# DnnWeaver: From High-Level Deep Network Models to FPGA Acceleration

**Hardik Sharma**   Jongse Park   Divya Mahajan   Emmanuel Amaro

Joon Kyung Kim   Chenkai Shao   Asit Mishra†   Hadi Esmaeilzadeh
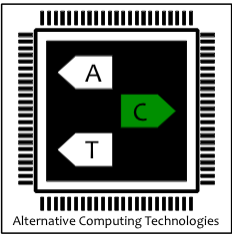
**Alternative Computing Technologies (ACT) Lab**
**Georgia Institute of Technology**

**†Intel Corporation**
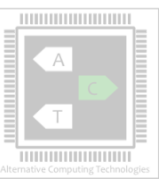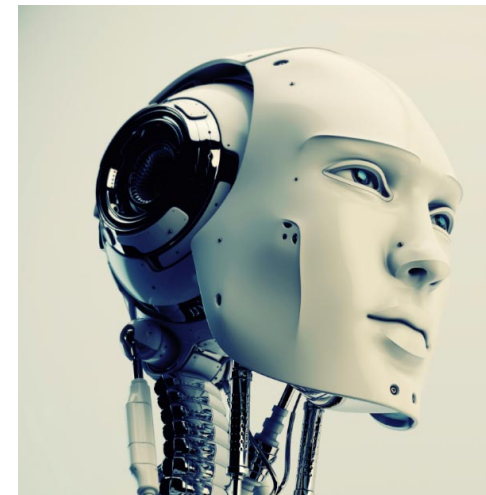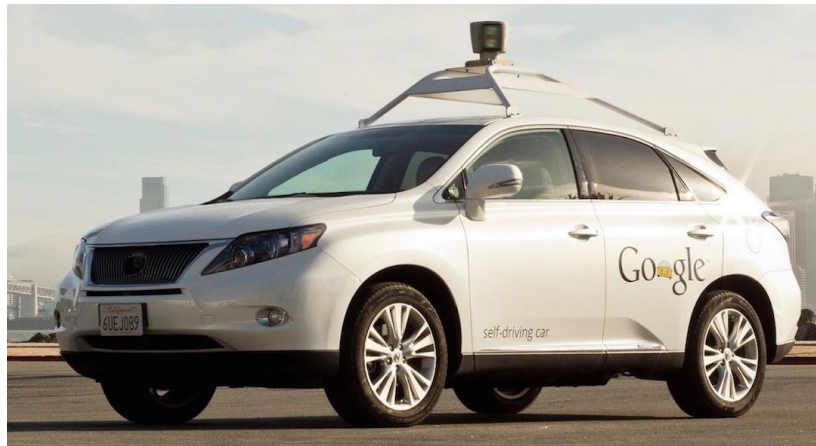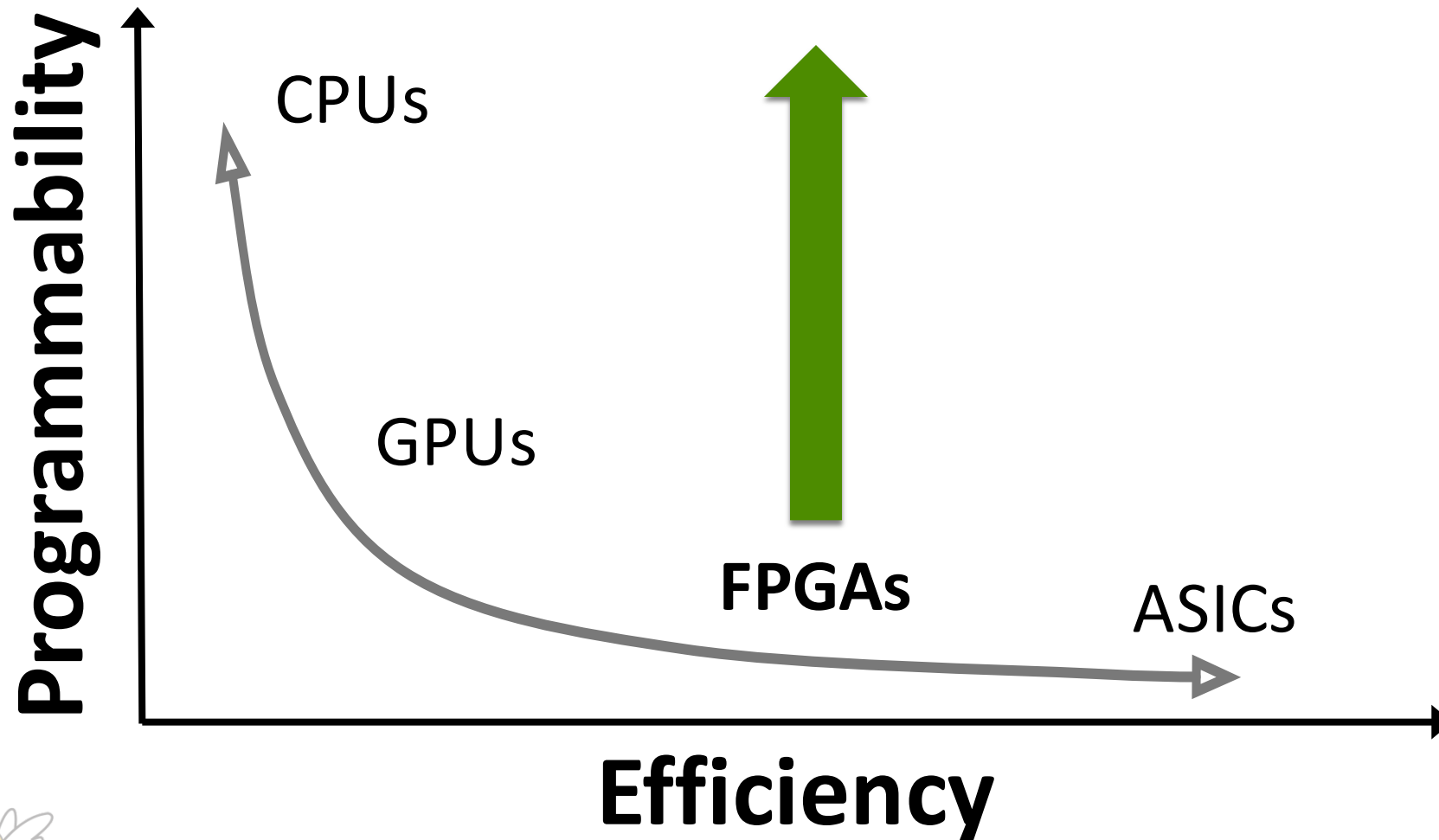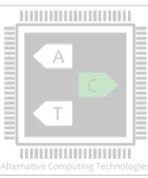
**MICRO'16**

# Deep Neural Networks: A Move Towards Artificial Intelligence

# Programmability is a First-Order Concern



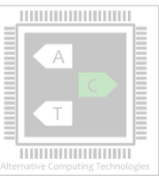*Different applications require different DNNs!*
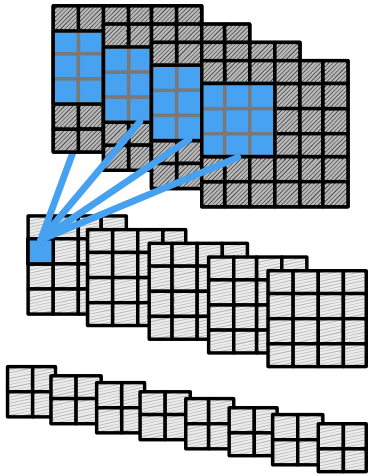
# FPGAs are Hard to Program!

Over 10,000 lines of code for DNN hardware templates in 14 months!

# DNNs are Big, on-chip FPGA storage is small!

AlexNet is 119 MBytes whereas Arria 10 has 5 MB!

# Bridging the Semantic Gap

**FPGA**

**Synthesizable Accelerator**

**Translator**

**Design Planner**

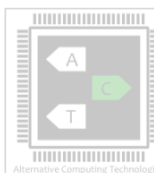**Design Weaver**

**Integrator**

High-level Model Specification

Macro Dataflow Graph

Execution Schedule

Accelerator Configuration

Accelerator Core

# 1. Translator

## High-Level DNN model – Caffe*
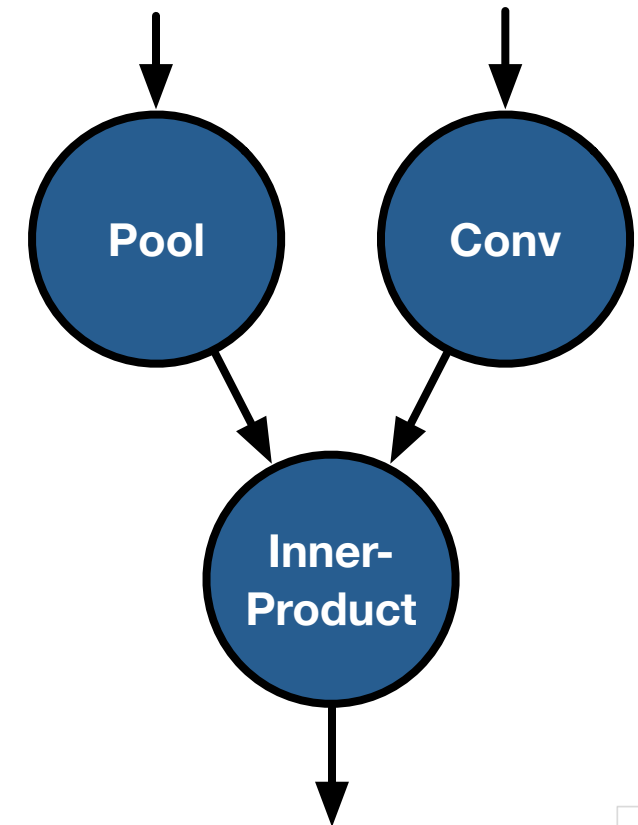
```
layer{
name: Pool,
type: POOLING,
params{...}}

    layer{
    name: Conv,
    type: CONVOLUTION,
    params{...}}

        layer{
        name: Inner-Product,
        type: INNER_PRODUCT,
        params{...}}
```
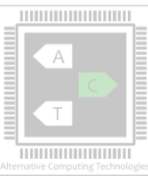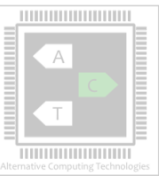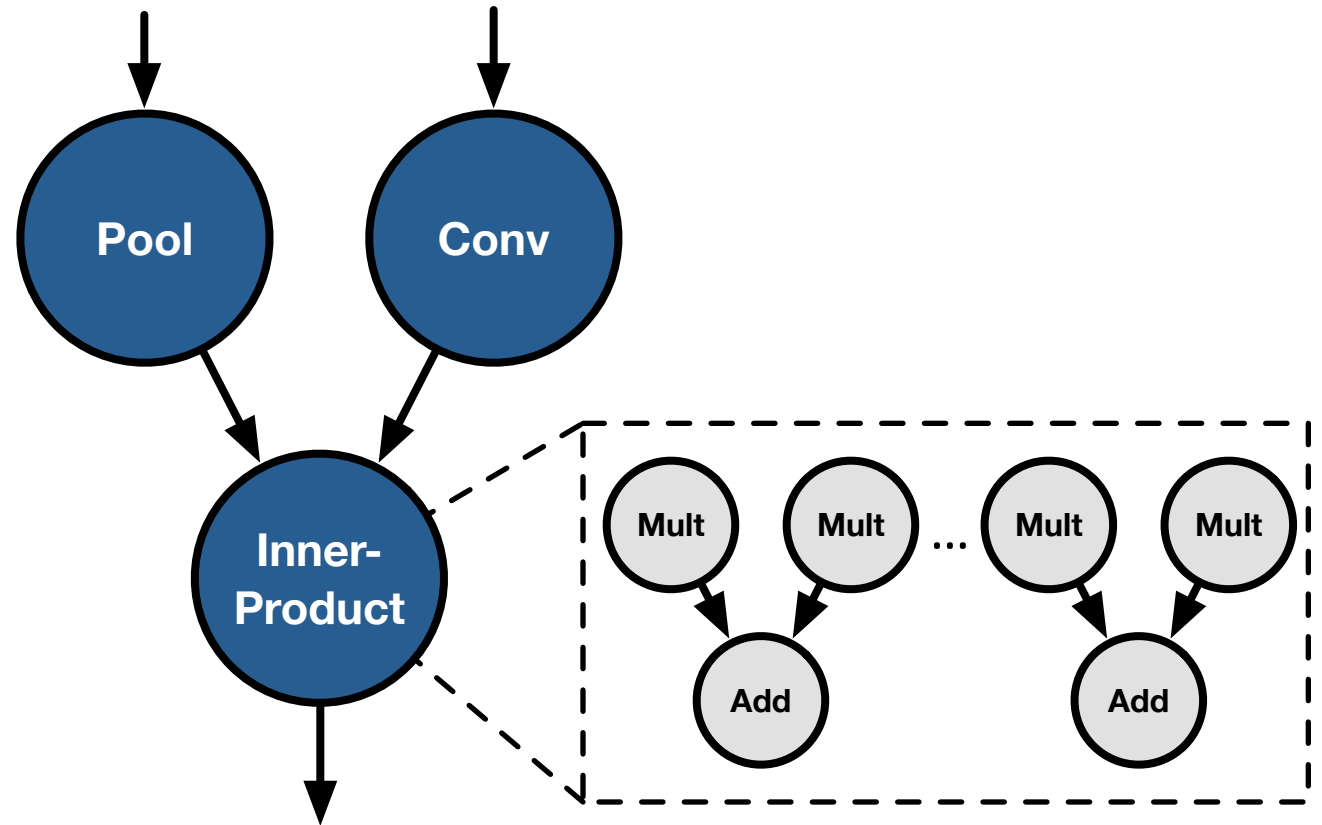
## DNNWEAVER ISA



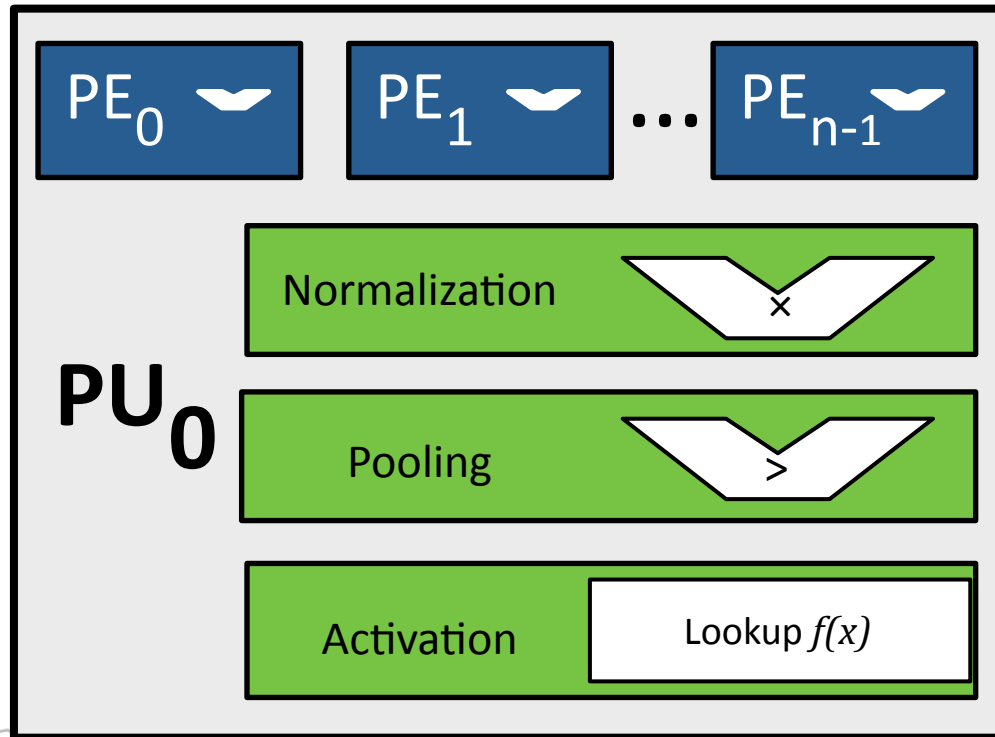**\*Caffe: Convolutional Architecture for Fast Feature Embedding. arXiv:1408.5093**

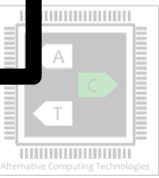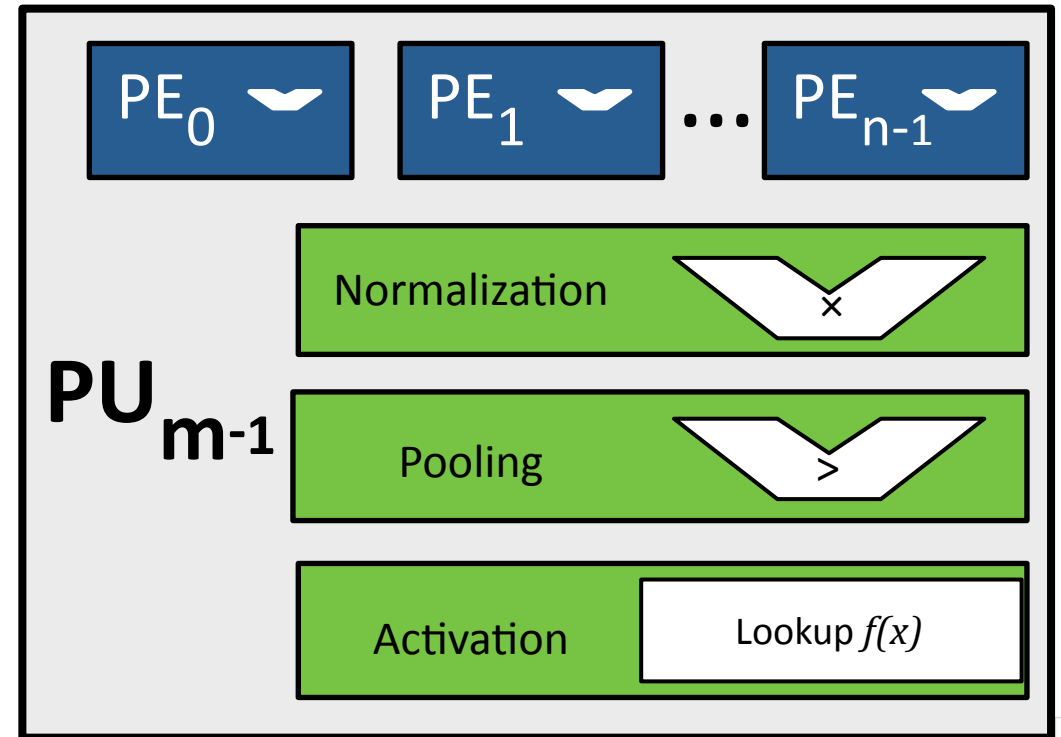# ISA: Abstracting DNNs

Abstract DNNs as a **macro-dataflow graph**

# Template Architecture

## Accelerator Core

PU$_0$

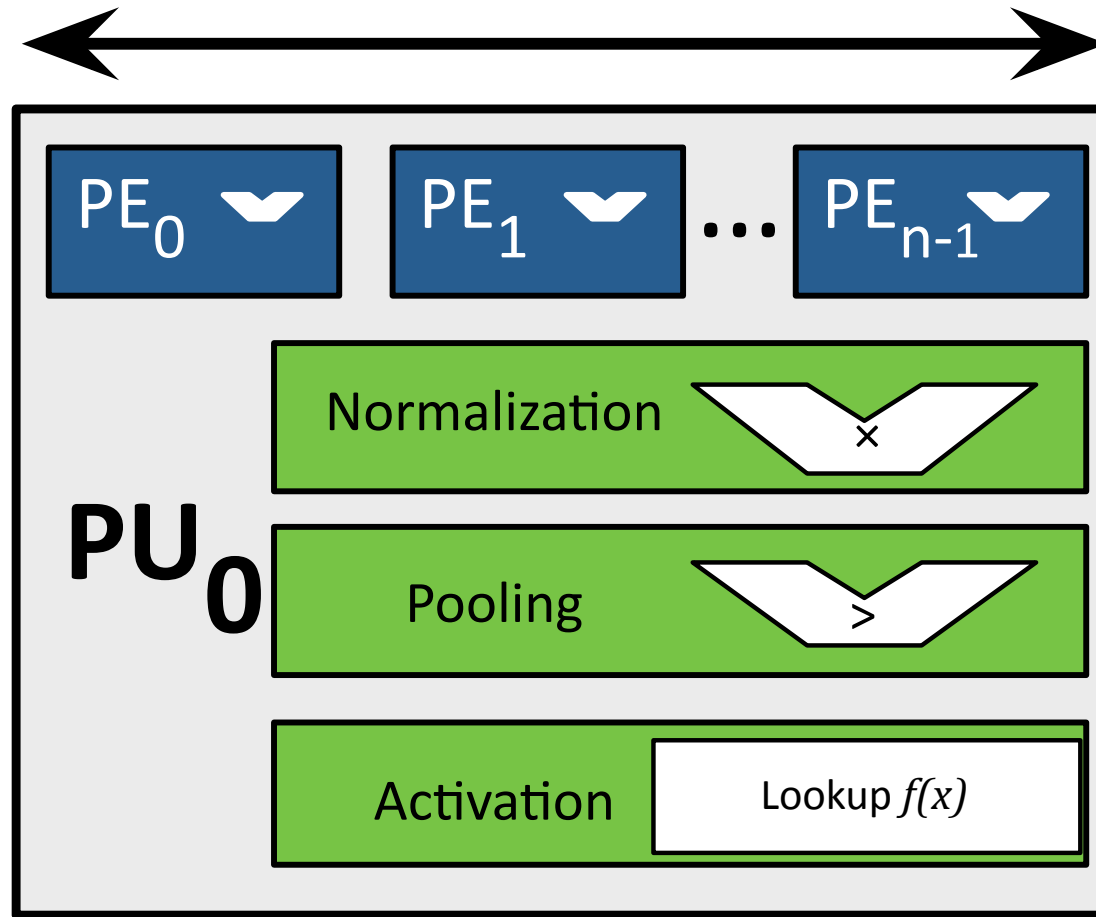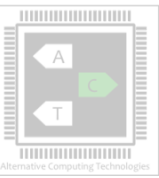PE$_0$  PE$_1$  ...  PE$_{n-1}$

Normalization

Pooling

Activation    Lookup $f(x)$

...  PU$_{m-1}$

PE$_0$  PE$_1$  ...  PE$_{n-1}$

Normalization

Pooling

Activation    Lookup $f(x)$

# Processing Unit
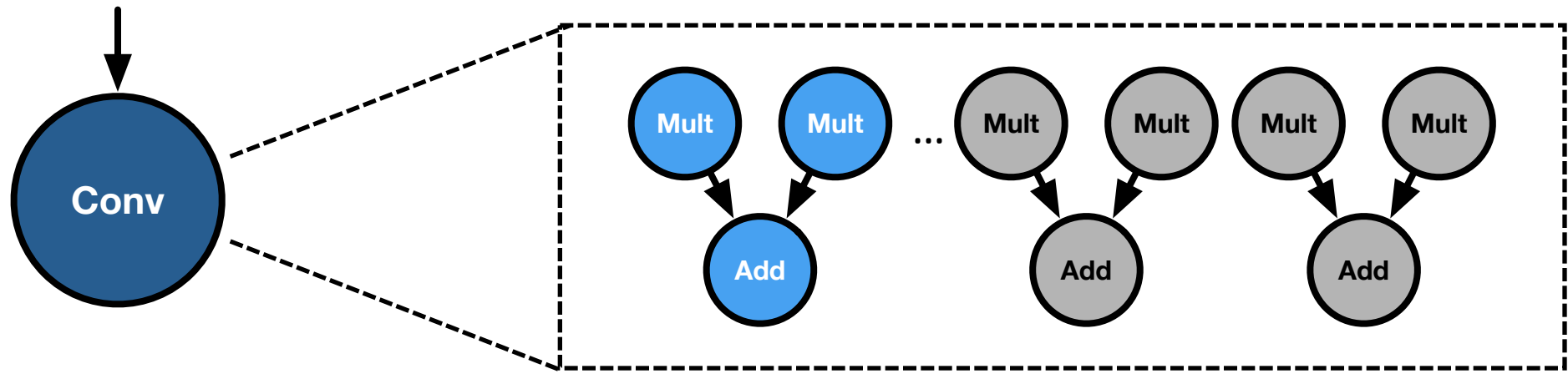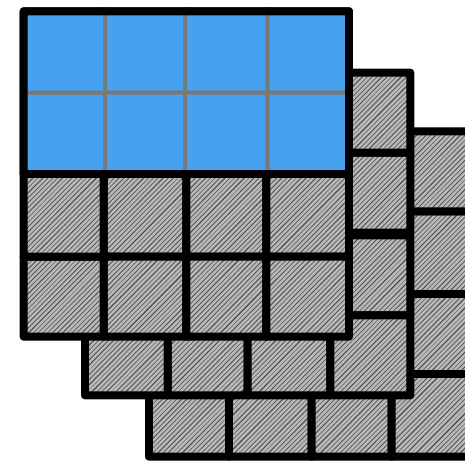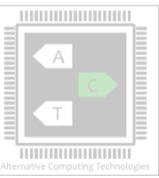
# Executing the Macro Dataflow Graph

# Slicing the Dataflow Graph



**Convolution Node Output**

**Sliced Output**

# 2. Design Planner

**Macro-Dataflow ISA**

**FPGA Specification**
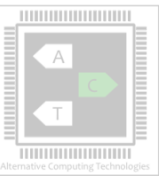
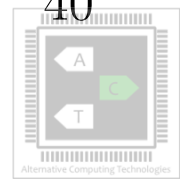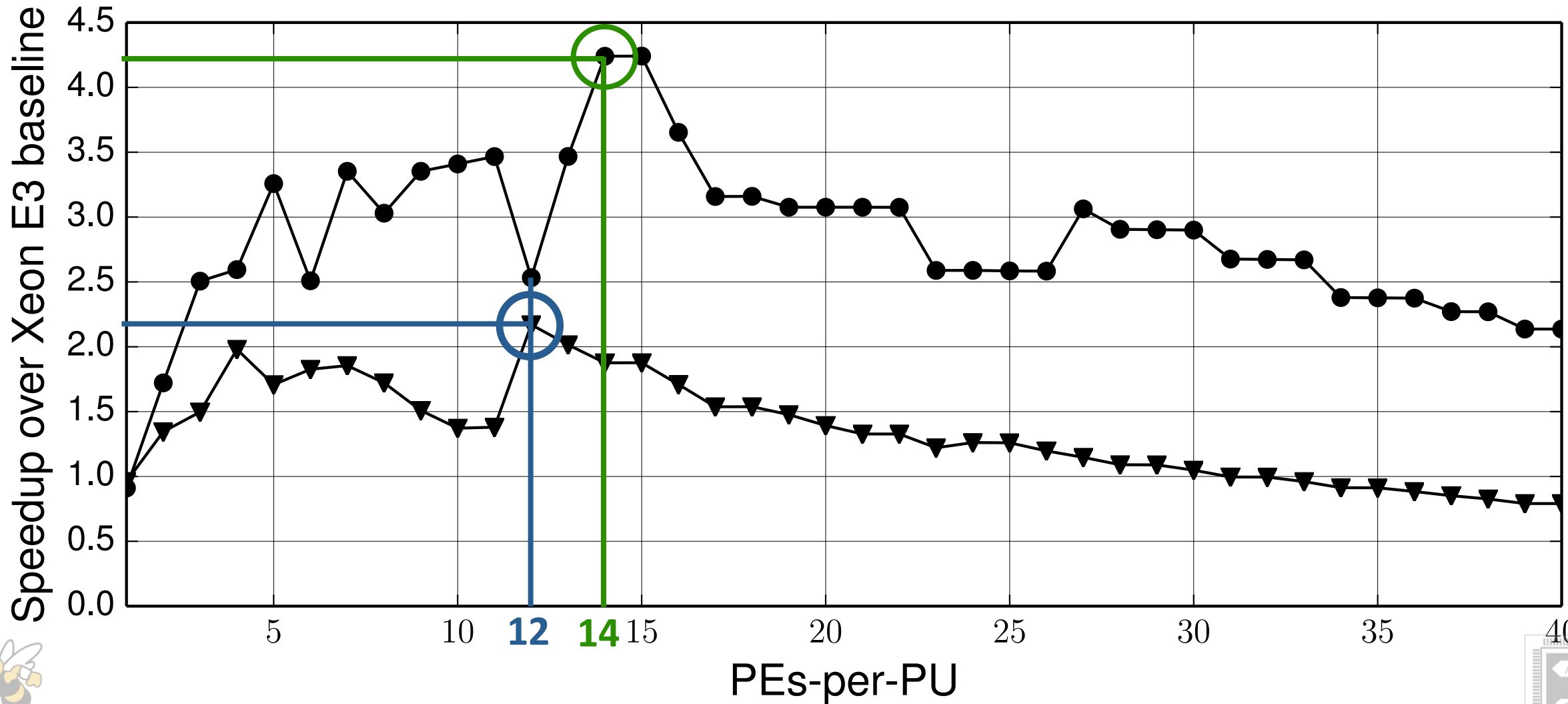**Co-optimize Hardware and Execution Schedule**

# Using FPGA's Programmability

Significant variations between different DNN models: # layers, Model Size, Operations, etc.

**Specialize for each DNN model!**

# Co-optimization
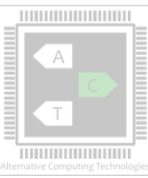
# 3. Design Weaver

**Design Planner**

**Accelerator Configuration** — *Verilog* →

**Execution Schedule** — *Decode* →

**Accelerator Core**

**PU$_0$**

PE$_0$   PE$_1$   ...   PE$_{n-1}$

Normalization    $\times$

Pooling    $>$

Activation    Lookup $f(x)$

...

**PU$_{m-1}$**

PE$_0$   PE$_1$   ...   PE$_{n-1}$

Normalization    $\times$

Pooling    $>$

Activation    Lookup $f(x)$

# 4. Integrator

# Benchmark DNNs

| Name | # Layers | Model Size(MB) | # Operations (Mops) | Lines of Code |
|---|---|---|---|---|
| VGG-16 | 36 | 324.0 MB | 16362 MOps | 347 |
| OverFeat | 16 | 278.0 MB | 2798 MOps | 196 |
| VGG-CNN-S | 19 | 196.0 MB | 2666 MOps | 200 |
| AlexNet | 20 | 119.0 MB | 1147 MOps | 278 |
| Djinn | 13 | 48.4 MB | 25 MOps | 105 |
| NiN | 28 | 14.5 MB | 1106 MOps | 516 |
| LeNet | 7 | 0.8 MB | 2 MOps | 128 |
| CIFAR-10Full | 12 | 0.2 MB | 12 MOps | 156 |

# Platforms Tested

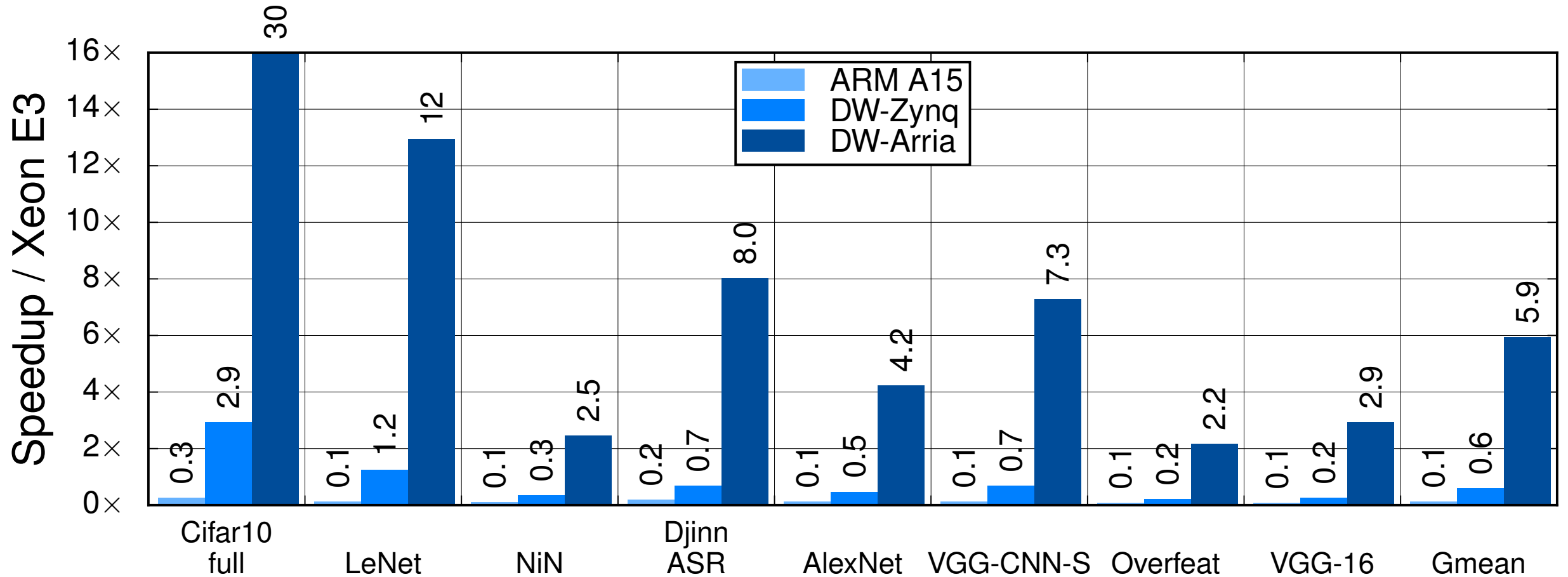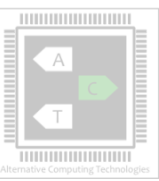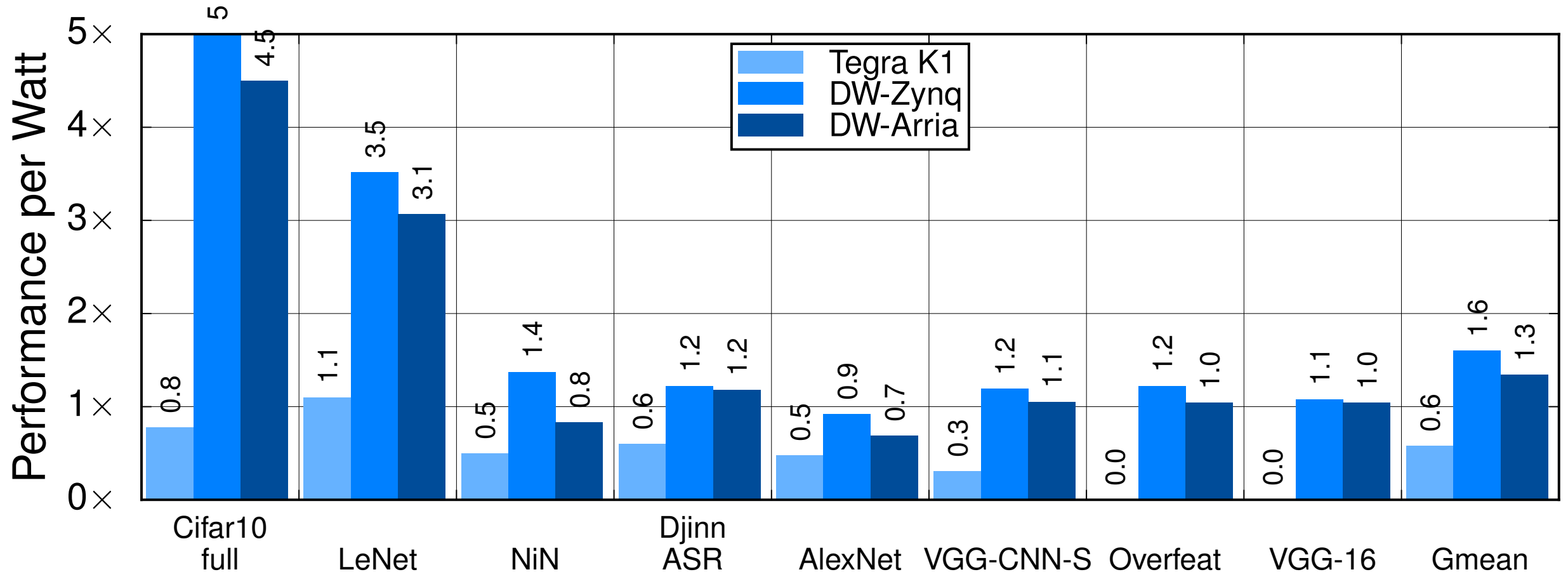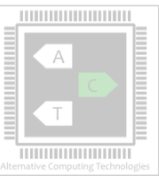|  | High Performance |  | Low Power |
|---|---|---|---|
| **FPGA** | Altera Arria 10 | Altera Stratix V | Xilinx ZC702 |
| **CPU** | Intel Xeon E3 |  | ARM Cortex 15 |
| **GPU** | Tesla K40 | GTX 650 Ti | Tegra K1 |

# Performance vs CPUs



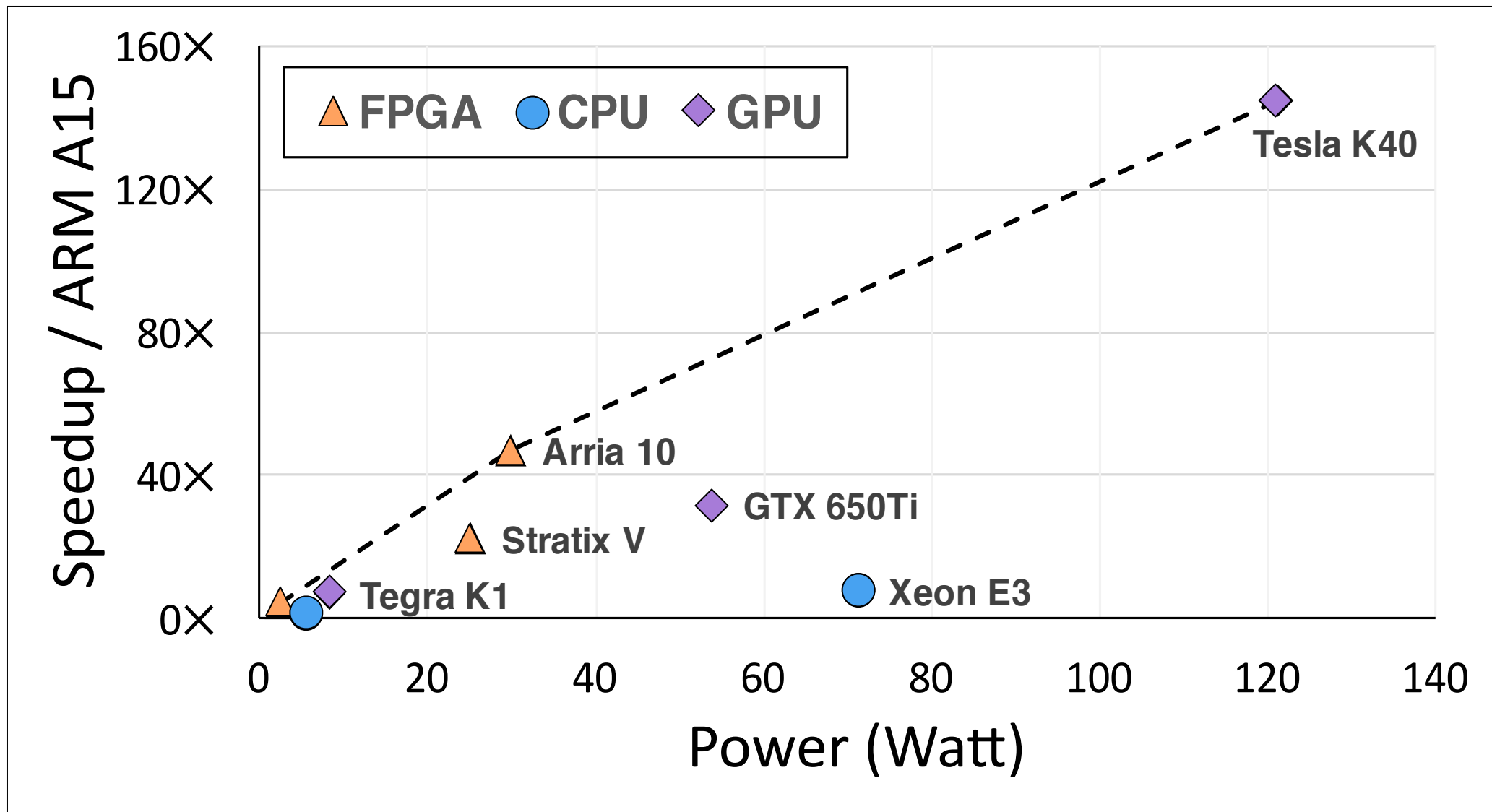Compared to Xeon, Arria10 is **5.9x** faster, and Zynq is **0.6x** faster.

# Performance-per-Watt vs GPUs



Compared to TeslaK40, Zynq is **1.6x**, and Arria10 is **1.3x** more power efficient.
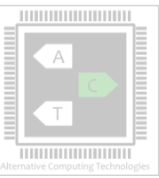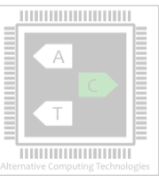
# Performance vs Power

# Conclusion

FPGAs are a promising option for low-mid power range

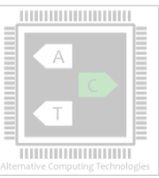However, there is a semantic gap between the high-level DNN models and FPGA acceleration

**DNNWEAVER is an initial step in making FPGAs more accessible to DNN programmers**
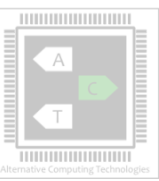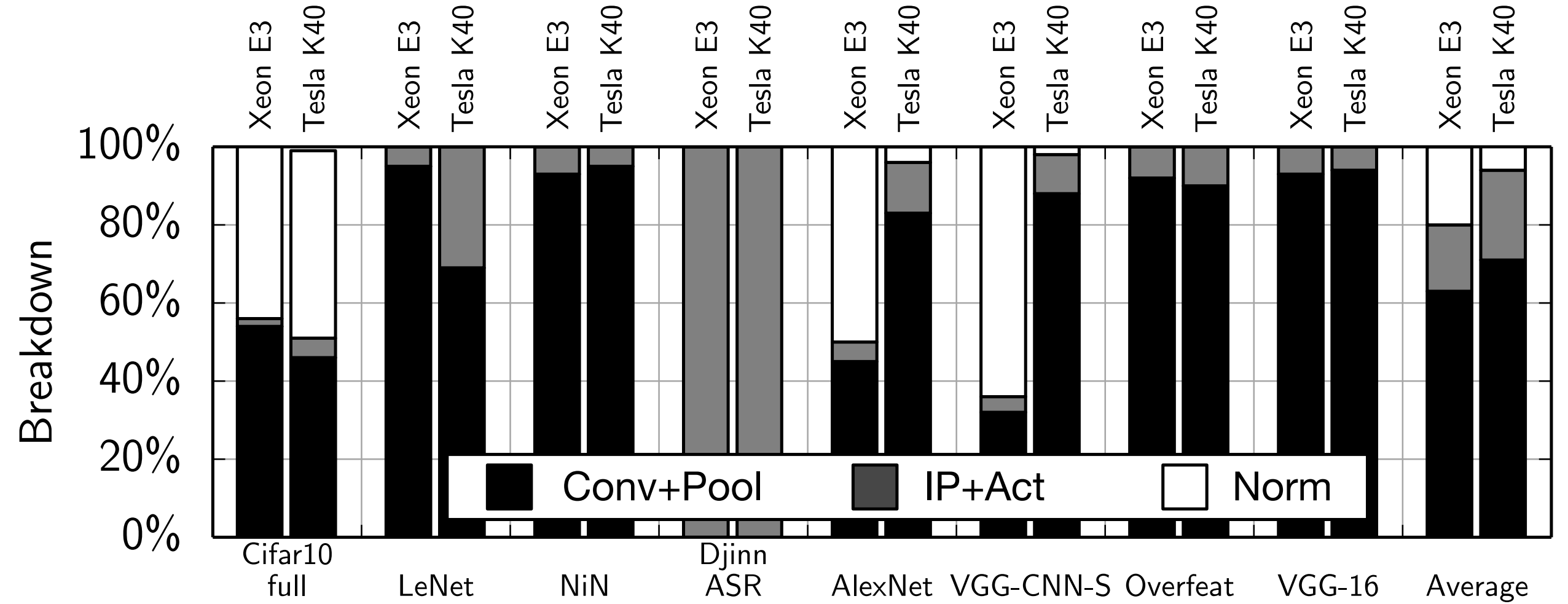
# Questions?

# Backup Slides

# Runtime Breakdown



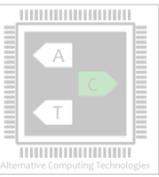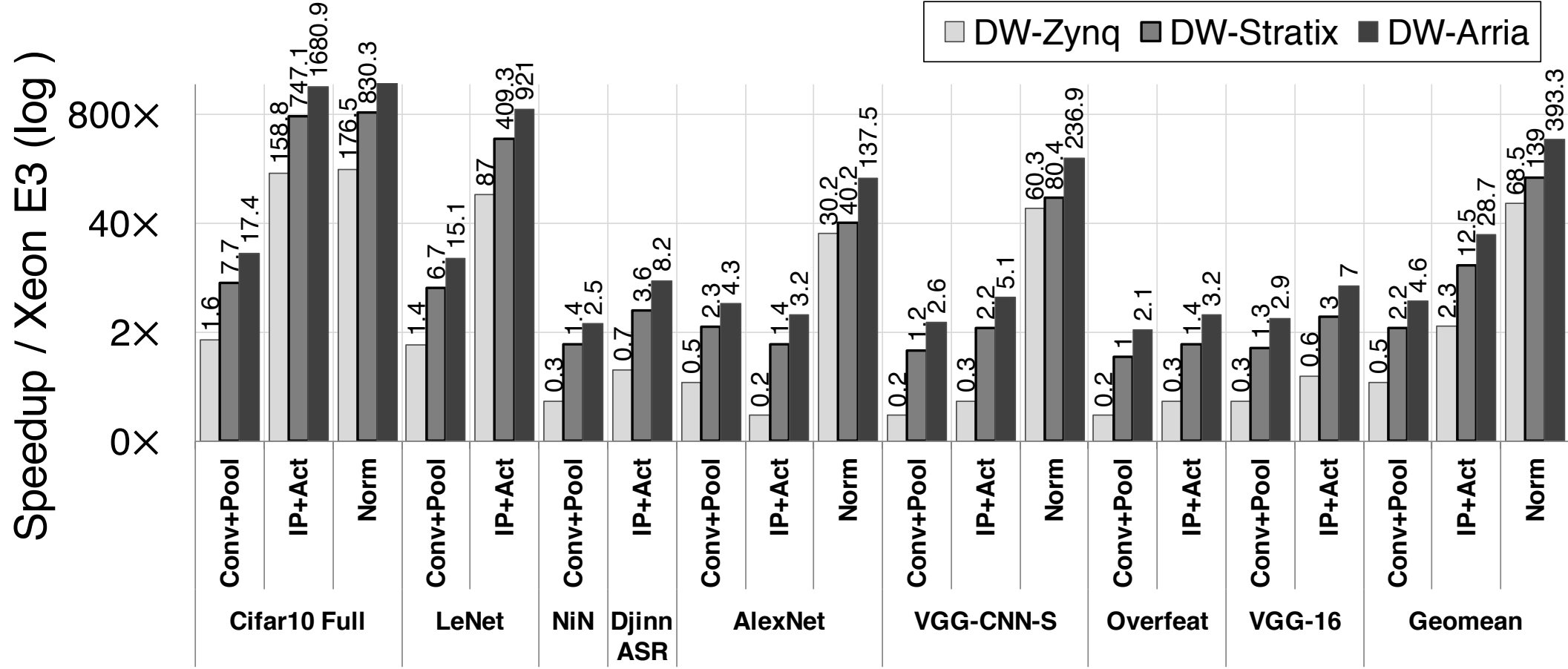Breakdown (%): 0%, 20%, 40%, 60%, 80%, 100%

Legend: Conv+Pool, IP+Act, Norm

Categories: Cifar10 full, LeNet, NiN, Djinn ASR, AlexNet, VGG-CNN-S, Overfeat, VGG-16, Average
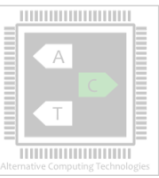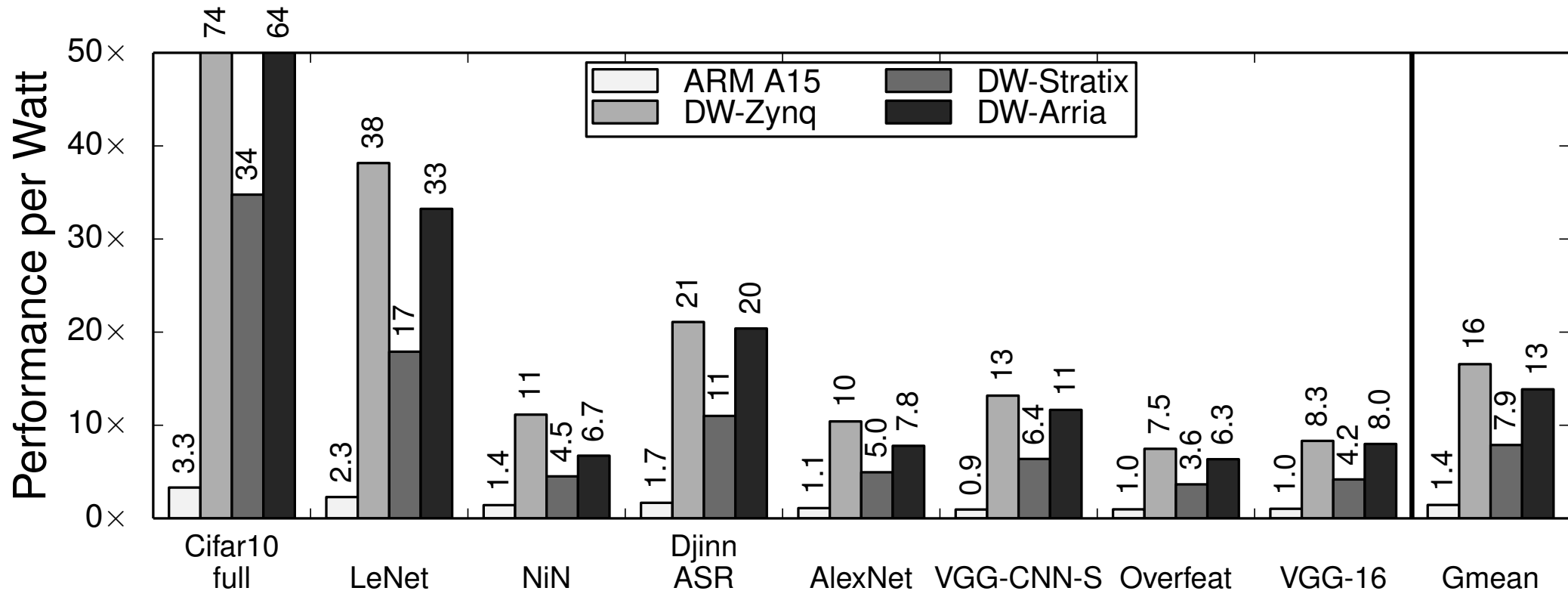
Each with: Xeon E3, Tesla K40

# Per Layer Speedup CPU

# Performance vs GPU
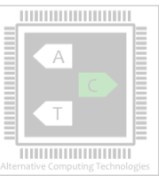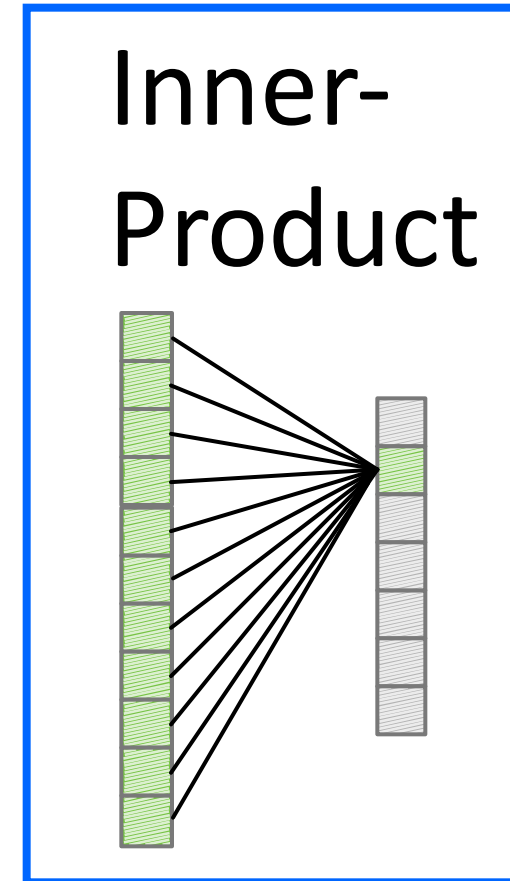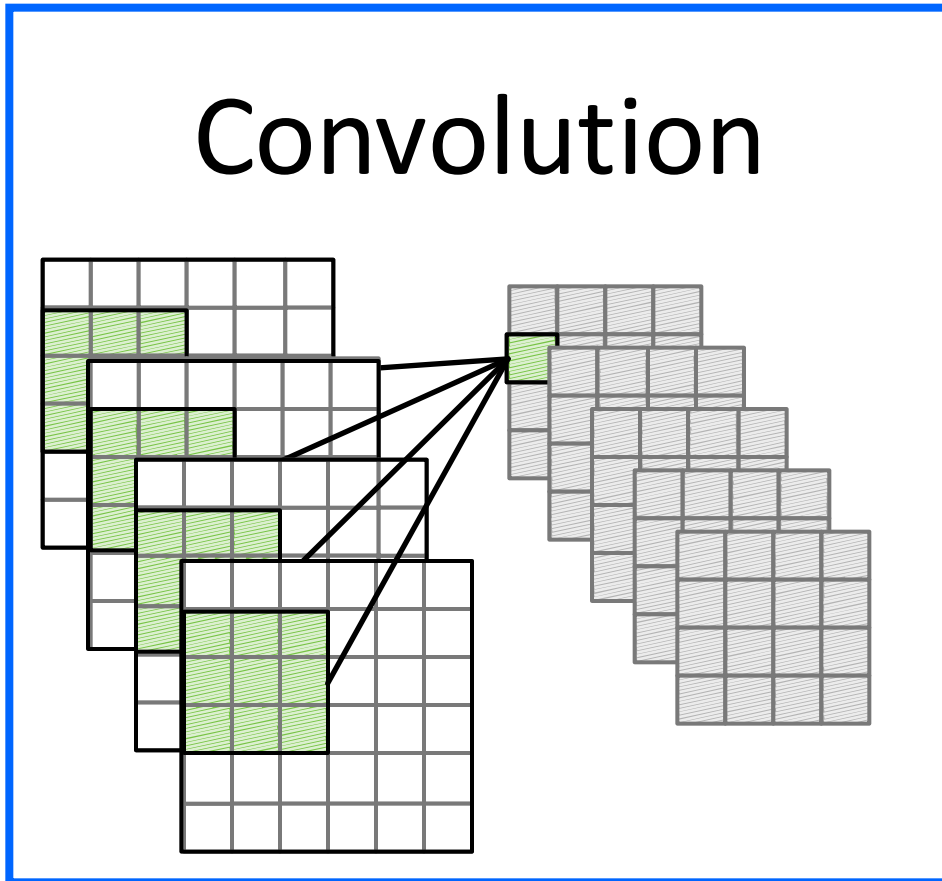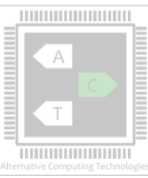
# Performance-per-Watt vs CPU

# Platforms Tested

| | | |
|---|---|---|
| **Altera Arria 10**<br>TDP:35W<br>$4495 | **Altera Stratix V**<br>TDP:25W<br>$6999 | **Xilinx Zynq ZC702**<br>TDP: 2W<br>$129 |
| **Intel Xeon E3-1276 V3**<br>TDP: 84W<br>$339 | **ARM Cortex 15**<br>TDP: 5W<br>$191 | |
| **Tegra K1 GPU**<br>TDP: 10 W<br>$191 | **GTX 650 Ti**<br>TDP: 110<br>$150 | **Tesla K40**<br>TDP: 235 W<br>$5499 |

# Deep Neural Networks (DNNs)
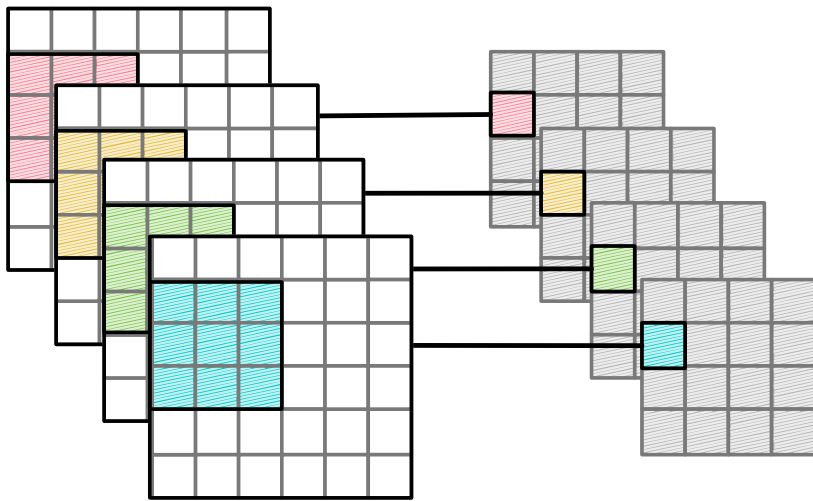
Convolution

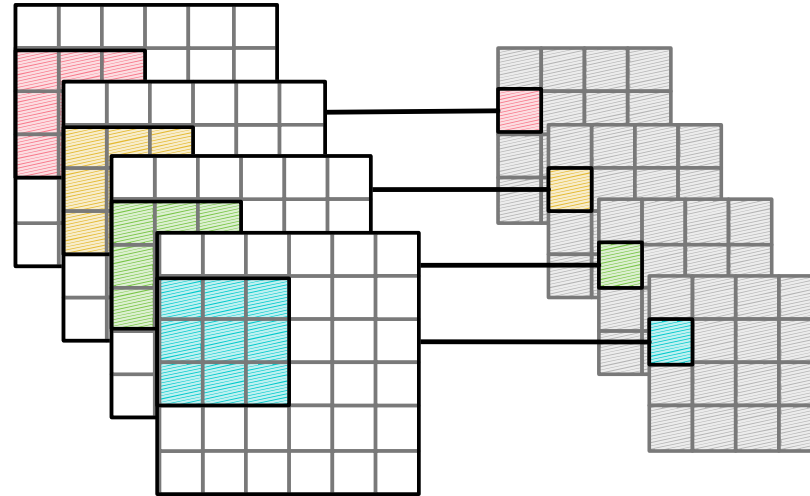Inner-Product

**Convolution and Inner-Product are the Learnable Layers**
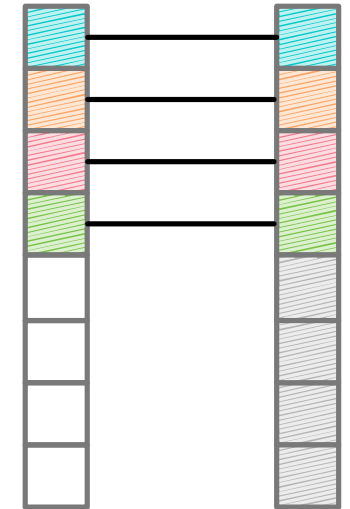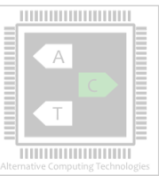
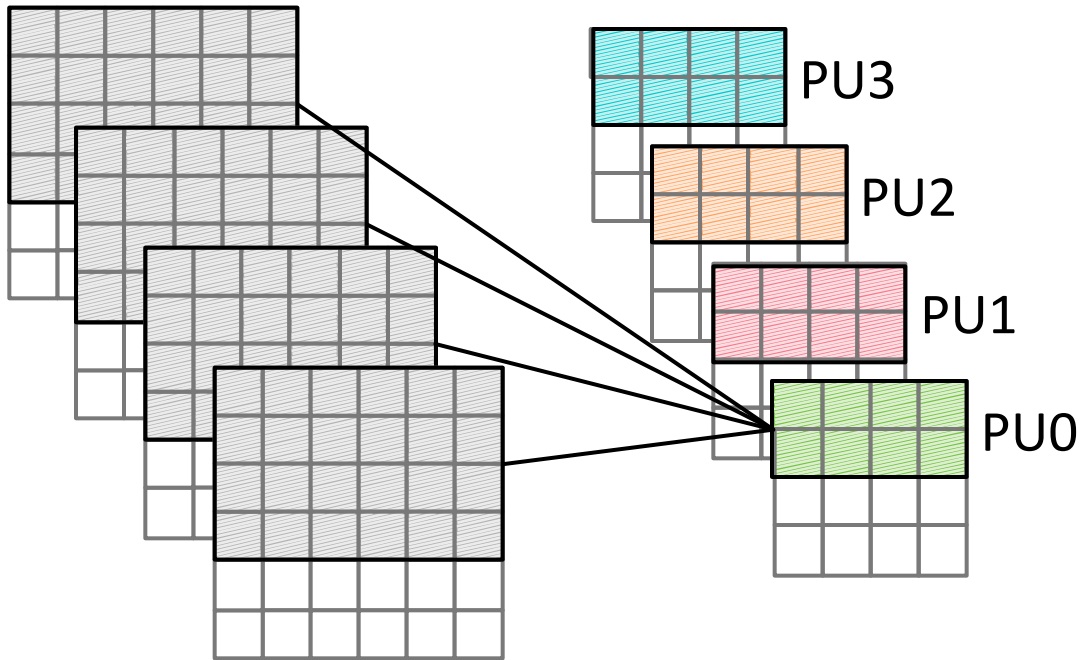# Deep Neural Networks (DNNs)



Pooling

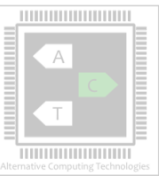Normalization

Activation

**DnnWeaver supports all these five layers.**

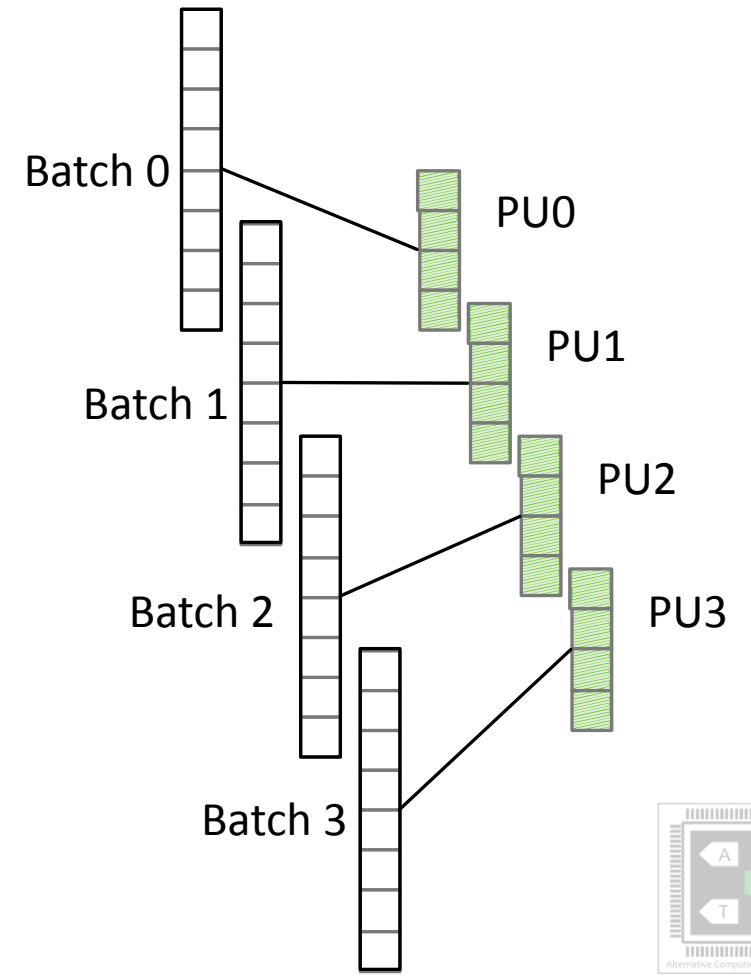# Scheduling Operations on Hardware

Convolution

Inner-product

# Performance vs GPU

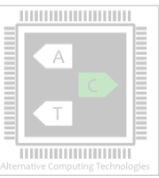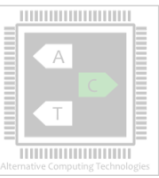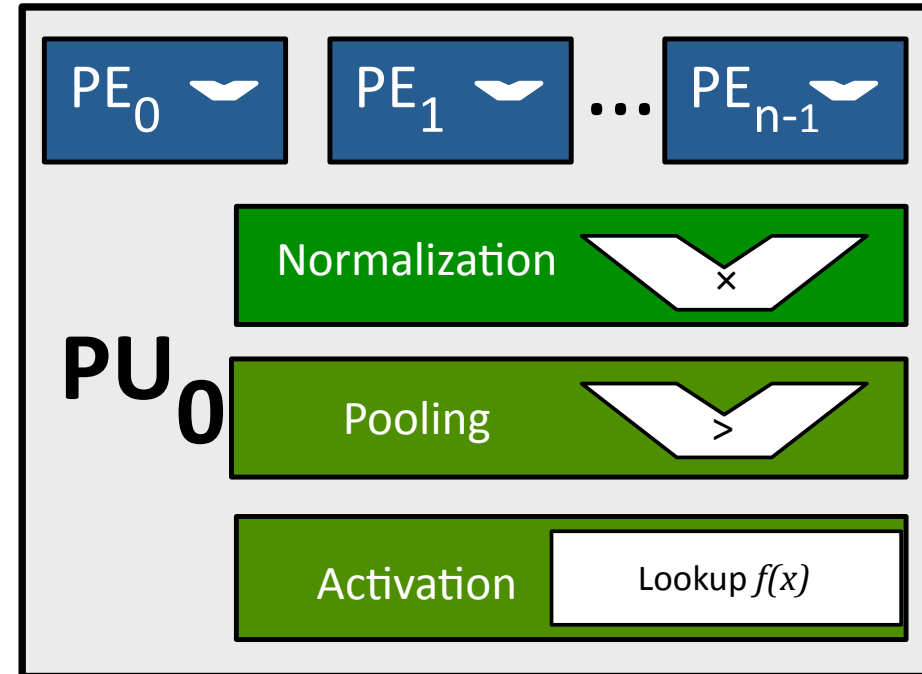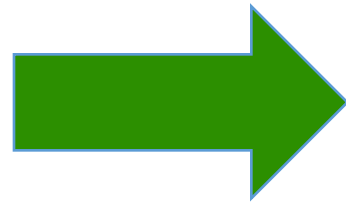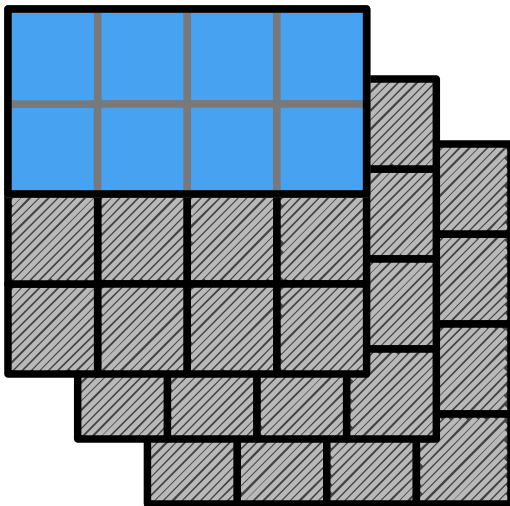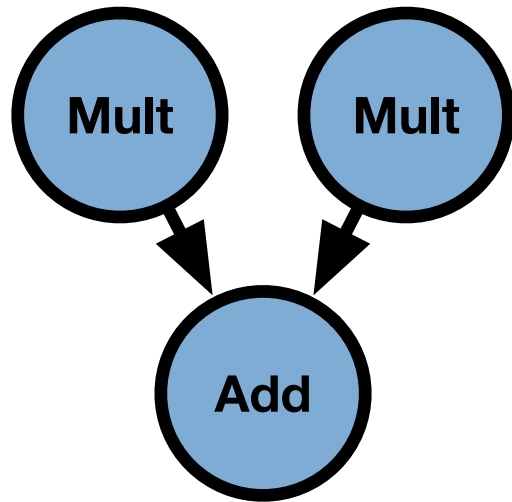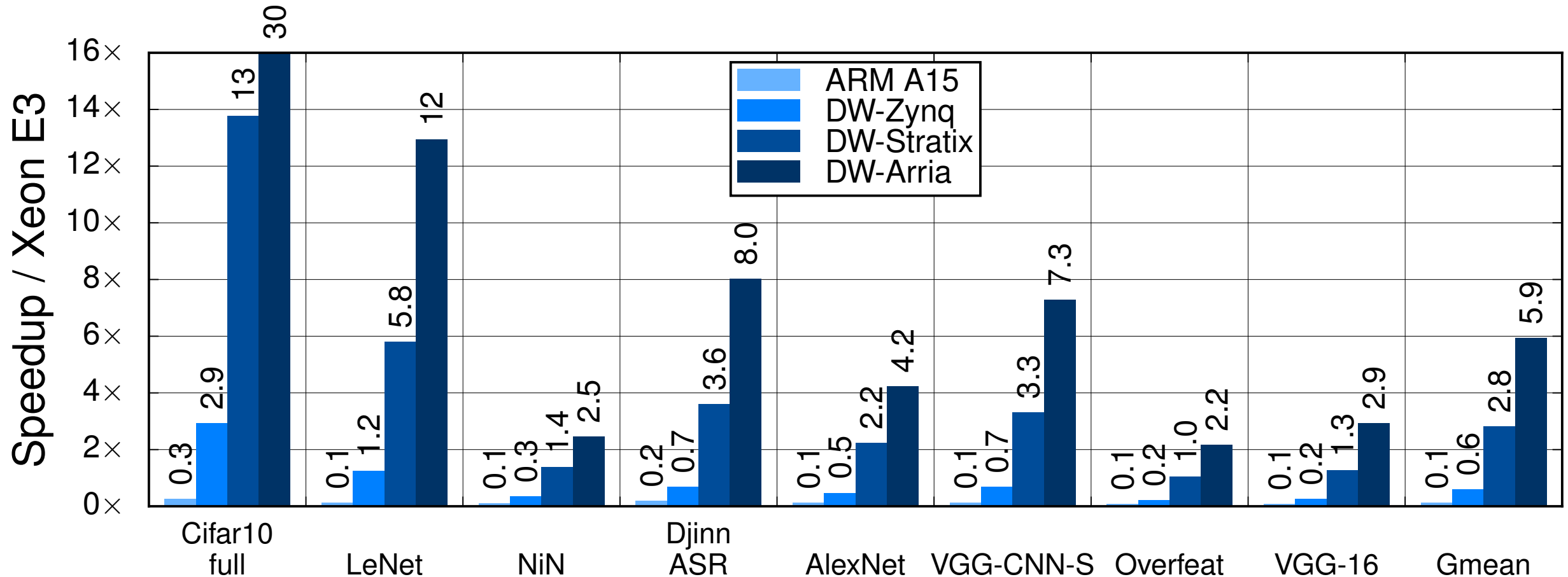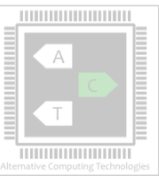# Performance-per-Watt vs CPU

# Executing the Slice

# Performance vs CPUs



Speedup / Xeon E3

Legend: ARM A15, DW-Zynq, DW-Stratix, DW-Arria

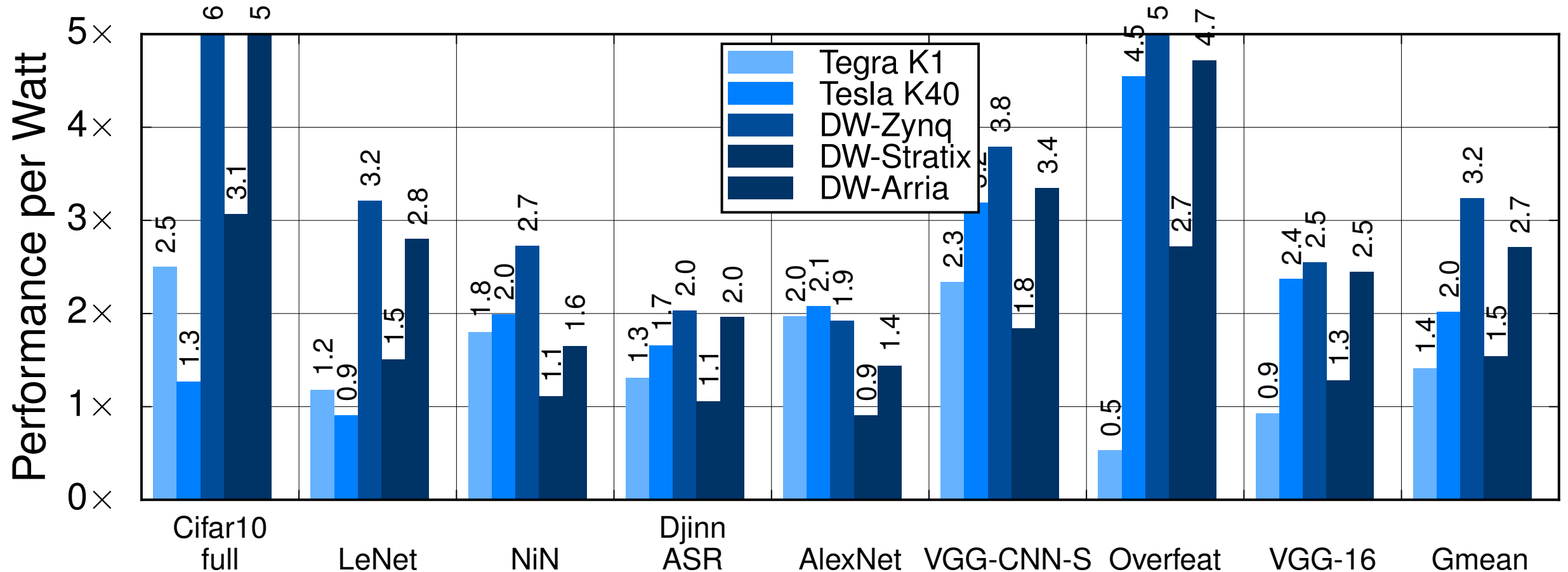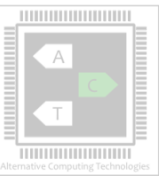| Benchmark | ARM A15 | DW-Zynq | DW-Stratix | DW-Arria |
|---|---|---|---|---|
| Cifar10 full | 0.3 | 2.9 | 13 | 30 |
| LeNet | 0.1 | 1.2 | 5.8 | 12 |
| NiN | 0.1 | 0.3 | 1.4 | 2.5 |
| Djinn ASR | 0.2 | 0.7 | 3.6 | 8.0 |
| AlexNet | 0.1 | 0.5 | 2.2 | 4.2 |
| VGG-CNN-S | 0.1 | 0.7 | 3.3 | 7.3 |
| Overfeat | 0.1 | 0.2 | 1.0 | 2.2 |
| VGG-16 | 0.1 | 0.2 | 1.3 | 2.9 |
| Gmean | 0.1 | 0.6 | 2.8 | 5.9 |

Compared to Xeon, Arria10 is **5.9x** faster, Stratix V is **2.8x** faster, and Zynq is **0.6x** faster.
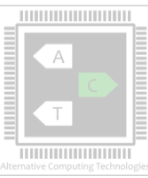
# Performance-per-Watt vs GPUs



Compared to GTX650, Zynq is **3.2x**, Stratix V is **1.5x**, and Arria10 is **2.7x** more power efficient.

# Co-optimize Hardware and Execution Schedule

**Hardware**

**Schedule**



PE$_0$ ... PE$_1$ ... PE$_{n-1}$

**PU$_0$**

Normalization ×

Pooling >

Activation | Lookup $f(x)$

# Processing Engine