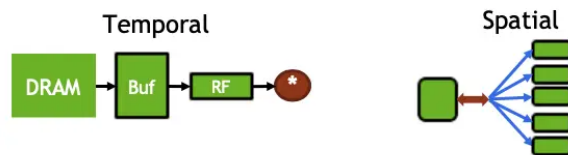# Dataflow



- **Temporal reuse:** the same data is used more than once over time by the same consumer.
- **Spatial reuse:** the same data is used by more than one consumer at different spatial locations of the hardware.

## Output-stationary

- **input:** $N, C_{in}, H_{in}, W_{in}$
- **kernel:** $C_{out}, C_{in}, H_k, W_k$
- **output:** $N, C_{out}, H_{out}, W_{out}$
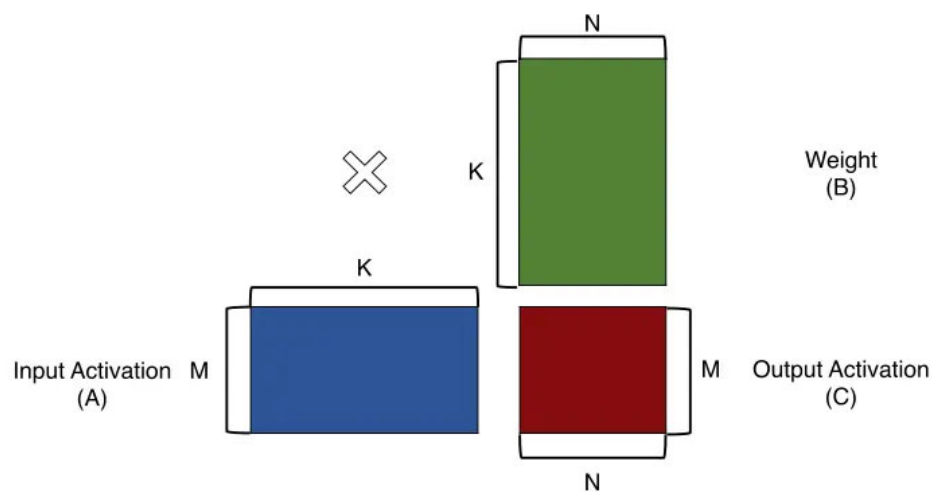
```
# padding=0, stride=1
for n in range(batch_size):
    for c_out in range(C_out):
        for h_out in range(H_out):
            for w_out in range(W_out):
                # per output element
                output[n,c_out,h_out,w_out] = 0
                for h_k in range(H_k):
                    for w_k in range(W_k):
                        for c_in in range(C_in):
                            output[n,c_out,h_out,w_out] += input[n,c_in,h_out+h_k,w_out+w_k]*
                                                           kernel[c_out,c_in,h_k,w_k]
```

## Weight-stationary

```
for n in range(batch_size):
  for h_k in range(H_k):
    for w_k in range(W_k):
      for c_in in range(C_in):
        for c_out in range(C_out):
          curr_k = kernel[c_out, c_in, h_k, w_k]
          for h_out in range(H_out):
            for w_out in range(H_out):
              output[n,c_out,h_out,w_out] += input[n,c_in,h_out+h_k,w_out+w_k]*curr_k
```
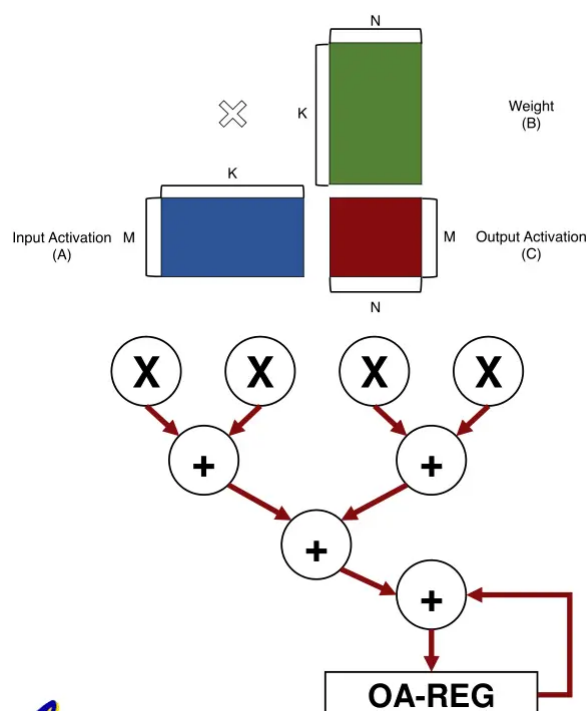
# Datapath Optimization

```
for m in range(M):
  for n in range(N):
    OA[m,n] = 0
    for k in range(K):
      OA[m,n] += IA[m,k] * W[k,n]
    OA[m,n] = Activation(OA[m,n])
```

# Spatial, K

```
for m in range(M):
  for n in range(N):
    OA[m,n] = 0
    spatial_for k in range(K):
      OA[m,n] += IA[m,k] * W[k,n]
    OA[m,n] = Activation(OA[m,n])
```
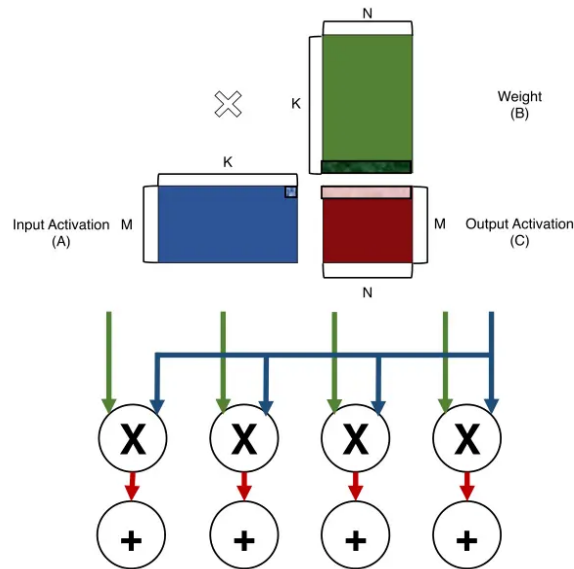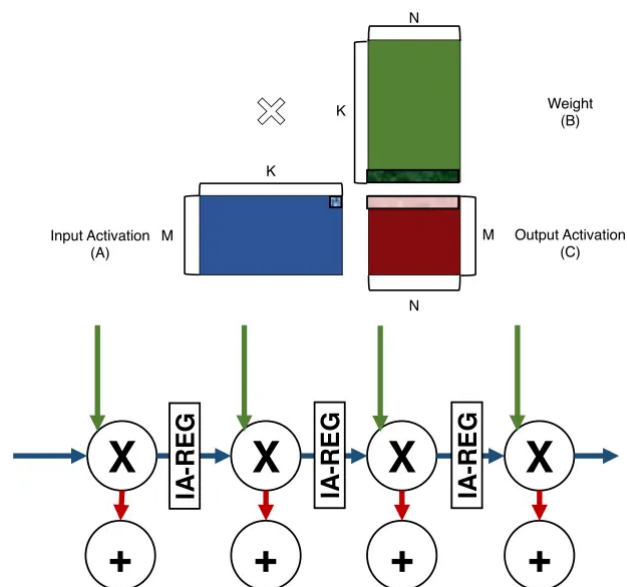


- **Adder Tree**

- Example: NVDLA, DianNao

# Spatial, N

```
for m in range(M):
  OA[m,:] = 0
  for k in range(K):
    OA[m,:]+= IA[m,k] * W[k,:]
```
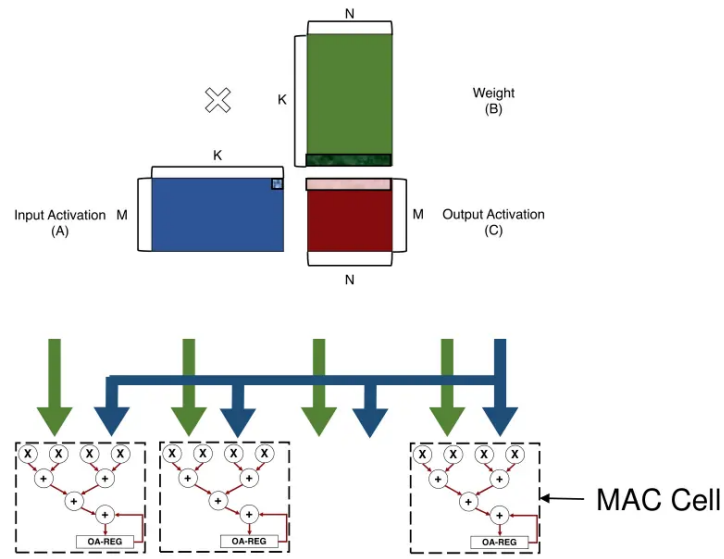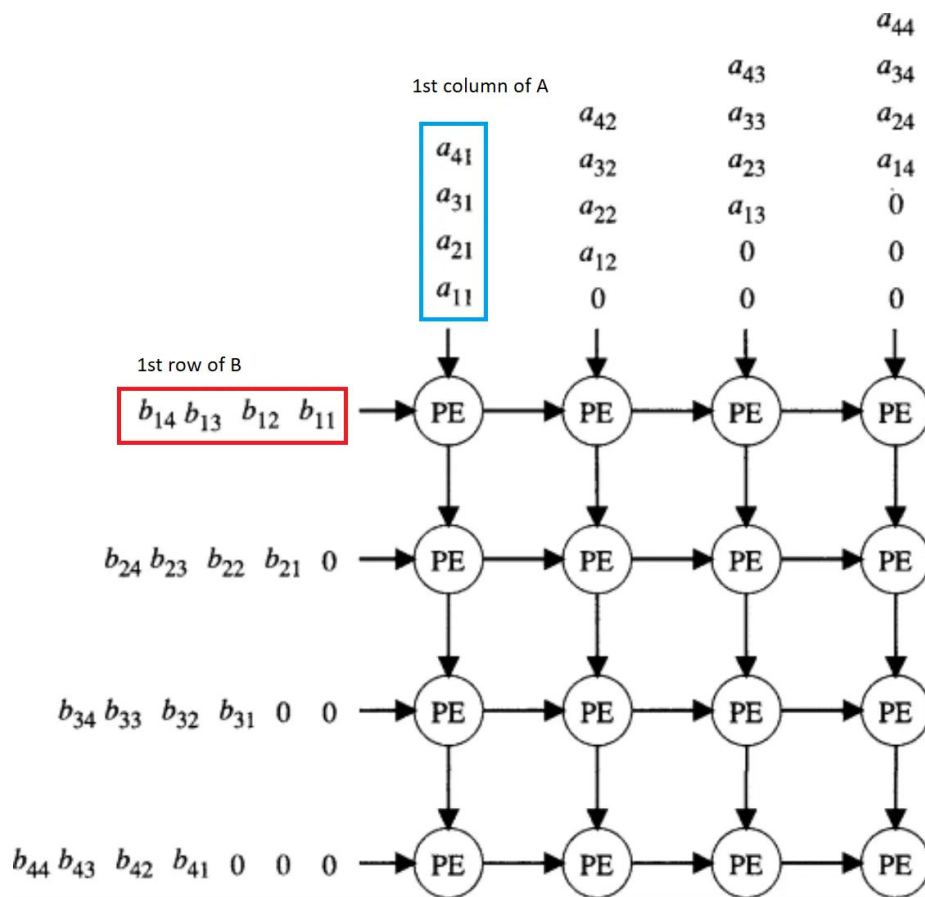


- Direct-wiring multicast
- Example: NVDLA, DianNao



- Systolic multicast
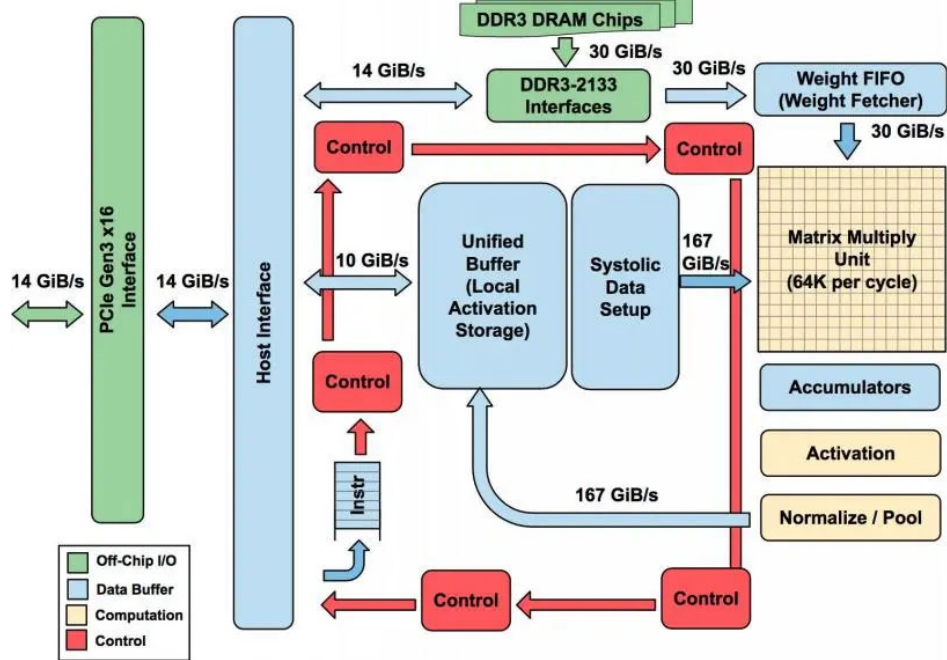- Example: TPU, Gemmini

# Combined : NVDLA

# Systolic Array

## Output-stationary



- 每个PE计算4次乘法, 4次累加
- 10 个周期完成计算

## Weight-stationary: TPU

$$A \times B \rightarrow C$$

- 脉动矩阵的第 $j$ 列固定 $B$ 的第 $j$ 列，计算 $C$ 的第 $j$ 列
- $C$ 中的每个值从脉动矩阵第一行落下
- 4个周期完成计算

https://inst.eecs.berkeley.edu/~ee290-2/sp21/