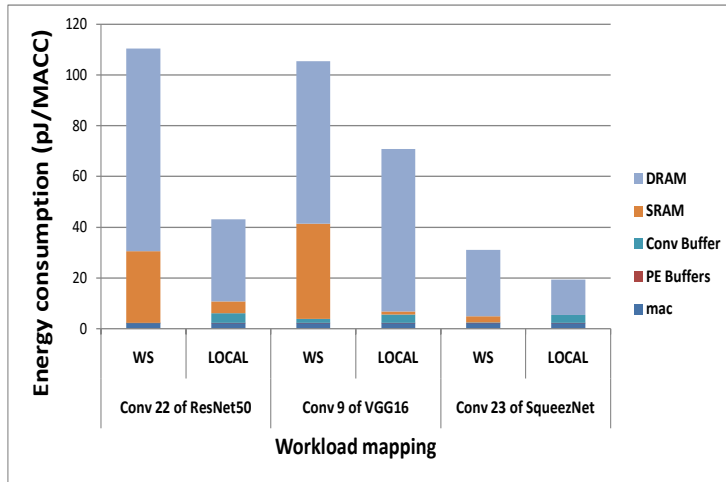


Energy consumption of output stationary and LOCAL mapping in Shi-diannao with High M value workload

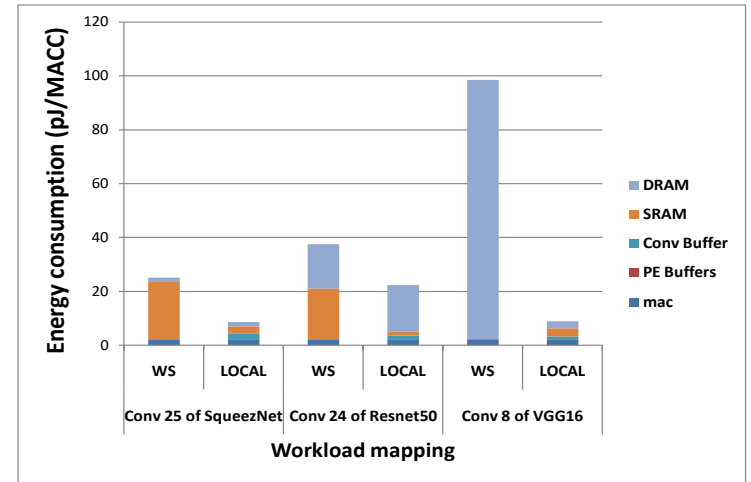


Energy consumption of output stationary and LOCAL mapping in Shi-diannao with High P and Q values workload

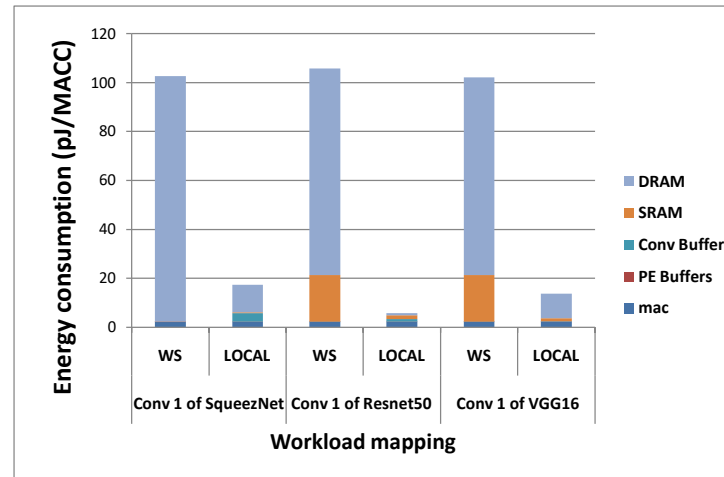
# Simulation Results



Energy consumption of weight stationary and LOCAL mapping in NVDLA with High C values workload



Energy consumption of weight stationary and LOCAL mapping in NVDLA with High M values workload



Energy consumption of weight stationary and LOCAL mapping in NVDLA with High P and Q values workload

Thank you for your attention