# ACCELERATION OPTIONS FOR AI AND HPEC
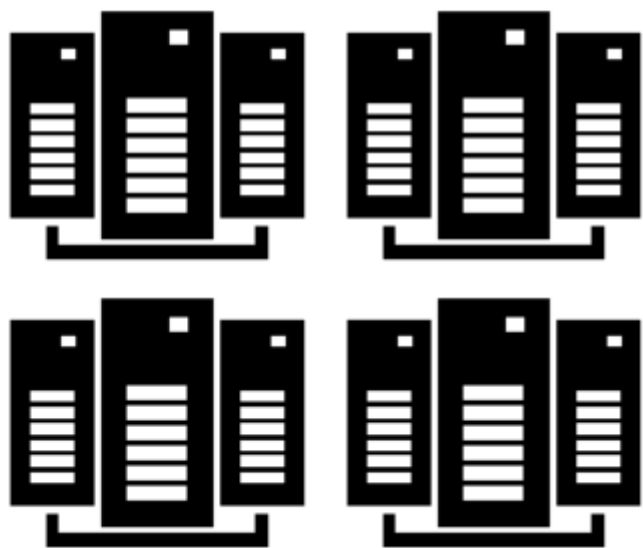
## Embedded Tech Trends, 2020
## Atlanta, GA

Dr. Mohamed BERGACH, Sr. Systems Architect – January 27, 2020
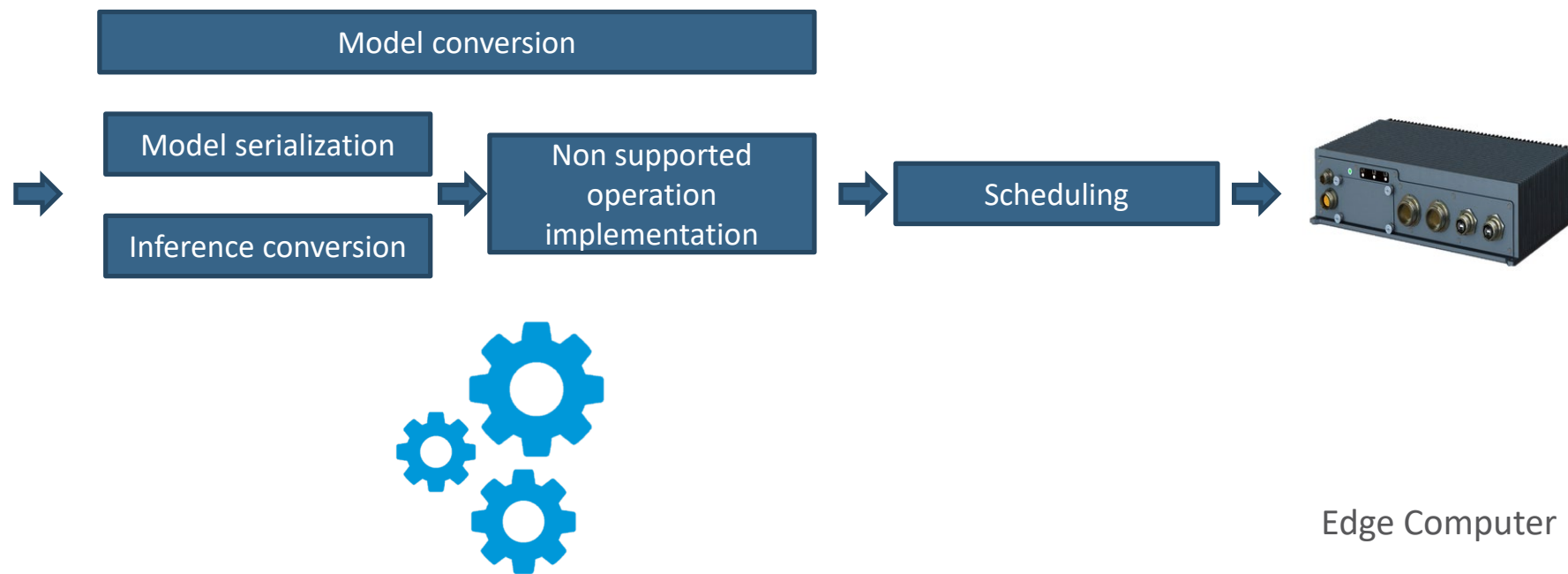
## AGENDA

- ❑ AI from datacenter to the edge
- ❑ Software acceleration
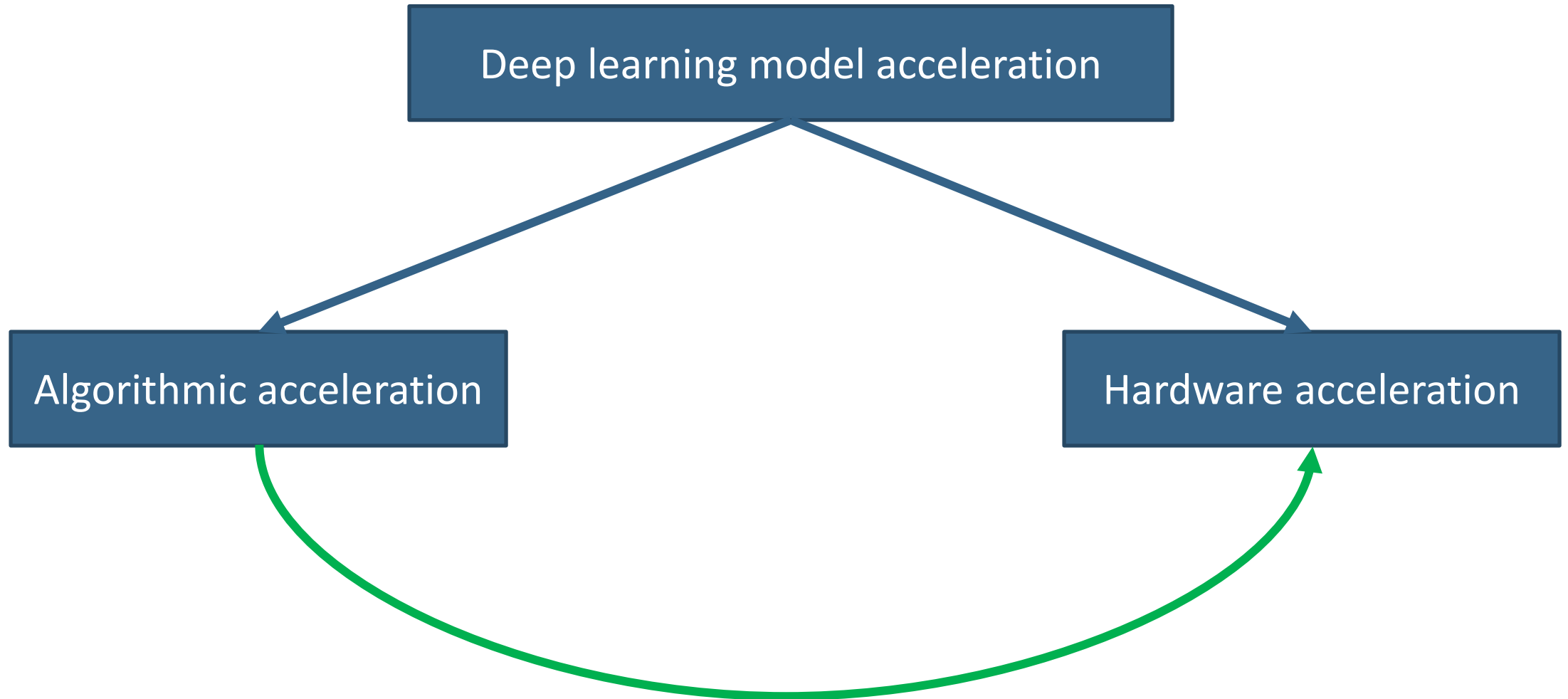- ❑ Hardware acceleration
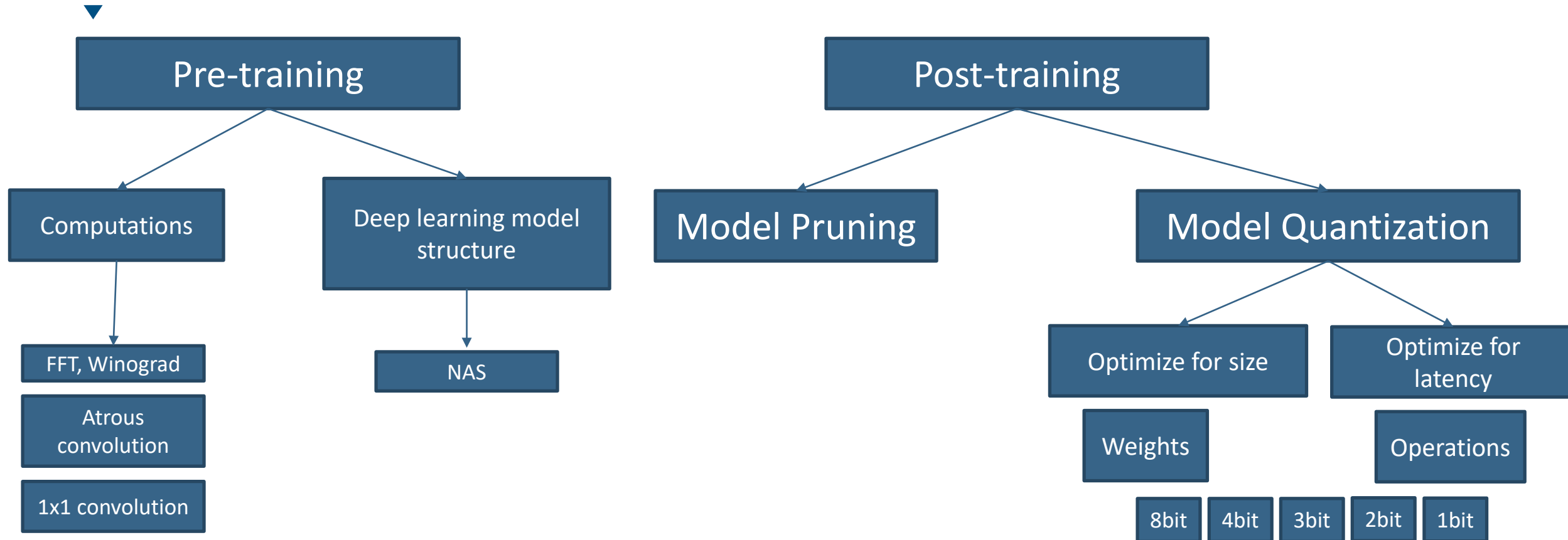- ❑ Conclusion and takeaway

# FROM TRAINED AI MODEL TO THE EDGE



| Model conversion |
|---|

| Model serialization | Non supported operation implementation |
|---|---|
| Inference conversion | |

| Scheduling |
|---|

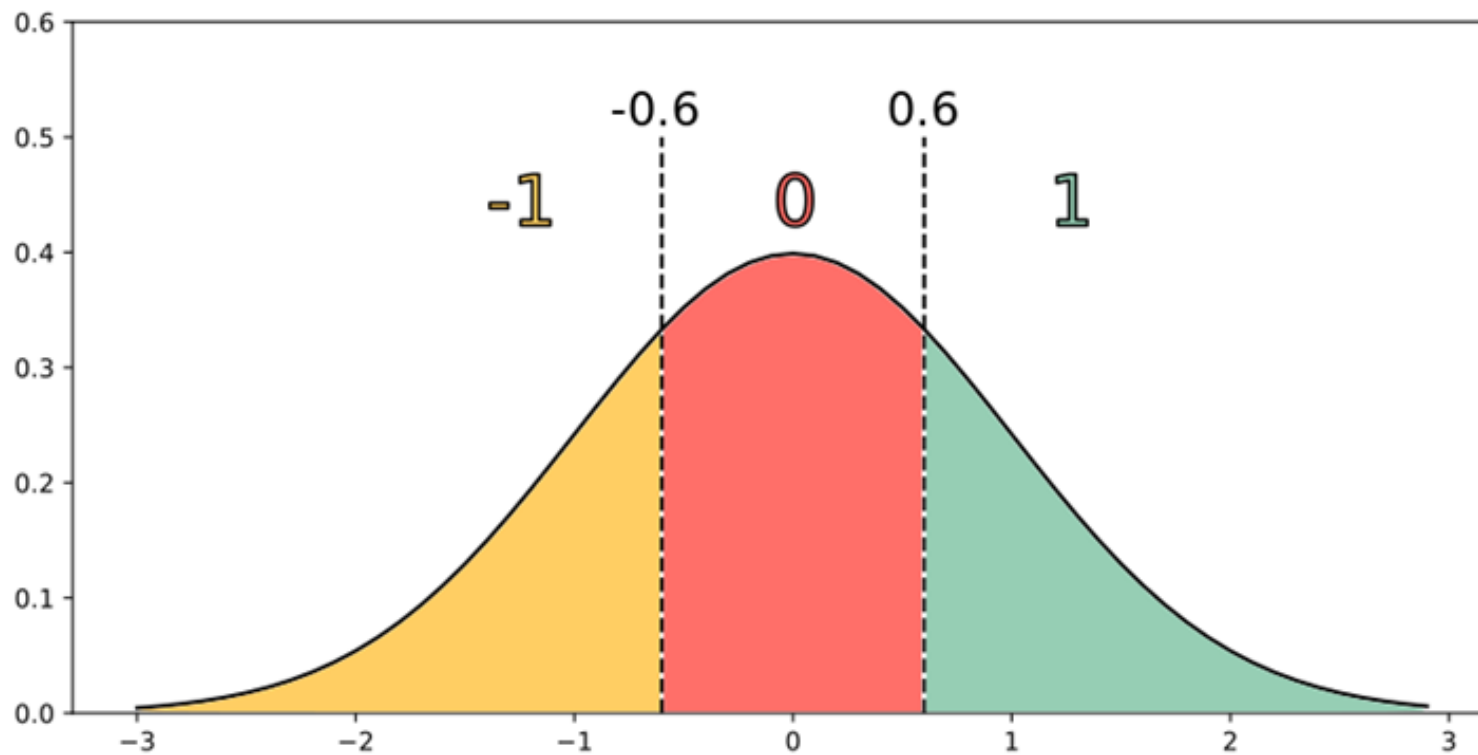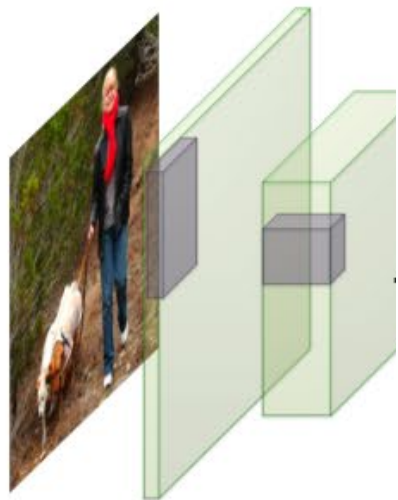AI model trained in the data center

Edge Computer

# ACCELERATION

# ALGORITHMIC ACCELERATION:

# ALGORITHMIC ACCELERATION: QUANTIZATION



| Computation Saving (Inference) | Accuracy on ImageNet (AlexNet) |
|---|---|
| 1x | %56.7 |
| ~2x | %56.8 |
| ~58x | %44.2 |

Source: [Rastegari et al. ECCV'16]

# ALGORITHMIC ACCELERATION: PRUNING



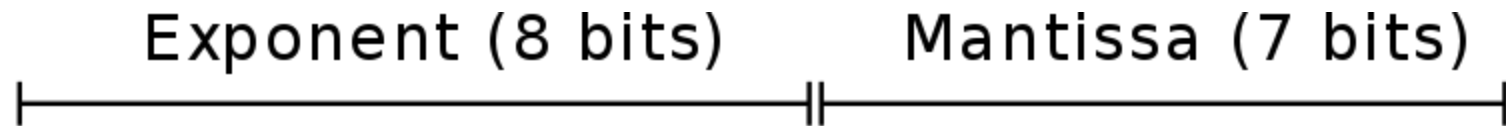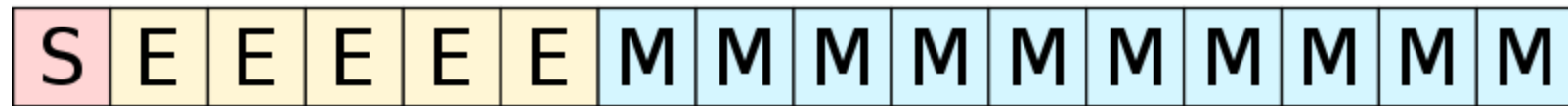before pruning     after pruning

pruning synapses
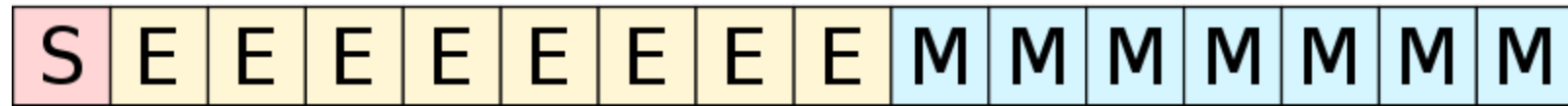
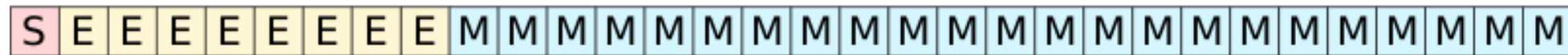pruning neurons

## 10X speedup

Source: [Han et al. NIPS'15]

# FLOAT VS BFLOAT

**Float16**
range (~5.96e-8 to 65,504)

Exponent (5 bits)  Mantissa (10 bits)
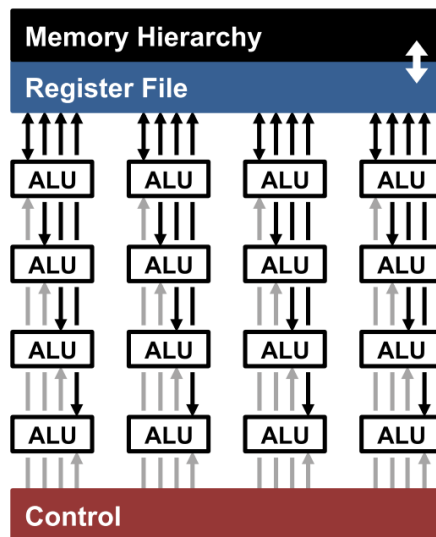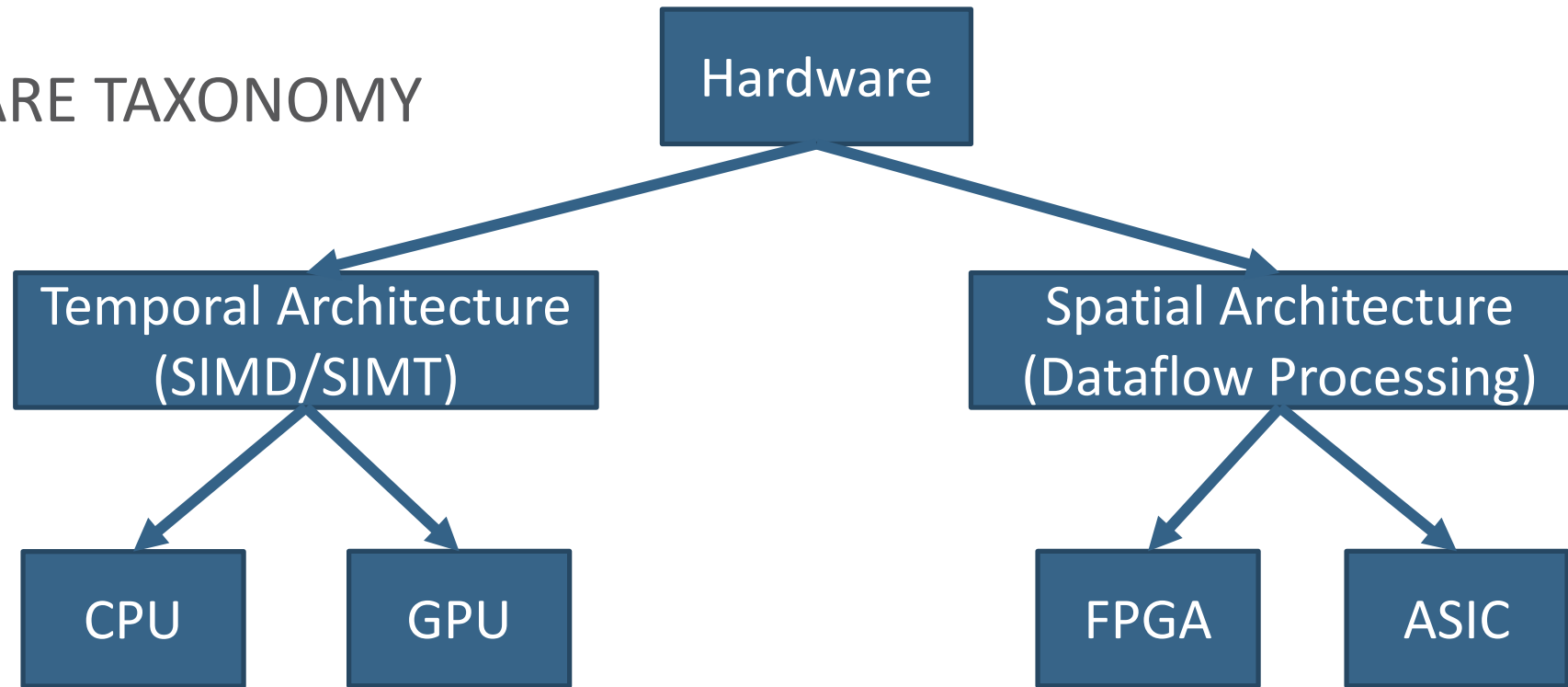
| S | E E E E E | M M M M M M M M M M |

**Bfloat16**
Range (~1e-38 to ~3e38)

Exponent (8 bits)  Mantissa (7 bits)

| S | E E E E E E E E | M M M M M M M |

**Float32**
Range (~1e-38 to ~3e38)

Exponent (8 bits)  Mantissa (23 bits)

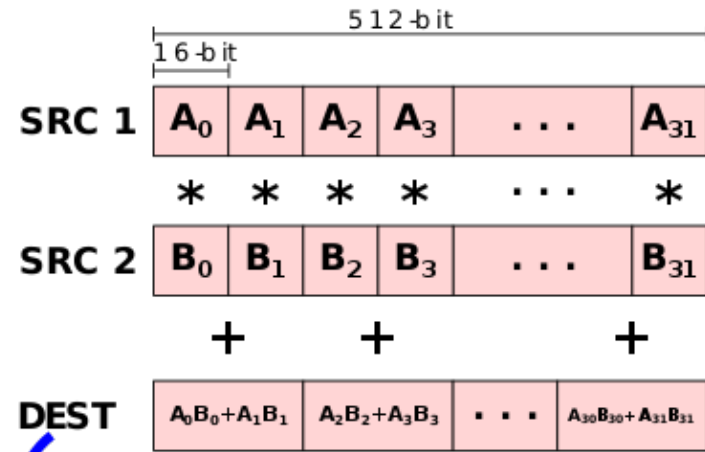| S | E E E E E E E E | M M M M M M M M M M M M M M M M M M M M M M M |

# HARDWARE TAXONOMY
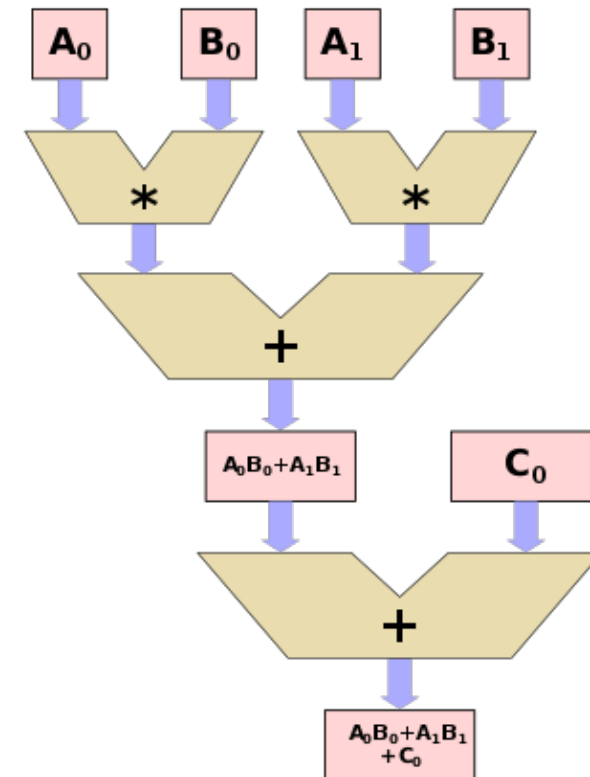


Source: [Sze et al. Proceedings of the IEEE 105(12): 2017]

9

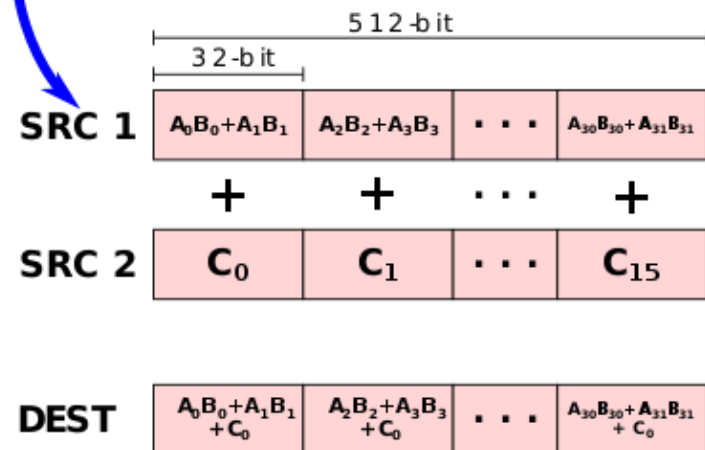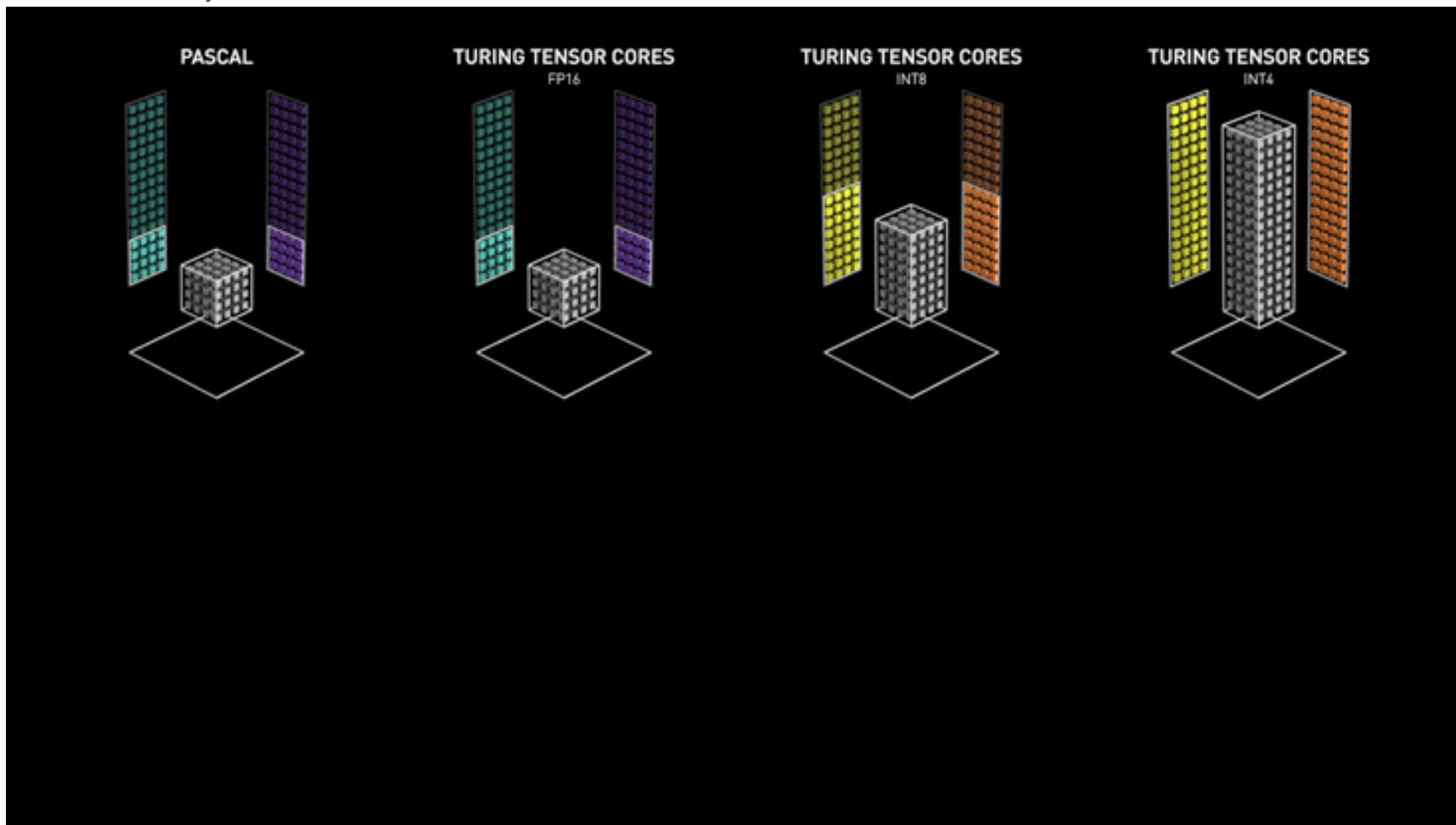# CPU (AVX512, VNNI)

2x64x2x2 = 0.5TOPS/core
4 core Intel Tiger lake CPU at 25W TDP
2TOPS => 80GOPS/W

# GPU (TENSOR CORES)
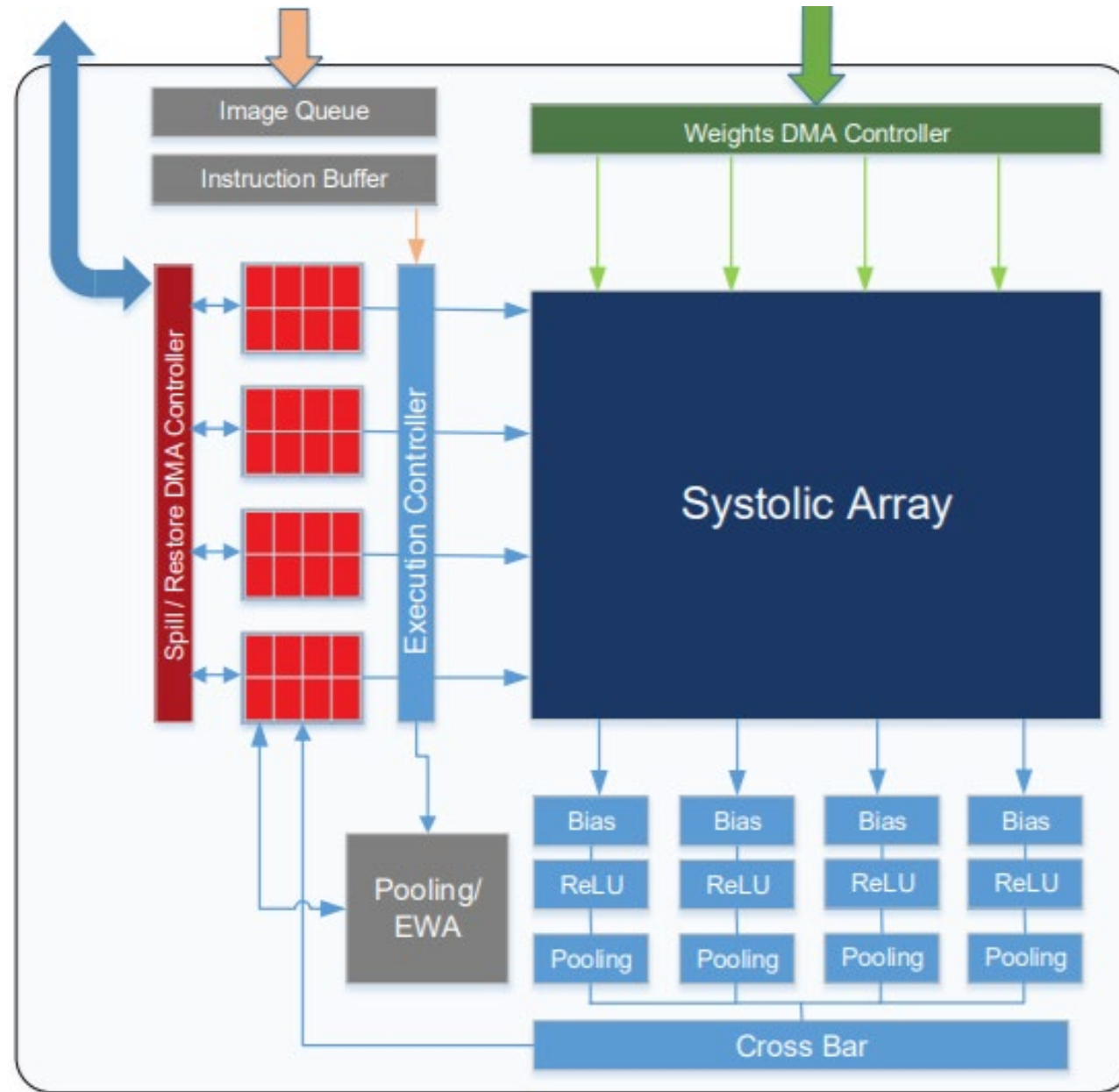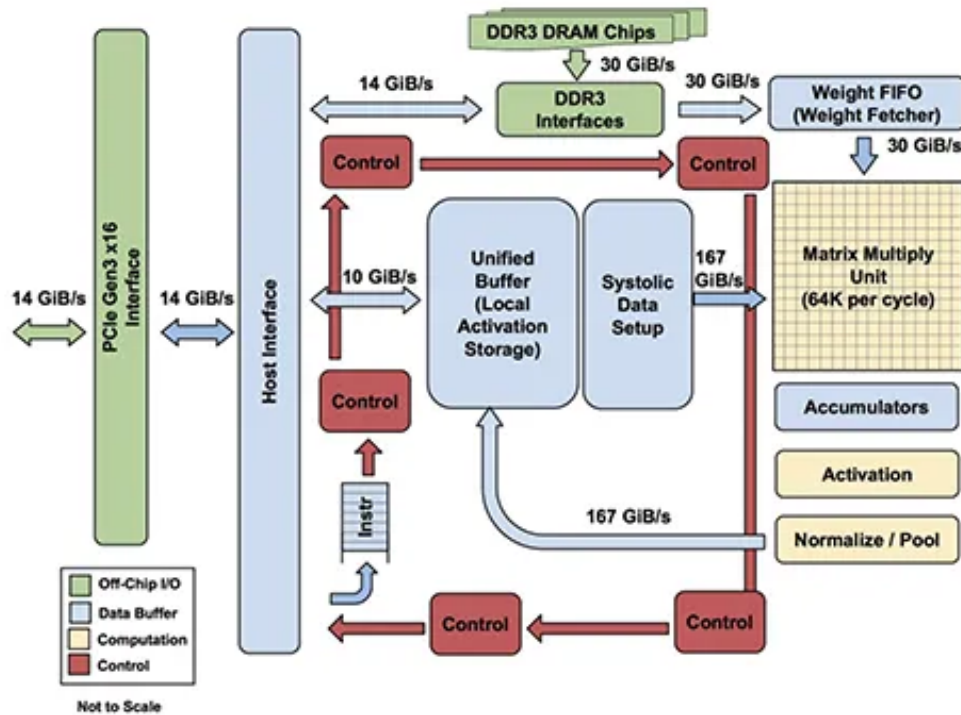
Nvidia Tesla T4
130TOPs for 75W TDP
.5TOPs/W

# FPGA

XILINX xDNN
75W
21TOPS
0.3TOPS/W
Very low latency

# ASIC



Google TPU
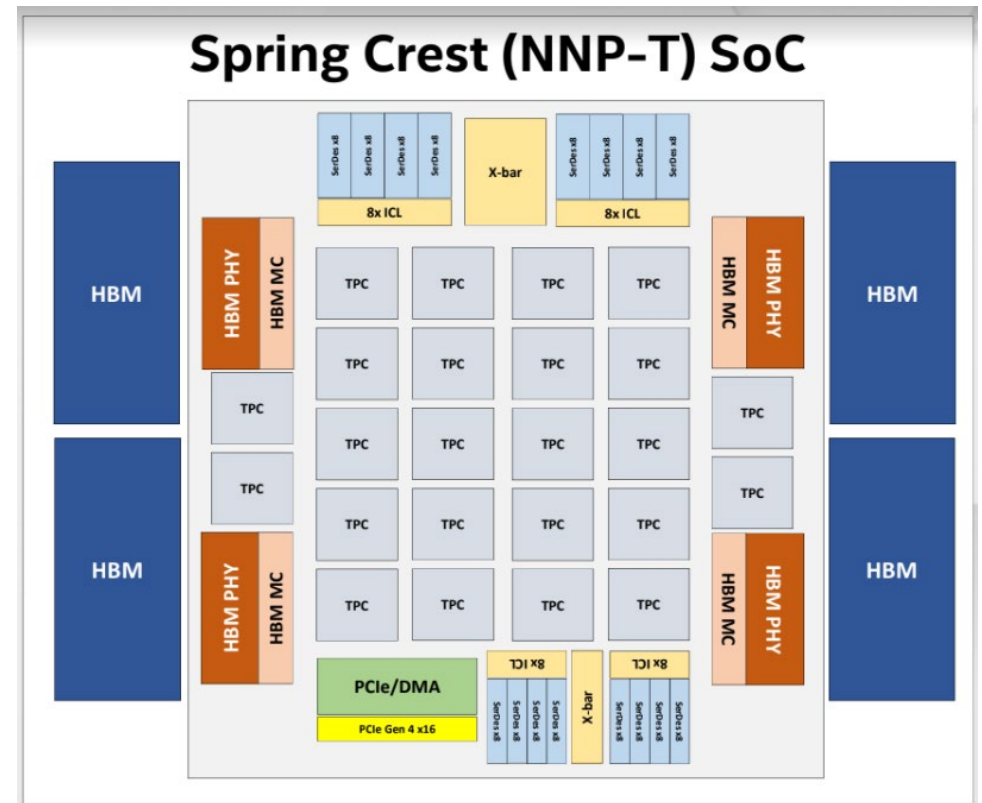
Edge TPU
4TOPS
2W
2TOPS/W



NNP-T

Intel Nervana Spring Crest
119TOPS
200W TDP
0.6TOPS/W

# HARDWARE COMPARISON

| HW Accelerator | TOPS/W | Latency | Efficiency |
|---|---|---|---|
| Intel Tiger lake CPU | 0.08 | ★★☆☆ | ★★★☆ |
| Nvidia T4 GPU | 0.5 | ★☆☆☆ | ★★☆☆ |
| XILINX xDNN FPGA | 0.3 | ★★★★ | ★★★★ |
| Intel NNP | 0.6 | ★★★☆ | ★★★☆ |
| Google Edge TPU | 2 | ★★☆☆ | ★★★☆ |

# TAKEAWAY

- ✓ AI on the Edge is a reality now, every major smartphone has an NPU
- ✓ GPUs are dominant right now but not in the future
- ✓ Specialized hardware are much more power efficient
- ✓ The software ecosystem is key ingredient to for best performance
- ✓ Quantization is a performance booster that needs to be considered

# Questions?

OUR CHANGING WORLD IS FILLED WITH
# BOUNDLESS OPPORTUNITIES

Dr. Mohamed BERGACH, Sr. Systems Architect – January 27, 2020