**Gigascale Reliable Energy Efficient Nanosystem (GREEN) Lab**
**School of Electrical and Computer Engineering, Georgia Tech**

*Exploring reliable, energy efficient computing solutions at nanometer nodes*
*— from devices to circuits to systems*

GIGASCALE RELIABLE ENERGY EFFICIENT NANOSYSTEMS LAB

# NeuroCube: A Programmable Digital Neuromorphic Architecture with High-Density 3D Memory

**Section 6 - 2**

**Duckhwan Kim**[1], Jaeha Kung[1], Sek Chai[2],
Sudhakar Yalamanchili[1], Saibal Mukhopadhyay[1]

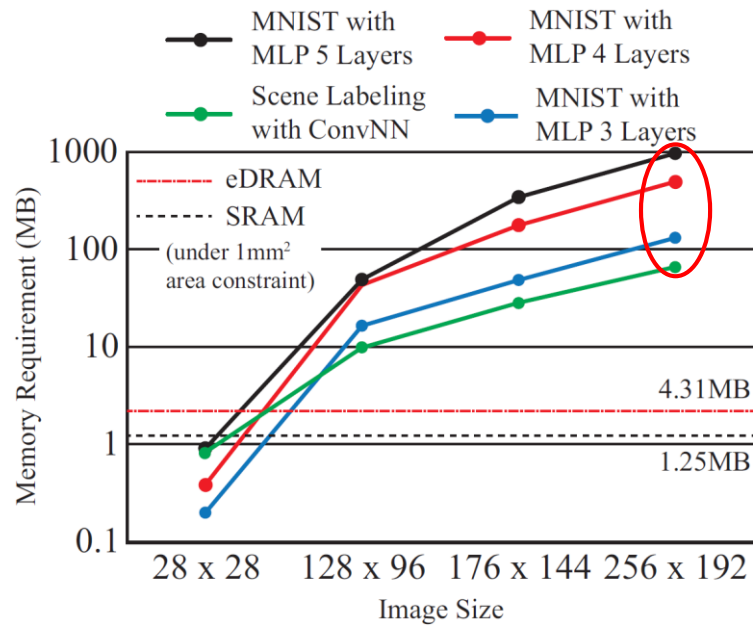[1]School of Electrical and Computer Engineering, Georgia Institute of Technology
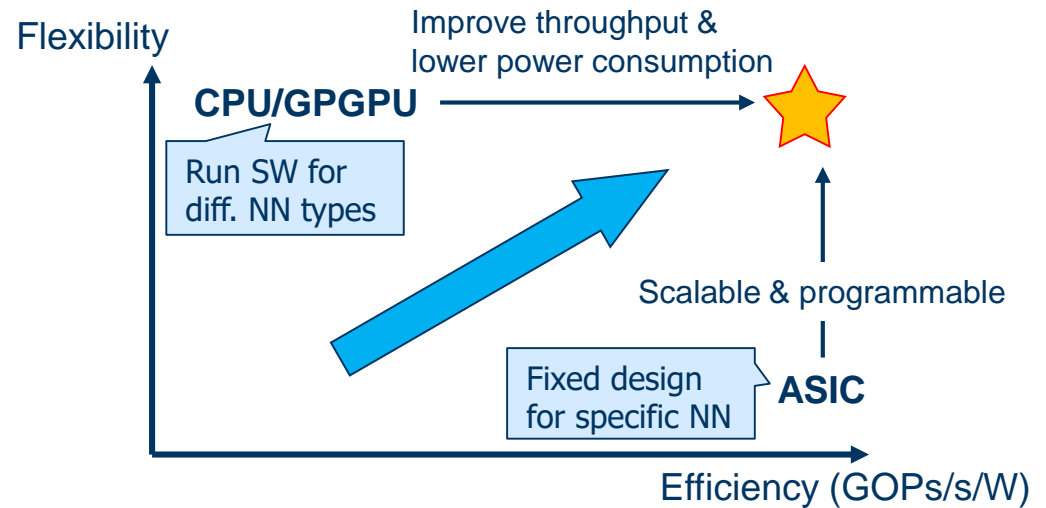[2]SRI International

1

# Outline

- Motivation
- Base Architecture of NeuroCube as PIM
- Programming NeuroCube
- Out-of-Order Packet Arrival
- Simulation
- Conclusion

# Digital Accelerator Design for Neuromorphic Algorithm

- Low operation density (ops/byte)
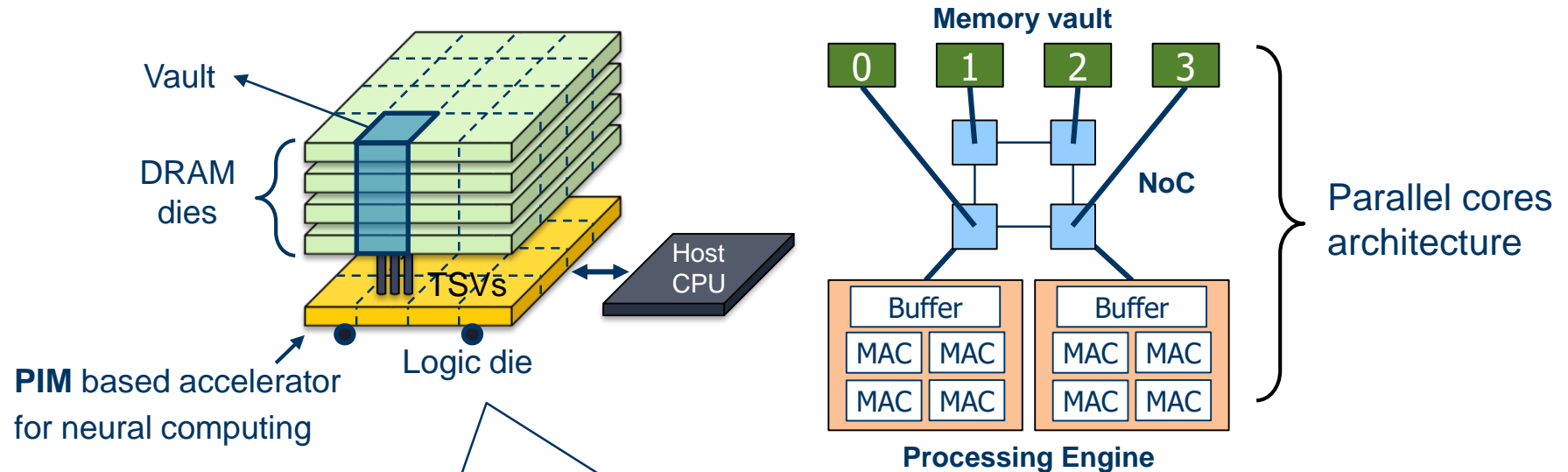
- Massive date required

Digital neuro-inspired architecture with **programmability** to cover different types of neural networks, **scalability**, and **high energy efficiency**

# NeuroCube: Process-in-Memory Architecture for Neural Computing

**Programmable, scalable platform as processor in memory**



Vault

DRAM dies

TSVs

Host CPU

**PIM** based accelerator for neural computing

Logic die

Memory vault

| 0 | 1 | 2 | 3 |

NoC

Parallel cores architecture

Buffer

| MAC | MAC |
| MAC | MAC |

Buffer

| MAC | MAC |
| MAC | MAC |

**Processing Engine**

### Hybrid Memory Cube (HMC)

- Heterogeneous integration
- Flexible logic die design
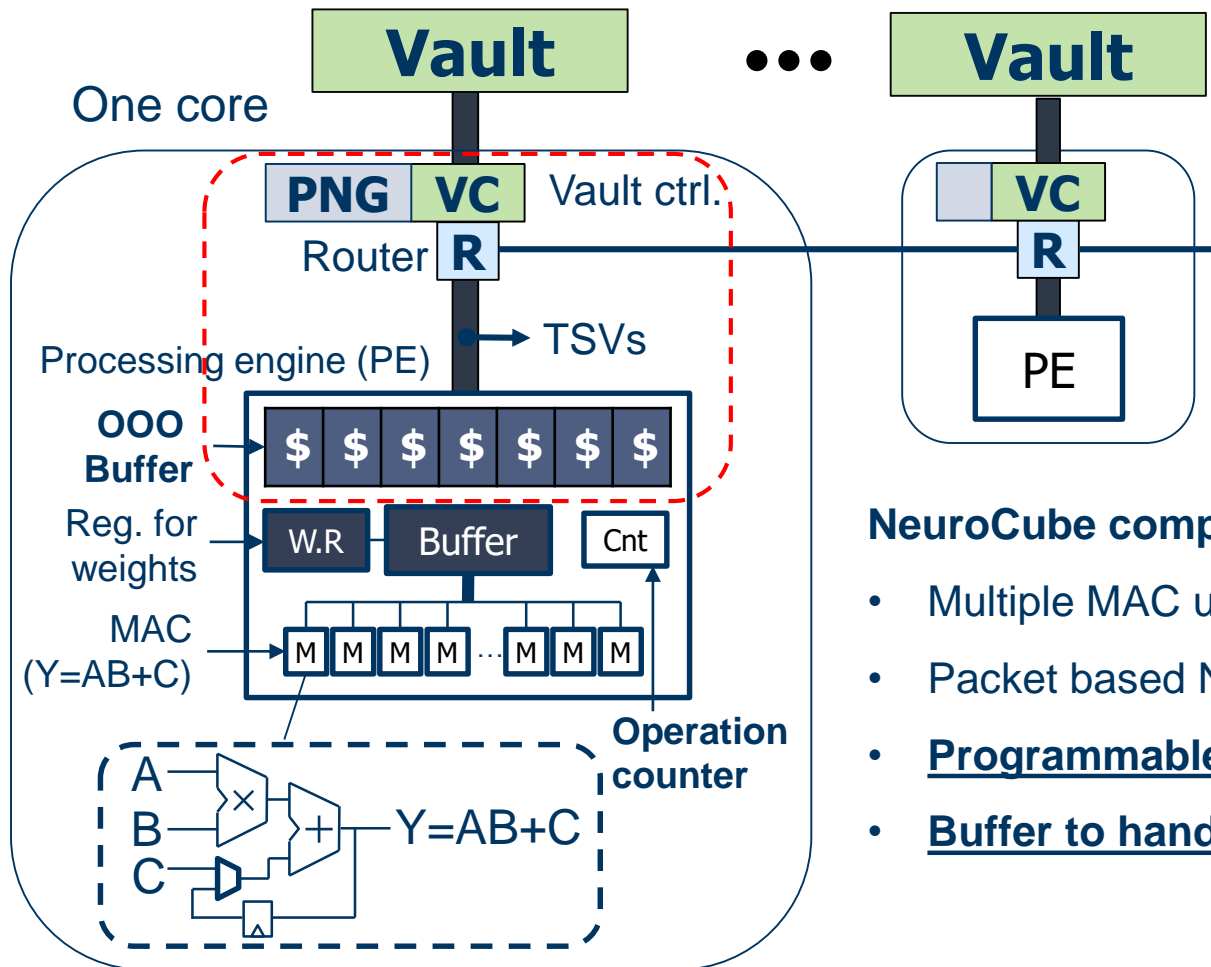
Fabricated by Micron
[J. Jeddeloh `12 TVLSI]

Q1. Neural computing layer should meet thermal and area constraint in 3D stacked DRAM

Q2. NeuroCube should be programmable to cover different types of neural network

# Basic NeuroCube Architecture

## Processor-in-memory + Parallelism



**NeuroCube computing core**

- Multiple MAC units and its temporal buffer

- Packet based NoC router

- **Programmable neurosequence generator**
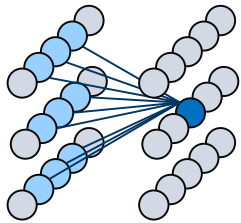
- **Buffer to handle out-of-order packet arrival**
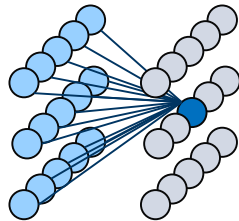
# Operational Model of NeuroCube

# Property of Neural Network: Deterministic Connections in Inference
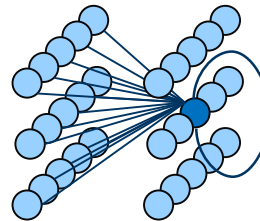
**2D convolutional**

**Fully connected**

**Recurrent connected**

**General expression of artificial neural network (ANN)**

locally
neighborhood

All neurons
in prev. layer

All neurons
in prev. layer + current layer

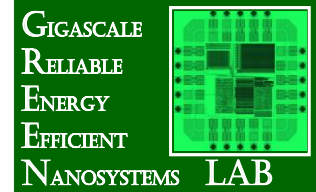$$y_j = \varphi \left( \sum_{i \in I} w_{(j,i)} \cdot y_i \right)$$

$I$: set of connected neurons

- Different NN layer can be mapped by changing set of connected neurons

  - Different data movements (memory address) can map different NN layers in NeuroCube

7

**Property of Neural Network:
Deterministic Connections in Training**

GIGASCALE
RELIABLE
ENERGY
EFFICIENT
NANOSYSTEMS  LAB

- Backpropagation with gradient descent
  - $\delta_i = \left(W_{i,i+1}^T \times \delta_{i+1}\right) * \varphi'(y_i)$ (hidden) or $= -(d - y_i) * \varphi'(y_i)$ (output)
    - $\varphi'$ : derivative of NL activation function
    - $*$: element-wise multiplication
    - $\gamma$: learning rate
  - $\Delta W_{i-1} = \delta_i \times y_{i-1}^T \,, W_{i-1} = W_{i-1} + \gamma \Delta W_{i-1}$

- It is composed of
  - Matrix-vector multiplication, element-wise multiplication, and outer product
  - Can be mapped to FC layer with <u>many zeros</u> in matrix

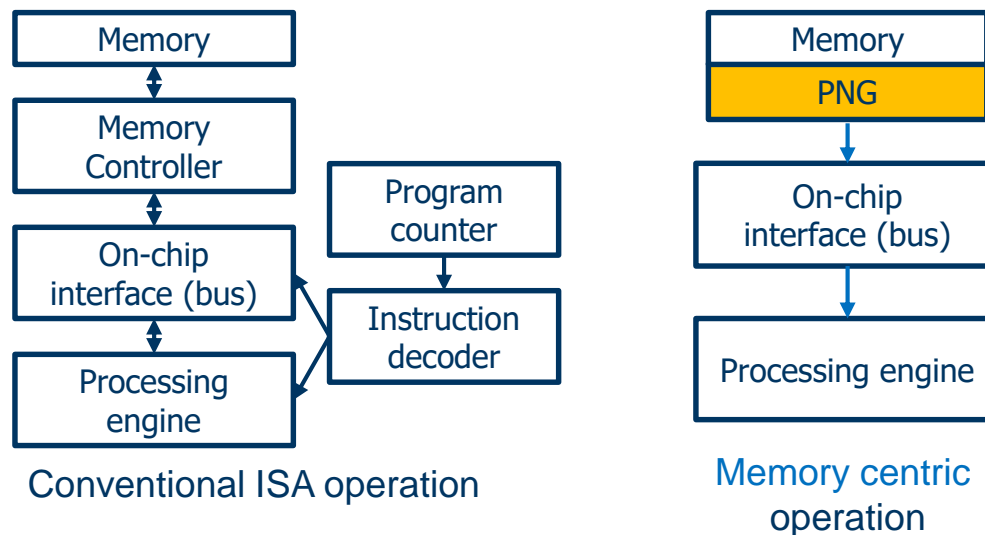- Training in neural network <u>still has deterministic connections</u>

# Memory Centric Neural Computing

**Programmable Neurosequence Generator (PNG)**

- **Sequence of operands** is predetermined
- Based on the sequence, memory **can push the data** without request
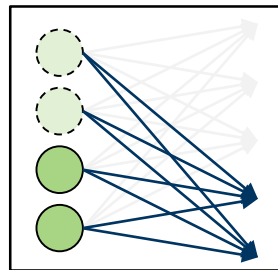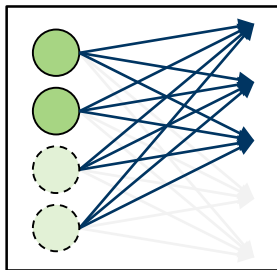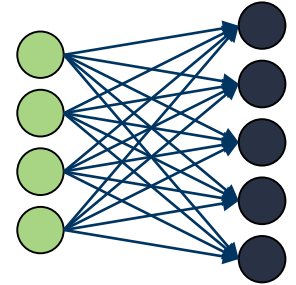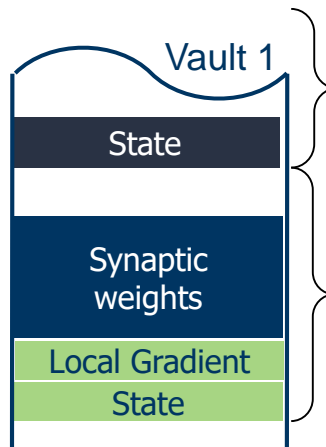- Data is delivered as packet through NoC
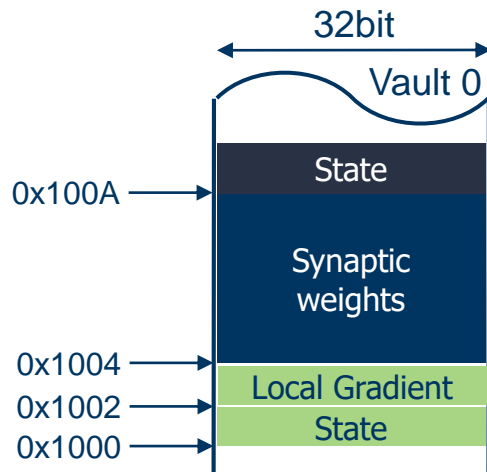
Conventional ISA operation

Memory centric operation

# MCNC Data Flow – Initial Memory Mapping

- Assume 2 vaults in NeuroCube, 3 MACs/PE
  - For synaptic weights
    - Divide synaptic weights matrix into 2 vaults
  - For previous layer
    - Divide input previous into 2 vaults or
    - **Duplicate prev. layer into all vaults** (reduce NoC traffic)

$W_{1,2}$= [5 by 4] matrix

- All neuron's states and synaptic weights are 16bit fixed point

32bit

Vault 0
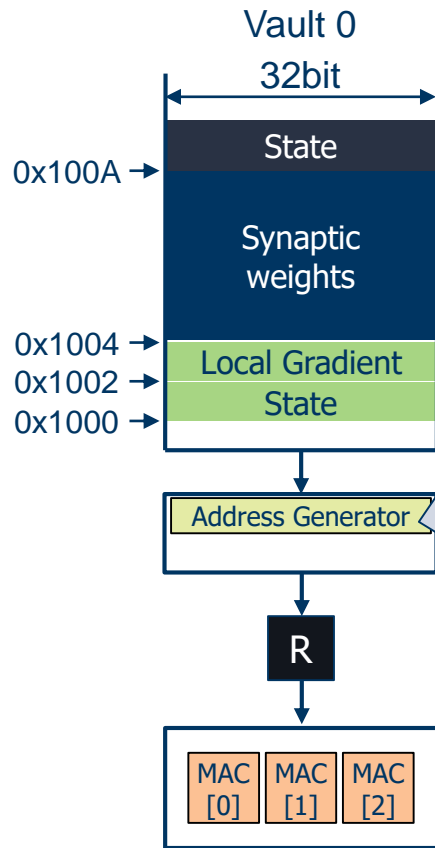
Vault 1

Memory range for next layer

State

State

0x100A

Synaptic weights

Synaptic weights

Memory range for this layer

0x1004

Local Gradient

Local Gradient

0x1002

State

State

0x1000

10

- Assume 2 vaults in NeuroCube, 3 MACs/PE
- For each vault, it will push data to NoC within packet form

Vault 0

32bit

| SRC | DST | DATA | MAC-ID | OP-ID |
|------|------|------|--------|-------|
| 4bit | 4bit | 16bit | 4bit | 8bit |

State

0x100A →

Synaptic weights

0x1004 →
0x1002 →   Local Gradient
0x1000 →   State

Address Generator

R

| MAC [0] | MAC [1] | MAC [2] |

```
for(j=0;j<1;j++) {//one MAC cover one neuron[j]
    for(i=0;i<4;i++) {//4 neurons[i] to one neuron[j]
        for(k=0;k<3;k++) {// state/weight/local grad
            for(m=0;m<3;m++) {//#MACs/PE = 3
                if k==0
                addr = addr_gen(i,j,m,state)
            else if k==1
                addr = addr_gen(i,j,m,weight)
            else
                addr = addr_gen(i,j,j,local_grad)
            }
        }
    }
}
```
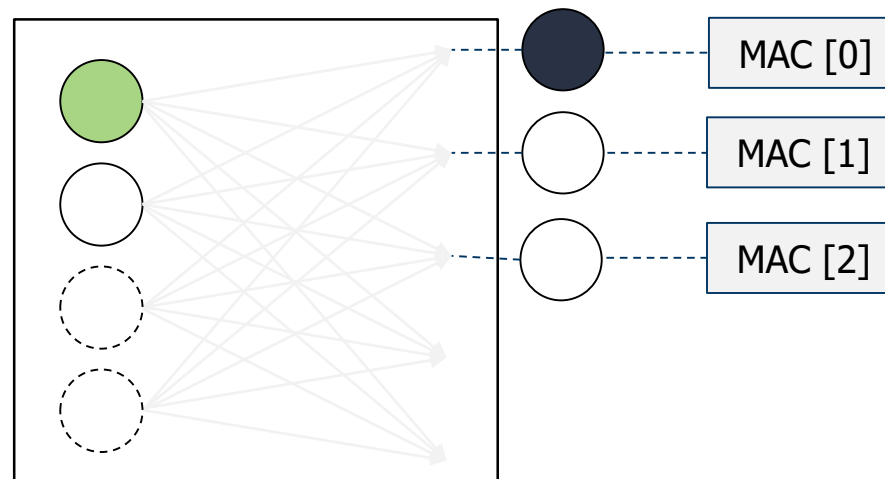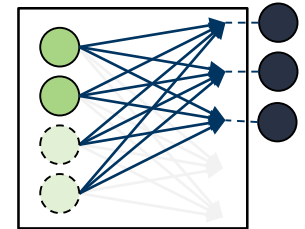
Three nested-counters

j,m

i

Push packet (state[0]) to MAC [0]

| SRC | DST | DATA | MAC-ID | OP-ID |
|-----|-----|------|--------|-------|
| 0 | 0 | State [0] | 0 | 0 |



MAC [0]

MAC [1]

MAC [2]

Push packet (state[0]) to MAC [1]

| SRC | DST | DATA | MAC-ID | OP-ID |
|-----|-----|----------|--------|-------|
| 0 | 0 | State [0] | 1 | 0 |

Push packet (state[0]) to MAC [2]

| SRC | DST | DATA | MAC-ID | OP-ID |
|:---:|:---:|:---:|:---:|:---:|
| 0 | 0 | State [0] | 2 | 0 |



MAC [0]

MAC [1]

MAC [2]

Push packet (weight[0,0]) to MAC [0]

| SRC | DST | DATA | MAC-ID | OP-ID |
|-----|-----|--------|--------|-------|
| 0 | 0 | W [0, 0] | 0 | 0 |



MAC [0]

MAC [1]

MAC [2]

# MCNC Data Flow (5)

Push packet (weight[0,1]) to MAC [1]

| SRC | DST | DATA | MAC-ID | OP-ID |
|-----|-----|--------|--------|-------|
| 0 | 0 | W [0, 1] | 1 | 0 |



MAC [0]

MAC [1]

MAC [2]

Push packet (weight[0,2]) to MAC [2]

| SRC | DST | DATA | MAC-ID | OP-ID |
|-----|-----|------|--------|-------|
| 0 | 0 | W [0, 2] | 2 | 0 |

Push packet (state[1]) to MAC [0]

| SRC | DST | DATA | MAC-ID | OP-ID |
|-----|-----|----------|--------|-------|
| 0 | 0 | State [1] | 0 | 1 |



MAC [0]

MAC [1]

MAC [2]

And so on …

# Out-of-order Packet Arrival Problem in NeuroCube

**Duplicating Previous Layer**

**Non Duplication**

$W_{1,2}$= [5 by 4] matrix

0    1

PE    PE

0    1

PE    PE

Although data access is sequential, data arrival can be **out-of-order** due to NoC congestion

# Out-of-Order data arrival (1)

R[2]

| 0 | 2 | 0x3F | 15 | 23 |

**PE[2]**

Weight Memory

OOO Buffer

| 0 |
| 1 |
| ... |
| 8 |
| ... |
| 15 |

Temp Buffer

0x7A ... | | | |
0    14   15

16 MACs

OP_CNT = 23

### Packet structure

| SRC | DST | DATA | MAC-ID | OP-ID |
|------|------|-------|--------|-------|
| 4bit | 4bit | 16bit | 4bit | 8bit |

- OP-ID == OP_CNT == 23
  - This packet is for current operation
  - It moves to temporal buffer [15] directly

# Out-of-Order data arrival (2)

R[2]

| 0 | 2 | 0x5A | 0 | 24 |

**Weight Memory**

**Temp Buffer**

| 0x7A | ... | 0x3F |
| 0 | 14 | 15 |

16 MACs

OP_CNT = 23

**PE[2]**

**OOO Buffer**

| 0 | |
| 1 | |
| ... | |
| 8 | |
| ... | |
| 15 | |

## Packet structure

| SRC | DST | DATA | MAC-ID | OP-ID |
|------|------|-------|--------|-------|
| 4bit | 4bit | 16bit | 4bit | 8bit |

- OP-ID != OP_CNT
  - This packet is for next operation (OP_CNT == 24)
  - It moves to OOO buffer [8]
    - 8 = mod(24,16)
    - OOO buffer [i] is FIFO with 64 depth

# Out-of-Order data arrival (3)



R[2]

| 1 | 2 | 0x32 | 14 | 23 |

**PE[2]**

Weight Memory

Temp Buffer

| 0x7A | ... | | 0x3F |
| 0 | | 14 | 15 |

16 MACs

OP_CNT = 23

OOO Buffer

| 0 | |
| 1 | |
| ... | |
| 8 | [0,0x5A] |
| ... | |
| 15 | |

## Packet structure

| SRC | DST | DATA | MAC-ID | OP-ID |
|------|------|-------|--------|-------|
| 4bit | 4bit | 16bit | 4bit | 8bit |

- OP-ID == OP_CNT == 23
  - This packet is for current operation
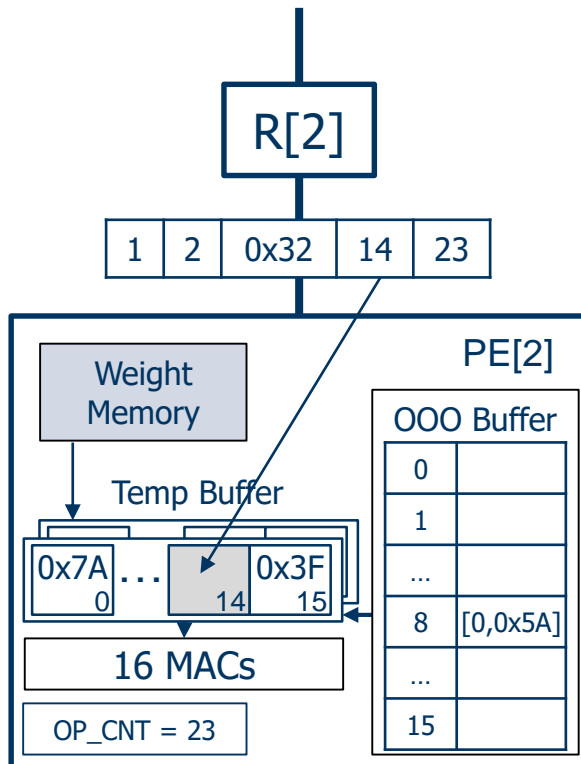  - It moves to temporal buffer [14] directly
  - Temp buffer (length 16) is full
    - It will trigger 16 MACs operation

# Out-of-Order data arrival (4)

R[2]

## Packet structure

| SRC | DST | DATA | MAC-ID | OP-ID |
|-----|-----|------|--------|-------|
| 4bit | 4bit | 16bit | 4bit | 8bit |

**PE[2]**

Weight Memory

Temp Buffer

| 0 | ... | 14 | 15 |

16 MACs

OP_CNT = 24

**OOO Buffer**

| 0 | |
| 1 | |
| ... | |
| 8 | [0,0x5A] |
| ... | |
| 15 | |

- Temp buffer trigger 16 MACs
- Increase OP_CNT
- Ready for 24th operation

GIGASCALE
RELIABLE
ENERGY
EFFICIENT
NANOSYSTEMS  LAB

R[2]

PE[2]

Weight Memory

Temp Buffer

... 

0    14    15

16 MACs

OP_CNT = 24

OOO Buffer

| 0 | |
| 1 | |
| ... | |
| 8 | [0,0x5A] |
| ... | |
| 15 | |

Packet structure

| SRC | DST | DATA | MAC-ID | OP-ID |
|------|------|-------|---------|--------|
| 4bit | 4bit | 16bit | 4bit | 8bit |

- Before capture new packet, check OOO buffer[8]
  - 8 = mod(24,16)
  - Full search OOO buffer [8] (64 depth)
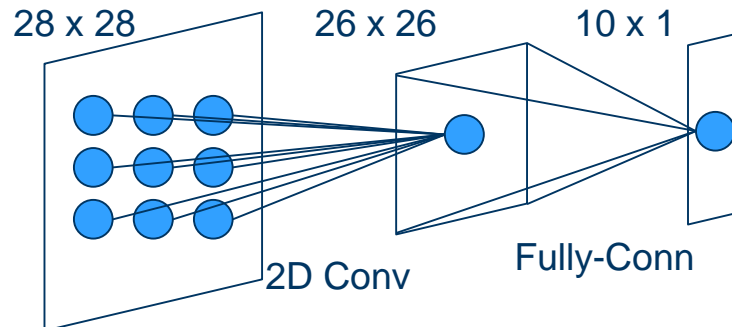  - Move [0;2;0x5A;0;24] to temp buffer [0]

# Layerwise Programming

Address generator can be designed using **three nested counters**



Counter | Comparator | Conf. Register

#N: number of neurons in current layer
#C: number of connections from previous layer to the neuron in current layer

For **each layer**, host program the NeuroCube by **writing conf. registers**



28 x 28      26 x 26      10 x 1

2D Conv      Fully-Conn

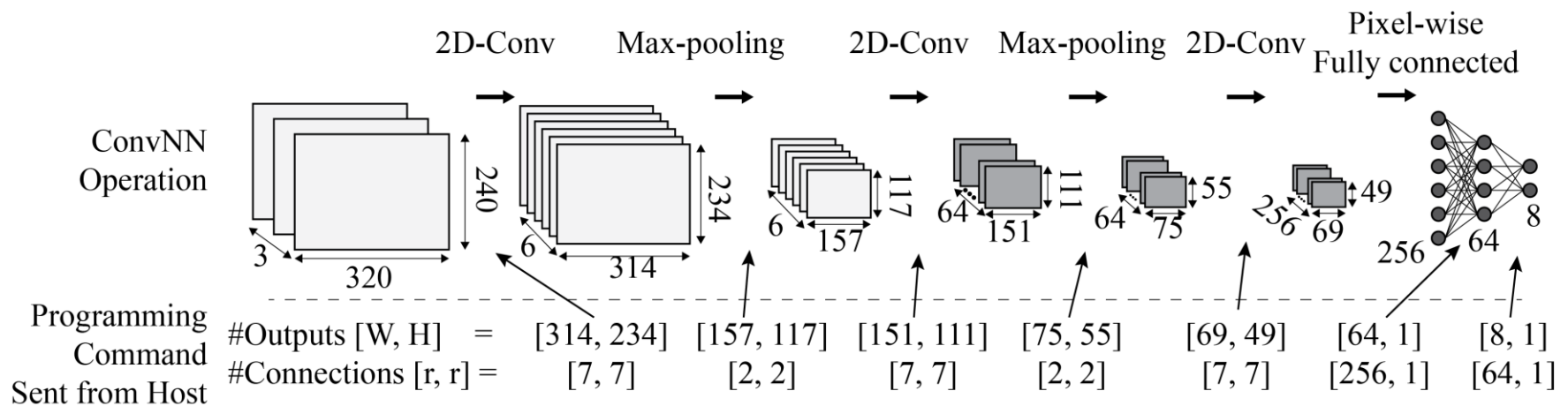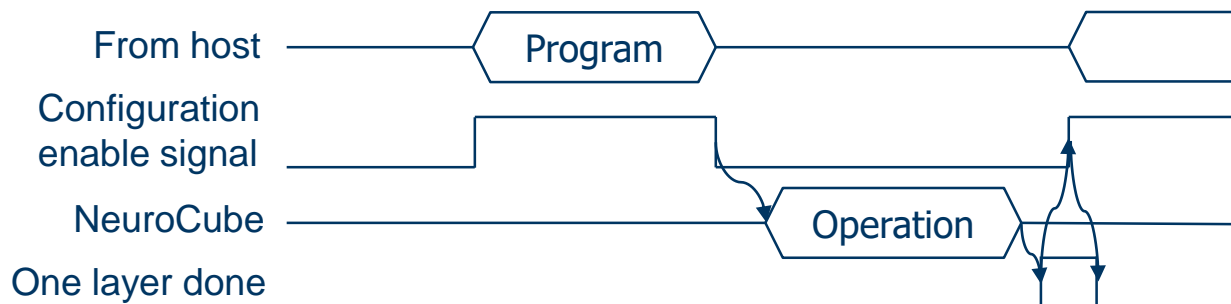For 2D-Conv: #N: 26 x 26, #C: 3 x 3
For FC:          #N: 10 x 1,   #C: 676 x 1

- Latency of writing configuration registers is negligible

- External interface is very rare

## Conv-NN for Scene labeling
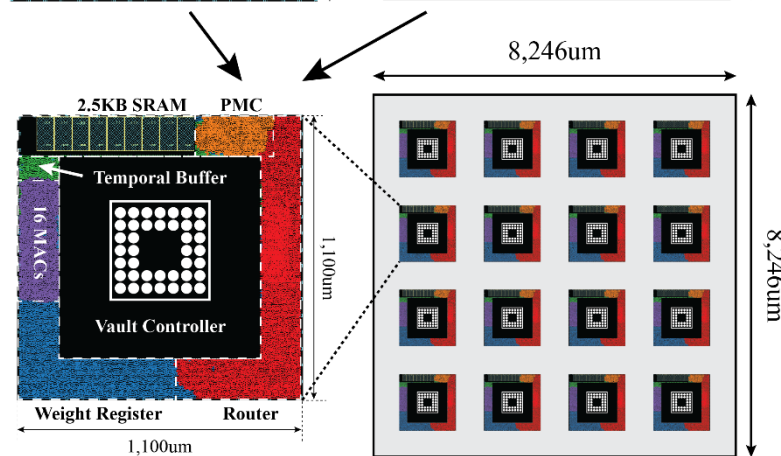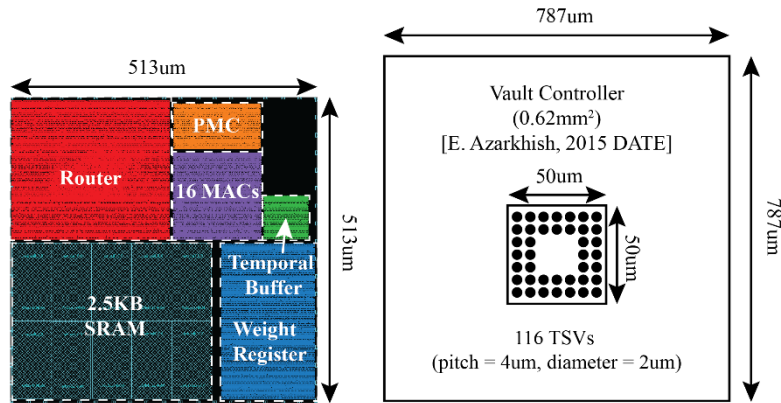
- Handshaking to start programming or main operation



[S. Goud, `09 ICCV]

# Simulation Results

# Synthesis Result

## 28nm CMOS process



513um

Router
PMC
16 MACs
513um
Temporal Buffer
2.5KB SRAM
Weight Register

787um

Vault Controller
(0.62mm$^2$)
[E. Azarkhish, 2015 DATE]

50um

50um

787um

116 TSVs
(pitch = 4um, diameter = 2um)

2.5KB SRAM    PMC
Temporal Buffer
16 MACs
Vault Controller
Weight Register    Router
1,100um

1,100um

One core in Neurocube

8,246um

8,246um

Neurocube on HMC logic die

Footprint of HMC 1.0: 68mm$^2$

## 15nm FinFet process

- To utilize HMC internal throughput fully, $f_{PE}$ = 5GHz
- Thermal analysis performed under 5GHz operating
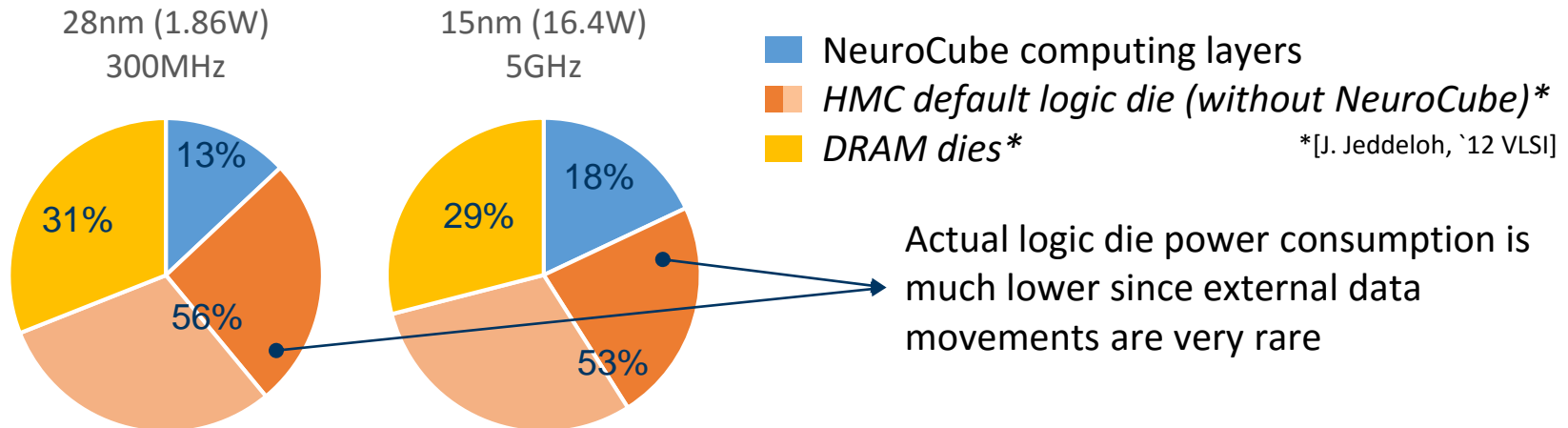- Max. allowable temp. of HMC: 378K



349K
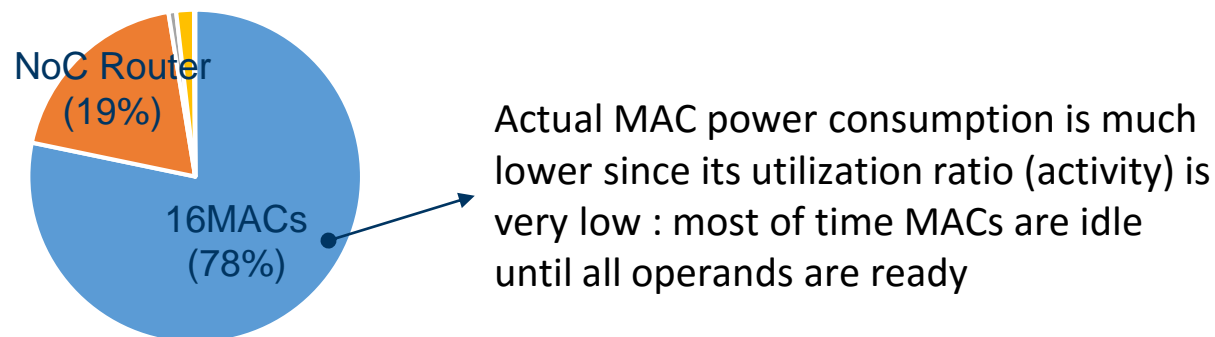
321K

[J. Jeddeloh `12 TVLSI] [HMC specification 1.0]

# Hardware Power Breakdown

- Measured energy consumption [J. Jeddeloh, `12 VLSI]
  - 3.7 pj/bit for the DRAM layers
  - 6.78 pj/bit for the logic layer, (most power hungry = ext. interface SERDES)

28nm (1.86W)
300MHz

15nm (16.4W)
5GHz

■ NeuroCube computing layers
■ *HMC default logic die (without NeuroCube)\**
■ *DRAM dies\**                    *[J. Jeddeloh, `12 VLSI]

**28nm (1.86W) 300MHz pie:** 13%, 56%, 31%

**15nm (16.4W) 5GHz pie:** 18%, 53%, 29%

Actual logic die power consumption is much lower since external data movements are very rare

- Single core power breakdown (15nm): 187mW

NoC Router (19%)

16MACs (78%)

Actual MAC power consumption is much lower since its utilization ratio (activity) is very low : most of time MACs are idle until all operands are ready
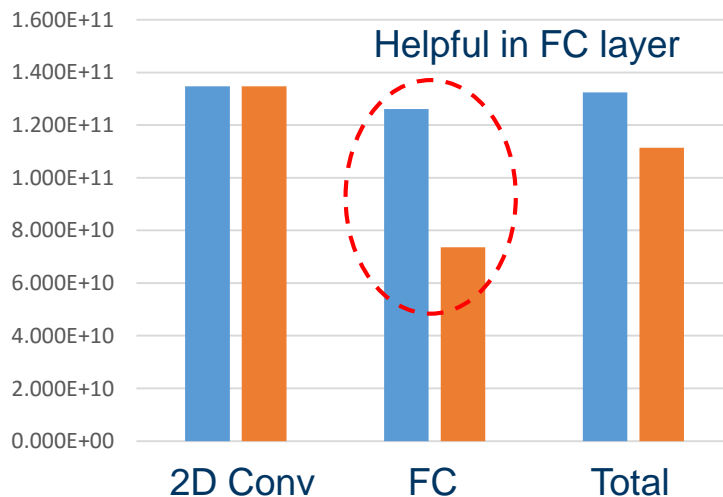
# Performance Simulation: Inference

**To see system throughput, cycle-level simulation is performed**

- Inference: scene labeling

  15nm Finfet design

- Two operating modes:

  1. Duplication prev. layer to reduce NoC traffic
  2. No duplication to reduce memory overhead



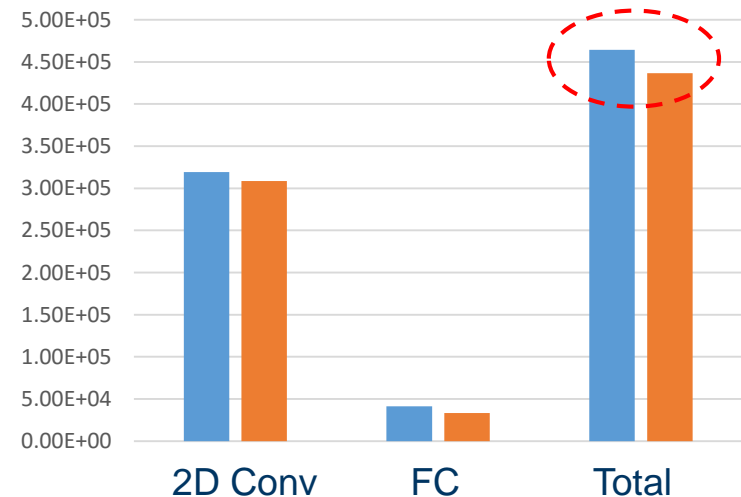Duplicating prev. layer       No duplicating

System throughput (OPs/s)

Helpful in FC layer

Memory requirement (byte)

Memory overhead is not significant because weights are dominant

2D Conv     FC     Total

2D Conv     FC     Total

31

15nm Finfet design

System throughput (OPs/s)

Memory requirement (byte)

Maintain constant system throughput

Memory requirement is high even for small image size (64 x 64)

# Performance Simulation: DDR3/Crossbar

15nm Finfet design



System throughput (OPs/s)

HMC   DDR3

FC layer

Duplicate   Non-duplicate

**Multiple channels in HMC improve system throughput**

System throughput (OPs/s)

Mesh grid 2D NoC   Crossbar

2D Conv   FC Layer

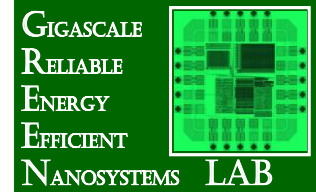**Crossbar improves system throughput for FC layer**

# Related Works

- **Most of prev. work focused on specific NN**
  - Shows high throughput based on optimized design (ASIC/FPGA)
  - Not programmable (not scalable)

- **Programmable + Scalable design**
  - General purposed architecture: mobile CPU or GPU
  - Integrated systems with external DRAM

|  | [L. Cabigelli `15 DAC] | | NeuroCube | |
| --- | --- | --- | --- | --- |
| Platform | Tegra K1 | GTX 780 | 28nm | 15nm |
| Bit precision ctrl. | N/A | N/A | 16bit | 16bit |
| Throughput (GOPs/s) | 76 | 1781 | 7.95 | 132.4 |
| Computing Power (W) | 11 | 206.8 | 0.249 | 3.41 |
| Efficiency (GOPs/s/W) | 6.91 | 8.61 | **31.92** | **38.82** |
| Inference/training | inference | inference | both | both |

# Conclusion

- **NeuroCube: Neuro-inspired architecture as PIM in HMC**
  - Utilize high memory bandwidth
  - Integrated with high density memory
  - Meet thermal/area constraints

- **Programmable architecture to cover diff. NN types**
  - Programming memory access pattern
  - Simple memory address generator (PNG) is embedded in memory
  - Memory centric neural computing (MCNC) scheme

- **System performance is simulated**
  - Network-on-chip traffic is next bottleneck
  - Optimized NoC design, data re-usage should be studied

# Thank you