# MAERI-FPGA: Enabling HW Design Space Exploration on Real FPGA Hardware Platform

## Tushar Krishna

Associate Professor
School of ECE
Georgia Institute of Technology

Email: tushar@ece.gatech.edu

**ICS 2022
Tutorial**

**June 27, 2022**

# Presenters

**Tushar Krishna**

*Associate Professor*
*Georgia Tech*

**Jianming Tong**

*PhD Student*
*Georgia Tech*

**Other Contributors**
- Yangyu Chen
- Yue Pan
- Abhimanyu Bambhaniya
- Taekyung Heo
- Hyoukjun Kwon

# Schedule (EST)

| Time slot | Topic | |
|-----------|-------|---|
| 14:00 to 14:30 | Introduction to DNN Accelerators | Tushar |
| 14:30 – 14:40 | Break | |
| 14:40: 15:10 | MAERI2.0 Architecture and Tool Flow | Jianming |
| 15:10 to 15:30 | Demo on FPGA | Jianming |

Brief Q/A at the end of each talk.

Please feel free to interrupt and ask questions or use chat

Attention: Tutorial is being recorded!

https://maeri-project.github.io/tutorials/ics-2022

# Deep Learning Applications
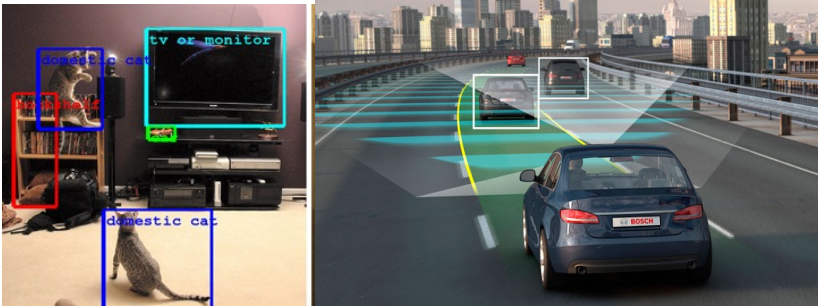
## "AI is the new electricity" – Andrew Ng

**Object Detection**

**Image Segmentation**

**Medical Imaging**

**Speech Recognition**

**Text to Speech**

**Recommendations**

**Games**

# Computation Platforms in Deep Learning

Training → Inference

ARM Trillum

NVDLA

Apple Neural Engine

CambriconX

ShiDianNao

Eyeriss

**Inference Accelerators**

**CPU / GPU / TPU clusters**

# Challenges in Design and Deployment



Training → Inference

DNN Model/Shape

Mapping (Dataflow)

Accelerator Microarchitecture

Energy

Runtime

Focus of this Tutorial

# Outline

- Background on DNNs

- DNN Accelerators

- Dataflow and Mapping

- Flexibility

# Outline

- <span style="color:red">Background on DNNs</span>

- DNN Accelerators

- Dataflow and Mapping

- Flexibility

# What is a Deep Neural Network?



**Neurons**

**Synapses**

$Y_j = \text{activation}\left(\sum_{i=1}^{3} W_{ij} \times X_i\right)$

Weighted Sum

$W_{11}$

$X_1$

$X_2$

$X_3$

$W_{34}$

$Y_1$

$Y_2$

$Y_3$

$Y_4$

input layer

hidden layer
[Image Source: Stanford]

output layer

*Each **synapse** has a **weight** for neuron **activation***

# Modern Deep Learning Landscape

Model Creation ⇄ Training → **Inference**

**Convolutional Layers (Feature Extraction)**

**Summarize features**

Conv. Layer → Conv. Layer → ... → Conv. Layer → Pool. Layer → FC Layer → **"Georgia Institute of Technology, Atlanta"**

**Intermediate features**

# Computations in a DNN → Linear Algebra

$x_1$

$w_1$

$x_2$

$w_2$

$w_k$

$x_k$

O

**Input**　　**Weight**　　**Output**

$n=1$

$k$

$m=1$

$k$

$x$

$n=1$

$= m$

**Neuron => Vector x Vector**

# Computations in a DNN → Linear Algebra



DNN Layer => Vector x Matrix

Data "Reuse"

# Computations in a DNN → Linear Algebra



Batching => Matrix x Matrix

Data "Reuse"

GEMM

# Convolutional Neural Networks

Input Neurons

First Hidden Layer



"Weight Filter" or
"Weight Kernel"

Even more
"Reuse"

"Input Feature Map"

**Shared Weights:**
All neurons use the *same* filter weights

# Convolution in CNN

**Input Image**

**Output Image**

**Filter**



**R**

**S**

$\otimes$

**dot
product**

**Y**

**X**

$\oplus$

**partial sum
accumulation**

**Y'**

**X'**

# Convolution in CNN

**Many Input Channels**

**Input Image**

**Output Image**



R

S

c

Y

X

c

Y'

X'

# Convolution in CNN



**Many Filters**

**Input Image**

**Output Image**

**Many Output Channels**

Tushar Krishna | School of ECE | Georgia Institute of Technology

# Convolution in CNN

**Filters**

**Many Input Images**

**Many Output Images**

**7 loops**

# Loop Nest Representation

7th (outermost) loop used during training

```
for(n=0; n<N; n++) { // Input feature maps (IFMaps)
  for(m=0; m<M; m++) { // Weight Filters
    for(c=0; c<C; c++) { // IFMap/Weight Channels
      for(y=0; y<H; y++) {  // Input feature map row
        for(x=0; x<H; x++) {  // Input feature map column
          for(j=0; j<R; j++) {  // Weight filter row
            for(i=0; i<R; i++) {  // Weight filter column
              O[n][m][x][y] += W[m][c][i][j] * I[n][c][y][x]}}}}}}}
```

# Challenges with DNN Computations

- **Millions of Parameters (i.e., weights)**
  - Billions of computations     ➡ **Need lots of parallel compute**

| DNN Topology | Number of Weights |
|---|---|
| AlexNet (2012) | 3.98M |
| VGGnet-16 (2014) | 28.25M |
| GoogleNet (2015) | 6.77M |
| Resnet-50 (2016) | 23M |
| DLRM (2019) | 540M |
| Megatron (2019) | 8.3B |

This makes CPUs inefficient

  - Heavy data movement     ➡ **Need to reduce energy**

Data Movement Energy Cost

DRAM → ALU    500×
Buffer → ALU    10×
PE → ALU    3×
RF → ALU    1×
ALU → ⊗⊕    1× (Reference)

This makes GPUs inefficient

# Outline

- Background on DNNs

- DNN Accelerators

- Dataflow and Mapping

- Flexibility

# The DL Inference Accelerator Zoo

**AI Chip Landscape**  *V0.5  August, 2019*  *S.T.*

**Tech Giants/System**
Google
Microsoft
facebook
aws
Apple
IBM
Alibaba Group 阿里巴巴集团
HUAWEI
Bai du 百度
Tesla
Hewlett Pack Enterprise
FUJITSU
DELL
Western Digital
NOKIA
LG

**IC Vender/Fabless**
intel
SAMSUNG
NVIDIA
QUALCOMM
AMD
XILINX
MEDIATEK
UNISOC
MARVELL
Rockchip 瑞芯微电子
Ambarella
NationalChip

**Automated Driving**
NXP
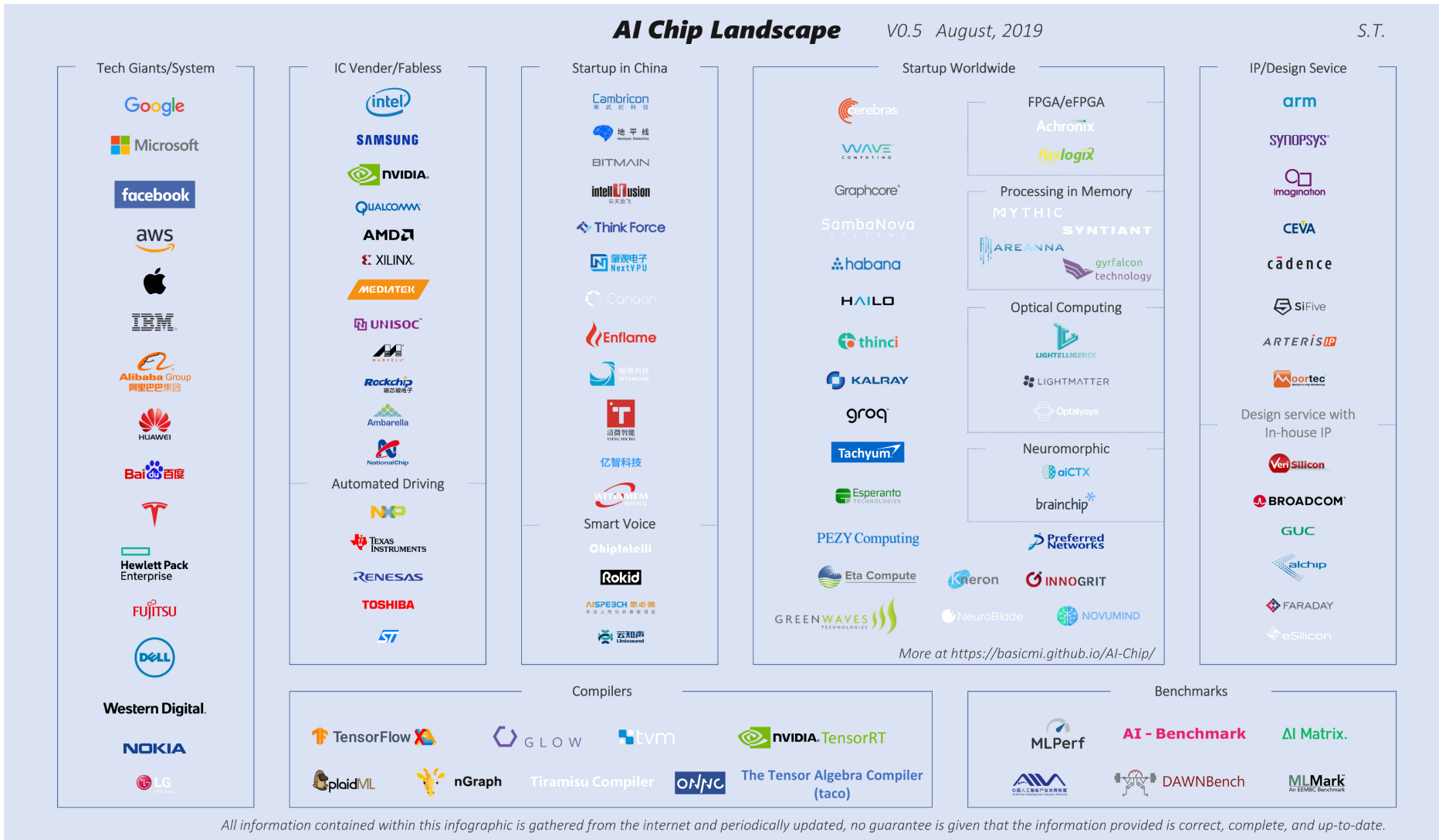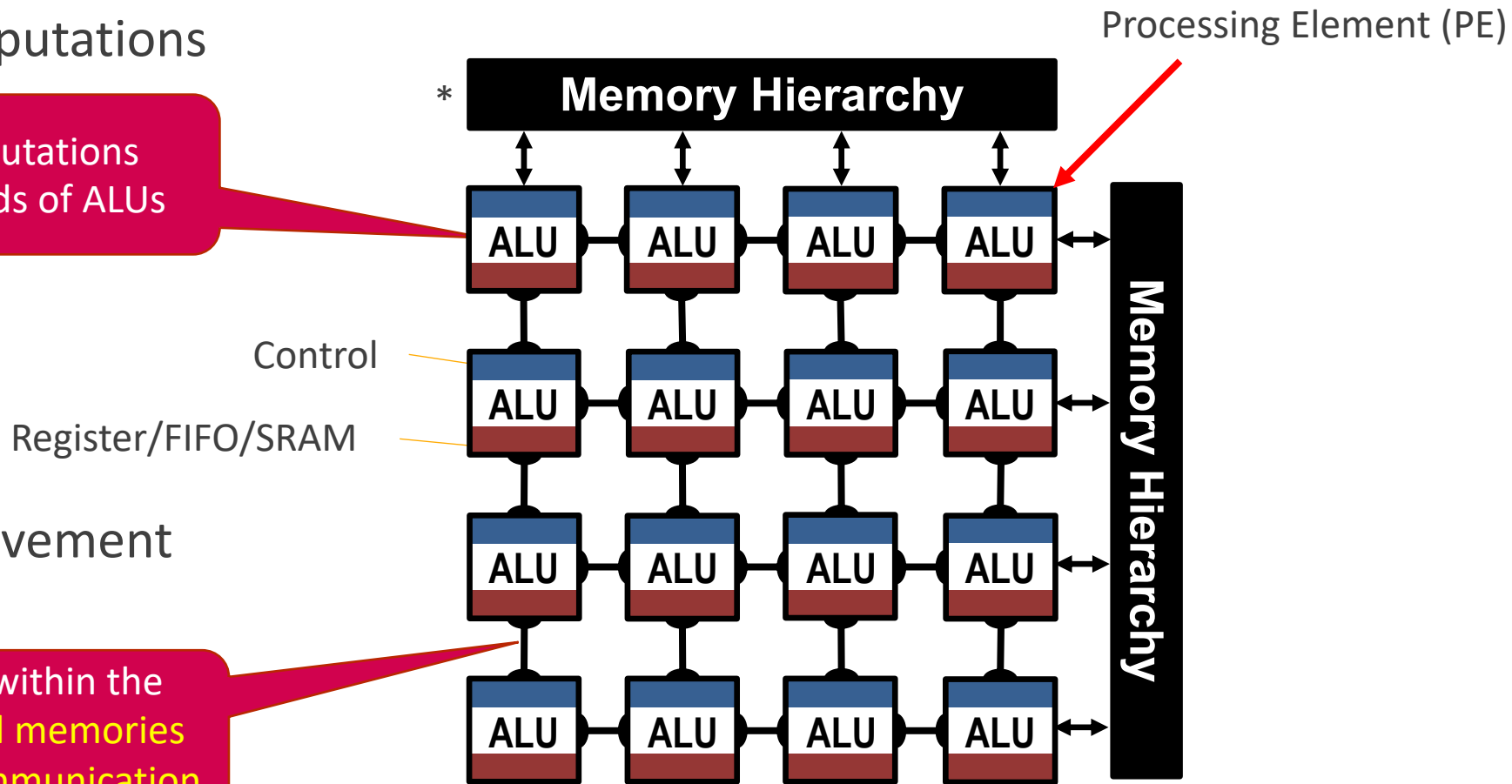TEXAS INSTRUMENTS
RENESAS
TOSHIBA
ST

**Startup in China**
Cambricon 寒武纪科技
地平线 Horizon Robotics
BITMAIN
intellifusion 云天励飞
Think Force
NextVPU 肇观电子
Canaan
Enflame
INTENGINE 深维科技
清微智能 TSING MICRO
亿智科技

WITMEM 知存科技

**Smart Voice**
ChipIntelli
Rokid
AISPEECH 思必驰
云知声 Unisound

**Startup Worldwide**
Cerebras
WAVE COMPUTING
Graphcore
SambaNova
habana
HAILO
thinci
KALRAY
groq
Tachyum
Esperanto TECHNOLOGIES
PEZY Computing
Eta Compute
GREENWAVES TECHNOLOGIES

**FPGA/eFPGA**
Achronix
flexlogix

**Processing in Memory**
MYTHIC
SYNTIANT
AREANNA
gyrfalcon technology

**Optical Computing**
LIGHTELLIGENCE
LIGHTMATTER
Optalysys

**Neuromorphic**
aiCTX
brainchip

Preferred Networks
Kneron
INNOGRIT
NeuroBlade
NOVUMIND

*More at https://basicmi.github.io/AI-Chip/*

**IP/Design Sevice**
arm
SYNOPSYS
Imagination
CEVA
cadence
SiFive
ARTERIS IP
Moortec

Design service with In-house IP
VeriSilicon
BROADCOM
GUC
alchip
FARADAY
eSilicon

**Compilers**
TensorFlow XLA
GLOW
tvm
NVIDIA TensorRT
plaidML
nGraph
Tiramisu Compiler
ONNC
The Tensor Algebra Compiler (taco)

**Benchmarks**
MLPerf
AI - Benchmark
AI Matrix.
中国人工智能产业发展联盟
DAWNBench
MLMark An EEMBC Benchmark

*All information contained within this infographic is gathered from the internet and periodically updated, no guarantee is given that the information provided is correct, complete, and up-to-date.*

# Spatial (or Dataflow) Accelerators

- **Millions of Parameters (i.e., weights)**
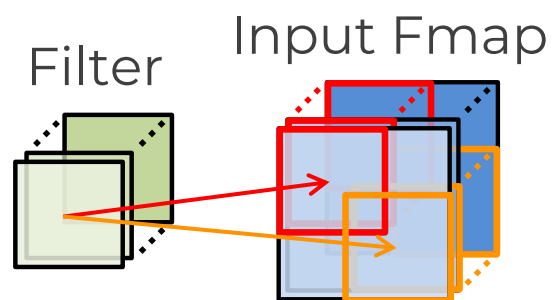  - Billions of computations

Processing Element (PE)

* **Memory Hierarchy**

Spread computations across hundreds of ALUs

Control

Register/FIFO/SRAM

- Heavy data movement

Reuse data within the array via local memories and direct communication

Memory Hierarchy

**\*Y. Chen et. al., "Eyeriss: A Spatial Architecture for Energy-Efficient Dataflow for Convolutional Neural Networks," ISCA, 2016.**

# Types of Algorithmic Data Reuse in DNNs

| **Convolutional Reuse** | **Fmap Reuse** | **Filter Reuse** |
|---|---|---|
| CONV layers only (sliding window) | CONV and FC layers | CONV and FC layers (batch size > 1) |

Filter    Input Fmap

Filters    Input Fmap

Filter    Input Fmaps

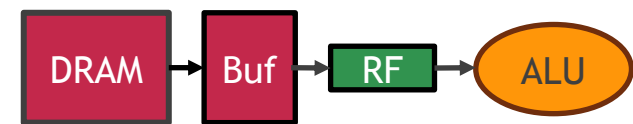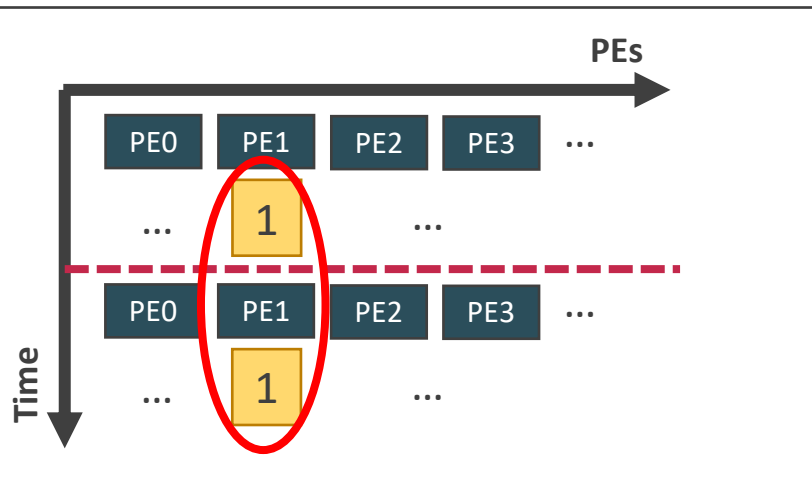Reuse: Activations / Filter weights

Reuse: Activations

Reuse: Filter weights

How to exploit reuse?

*Slide Acknowledgment: Yu-Hsin Chen, Vivenne Sze, Joel Emer (MIT)*
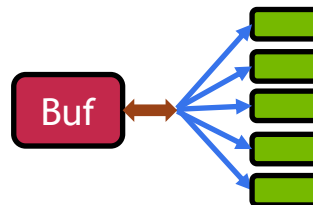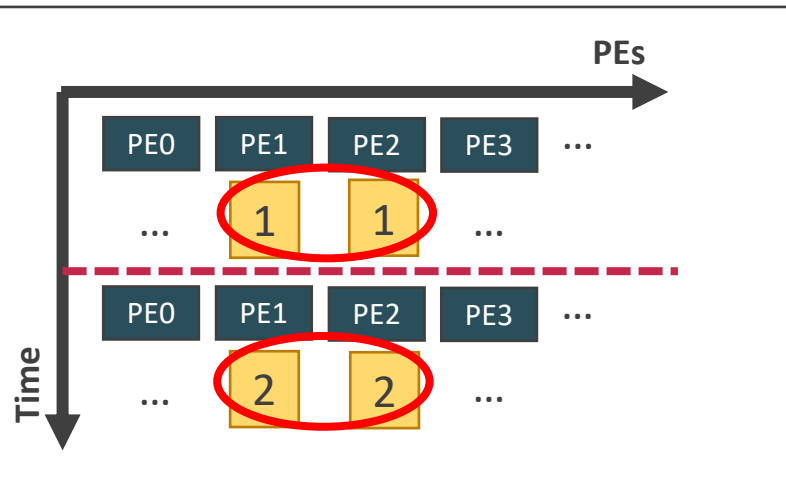
# Hardware structures to exploit reuse

| **Temporal Reuse** | **Spatial Reuse** | **Spatio-Temporal Reuse** |
|---|---|---|
|  |  |  |
|  Memory Hierarchy / Staging Buffers |  Multicasting-support NoCs |  Neighbor-to-Neighbor Connections |

E.g., Custom memory hierarchies in accelerators.

E.g., Hierarchical Bus in Eyeriss (ISCA 2016), Tree in MAERI (ASPLOS 2018)

E.g., TPU (ISCA 2017), local network in Eyeriss (ISCA 2016)

# Mapping and Dataflow

**7-dimensional network layer**



**7D Computation Space:** R * S * X * Y * C * K * N

- **Goal of Mapping**: *translate algorithmic data reuse to HW data reuse*

- **Precise Definition of Mapping:** Fine-grained schedule of computations within DNN accelerators
  - **Computation Order** *(slowest tensor dimension often called "stationary")*
  - **Parallelization Strategy** *(which loops to unroll spatially)*       *Dataflow*
  - **Tiling Strategy** *(number of levels of memory hierarchy)*
  - **Tile Sizes**

# Architectural Components of a DNN Accelerator

**Dataflow**

Accelerator

Mapping

HW Resources

Model

Tiling | Ordering | Parallelism | Shape

Latency
Energy
Power
Area

Num of PEs

Buffer Sizes

NoCs

# Architectural Components of a DNN Accelerator

# Architectural Components of a DNN Accelerator

**Workload (Resnet50)**

CONV

Filter Weight · Input Activation · Output Activation

**Mapping**

**PE**
ALU
L1 Spad (S1)
Controller

**Target Accelerator**

PE — Spad (array of PEs with Spad, connected to Global Spad)

*Number of PEs*

*Buffer sizes (global/ local)*

*NoCs*

**Data / Computation Tile Sizing**

Number: Tile IDs

Filter Tiles · Input Tiles · Output Tiles

**Dataflow**

Mapping on entire accelerator at time = 1

**Tile Scheduling** · **Spatial Partitioning**

Number: Tile IDs

**HW Design-Space**

**Mapping Design-Space** *aka* **Map-Space**

*Tiling*

*Ordering*

*Parallelism Dimension*

*Mapping Shape (Level of Tiling)*

# GEMM vs CONV2D Accelerators

**GEMM Operation**



**CONV2D Operation**



Filter Weight     Input Activation     Output Activation

**3 Loops**

- *Less Opportunities for Reuse*

- *More general: any DNN layer (including convolutions) can be lowered to GEMM (e.g, Im2Col)*

- *E.g., NVIDIA Tensor Core, Google TPU*

**7 Loops**

- *More Opportunities for Reuse*

- *Only applicable for convolution layers*

- *E.g., NVDLA, MAERI (this work)*

# Outline

- Background on DNNs

- DNN Accelerators

- Dataflow and Mapping

- Flexibility

# Dataflow and Mapping

7-dimensional network layer



**7D Computation Space:** R * S * X * Y * C * K * N

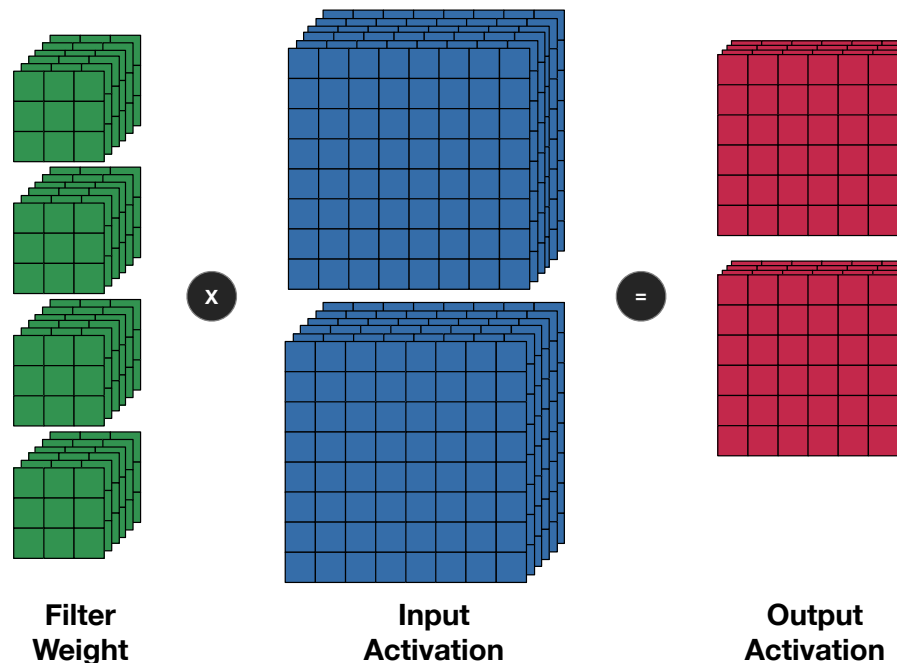- **Number of PEs**
- **Memory Hierarchy**
- **Interconnect Bandwidth**
- **...**

▪ **Goal of Mapping**: *translate algorithmic data reuse to HW data reuse*

▪ **Precise Definition of Mapping:** Fine-grained schedule of computations within DNN accelerators

  - **Computation Order** *(slowest tensor dimension often called "stationary")*
  - **Parallelization Strategy** *(which loops to unroll spatially)*
  - **Tiling Strategy** *(number of levels of memory hierarchy)*
  - **Tile Sizes**

  *Dataflow*

# Impact of Computation Order

*"Output Stationary" Dataflow*

Weights * Inputs = Outputs*  CONV1D

S    X    X' = X-S

**Computation**
```
for(int x = 0; x < X'; x++)
  for(int s = 0; s < S; s++)
Output[x] += Weight[s] * Input[x+s]
```

**Data**
PartialSum[X'][S] needs to access:
- Weight[s]
- Output[x']
- Input[x'+s]

**Time = 0**

*Spatial multicast opportunity for weights*

**Data Space**

*Suppose we map this computation over three PEs*

**PE2**
**PE1**
**PE0**

**Computation Space**

*Each point is a partial sum*

*Output does not change over time => Temporal reuse opportunity*

S
2
1
0
0 1 2 3 4 5  X'

Weight
0 1 2  **S**

*Each point is a data point*

Output
0 1 2 3 4 5  **X'**

Input
0 1 2 3 4 5  **X**

# Impact of Computation Order

**"Weight Stationary" Dataflow**

Weights  *  Inputs  =  Output*  CONV1D

S    X    X' = X-S

**Computation**
```
for(int s = 0; s < S; s++)
    for(int x = 0; x < X'; x++)
Output[x] += Weight[s] * Input[x+s]
```

**Data**
PartialSum[X'][S] needs to access:
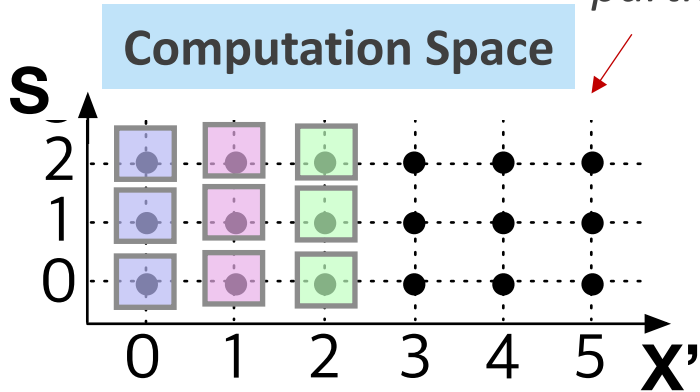- Weight[s]
- Output[x']
- Input[x'+s]

**Time = 0**

*Suppose we **map** this computation over three PEs*

- PE2
- PE1
- PE0

*Need **Spatial reduction** for output*

**Computation Space**

*Each point is a partial sum*

S
2
1
0
  0  1  2  3  4  5  X'

*Weight does not change over time => **Temporal reuse** opportunity*

**Data Space**

*Each point is a data point*

Weight
  0  1  2  S

Output
  0  1  2  3  4  5  X'

Input
  0  1  2  3  4  5  X

# Takeaways: Data Reuse + Hardware Support

- Dataflow exposes data reuse opportunities
- Hardware support is needed to leverage reuse opportunity

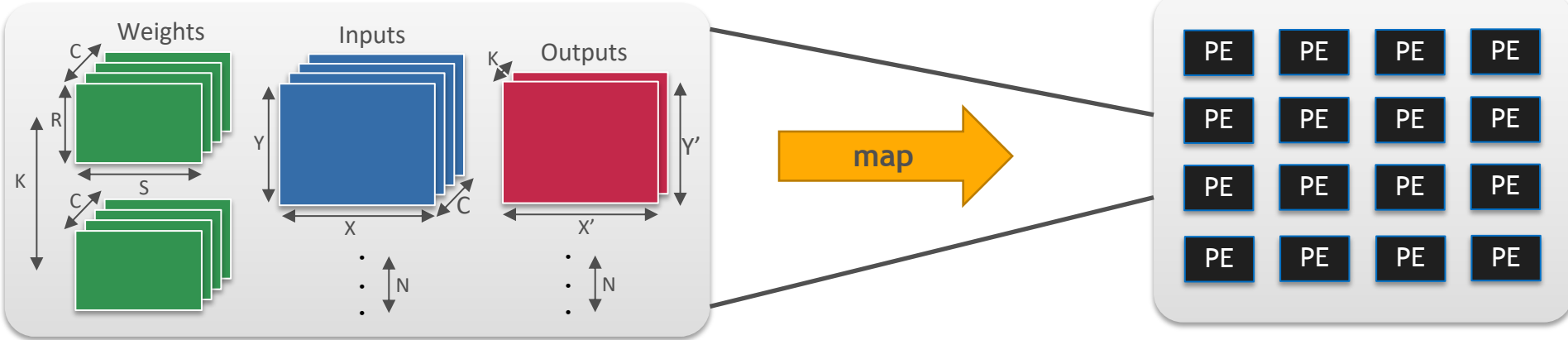| Hardware Structure | Per Data Type | Weight Stationary Dataflow Implication | Output Stationary Dataflow Implication |
|---|---|---|---|
| Bandwidth to MAC | Weight Fetch Rate | Every S Cycles | Every Cycle |
| | Input Fetch Rate | Every Cycle | Every Cycle |
| | Output Fetch Rate | Every Cycle | Every S Cycles |
| Local Buffer Sizes for Temporal Reuse | Weight Buffer Size | 1 | 3 |
| | Input Buffer Size | 3 | 3 |
| | Output Buffer Size | 3 | 1 |
| Network-on-Chip for Spatial Reuse | Weight Distribution | Unicast | Spatial Multicast |
| | Input Distribution | Spatial Multicast | Unicast |
| | Output Collection | Spatial Reduction | Temporal Reduction |

*Note: for full 6D conv, trillions of valid dataflow choices → Huge Design Space*

# Dataflow and Mapping

7-dimensional network layer



**7D Computation Space:** R * S * X * Y * C * K * N

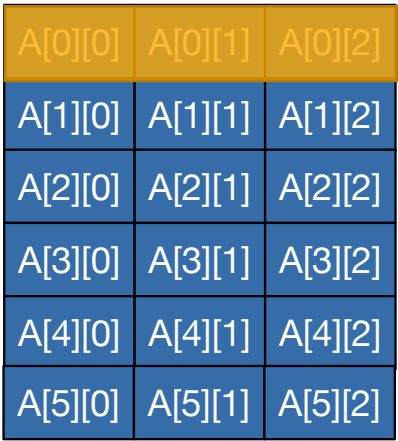- **Goal of Mapping**: *translate algorithmic data reuse to HW data reuse*

- **Precise Definition of Mapping:** Fine-grained schedule of computations within DNN accelerators

  - **Computation Order** *(slowest tensor dimension often called "stationary")*
  - **Parallelization Strategy** *(which loops to unroll spatially)*
  - **Tiling Strategy** *(number of levels of memory hierarchy)*
  - **Tile Sizes**

*Dataflow*

# Impact of Parallelization

**Example Model A: Matrix-Vector Multiplication (i.e., Simplified Fully-connected layer)**

**Parallelize**

| A[0][0] | A[0][1] | A[0][2] |
|---------|---------|---------|
| A[1][0] | A[1][1] | A[1][2] |
| A[2][0] | A[2][1] | A[2][2] |
| A[3][0] | A[3][1] | A[3][2] |
| A[4][0] | A[4][1] | A[4][2] |
| A[5][0] | A[5][1] | A[5][2] |

**Matrix A**

**X**

**Parallelize**

| B[0] | B[1] | B[2] |
|------|------|------|

**Vector B**

**=**

| C[0] |
|------|
| C[1] |
| C[2] |
| C[3] |
| C[4] |
| C[5] |

**Vector C**

C[0] = A[0][0] * B[0]
      + A[0][1] * B[1]
      + A[0][2] * B[2]

## Example Accelerator

Global Scratchpad

PE0 — PE1 — PE2

**Map**

Matrix A:  A[0][0]   A[0][1]   A[0][2]

Vector B:  B[0]   B[1]   B[2]

**Avg. Utilization: 100%**

# Impact of Parallelization

**Example Model B: Matrix-Vector Multiplication (i.e., Simplified Fully-connected layer)**



**Parallelized**

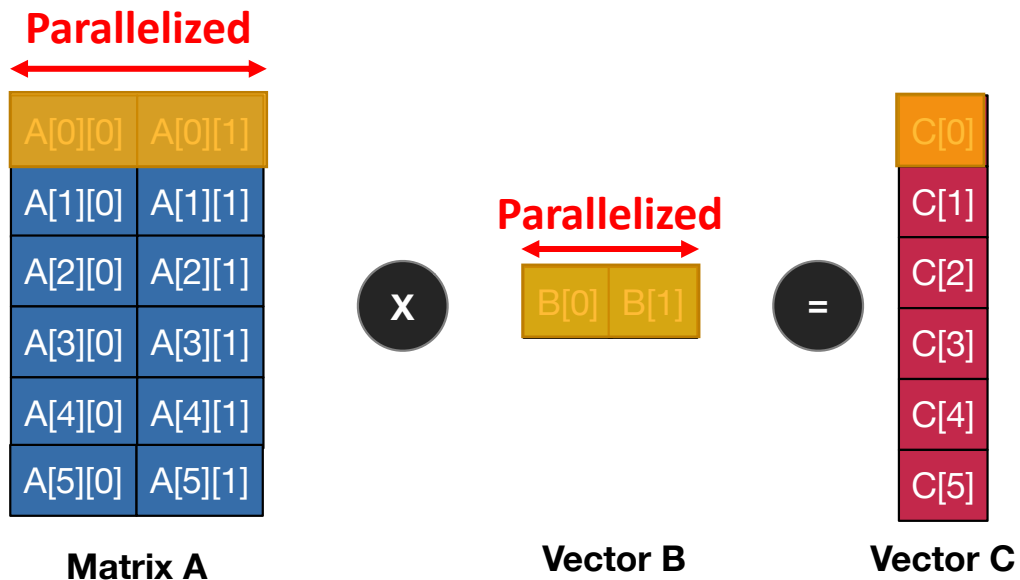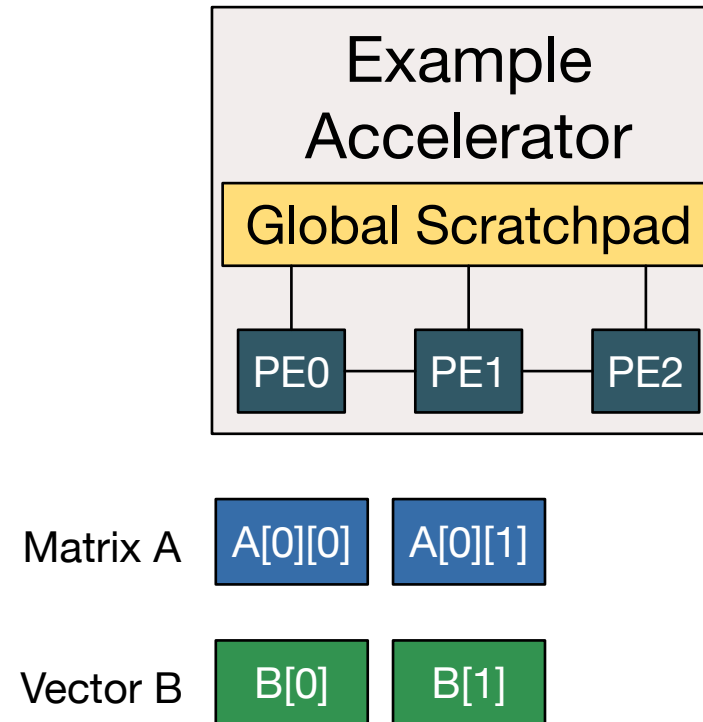| A[0][0] | A[0][1] |
| A[1][0] | A[1][1] |
| A[2][0] | A[2][1] |
| A[3][0] | A[3][1] |
| A[4][0] | A[4][1] |
| A[5][0] | A[5][1] |

**Matrix A**

X

**Parallelized**

| B[0] | B[1] |

**Vector B**

=

| C[0] |
| C[1] |
| C[2] |
| C[3] |
| C[4] |
| C[5] |

**Vector C**

C[0] = A[0][0] * B[0]
       + A[0][1] * B[1]

**Example Accelerator**

Global Scratchpad

PE0   PE1   PE2

**Map**

Matrix A   A[0][0]   A[0][1]
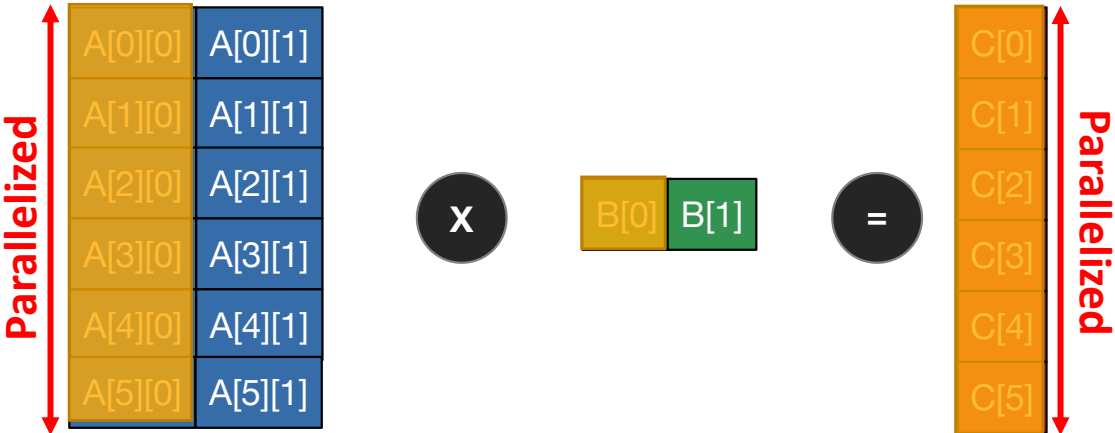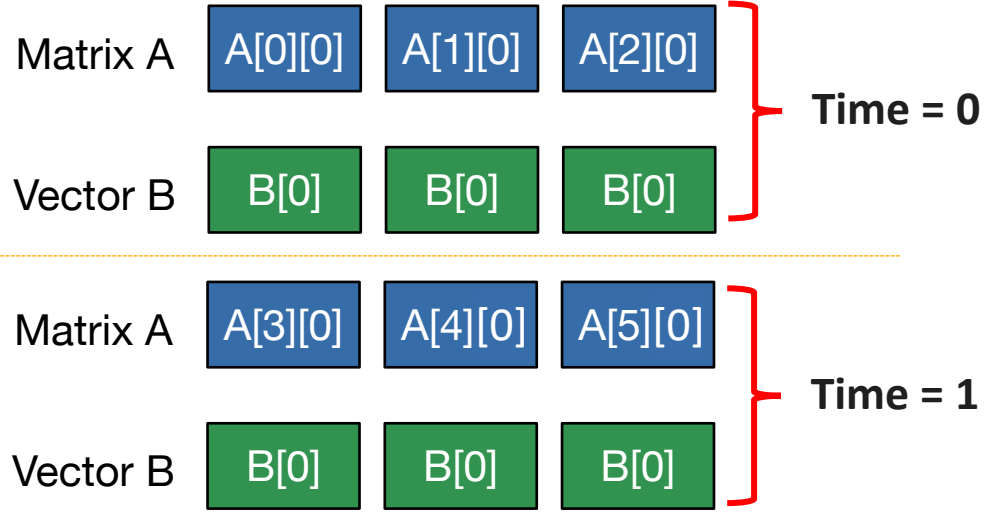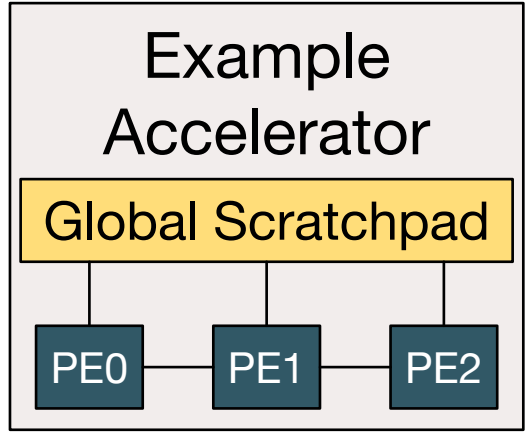
Vector B   B[0]   B[1]

**Avg. Utilization: 66%**

**Can we map it in a better way?**

# Impact of Parallelization

**Example Model B: Matrix-Vector Multiplication (i.e., Simplified Fully-connected layer)**



Matrix A

C[0] = A[0][0] * B[0]
    + A[0][1] * B[1]
    + A[0][2] * B[2]

**Example Accelerator**

Global Scratchpad

PE0    PE1    PE2

**Map**

Matrix A    A[0][0]    A[1][0]    A[2][0]

Vector B    B[0]    B[0]    B[0]

**Time = 0**

Matrix A    A[3][0]    A[4][0]    A[5][0]

Vector B    B[0]    B[0]    B[0]

**Time = 1**

**Avg. Utilization: 100%**

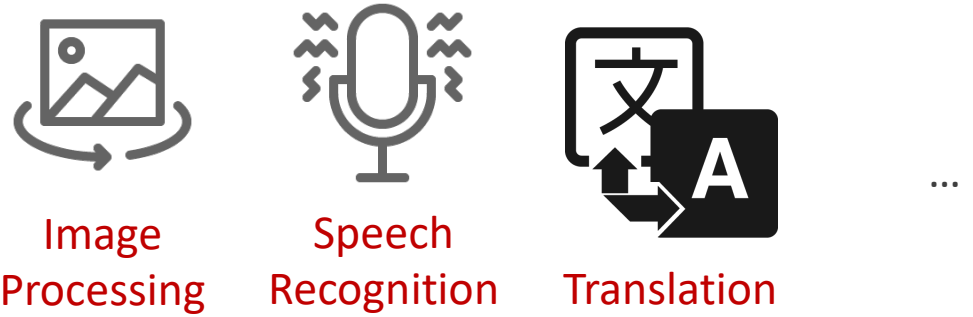**The more dimensions, the more optimization opportunities**

# Outline

- Background on DNNs
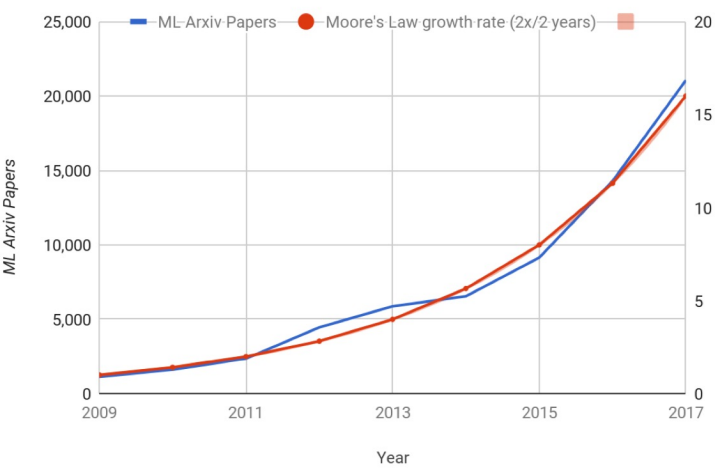
- DNN Accelerators

- Dataflow and Mapping

- Flexibility

# Why do we need *flexible* DNN accelerators?

- **Trend 1: Diversity in DNN Models**
  - Layer Sizes
  - Layer Shapes
  - Layer Types

**<Number of new ML papers in Arxiv>**



Image Processing     Speech Recognition     Translation     ...

**Evolution of DNN Applications**

AlexNet     VGGNet     ResNet     Transformer     MobileNetv2     EfficientNet

2012     2014     2015     2016     2017     2018     2019

**Evolution of DNN Models**

# Why do we need *flexible* DNN accelerators?

- **Trend 1: Diversity in DNN Models**
  - Layer Sizes
  - Layer Shapes
  - Layer Types

- **Trend 2: Diversity in Implementations**
  - Depth-wise/Point-wise Convolutions
  - Pruning → Sparsity



**e.g. of Depth-wise Separable CONV**



**e.g. of Pruning**

# Why do we need *flexible* DNN accelerators?

- **Trend 1: Diversity in DNN Models**
  - Layer Sizes
  - Layer Shapes
  - Layer Types

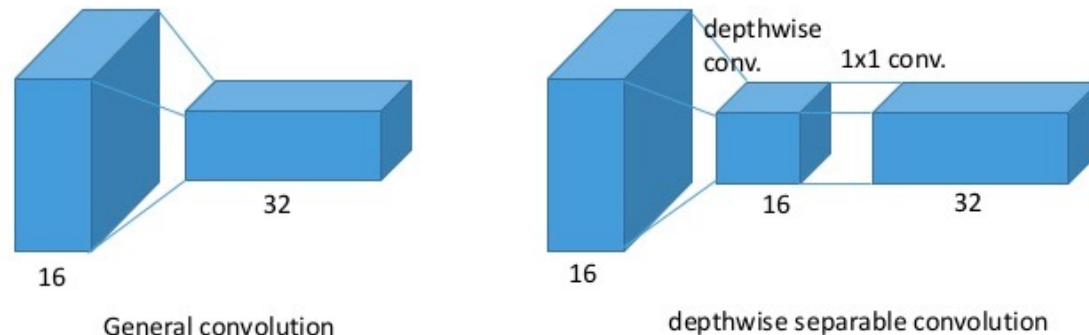- **Trend 2: Diversity in Implementations**
  - Depth-wise/Point-wise Convolutions
  - Pruning → Sparsity

- **Trend 3: Diversity in Mapping/Dataflow**
  - Loop Transformations ("Dataflow")
    - Order, Parallelization, Tiling
    - "Weight Stationary", "Row Stationary"
  - Partitioning Strategies – Per Layer, Cross Layer, ..



Computation Order
Parallelization
Tiling

*Dataflow*

*Data Reuse*

*Data Movement*

# Why do we need *flexible* DNN accelerators?
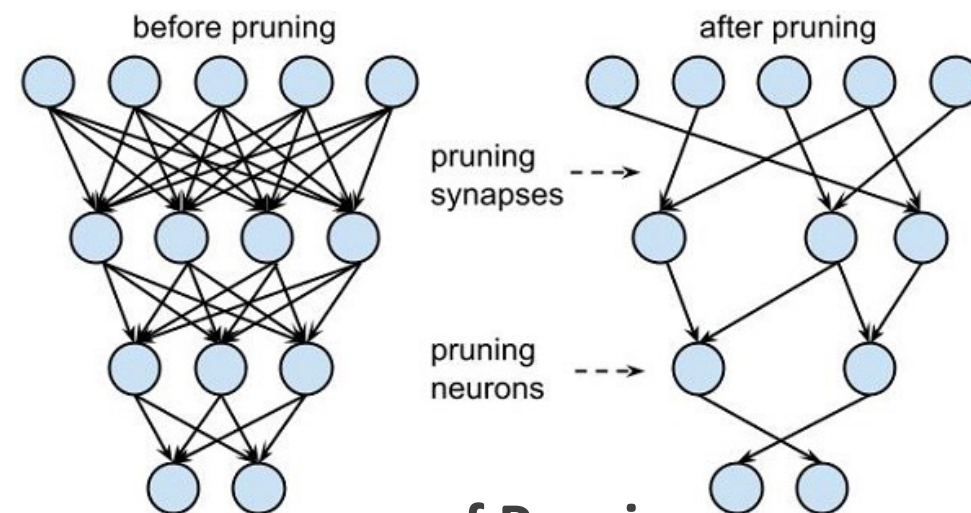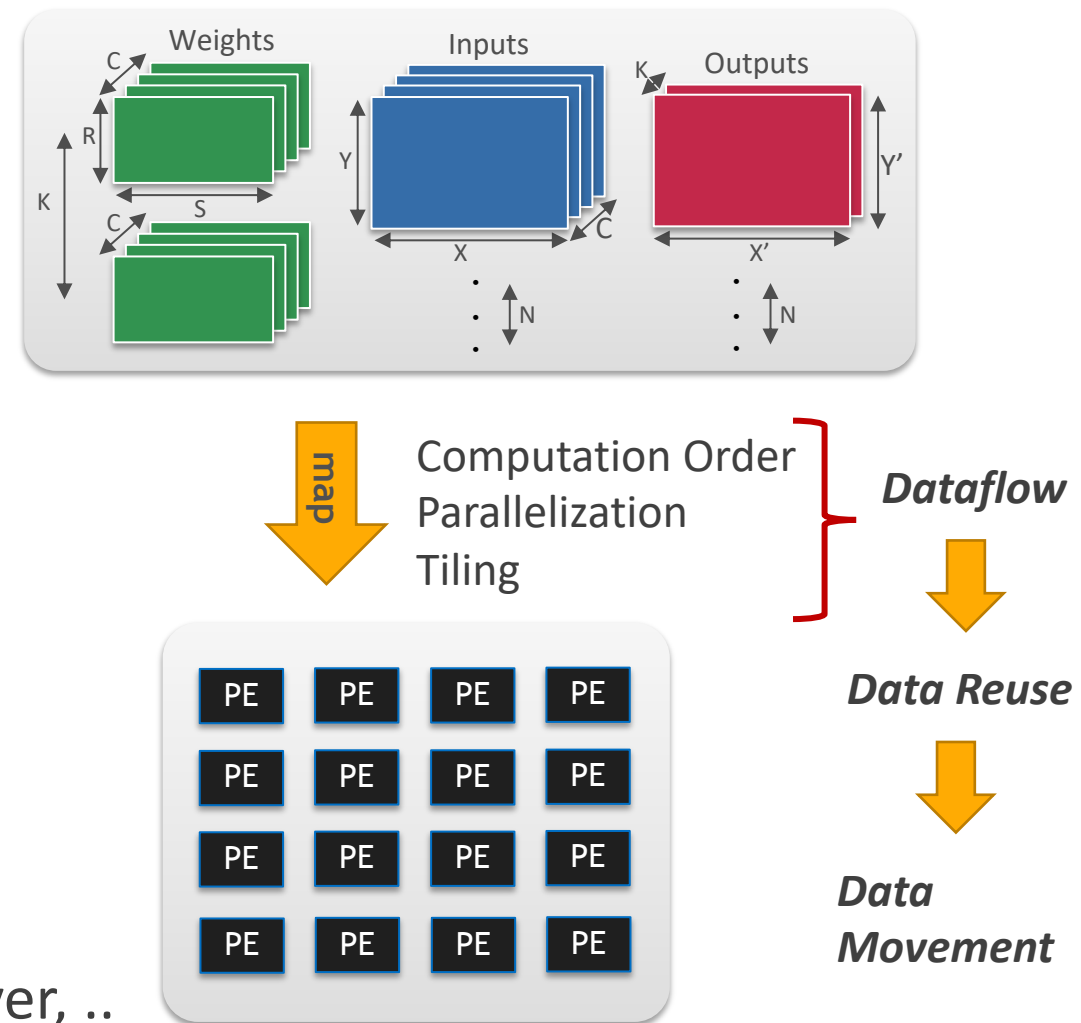
- **Trend 1: Diversity in DNN Models**
  - Layer Sizes
  - Layer Shapes
  - Layer Types

- **Trend 2: Diversity in Implementations**
  - Depth-wise/Point-wise Convolutions
  - Pruning → Sparsity

- **Trend 3: Diversity in Mapping/Dataflow**
  - Loop Transformations ("Dataflow")
    - Order, Parallelization, Tiling
    - "Weight Stationary", "Row Stationary"
  - Partitioning Strategies – Per Layer, Cross Layer, ..

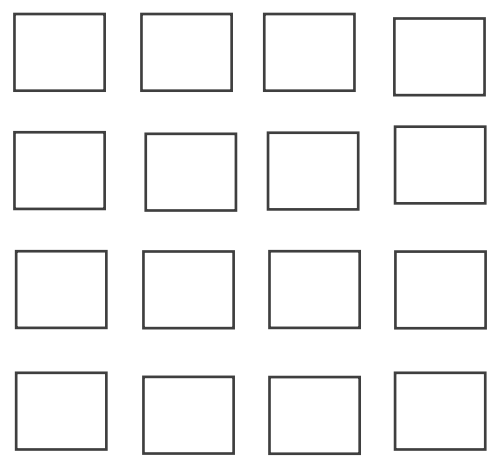→ **Myriad "irregular" shapes, sizes, accesses**

**Challenge:**
Getting high-utilization from accelerator for all cases.

*Why?*
*Aren't DNNs essentially Matrix-Matrix multiplications?*

# Example of GEMM Operation

K

N

N

M  ✖  K  ⬛  =  ⬛  M

⬇

**Metric 1: Mapping Efficiency**
*Percentage of PEs with useful computations mapped over them*

**Metric 2: Utilization**

*Mapping Efficiency x Activity*

e.g., TPU

# Mapping Examples

**Communication**

→ Distribute

→ Collect

**Regular**

4

8

4

4

*Physical Array: 4x4*

**Map Effic. = 100%**

**Distribute**  Row multicast

**Collect**  Column Reduce

# Mapping Examples

**Communication**

→ Distribute

→ Collect



**Regular**

4

4

X

8

4

=

*Physical Array: 4x4*

Map Effic. = 100%

**Distribute**  Row multicast

**Collect**  Column Reduce

**Irregular**

2

8

X

8

2

=

Map Effic. = 50%

Map Effic. = 50%

**Irregular**

5

3

X

8

5

=

Map Effic. = 75%

Map Effic. = 19%

**Sparse**

4

6

X

8

4

=

Map Efffic. = 62%

Map Effic. = 38%

# Mapping Efficiency needs Mapping Flexibility



**Regular**

4, 8, 4, 4

*Physical Array: 4x4*

Map Effic. = 100%

**Irregular**

2, 8, 2, 8

*Logical: 2x8*

Map Effic. = 100%

**Irregular**

5, 8, 3, 5

*Logical: 5x3*

Map Effic. = 94%

**Sparse**

4, 8, 6, 4

*Logical: {3x1, 2x1, 4x1, 1x1, 4x1, 2x1}*

Map Effic. = 100%

**How to support Mapping Flexibility?**

| | Regular | Irregular | Irregular | Sparse |
|---|---|---|---|---|
| **Distribute** | Row multicast | Spatial Multicast | Multicast to non-neighbors | Only send non-zeros |
| **Collect** | Column Reduce | Multiple Parallel | Variable Length | Variable Non-Uniform Length |

**Flexible data distribution and reduction**

# Levels of Flexibility



**Regular**     **Irregular**     **Irregular**     **Sparse**

*Physical Array: 4x4*

*Logical: 2x8*

*Logical: 5x3*

*Logical: {3x1, 2x1, 4x1, 1x1, 4x1, 2x1}*

Map Effic. = 100%

Map Effic. = 100%

Map Effic. = 94%

Map Effic. = 100%

*"Cluster"*

**Fixed Homogeneous Clusters**
*(i.e., fixed cluster size => fixed aspect ratio)*

**Partially-Flexible Homogeneous Clusters**
*(configurable (limited choices) number of PEs per cluster)*

**Fully-Flexible Homogeneous Clusters**
*(configurable (any choice) number of PEs per cluster)*

**Fully-Flexible Heterogeneous Clusters**
*(configurable (any choice) unequal sized clusters)*

# Introducing MAERI2.0 – A Flexible DNN Accelerator



**Spatio-Temporal Reuse** via forwarding of inputs

**Temporal Reuse** via memory hierarchy

**Spatial Reuse** via Broadcasts

**Temporal Reuse** i.e. "stationary" via local buffers

**Dot Product Engine (DPE)**

*MAERO 2.0 builds upon:*
*MAERI: Enabling Flexible Dataflow Mapping over DNN Accelerators via Reconfigurable Interconnects:*
*Hyoukjun Kwon, Ananda Samajdar, and Tushar Krishna*
*ASPLOS 2018, IEEE Micro Top Picks 2019 Honorable Mention*

# Focus of Today's Tutorial

- Supported Neural Network Model

- Quantization Flow

- Memory Layout

- Heterogeneous Scheduling

- MAERI 2.0 Microarchitecture

- FPGA DEMO

**Future Work:**
- Support for Sparsity
- Support for Multi-layer Mapping
- Compiler support

# Schedule (EST)

| Time slot | Topic | |
|-----------|-------|---|
| 14:00 to 14:30 | Introduction to DNN Accelerators | Tushar |
| 14:30 – 14:40 | Break | |
| 14:40: 15:10 | MAERI2.0 Architecture and Tool Flow | Jianming |
| 15:10 to 15:30 | Demo on FPGA | Jianming |

Brief Q/A at the end of each talk.

Please feel free to interrupt and ask questions or use chat

Attention: Tutorial is being recorded!

https://maeri-project.github.io/tutorials/ics-2022