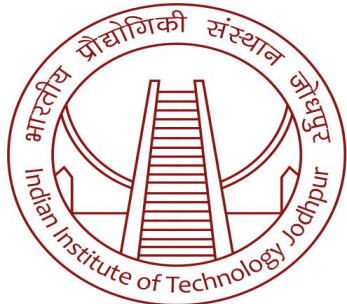


# Hardware Architectures and Design for AI

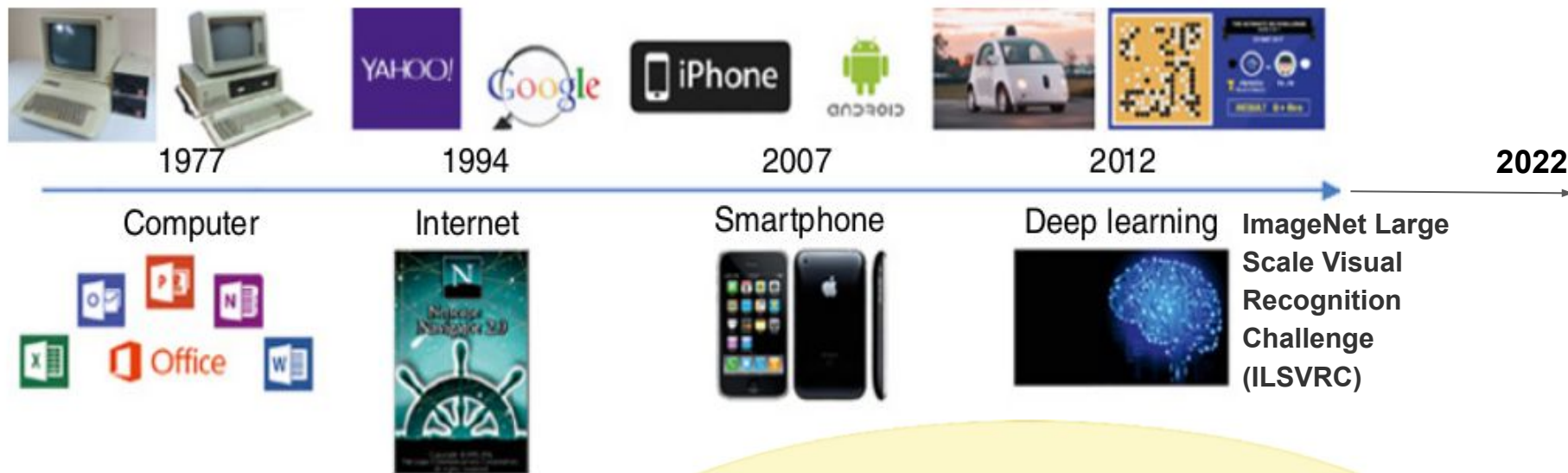
**Binod Kumar**  
**Electrical Engineering Department**  
**IIT Jodhpur**



॥ त्वं ज्ञानमयो विज्ञानमयोऽसि ॥

***Talk at IIT Tirupati (6th April 2022)***

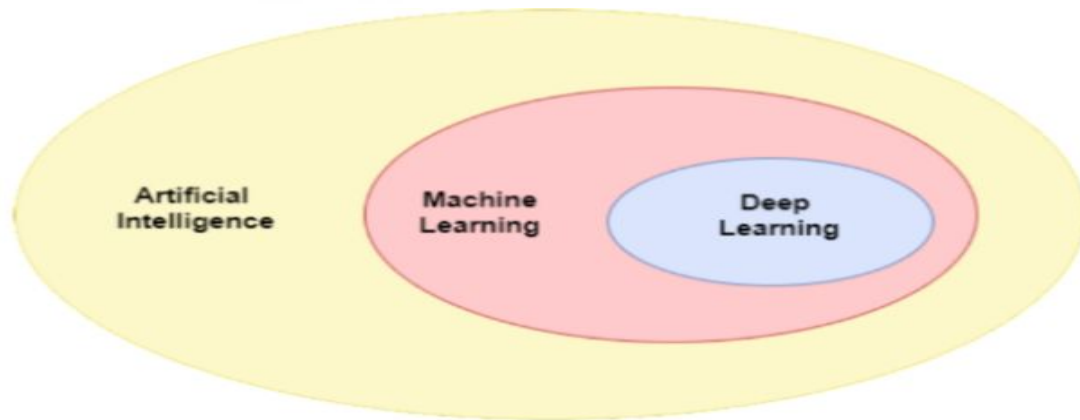
# The Technology Evolution

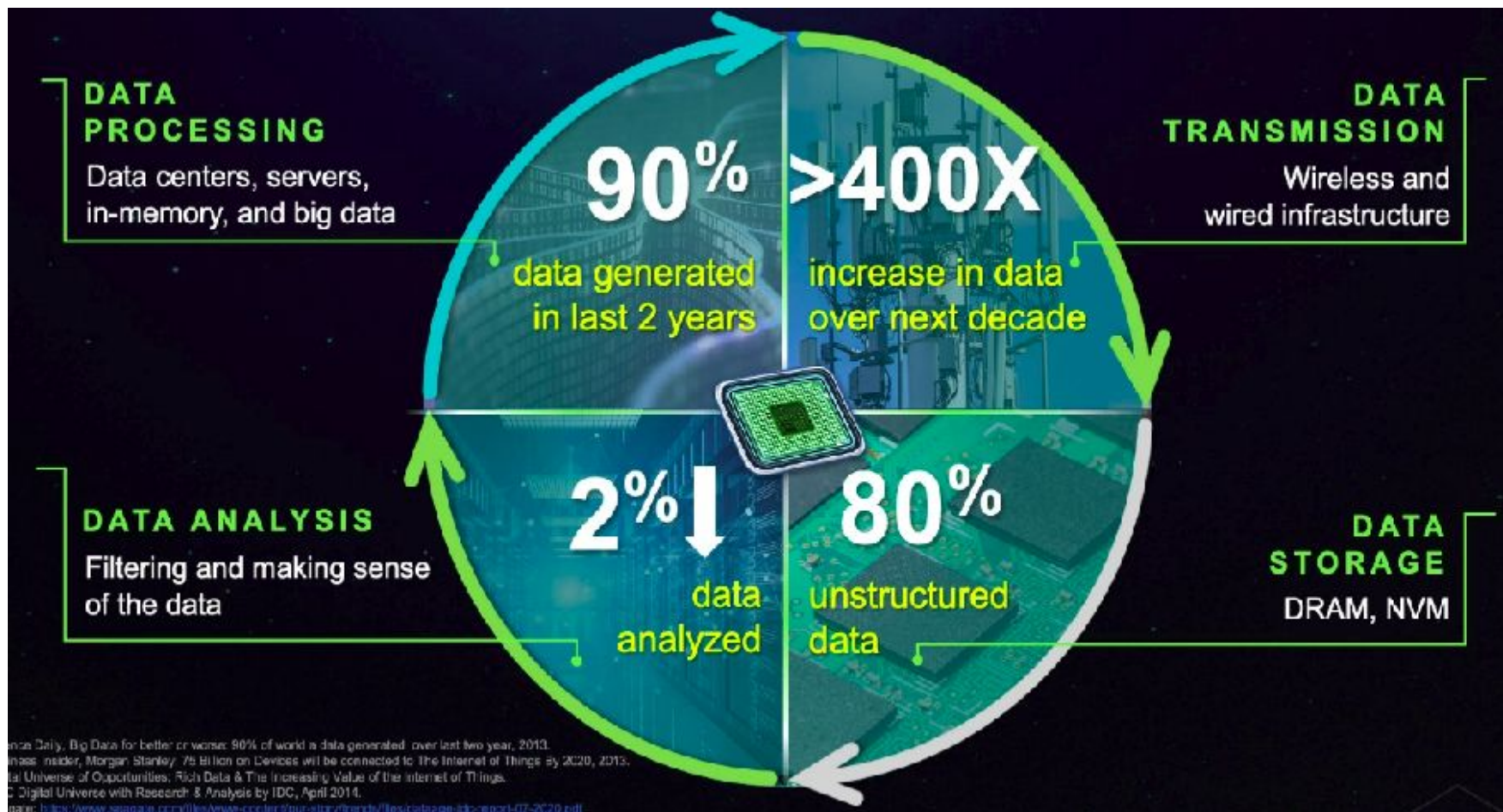


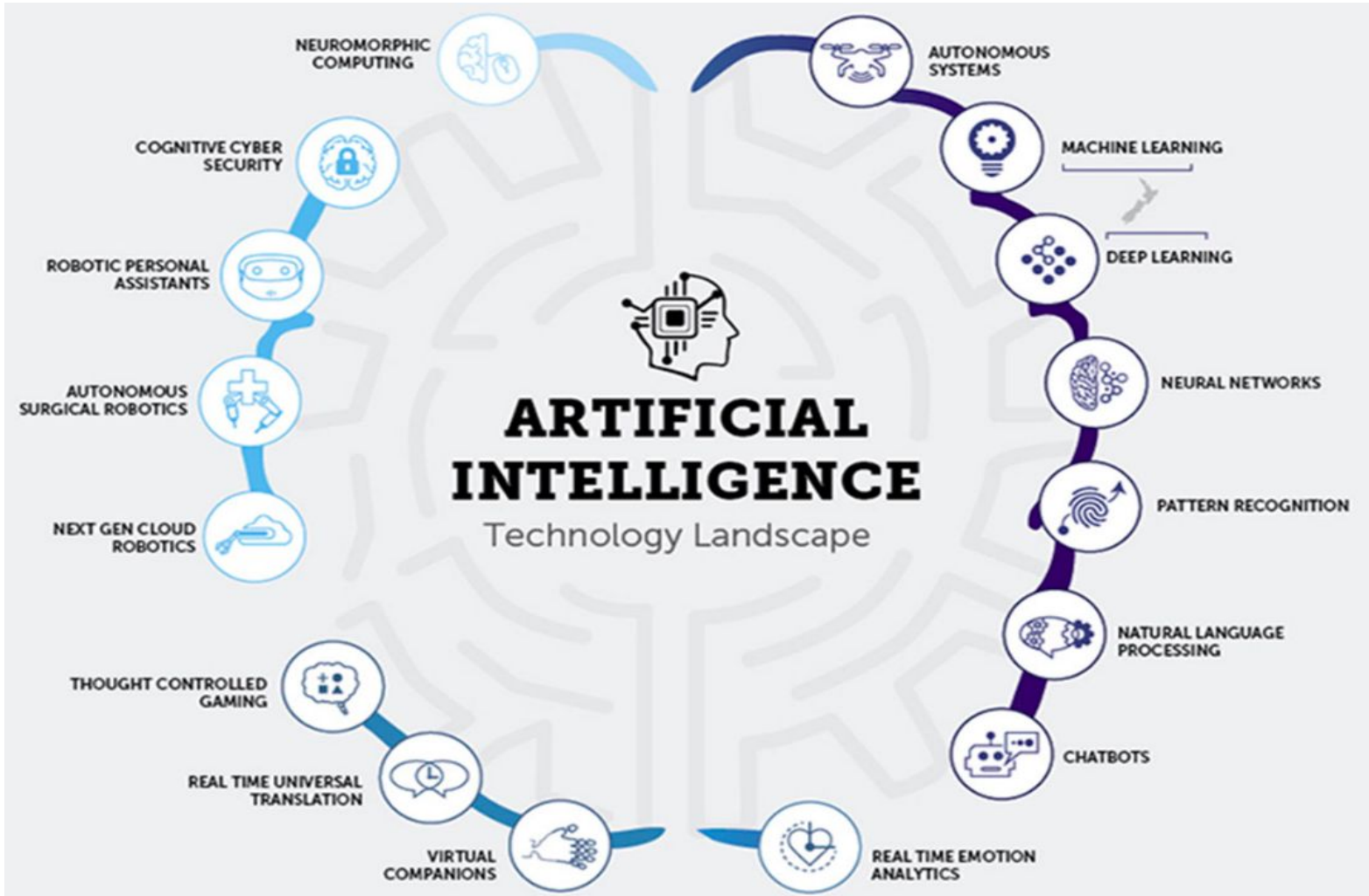
**AI:** Artificial Intelligence

**DL:** Deep Learning

**ML:** Machine Learning

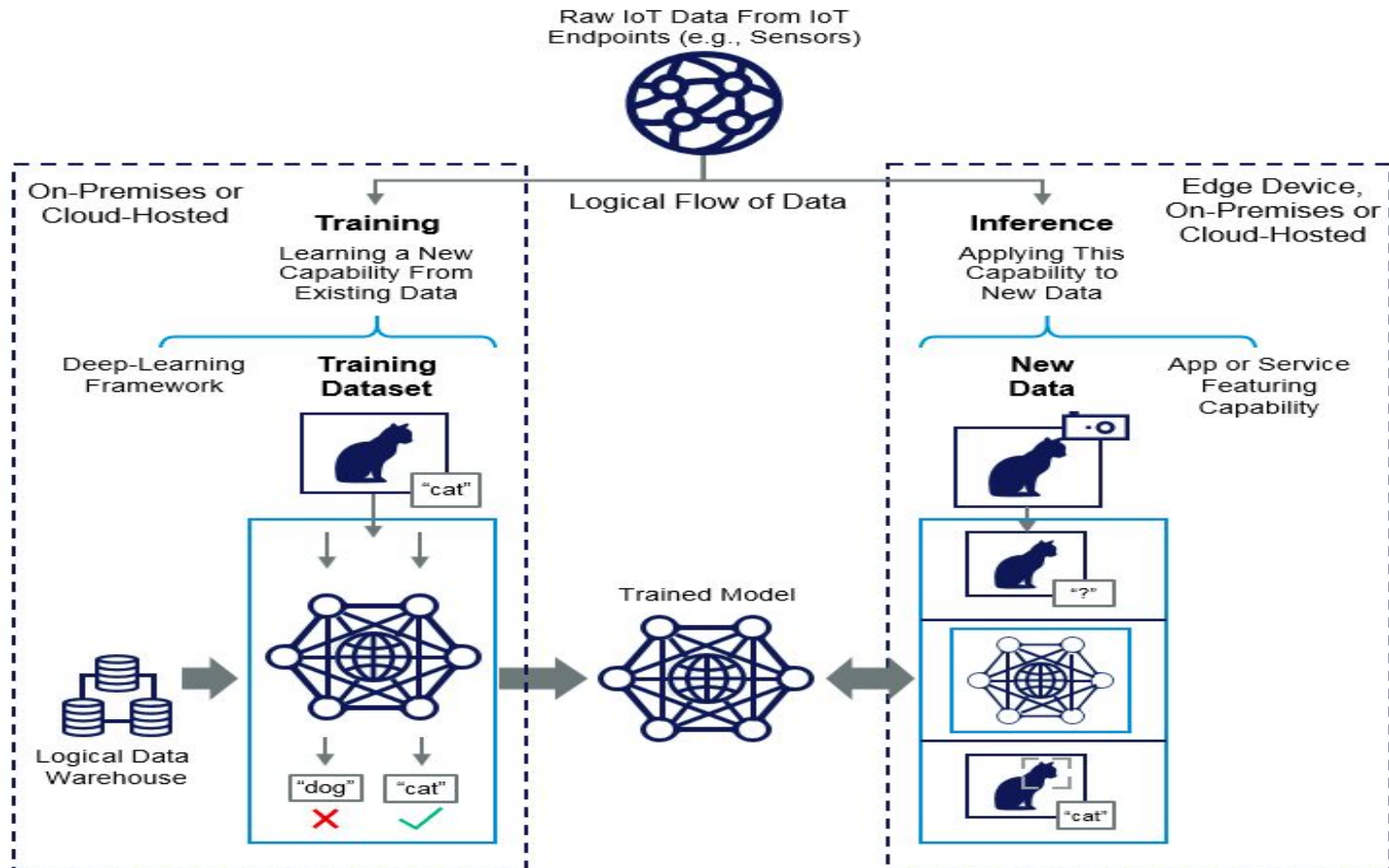






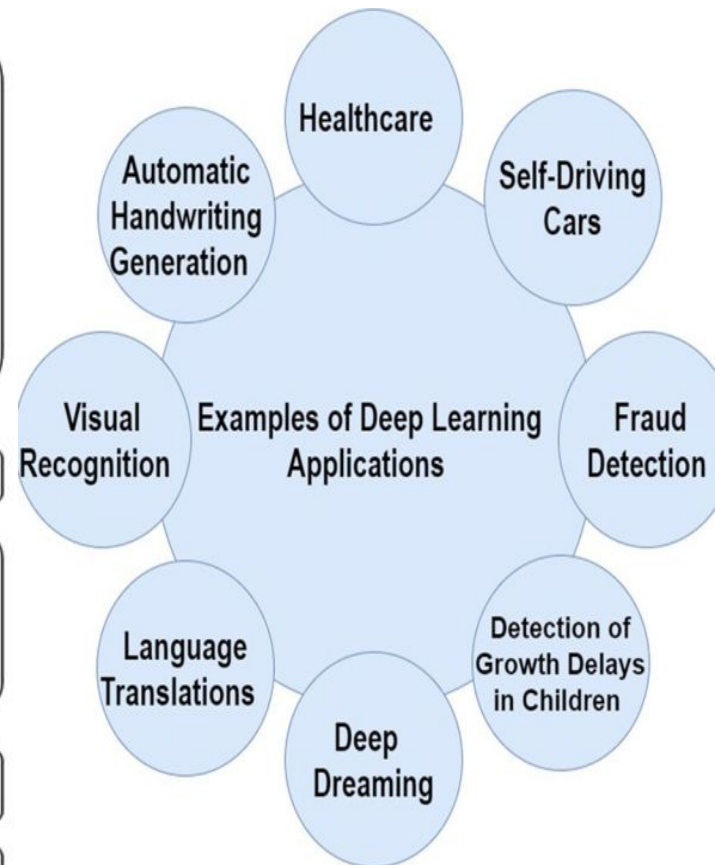
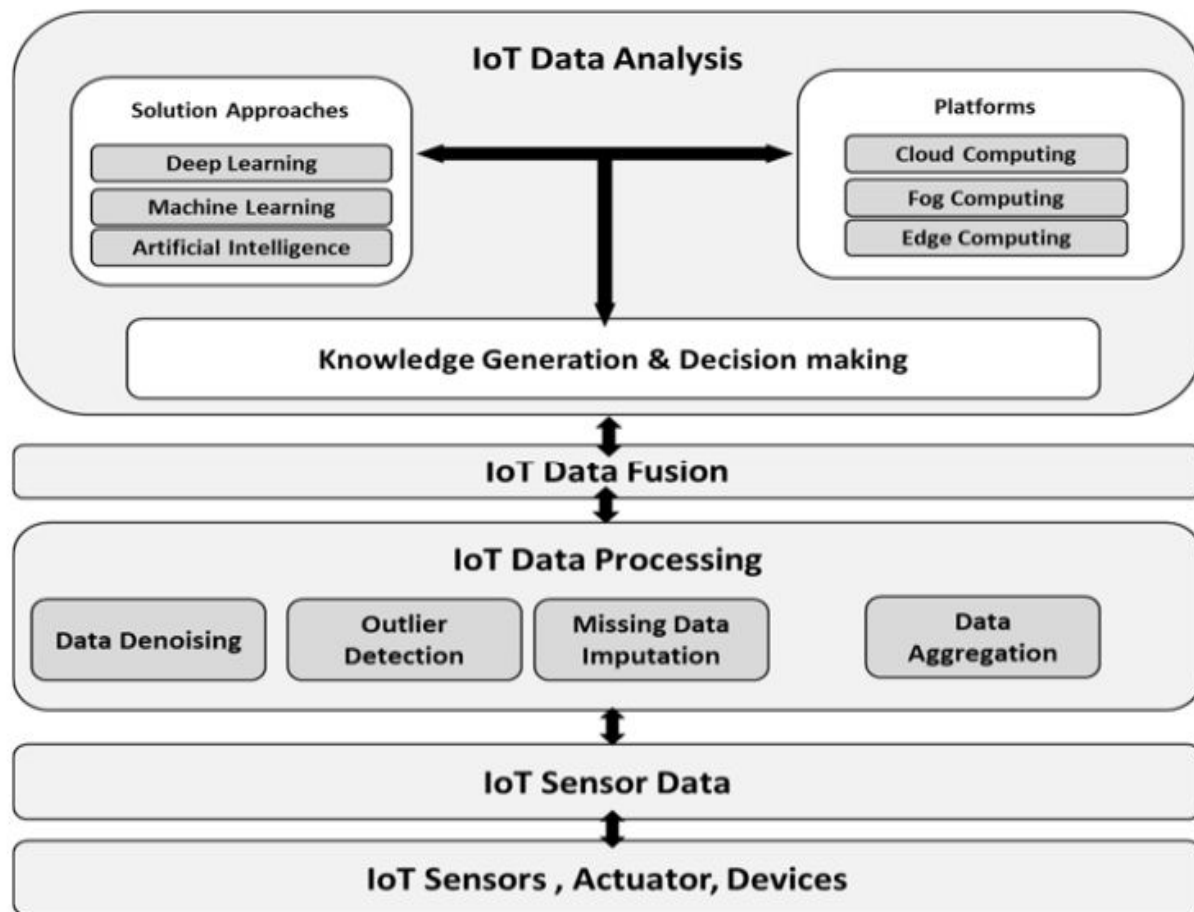
## IoT Data Input to ML Models (Training vs. Inference)

# AI Training & Inference

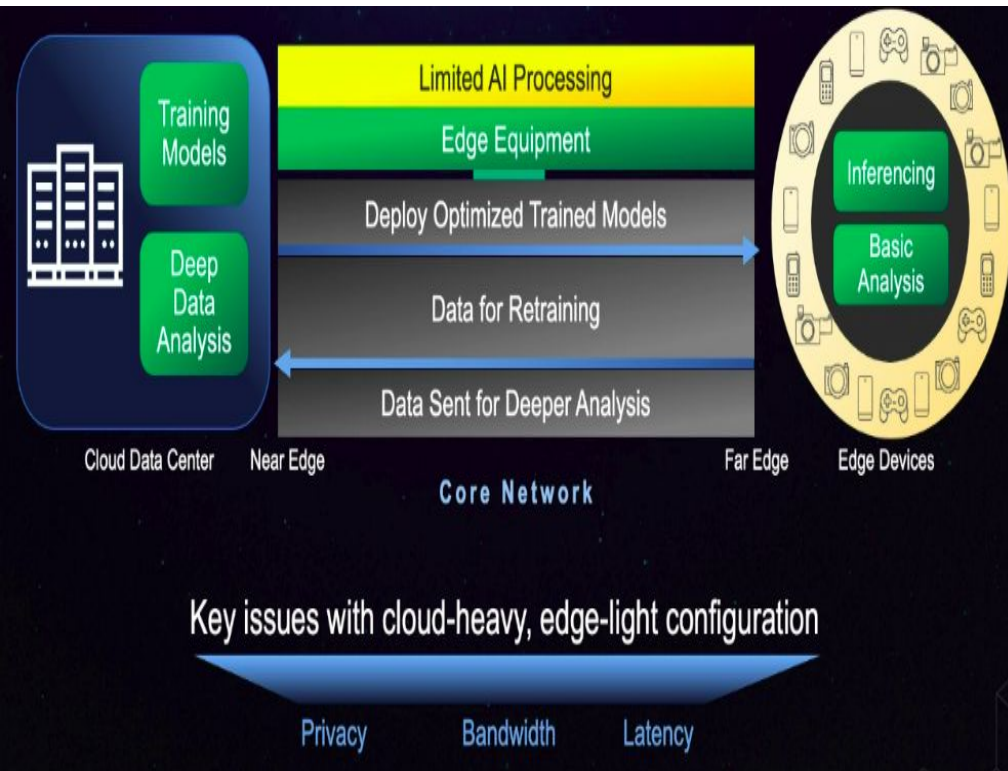




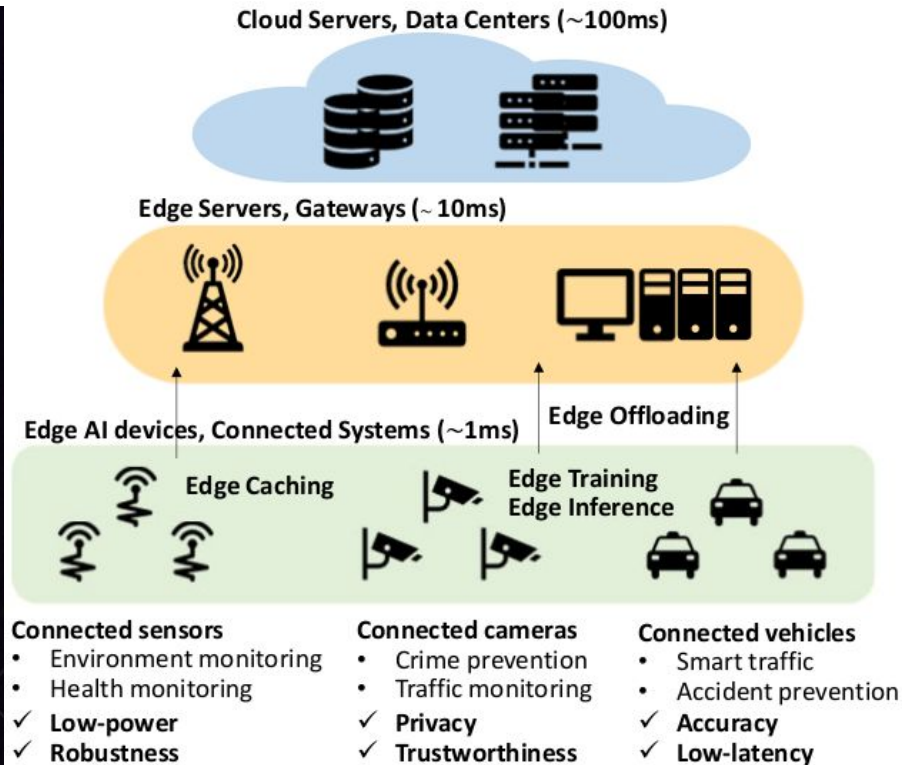
# AI use case: IoT Framework



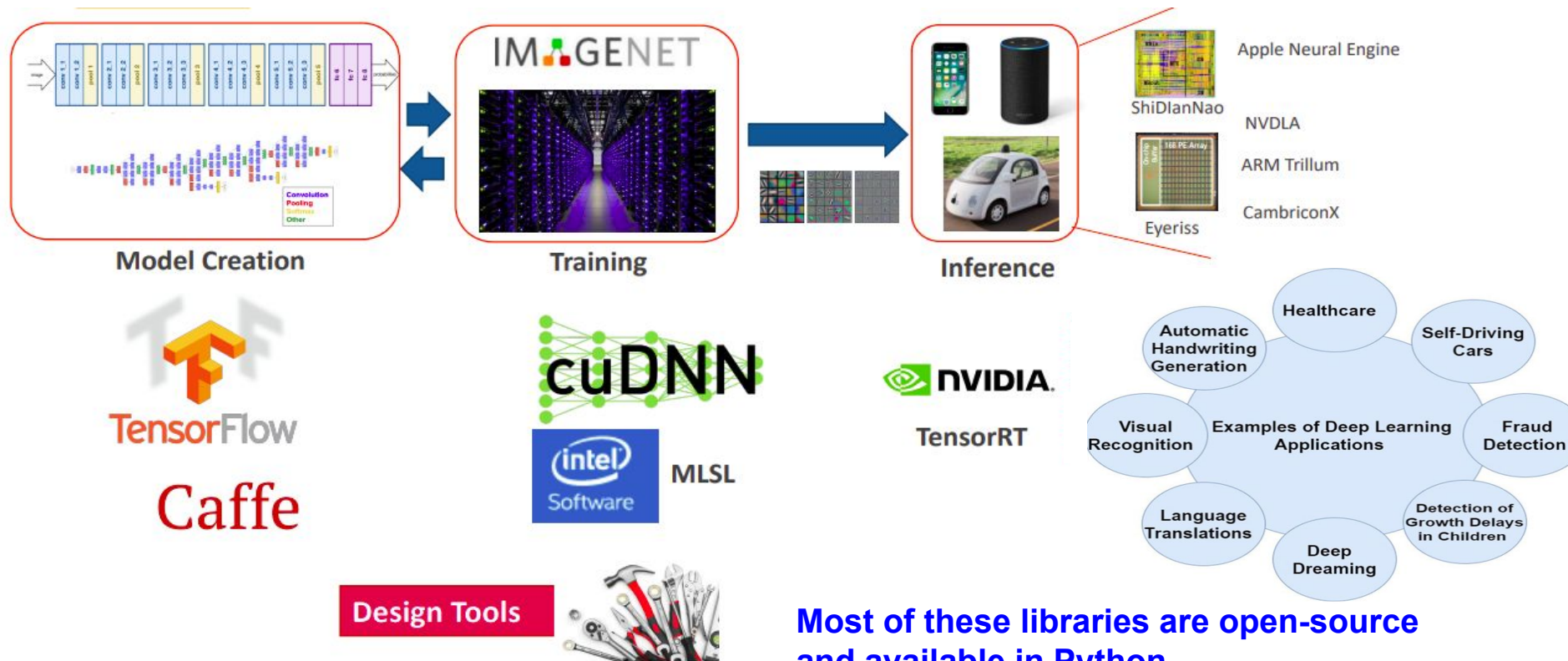
# AI Training & Inference at Different Levels



Source: CadenceLive 2021

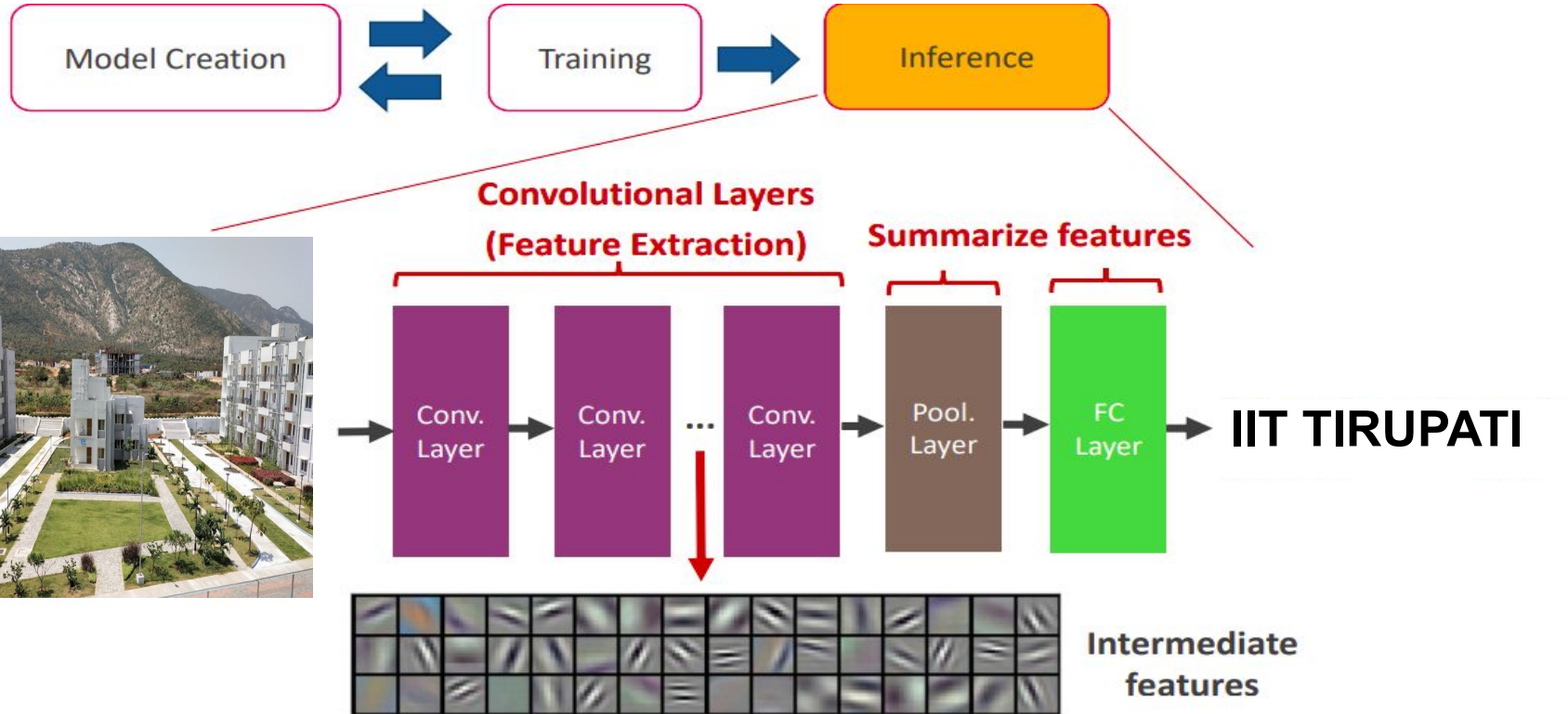


# Deep Learning Landscape

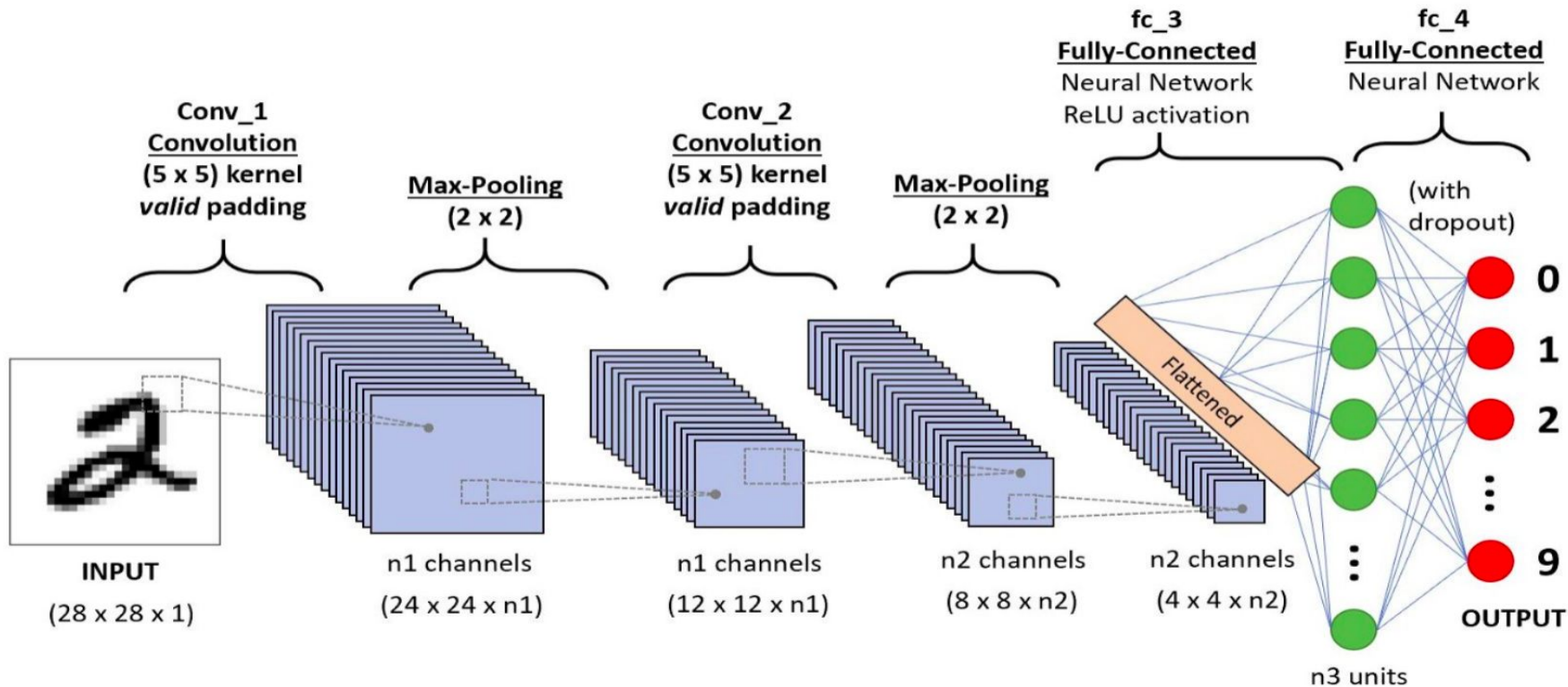




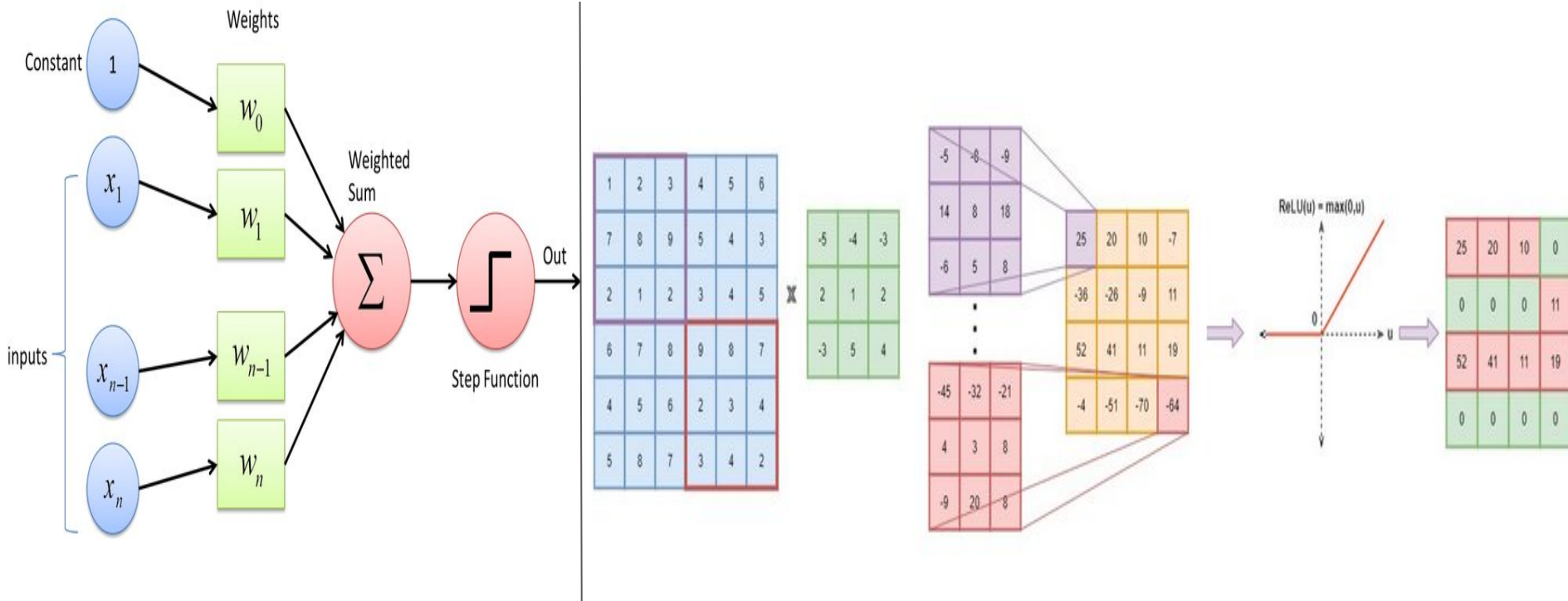
# AI Task Example: Image Recognition



# Convolution Neural Network (CNN)



# Understanding Convolution Neural Network Operations



**A large number of computations are involved!**

# Understanding Deep Learning Operations

0	0	0	0	0	0	...
0	156	155	156	158	158	...
0	153	154	157	159	159	...
0	149	151	155	158	159	...
0	146	146	149	153	158	...
0	145	143	143	148	158	...
...	...	...	...	...	...	...

Input Channel R

0	0	0	0	0	0	...
0	156	155	156	158	158	...
0	153	154	157	159	159	...
0	149	151	155	158	159	...
0	146	146	149	153	158	...
0	145	143	143	148	158	...
...	...	...	...	...	...	...

Input Channel G

0	0	0	0	0	0	...
0	156	155	156	158	158	...
0	153	154	157	159	159	...
0	149	151	155	158	159	...
0	146	146	149	153	158	...
0	145	143	143	148	158	...
...	...	...	...	...	...	...

Input Channel B

-0.3	-0.7	1
0	1.8	-1
0	1	1

Filter Channel #1

-1	-1	1
0	1	-1
0	1	1

Filter Channel #2

-1	-1	1
0	1	-1
0	1	1

Filter Channel #3

$$\mathbf{F} \begin{pmatrix} k_1 \\ k_2 \\ k_3 \end{pmatrix} + b = p_{ij}$$

$p_{ij} \rightarrow$

2	4	3	6	...
6	8	9	4	...
2	5	6	9	...
7	0	10	7	...
...	...	...	...	...

Feature Map  $a_{i,l}$



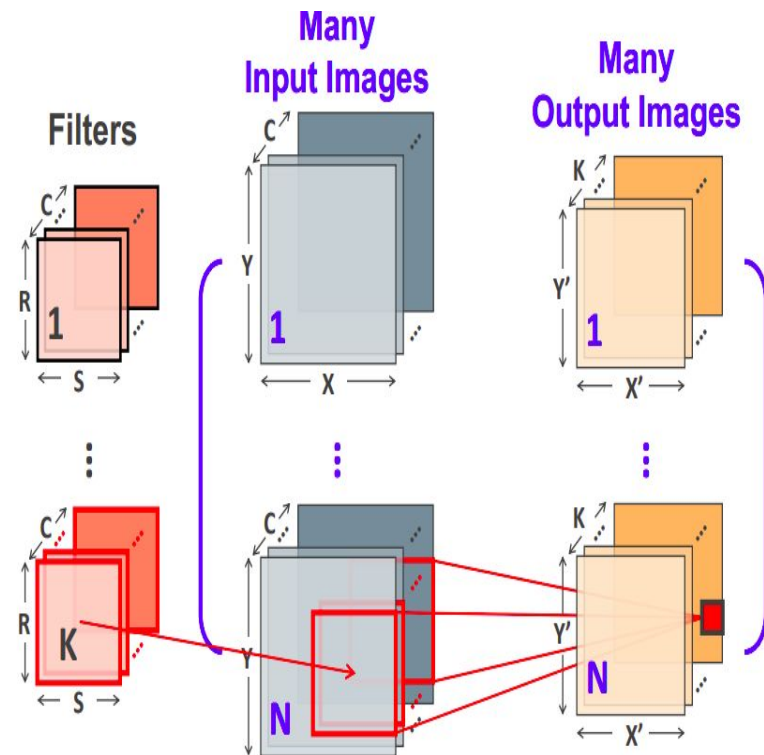
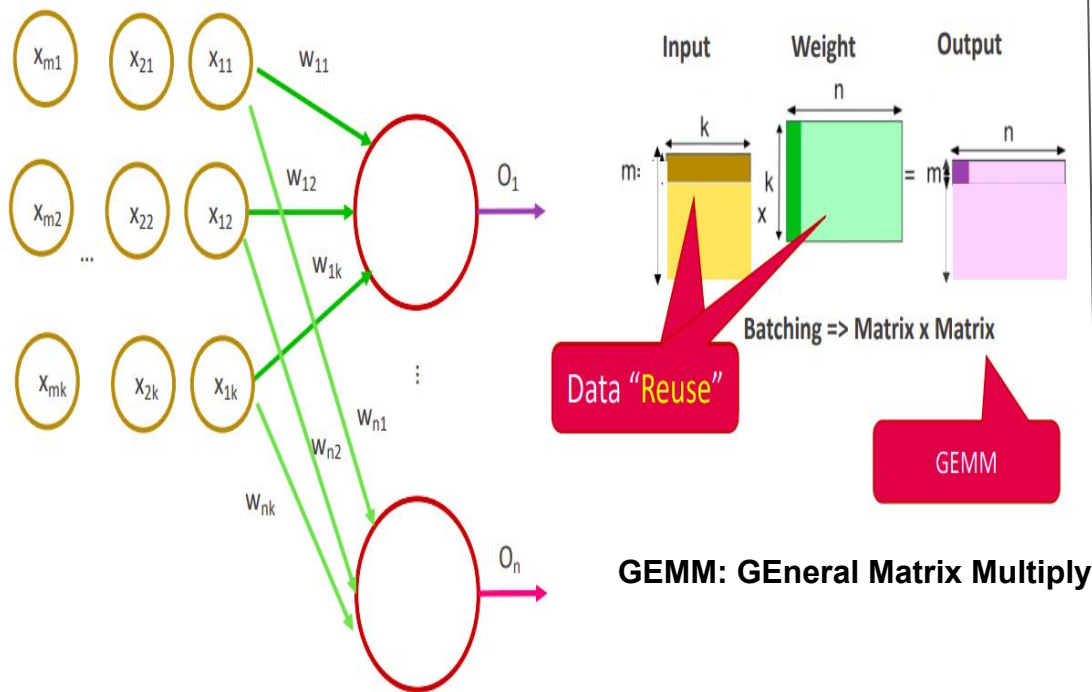
2\*2 filter



8	9	...
7	10	...
...	...	...

Pooled Feature Map

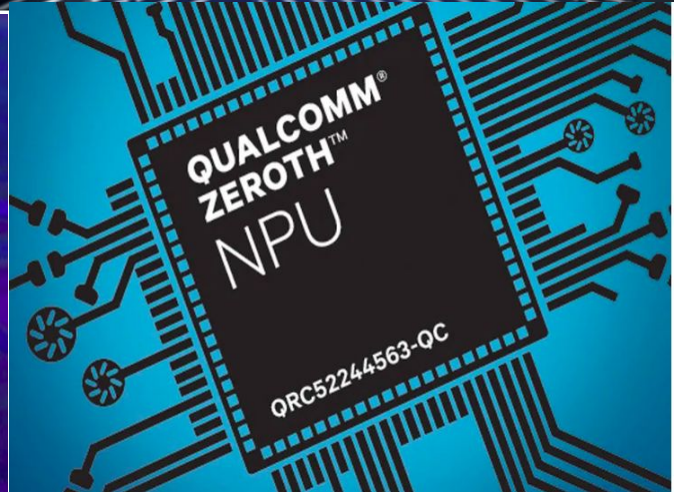
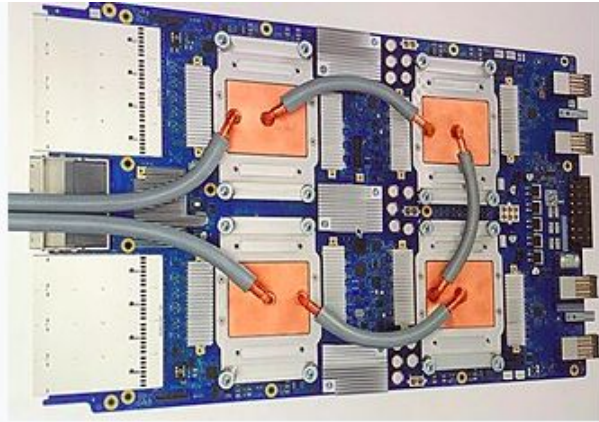
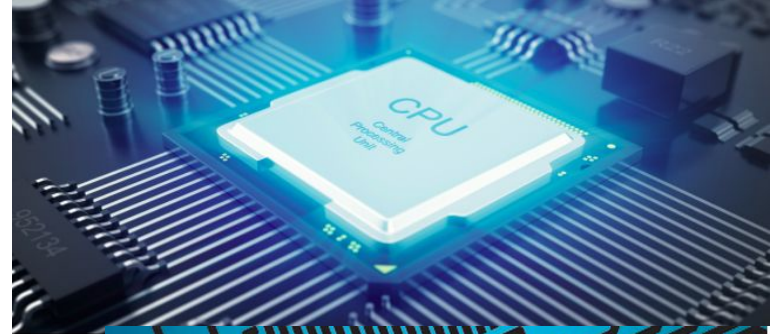
# Computations in a CNN: Linear Algebra





# Different Hardware Platforms

1. *CPU (Central Processing Units)*
2. *GPU (Graphics Processing Units)*
3. *TPU (Tensor Processing Units)*
4. *NPU (Neural Processing Units)*
5. *FPGAs*

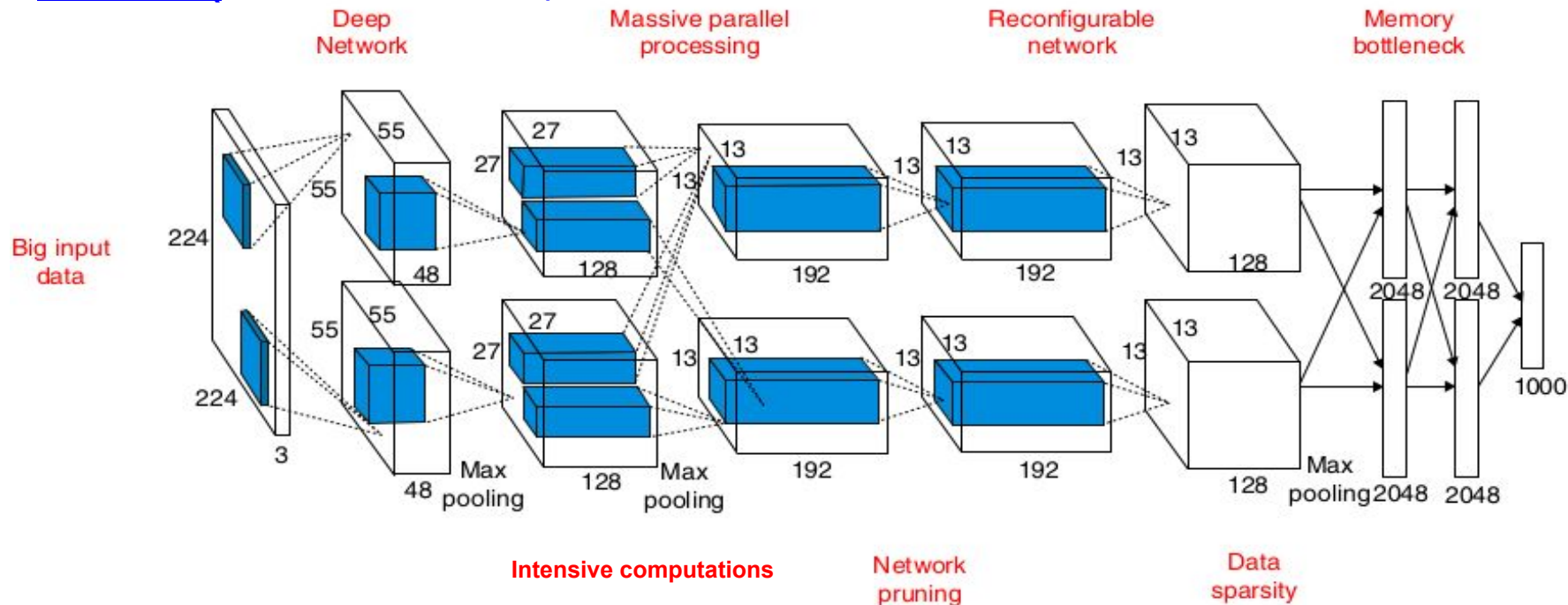


# Comparative Evaluation of Different Devices

	CPU-based	GPU-based	FPGA-based	ASIC-based
Advantages	Good versatility Lowest price Multitasking High programmability	Medium versatility Massive parallelism Moderate programmability	Customized designs Low latency High performance/watt	Extremely low power Highest performance
Limitations	Limited parallelism	Power hungry	Limited on-chip memory Requires design expertise	High development cost Long Time-To-Market Low flexibility
Example Devices	Arm Cortex-M Series Raspberry Pi Series NanoPi Series Sipeed MAIX Series	Nvidia Jetson Series AMD Ryzen Family Arm Mali GPUs	Xilinx Zynq FPGAs Intel Arria 10 FPGAs Lattice iCE40 FPGAs	Google Edge TPU Ascend 310 processor In-memory chips Neuromorphic chips
Development Tools	Arm NN TensorFlow Lite	TensorRT Intel OpenVino	Intel OpenVino Xilinx Edge AI platform	Apache TVM

# Computation Challenges: DL Hardware

## AlexNet (ILSVRC 2012 Winner)



# Challenges with DNN Computations

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

- Billions of computations → **Need lots of parallel computations**

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

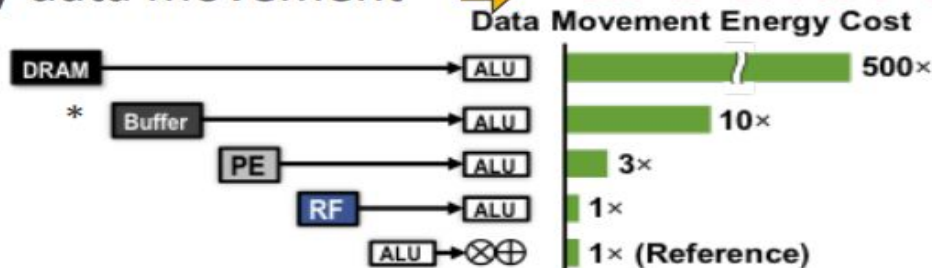
This makes CPUs inefficient

PE: Processing Element

RF: Register File

DRAM: Dynamic Random Access Memory


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



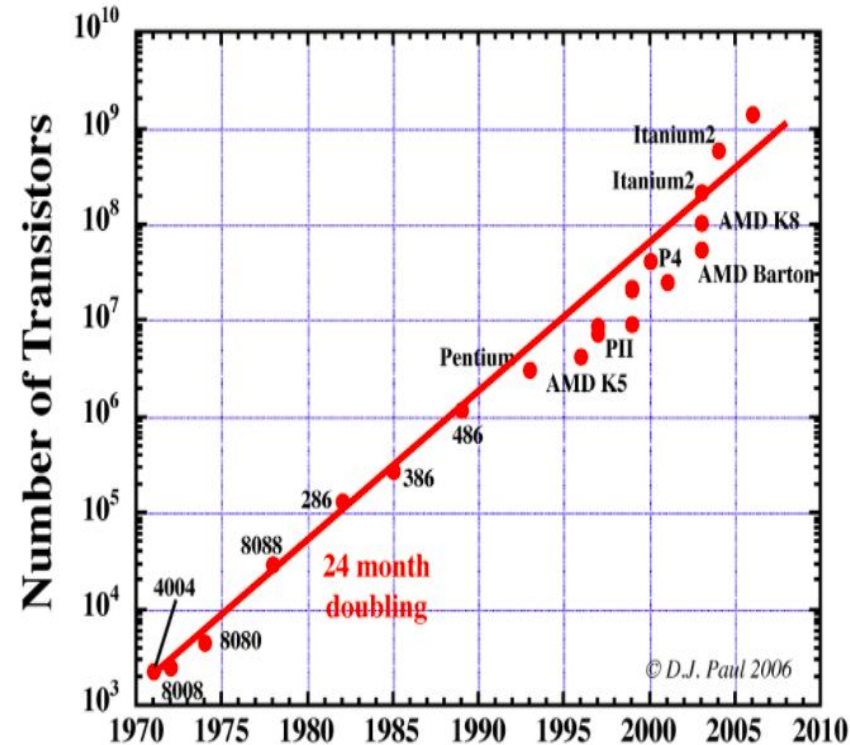
This makes GPUs inefficient



# Technology challenges: There are 4 types of people!

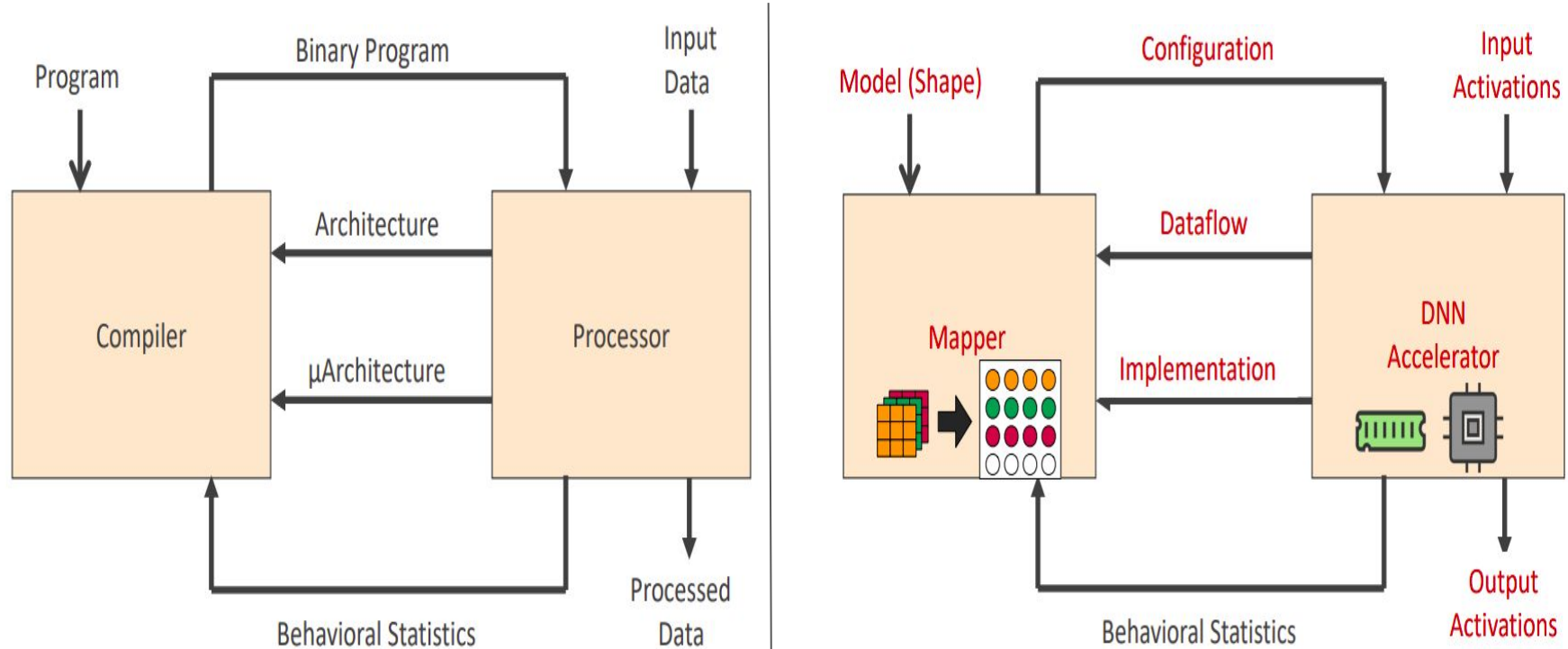


- Those who say Moore's dead → find another type of job .
- Those who say Moore's dead → find another type of transistor
- Those who say Moore's dead → use other technologies (photonics, spintronics, quantum computing, ...)
- Those who say Moore's dead → now our job has begun

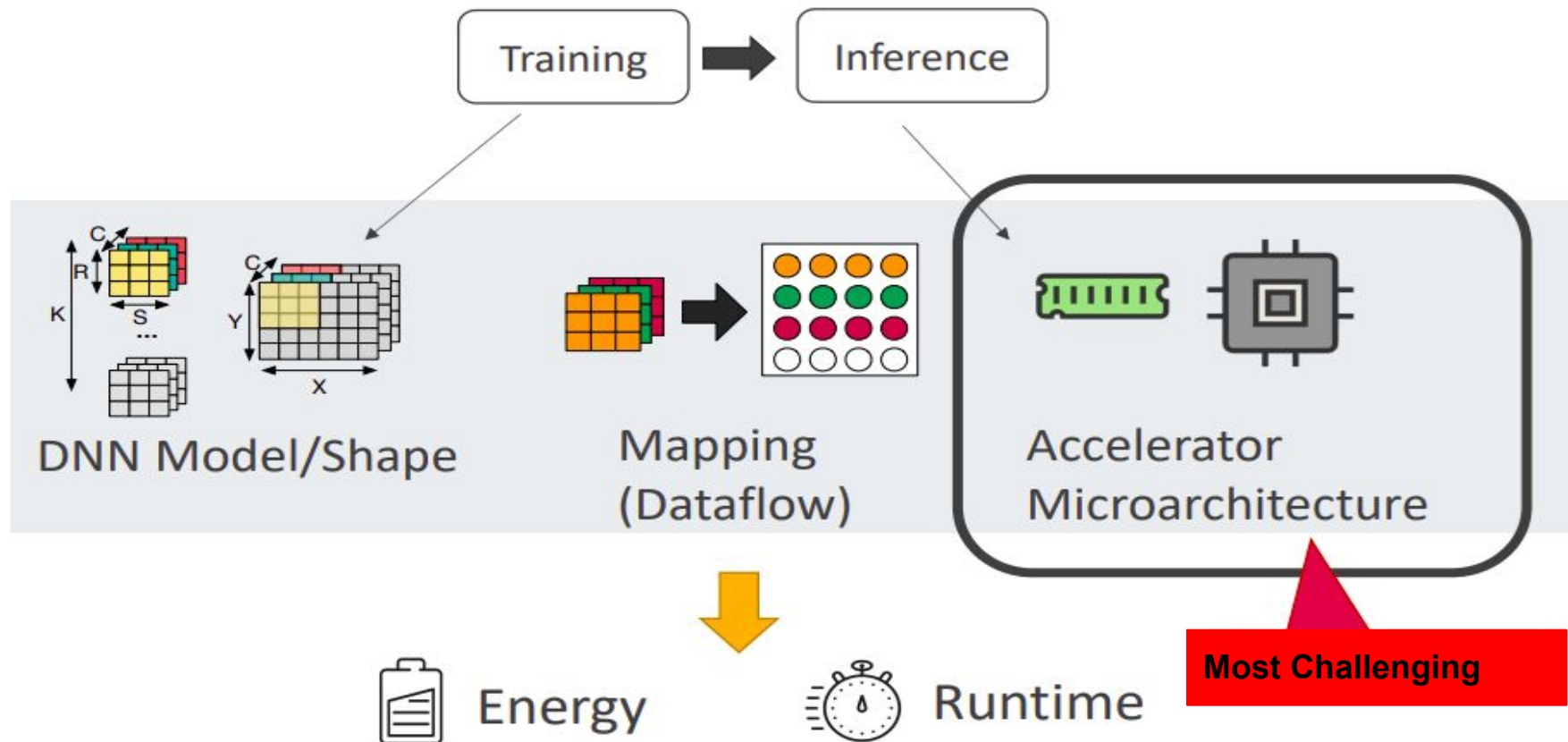




# Compute Architectures: CPU v/s Accelerator



# Challenges in AI Hardware Design & Deployment

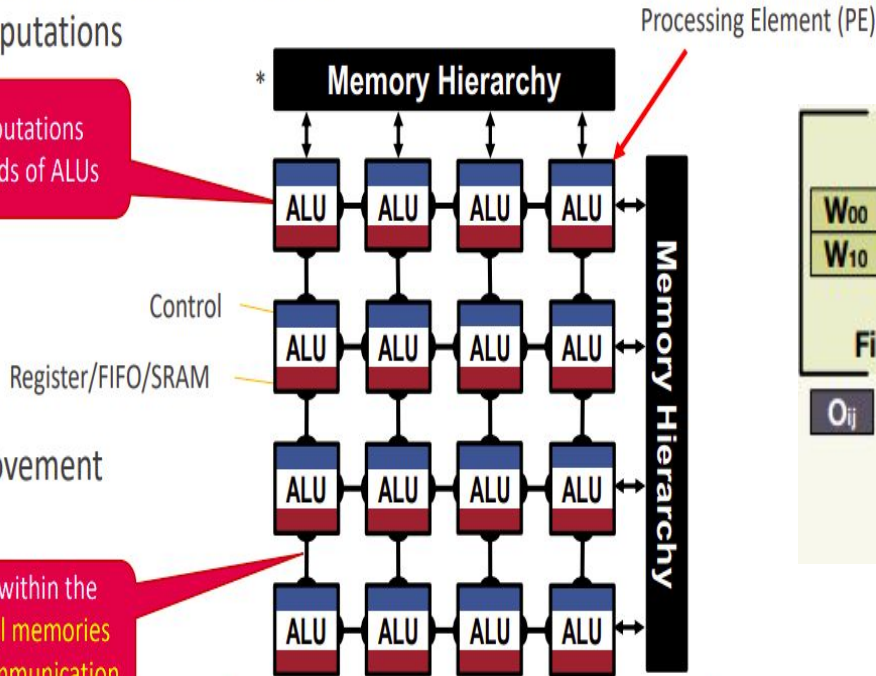


# Customised Architectures (Accelerators)

- Millions of Parameters (i.e., weights)

- Billions of computations

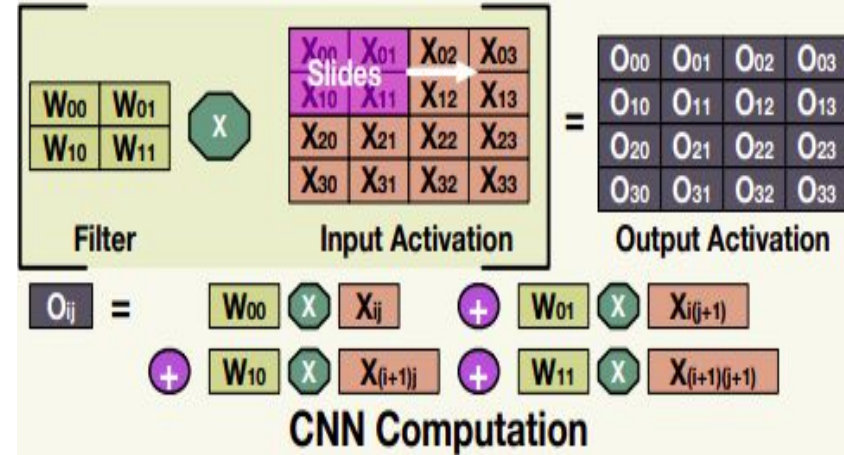
Spread computations  
across hundreds of ALUs



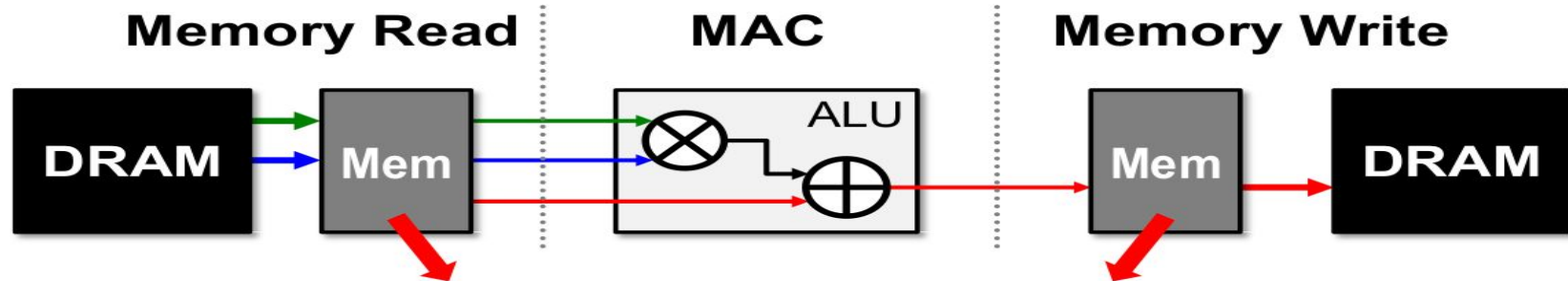
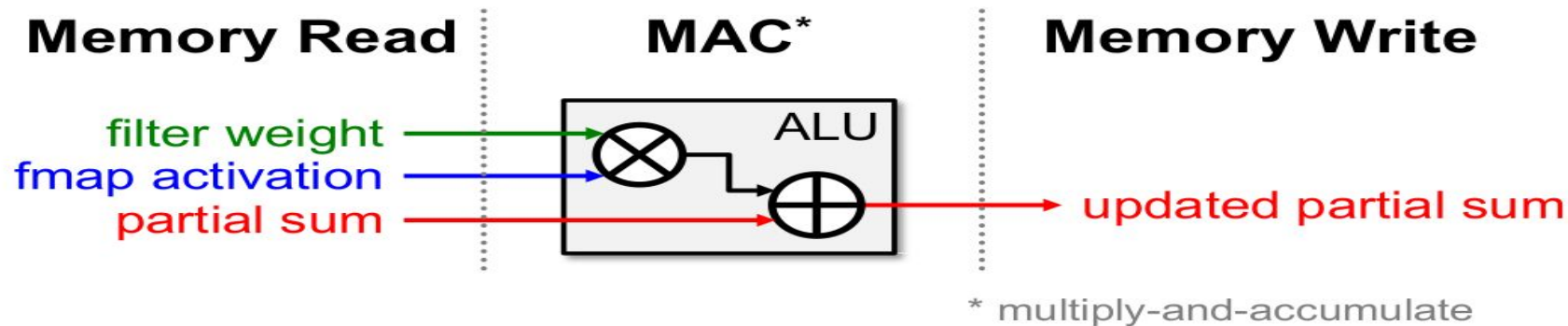
- Heavy data movement

Reuse data within the  
array via **local memories**  
and **direct communication**

Examples: MIT Eyeriss, Google TPU, ...

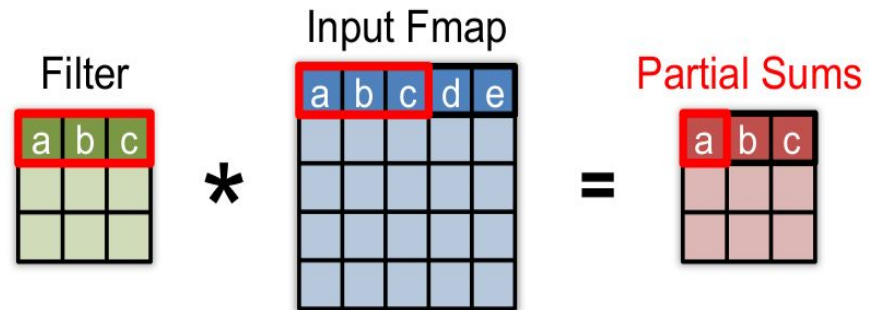
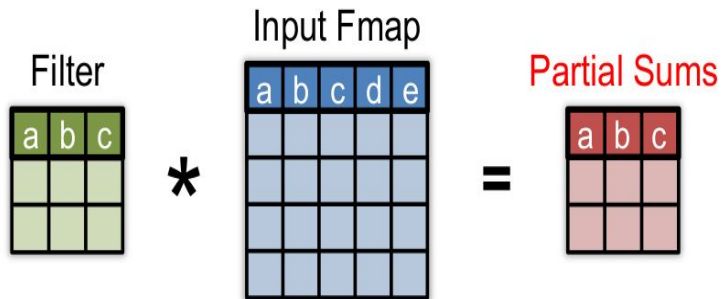


# Computations in DNN/CNN

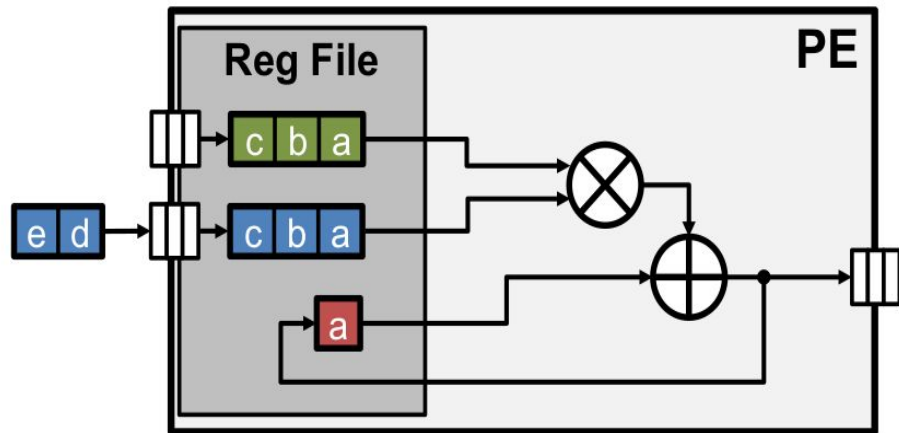
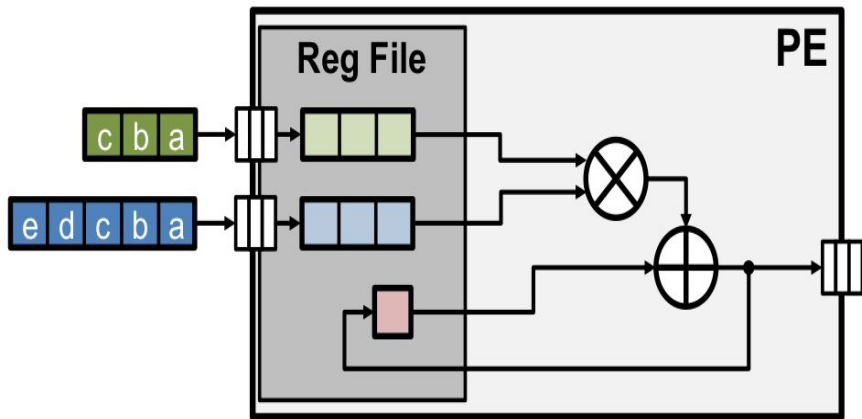


Extra levels of local memory hierarchy  
**Smaller, but Faster and more Energy-Efficient**

# Implementing Convolution

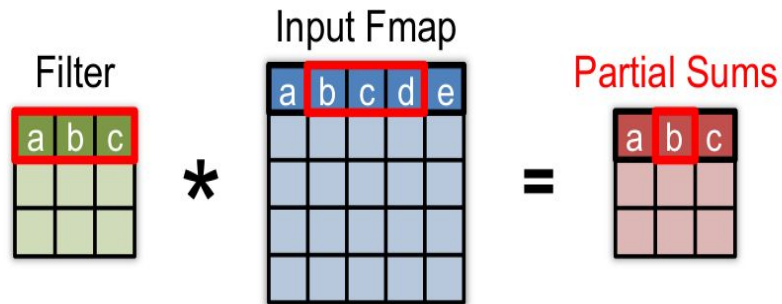


Step 1

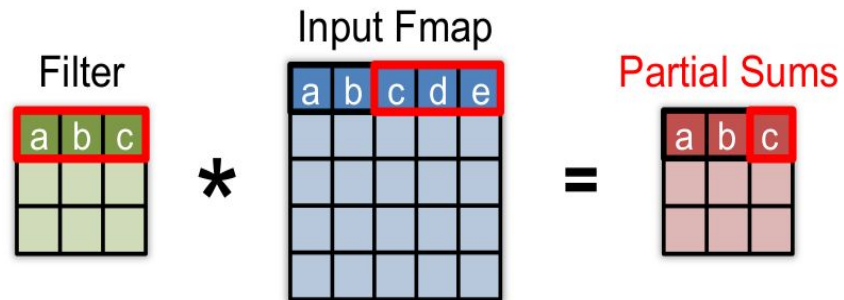
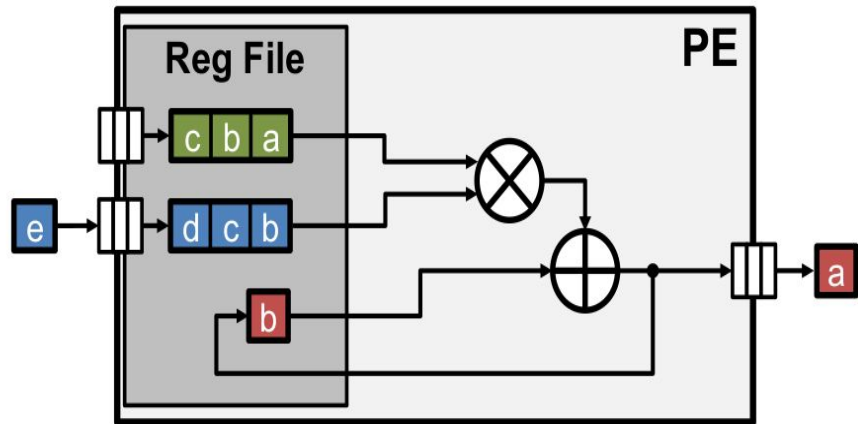




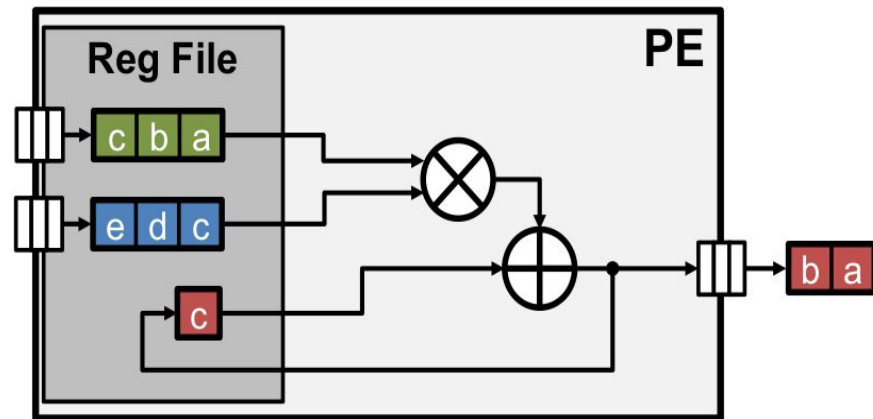
# Implementing Convolution



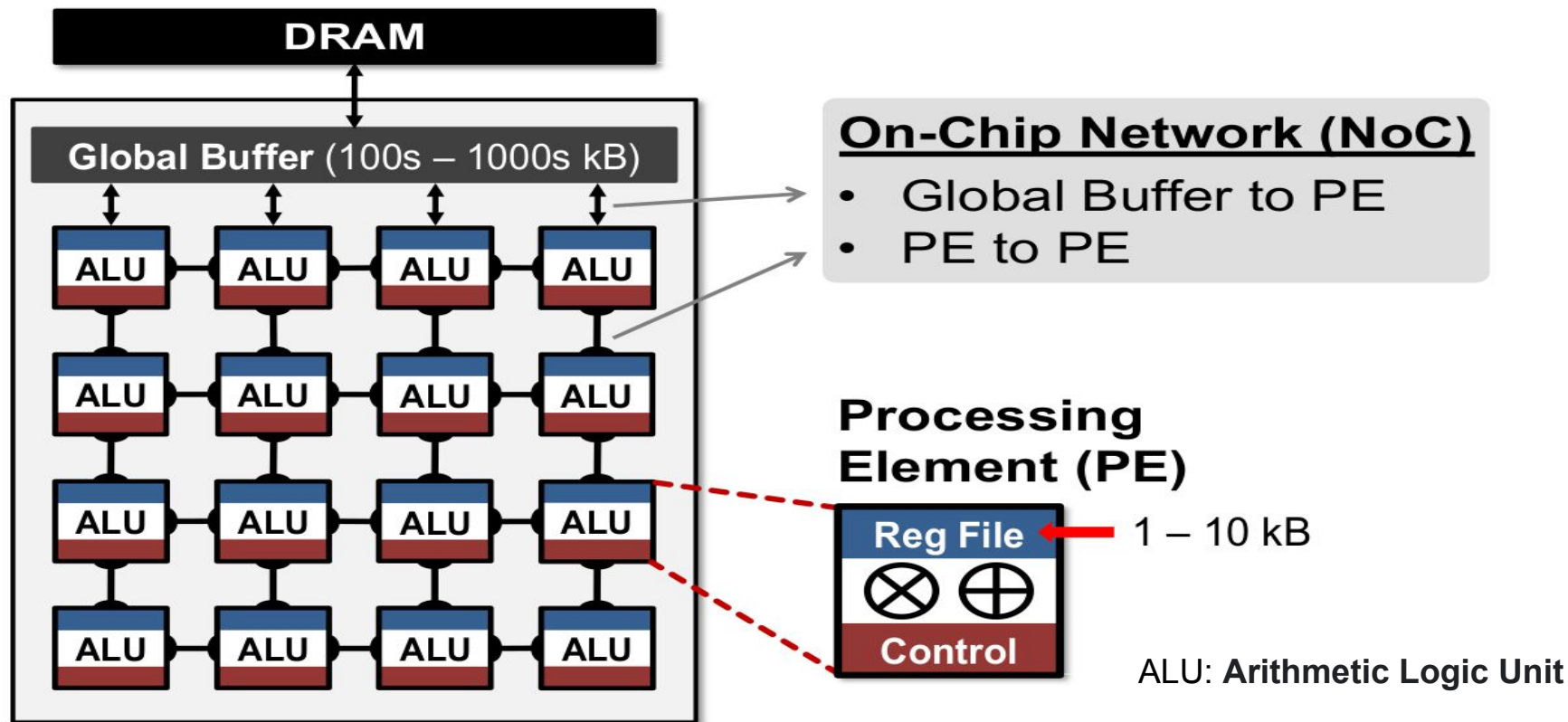
Step 2



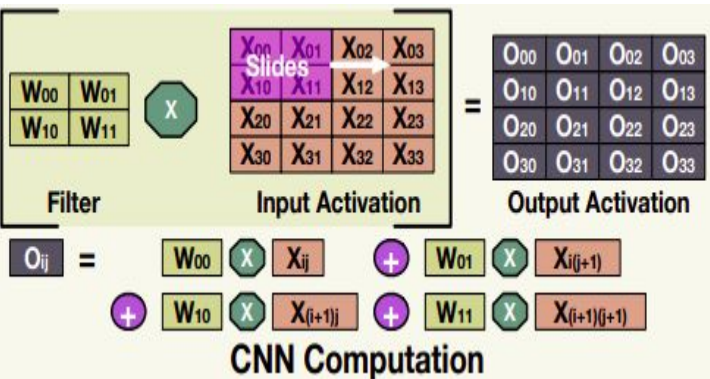
Step 3



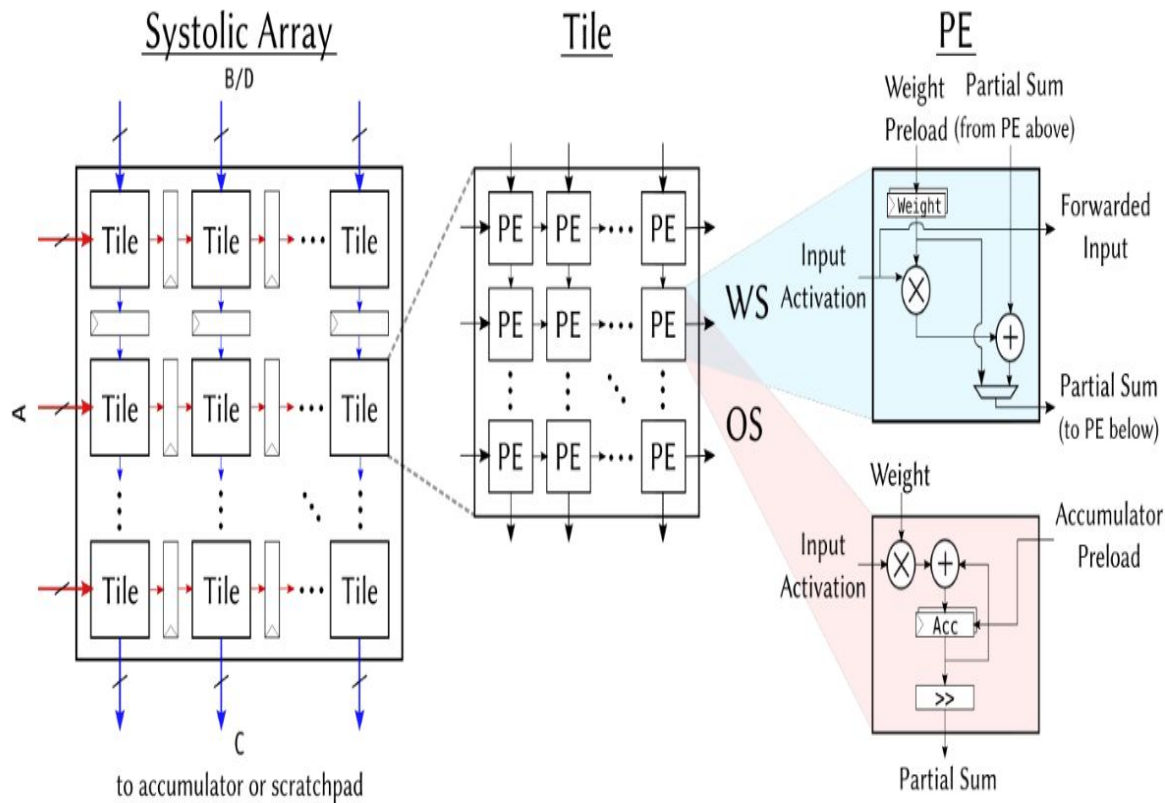
# Accelerators for Computations in DNN/CNN



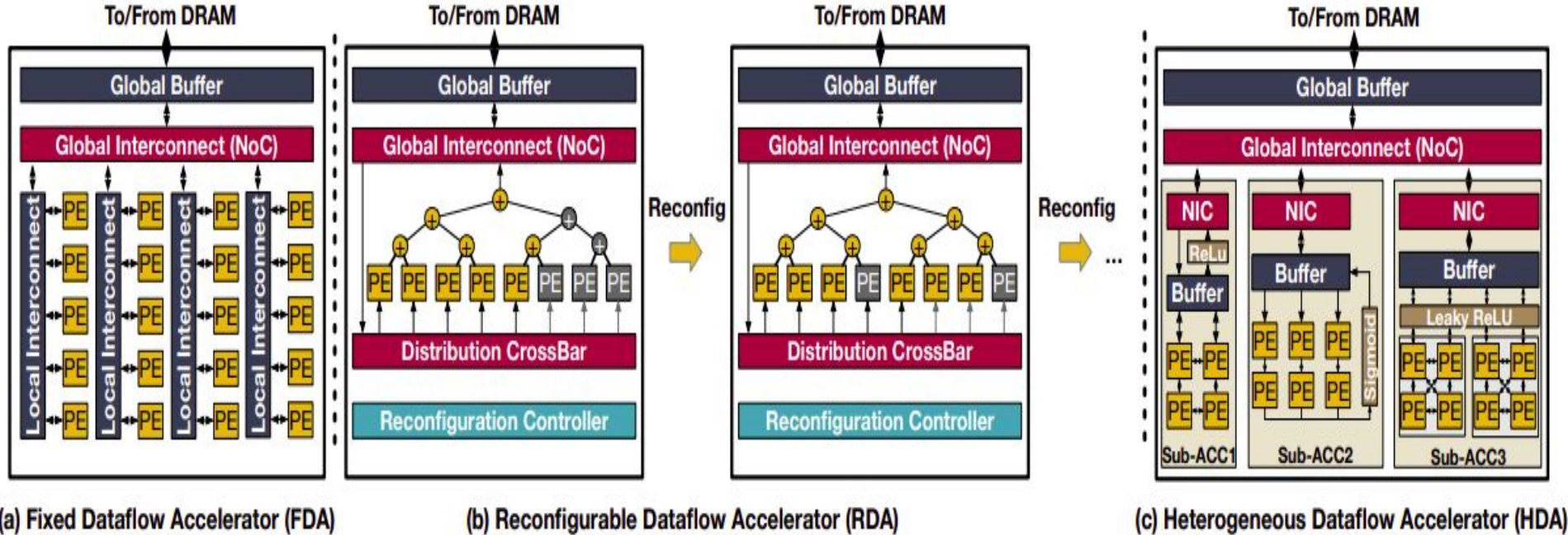
# Different Dataflow Architectures



***Each dataflow architecture has a set of benefits in terms of throughput and resource utilization.***

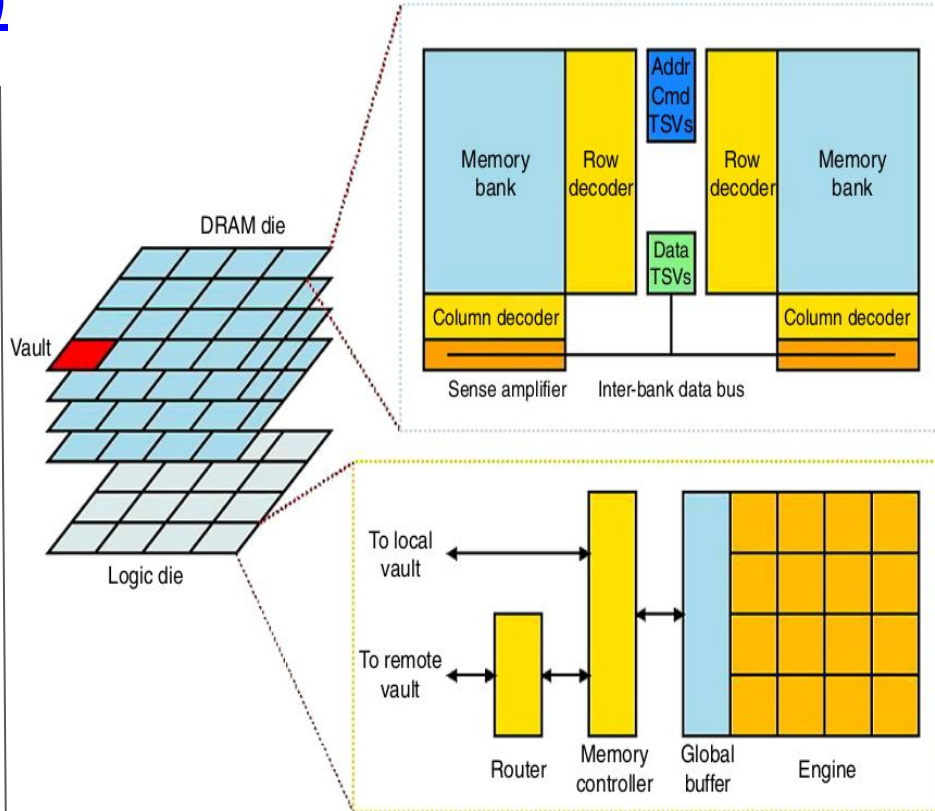
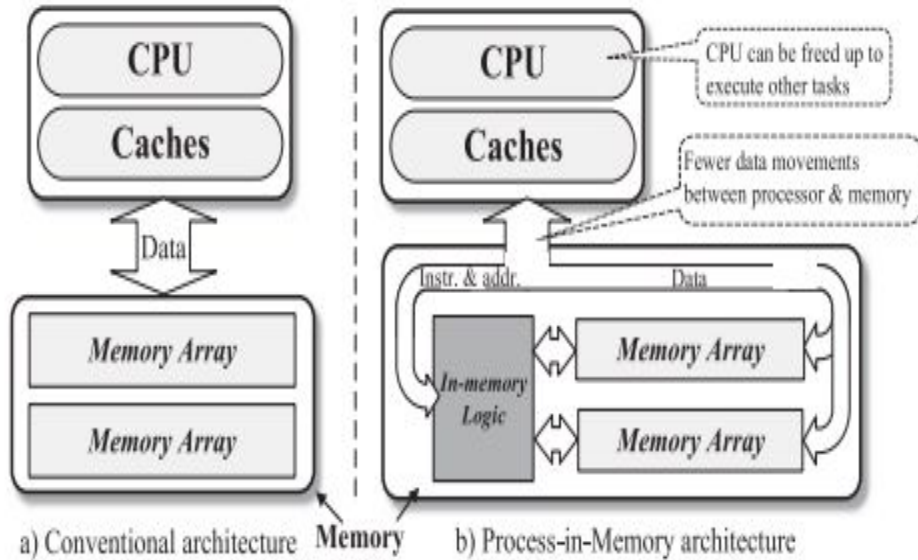


# Different Dataflow Architectures



# Processing-In-Memory (PIM)

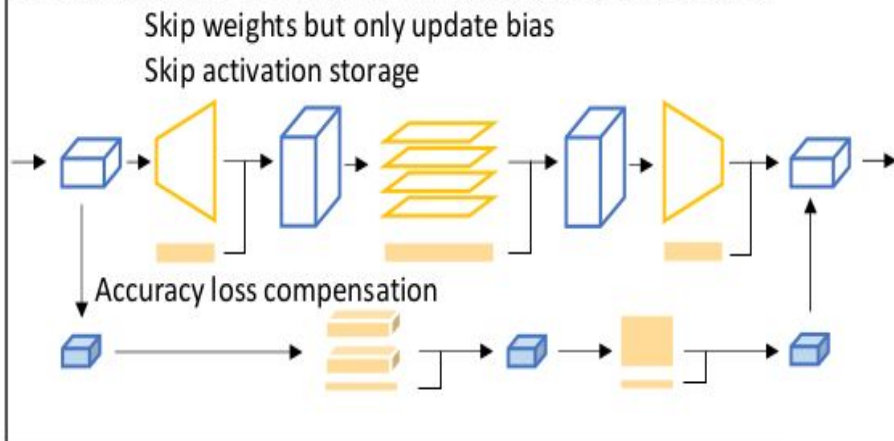
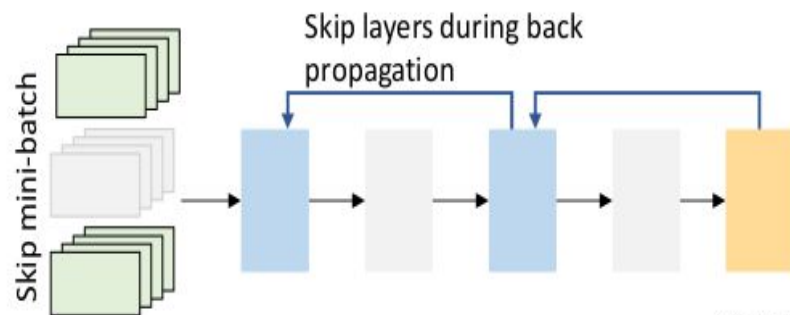
## Novel architectural considerations:



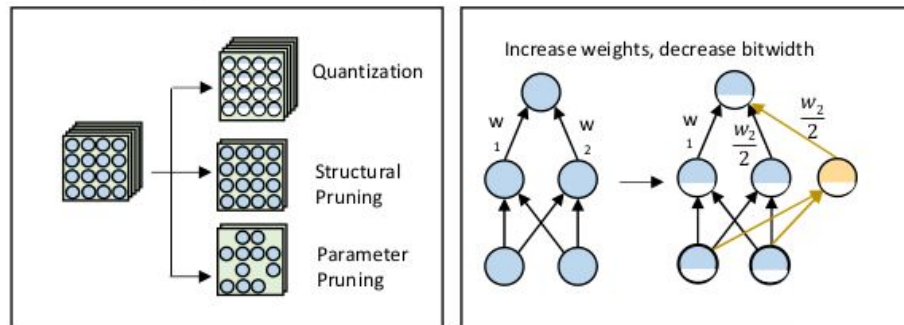


# Specialized Techniques

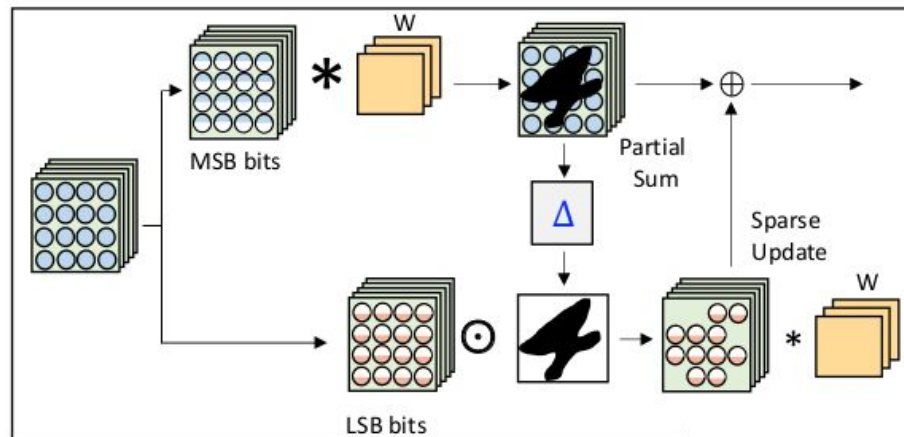
## On-device Training



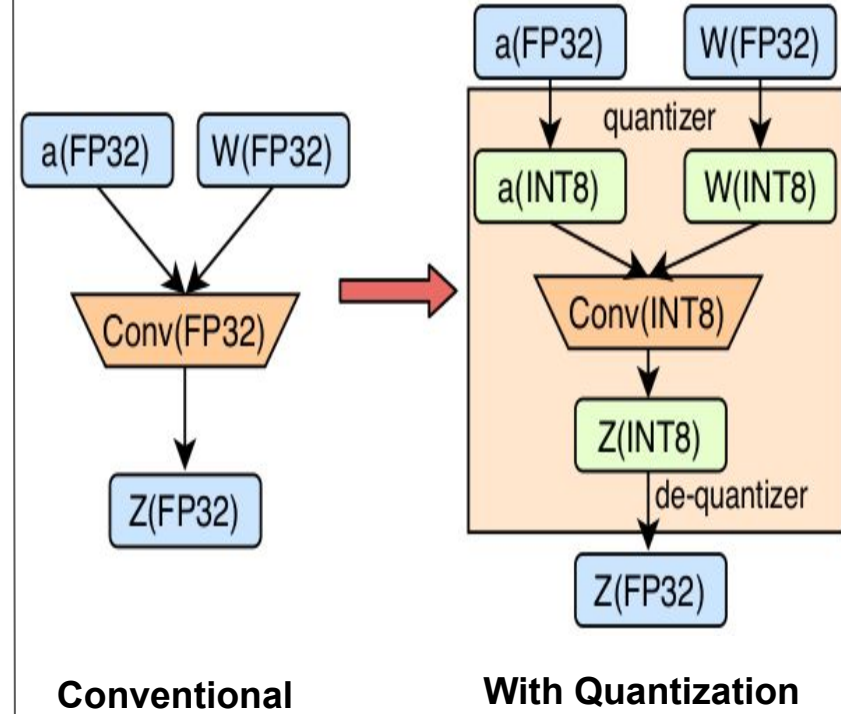
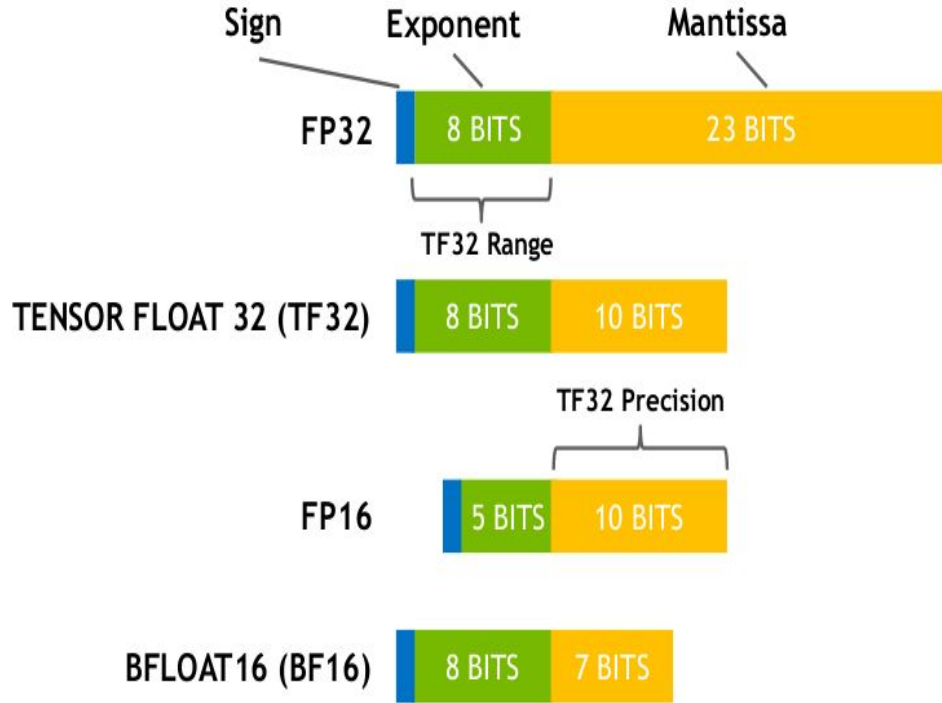
## Model Compression



## Adaptive Inference

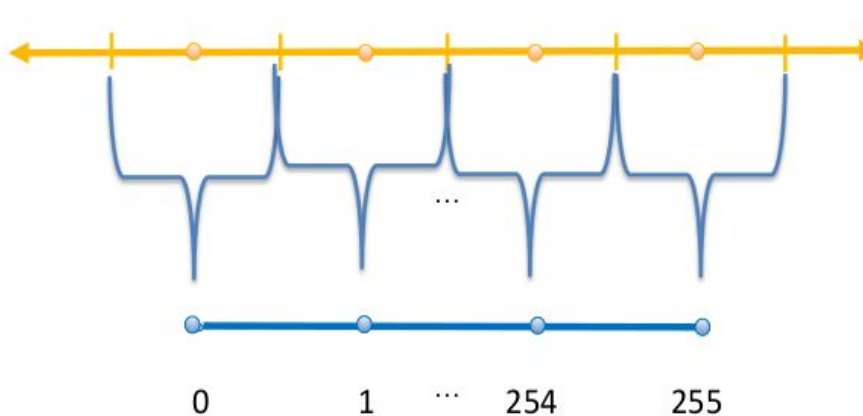


# Neural Network Quantization



# Quantization Example

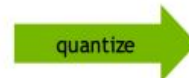
Floating Point Values



8-bit Quantized Values

0.34	3.75	5.64
1.12	2.7	-0.9
-4.7	0.68	1.43

FP32  
(pre-quantized)

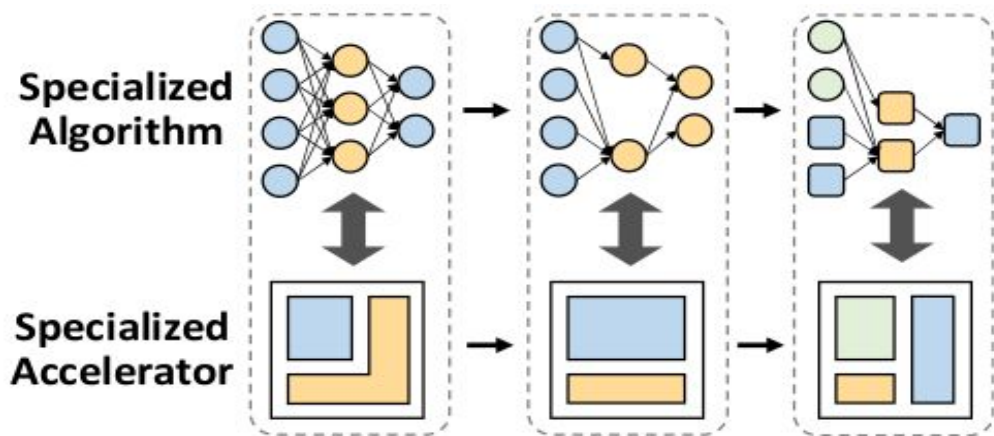


64	134	217
76	119	21
3	81	99

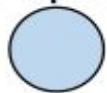
INT8  
(quantized)

**How can we find the right trade-off for a given latency/power constraint for a given hardware platform?**

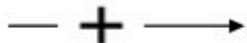
# Co-design Methodology



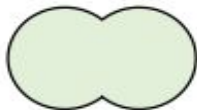
Neural Architecture  
Search Space  $\{A\}$



Implementation  
Search Space  $\{I\}$



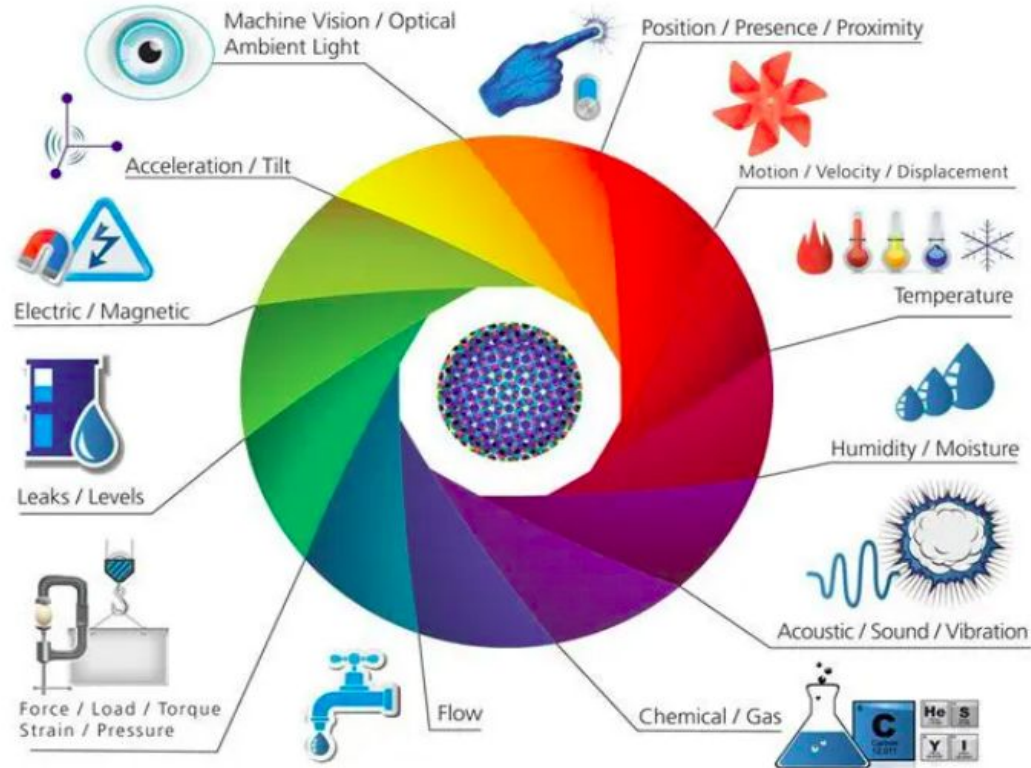
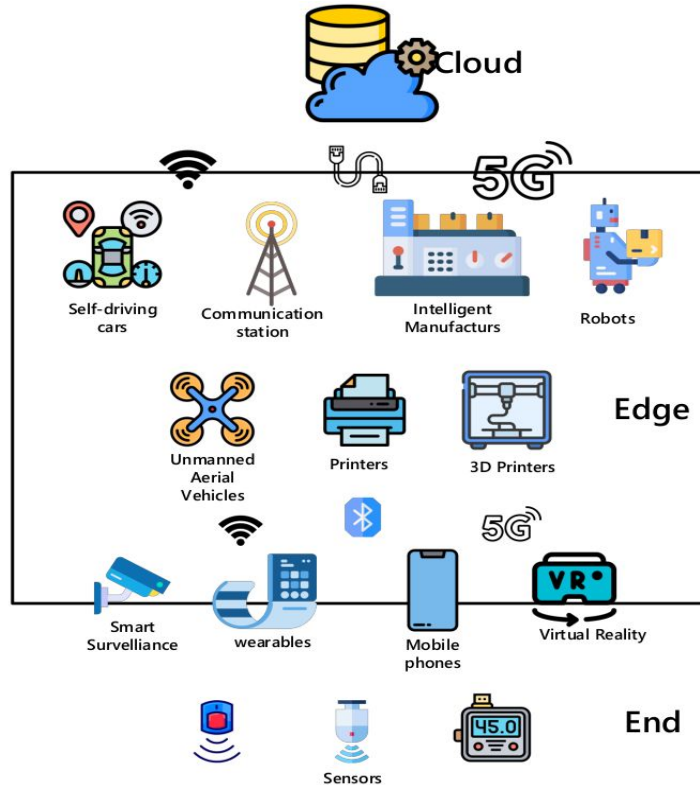
Merged Neural Architecture and  
Implementation Search Space  $\{A, I\}$



- Formulate  $\{A, I\}$  as differentiable
- Solve  $\{A, I\}$  simultaneously using Gradient Descent

- **Data reuse** is the key to achieving **high energy efficiency**.
- High PE utilization with **adaptive on-chip networks** is the key to achieving **high performance**
- Co-design of **dataflow** and **hardware** is critical for the optimization of **performance**, **energy efficiency** and **flexibility** for DNN accelerators.

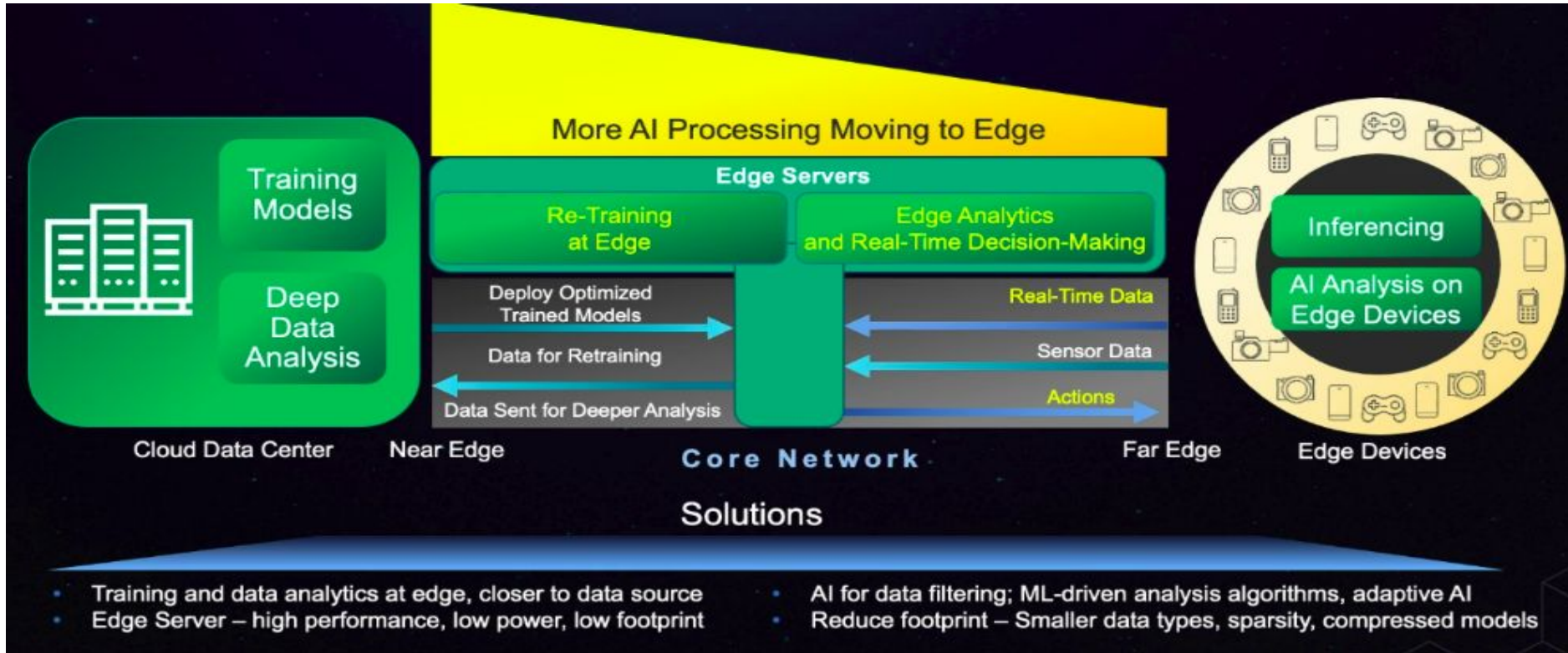
# AI Moving Close to Edge



Sensors are everywhere!



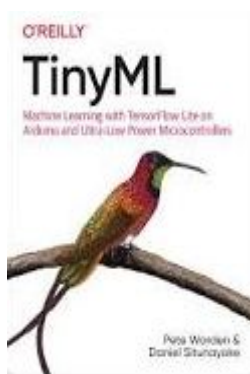
# AI Moving Close to Edge



# Power Concerns



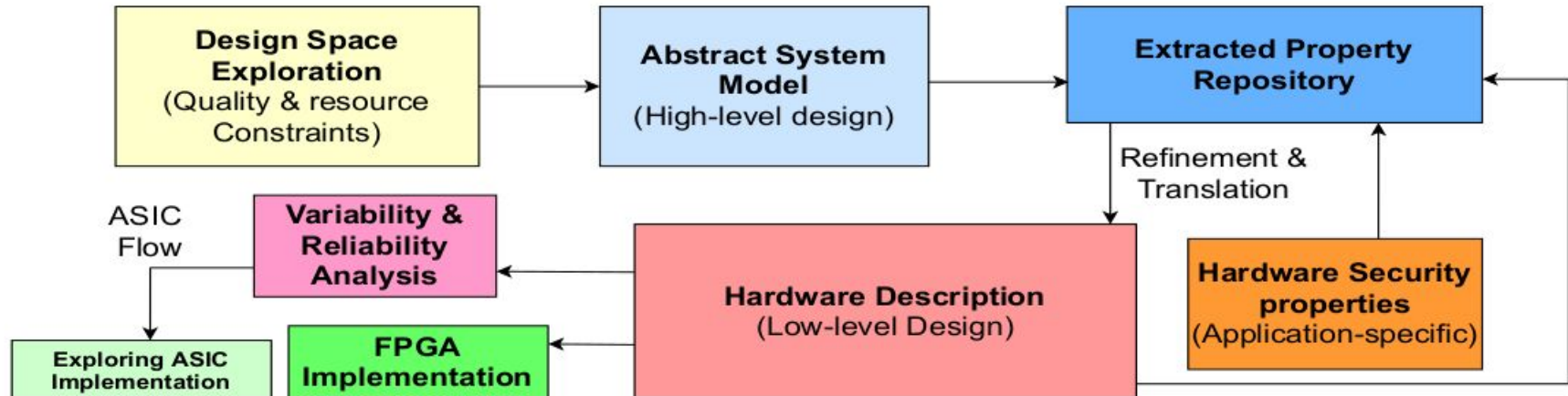
<https://www.tinyml.org/>



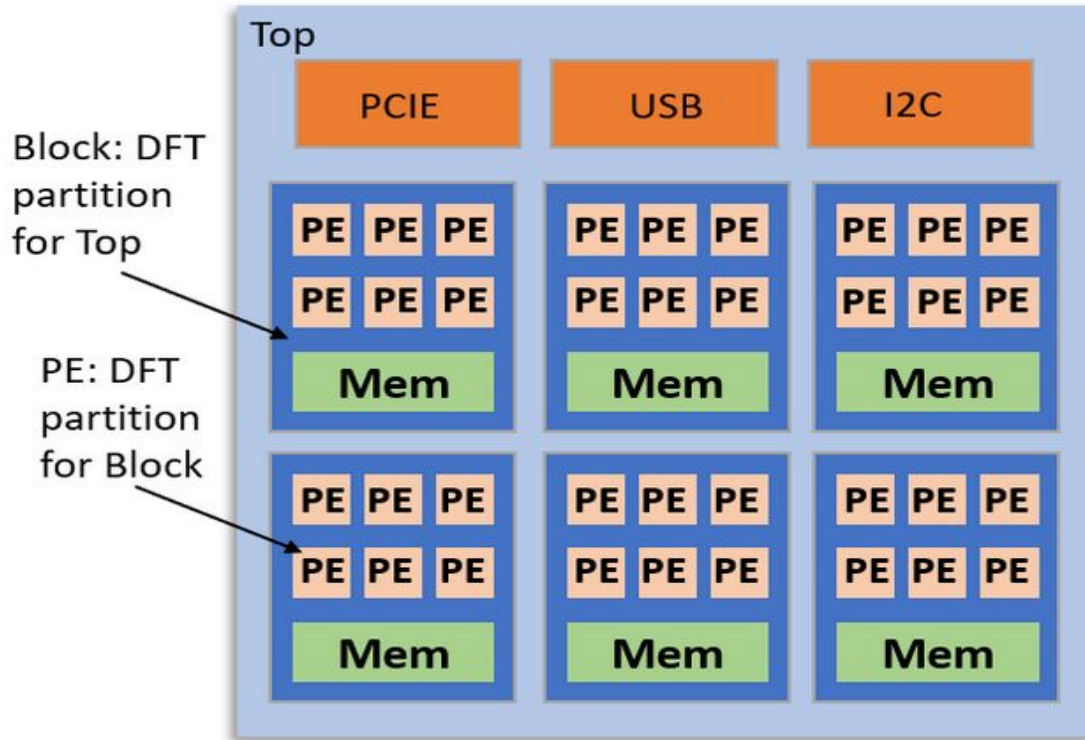
Slide: Courtesy of Xiangyu Yu, and Kurt Keutzer

# Verification of AI Designs/Architectures

- Property-driven methodology development for hardware designs
- Design size can be relatively small with mix of datapath and control mechanisms
- Verification goals: security and dependability
- Target systems: safety-critical applications (e.g. healthcare)
- **Formal guarantees to ensure correct operation under all conditions**



# Test Techniques (DFT) for AI Architectures



1. Incremental test-pattern generation (reuse of test patterns for new/optimized version of the design).
2. Core-level pattern generation and retargeting at top-level.

# Some Available Resources (Open-Source)

<b>High-level AI Accelerator Modeling</b>	<a href="https://github.com/harvard-acc/ALADDIN">https://github.com/harvard-acc/ALADDIN</a>
<b>DNN Accelerators Design Space Exploration</b>	<a href="https://github.com/ARM-software/SCALE-Sim">https://github.com/ARM-software/SCALE-Sim</a>
<b>CPU/Memory Architecture Modeling</b>	<a href="https://www.gem5.org/">https://www.gem5.org/</a>
<b>Accelerator Modelling (Eyeriss Design, MIT)</b>	<a href="https://github.com/taoyilee/clacc">https://github.com/taoyilee/clacc</a>
<b>Python-based Modeling of Hardware</b>	<a href="https://www.myhdl.org/">https://www.myhdl.org/</a>
<b>DRAM Power &amp; Energy Estimation</b>	<a href="https://www.es.ele.tue.nl/drampower/">https://www.es.ele.tue.nl/drampower/</a>
<b>Running ML Algorithms (in Python)</b>	<a href="https://scikit-learn.org/stable/">https://scikit-learn.org/stable/</a>
<b>Open-source ASIC flow</b>	<a href="https://theopenroadproject.org/">https://theopenroadproject.org/</a>
<b>RTL simulation (Verilog/VHDL)</b>	<a href="https://www.edaplayground.com/">https://www.edaplayground.com/</a>
<b>Online Platform for ML/DL Training/Inference</b>	<a href="https://colab.research.google.com/?utm_source=scs-index">https://colab.research.google.com/?utm_source=scs-index</a>
<b>Processing-in-Memory simulator</b>	<a href="https://github.com/CMU-SAFARI/ramulator-pim">https://github.com/CMU-SAFARI/ramulator-pim</a>

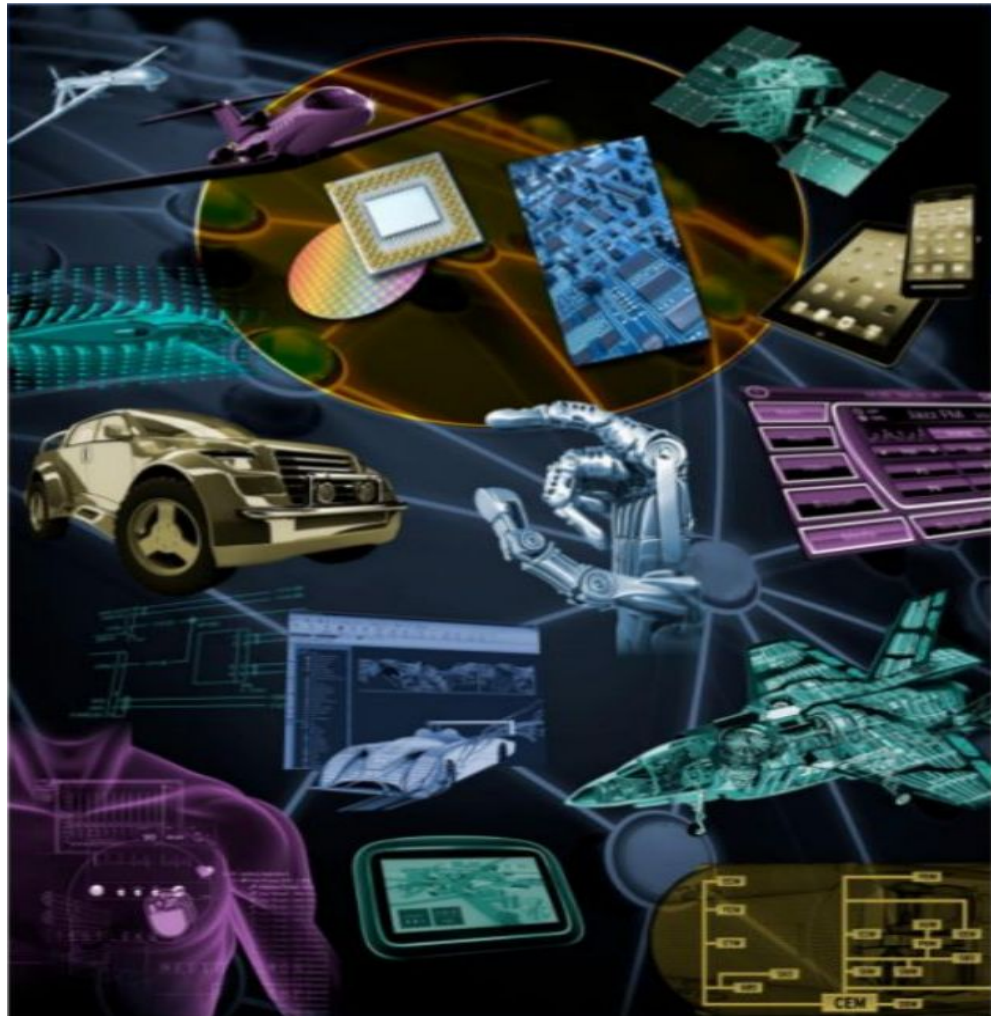


< UBUNTU = “I am  
because  
you are.” >

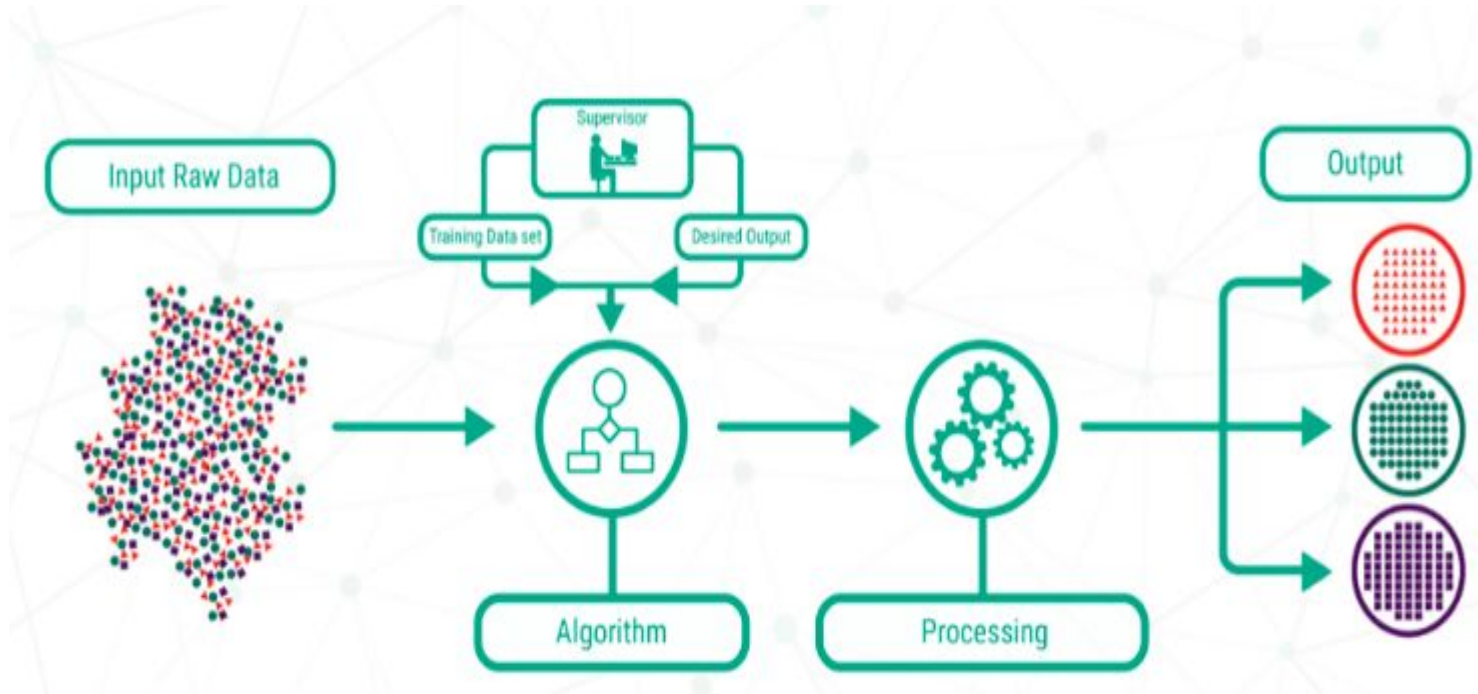


# THANK YOU

All queries be addressed to [binod@iitj.ac.in](mailto:binod@iitj.ac.in)



# Supervised Learning



# Unsupervised Learning

