# A Systolic Dataflow Based Accelerator for CNNs

Saptarsi Das, Arnab Roy, Kiran Kolar Chandrasekharan, Ankur Deshwal
SAIT, Samsung R&D Institute India-Bangalore, India
Email: saptarsi.das, arnab.roy, kiran.kc, a.deshwal@samsung.com
Sehwan Lee
SAIT, Samsung Electronics, South Korea
Email: sehwan.b.lee@samsung.com

**2020 IEEE International Symposium on Circuits and Systems**
**Virtual, October 10-21, 2020**

# Content

- Motivation & Background
- Design & Implementation of the Proposed Accelerator
- Experimental Results
- Conclusions

# Motivation & Background

**Why Accelerate CONV?**

- Convolution Neural Networks (CNN) are ubiquitous in modern AI systems.
- Power efficient execution of CNN models is crucial especially in mobile phones and other hand held devices.
- In CNN models, the most compute intensive parts are convolution operations (CONV).

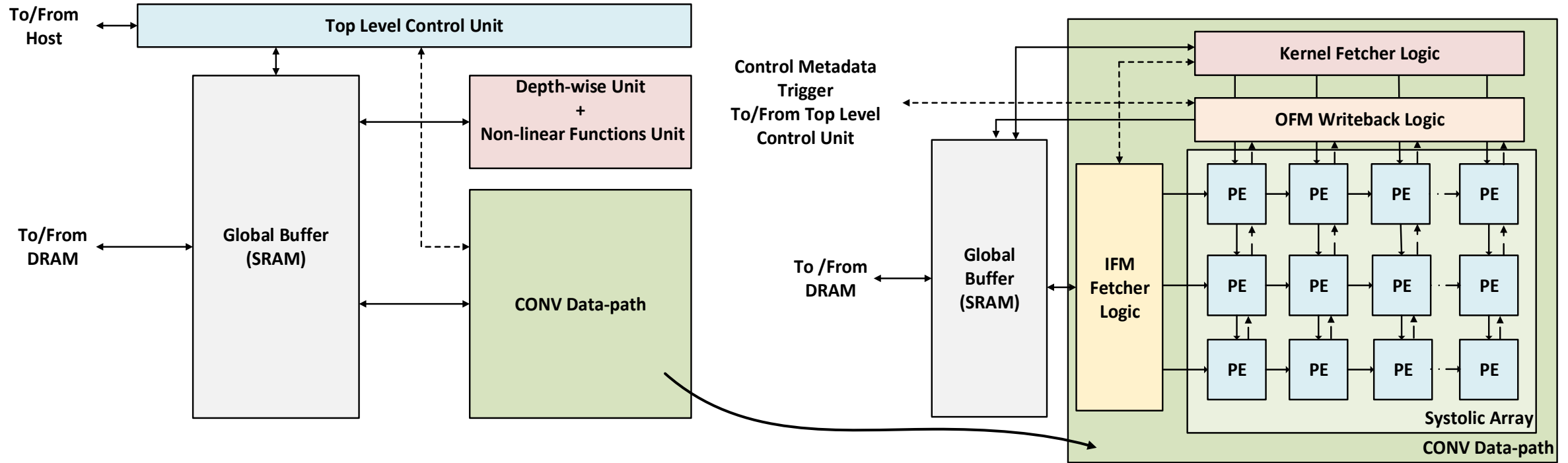SAMSUNG

# Motivation & Background

**Why Use Systolic Array?**

- Systolic arrays have a long history
- Have largely been off the mainstream
- Offers scalability, structural simplicity, ability to exploit parallelism and data reuse
- Renewed interest as candidate architecture for neural net accelerators

**Why A New Systolic Array?**

- Typically systolic arrays use scalar Processing Elements (PE)
- Cost of accumulation is high
- To mitigate this we propose a systolic array with inner-product units as the PEs

SAMSUNG ADVANCED INSTITUTE OF TECHNOLOGY

SAMSUNG

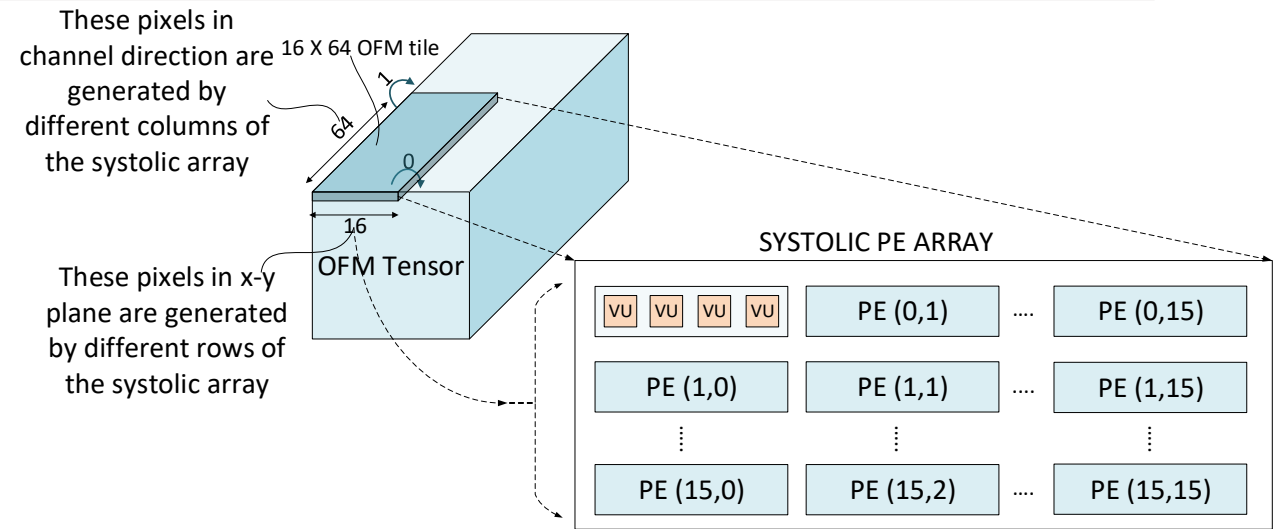# High Level View of the Accelerator



**Primary Components of the Accelerator**

- A Top Level Control Unit
- A Depth-wise Unit + Non-Linear Functions Unit
- A CONV Data-path
- A Global SRAM Buffer

# Dataflow of CONV

- We follow an Output Stationary Traversal (OS)
- The Output Feature Map (OFM) 3D tensor is tiled as shown in Fig 1.
- Fig 2. shows delivery of Input Feature Map (IFM) & kernel into the systolic array

These pixels in channel direction are generated by different columns of the systolic array

16 X 64 OFM tile

64
16
1
0

OFM Tensor

These pixels in x-y plane are generated by different rows of the systolic array

SYSTOLIC PE ARRAY

| VU VU VU VU | PE (0,1) | .... | PE (0,15) |
| PE (1,0) | PE (1,1) | .... | PE (1,15) |
| PE (15,0) | PE (15,2) | .... | PE (15,15) |

```
for ofm_ch = 0 to C, stride = 64        //Loop level 1 – traversing OFM channels
    for ofm_px = 0 to HxW, stride = 16      //Loop level 0 – traversing OFM pixels
        systolic_execution()                    //Generation of 64x16 OFM tile
```

Fig. 1

- The IFM streams flow through a rows and get reused
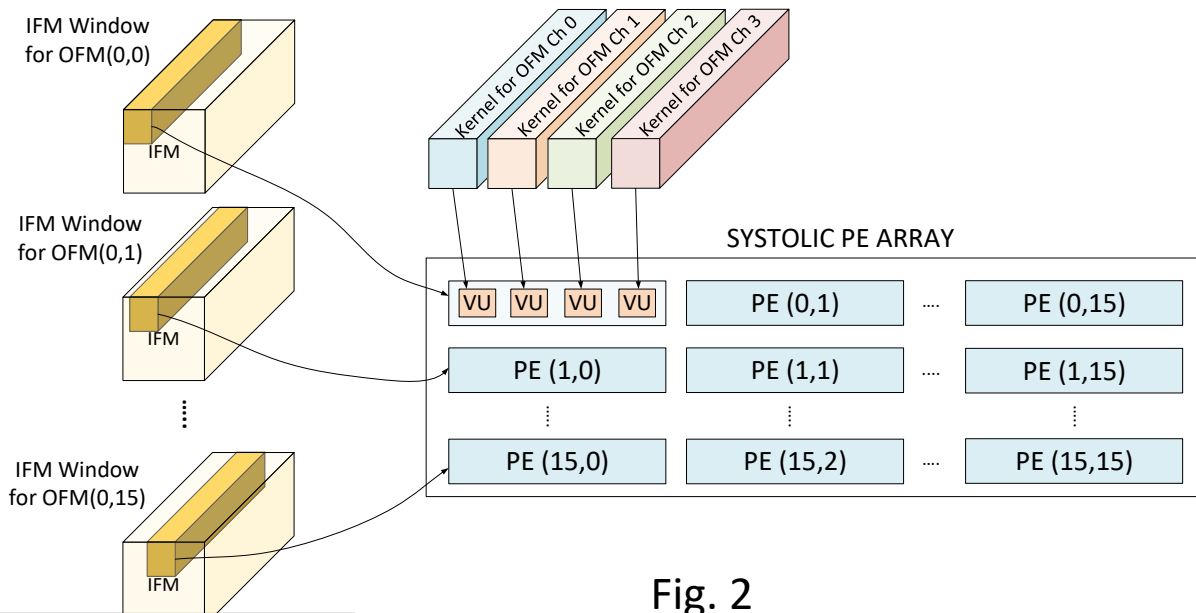- The kernel streams flow through columns and get reused

IFM Window for OFM(0,0)

IFM

Kernel for OFM Ch 0
Kernel for OFM Ch 1
Kernel for OFM Ch 2
Kernel for OFM Ch 3

IFM Window for OFM(0,1)

IFM

SYSTOLIC PE ARRAY

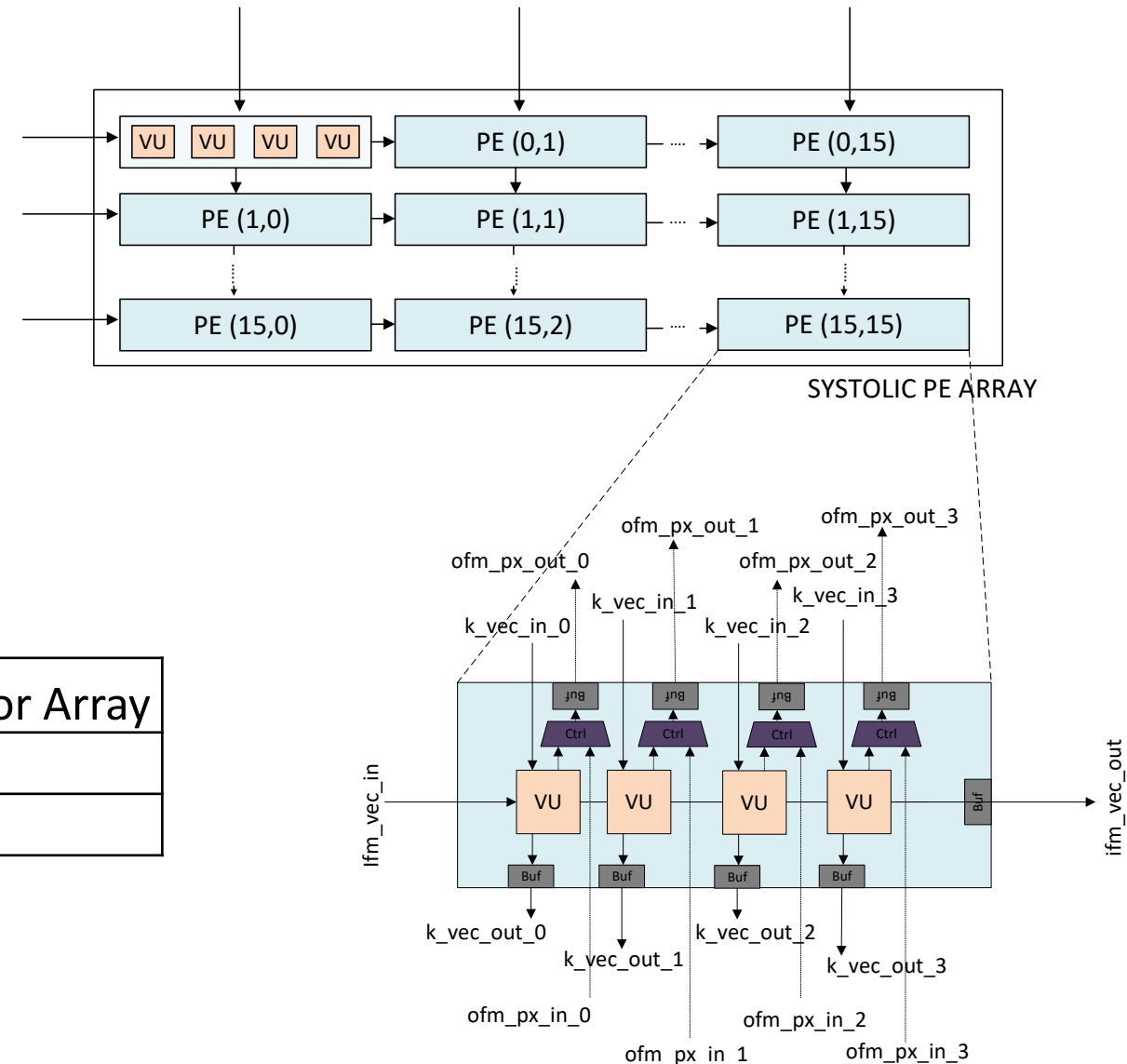| VU VU VU VU | PE (0,1) | .... | PE (0,15) |
| PE (1,0) | PE (1,1) | .... | PE (1,15) |
| PE (15,0) | PE (15,2) | .... | PE (15,15) |

IFM Window for OFM(0,15)

IFM

Fig. 2

# Features of the Accelerator

- Hierarchical vector units in PEs
- Channel-major data layout
- Overlapping of computation and communication
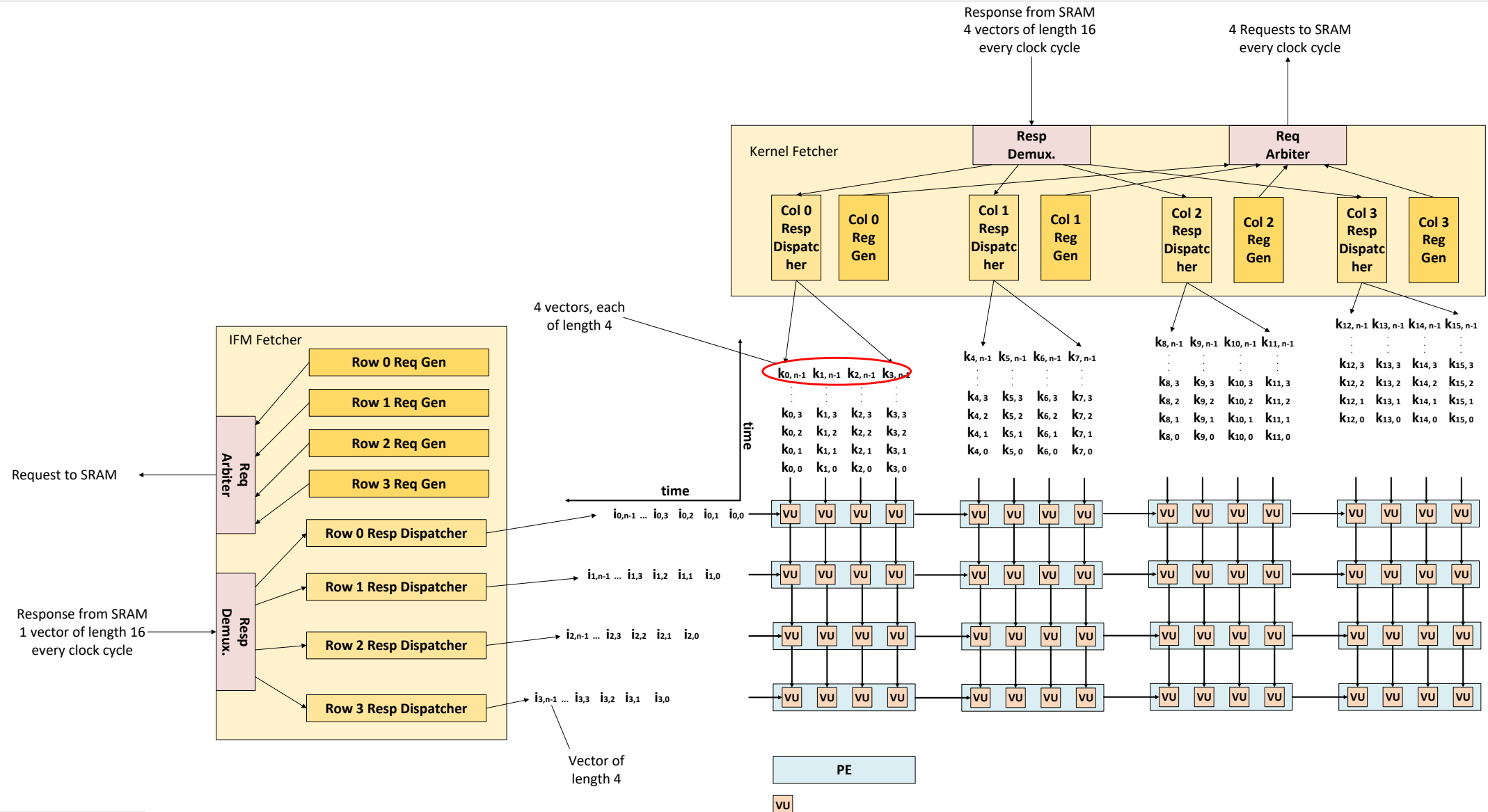- Software-controlled scratch-pad SRAM with configurable logical buffers

# Hierarchical vector units in PEs

- Each PE houses a multiplicity of vector units (VU)
- We implement two systolic arrays using Samsung 10 nm library (at 800 MHz)
  - One with scalar PEs
  - One with PEs consisting of 4 VUs. Each VU performs inner-product on vectors of length 4
  - Peak throughput of both arrays = 4096 INT8 Multiply-Accumulate (MAC)/cycle

|  | Scalar Array | Hierarchical Vector Array |
|---|---|---|
| Power-Efficiency (TOPs/W) | 7.1 | 14.2 |
| Area-Efficiency (TOPs/mm²) | 10.7 | 21.8 |



SYSTOLIC PE ARRAY

# IFM & Kernel Dataflow through a 4x4 Array

# Channel-major Data Layout

- We employ a channel-major data layout
- The IFM and kernel tensors are fetched as 1x1x64 vectors (in a 16x16 array)
- We measure memory read efficiency as #Bytes Used/#Bytes Read
- We set the length of channel vector in channel-major storage to 64 and dimension of 2D tile to 4x4 in the XY-major storage.

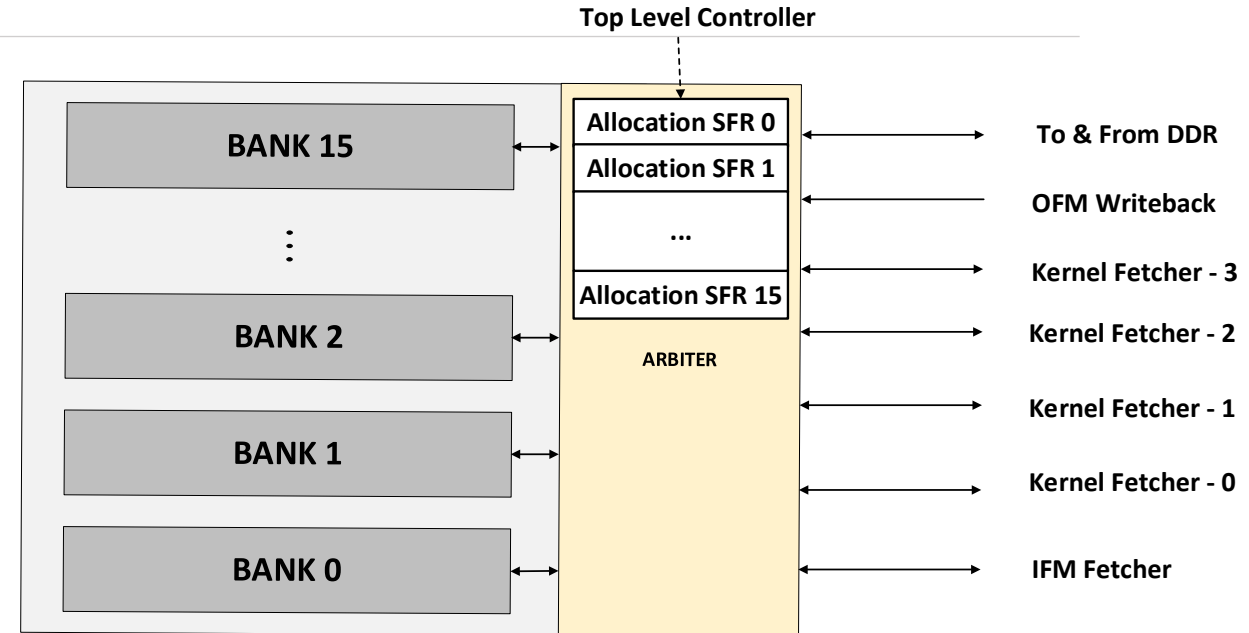| Network Models | Memory read efficiency | |
| --- | --- | --- |
| | Channel-major | XY-major |
| InceptionV3 | 78.89% | 17.79% |
| Resnet50 | 99.13% | 10.88% |
| Inception-Resnet | 90.63% | 15.96% |

# Overlapping of Computation & Communication

- Each OFM tile generation involves computation & write-back
- To improve performance we overlap next tile computation with current tile write-back
- Data hazard is prevented through software-controlled mode selection

| Network Models | TOPs | | |
|---|---|---|---|
| | Non-overlapped | Overlapped | % Increase |
| InceptionV3 | 3.2974 | 4.16 | 26.1 |
| Inception-Resnet | 3.8859 | 4.84 | 24.5 |
| Resnet50 | 3.3467 | 3.77 | 12.6 |

# Software-controlled Scratchpad SRAM

- Software-controlled scratchpad memory
- Configurable logical buffers
  - Minimal hardware overhead
  - Support for variable size IFM, kernel, OFM tensors
  - Logic synthesis performed using Samsung 10 nm library (at 800 MHz)
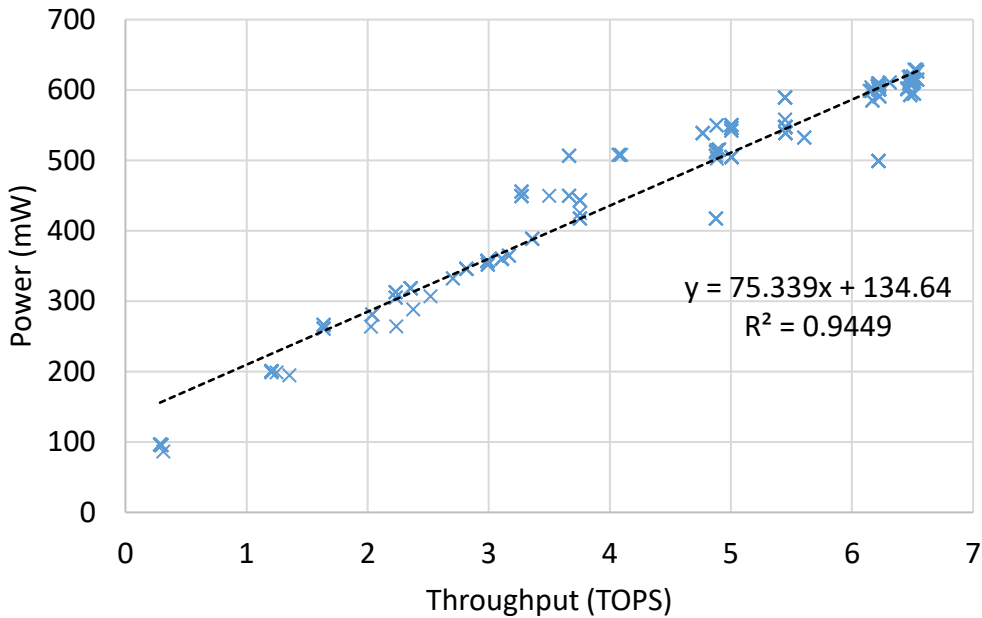
**Top Level Controller**

BANK 15

⋮

BANK 2

BANK 1

BANK 0

Allocation SFR 0
Allocation SFR 1
...
Allocation SFR 15

ARBITER

To & From DDR
OFM Writeback
Kernel Fetcher - 3
Kernel Fetcher - 2
Kernel Fetcher - 1
Kernel Fetcher - 0
IFM Fetcher

| | Area (mm²) | | |
|---|---|---|---|
| Total Area | Combinational Area | Sequential Area | SRAM Bank Area |
| 0.785442 | 0.01068 | 0.001243 | 0.774739 |
| 100.00% | 1.35% | 0.15% | 98.50% |

# Area & Power Results

| Module | Area | | Peak Power | |
|---|---|---|---|---|
| | mm² | Percentage | W | Percentage |
| SRAM | 0.785442 | 67.72% | 0.1147 | 19.56% |
| Systolic Array | 0.3008 | 25.93% | 0.416 | 70.97% |
| Support Logic | 0.073505 | 6.33% | 0.0555 | 9.47% |
| Total | 1.159747 | | 0.5864 | |

Table 1

| Module | Module-wise Power Breakup | | |
|---|---|---|---|
| | InceptionV3 | Resnet50 | Inception-Resnet |
| SRAM | 19.60% | 17.05% | 18.34% |
| Systolic | 71.00% | 73.53% | 72.78% |
| Support | 9.40% | 9.42% | 8.88% |

Table 2



$y = 75.339x + 134.64$
$R^2 = 0.9449$

# Comparison with Other Accelerators

| | Google TPU | Samsung NPU | Our Proposal |
|---|---|---|---|
| Technology | 28 nm | 8 nm | 10 nm |
| Architecture | Weight-stationary Systolic Array | Wide-SIMD with zero-skipping | Output-stationary systolic array |
| Network | InceptionV1 | InceptionV3 | InceptionV3 |
| TOPs/mm² | 0.28 | 1.25 | 3.6 |
| TOPs/W | 2.3 (peak) | 3.4 (avg) | 8.95 (avg) |

# Concluding Remarks

- In this presentation we presented a systolic dataflow based CNN accelerator
- Features like hierarchical vector units in PEs and channel major data layout are effective for improving energy efficiency
- Computation-communication overlap improves performance
- Software controlled scratch-pad memory makes SRAM more area efficient
- Effective clock gating further improves energy efficiency

# Thank You!