# *System Driven Technology Optimization*

Kaushik Roy, Edward G. Tiedemann Distinguished Professor of ECE
kaushik@purdue.edu
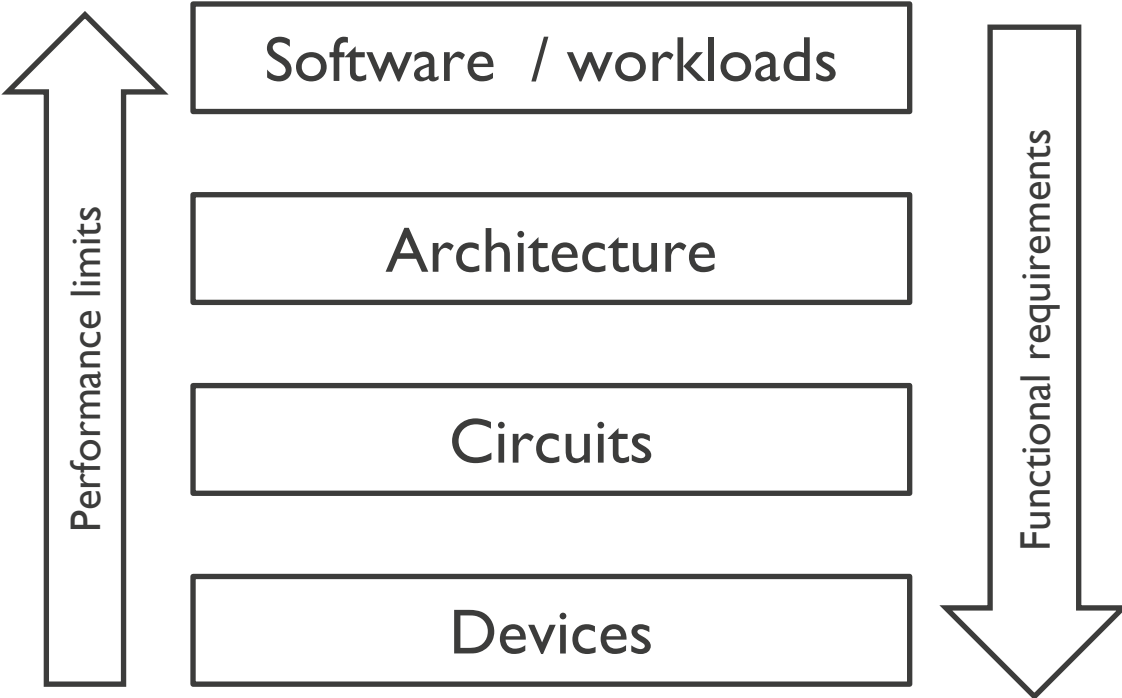Arindam Mallik, IMEC
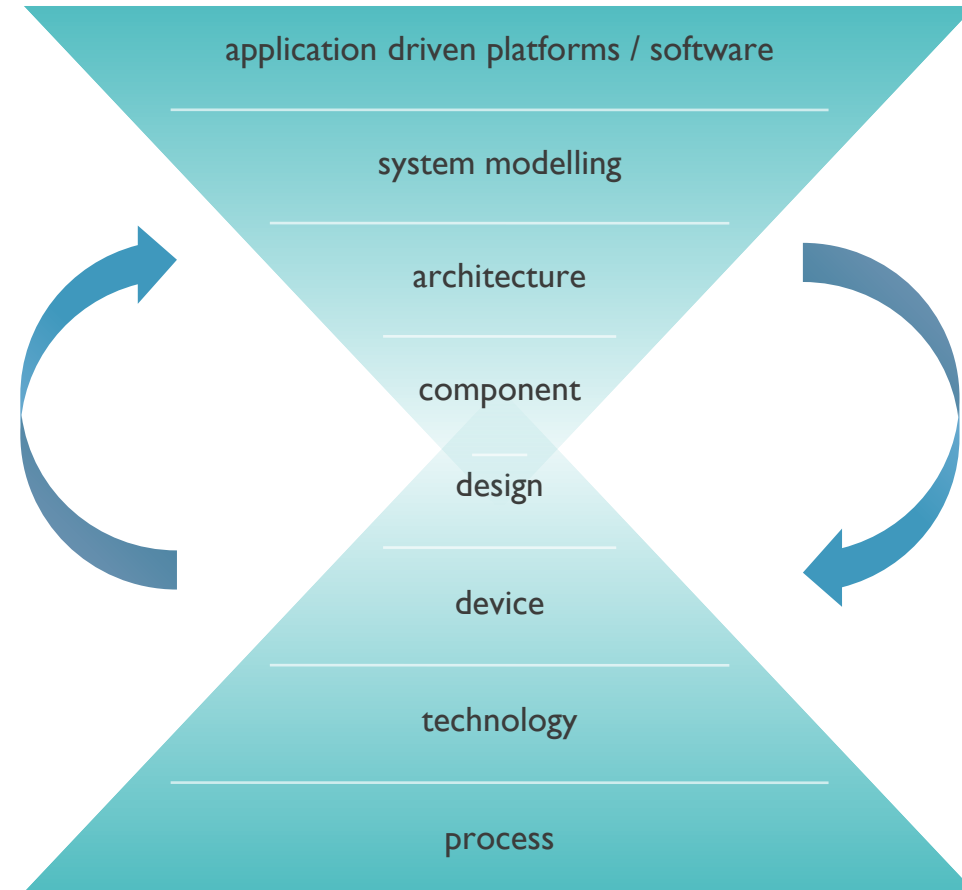Anand Raghunathan, Silicon Valley Professor of ECE
Sumeet Gupta, Elmore Associate Professor of ECE

**PURDUE**
UNIVERSITY®

# Tools of the Trade

Software / workloads

Architecture

Circuits

Devices

Performance limits

Functional requirements

**Need for a full stack approach**

application driven platforms / software

system modelling

architecture

component

design

device

technology

process

Technology impact might not be observable or significant in a specific architecture, but extremely important in another.
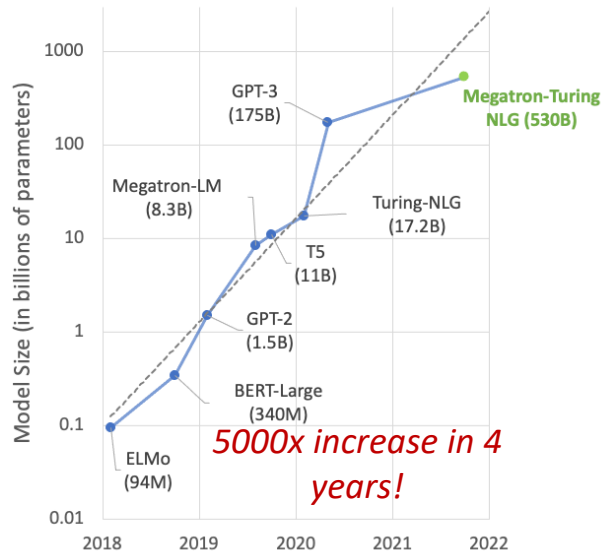
Technology: What, When, Where?

To translate technology improvements to *workload* performance, a HW-SW codesign methodology is needed.

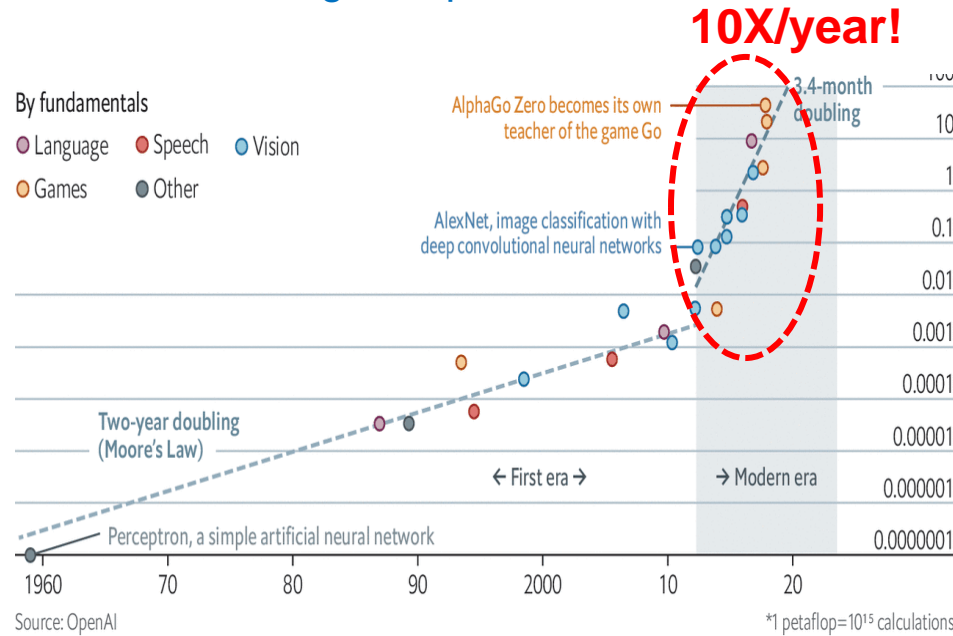# Specific Workloads will Drive Demand for Computing and Devices

■ AI will continue to set the pace in demand for compute efficiency

**NLP model growth**



*5000x increase in 4 years!*

Source: NVIDIA

**AI Training Compute Demands**

**10X/year!**



Source: The Economist, OpenAI
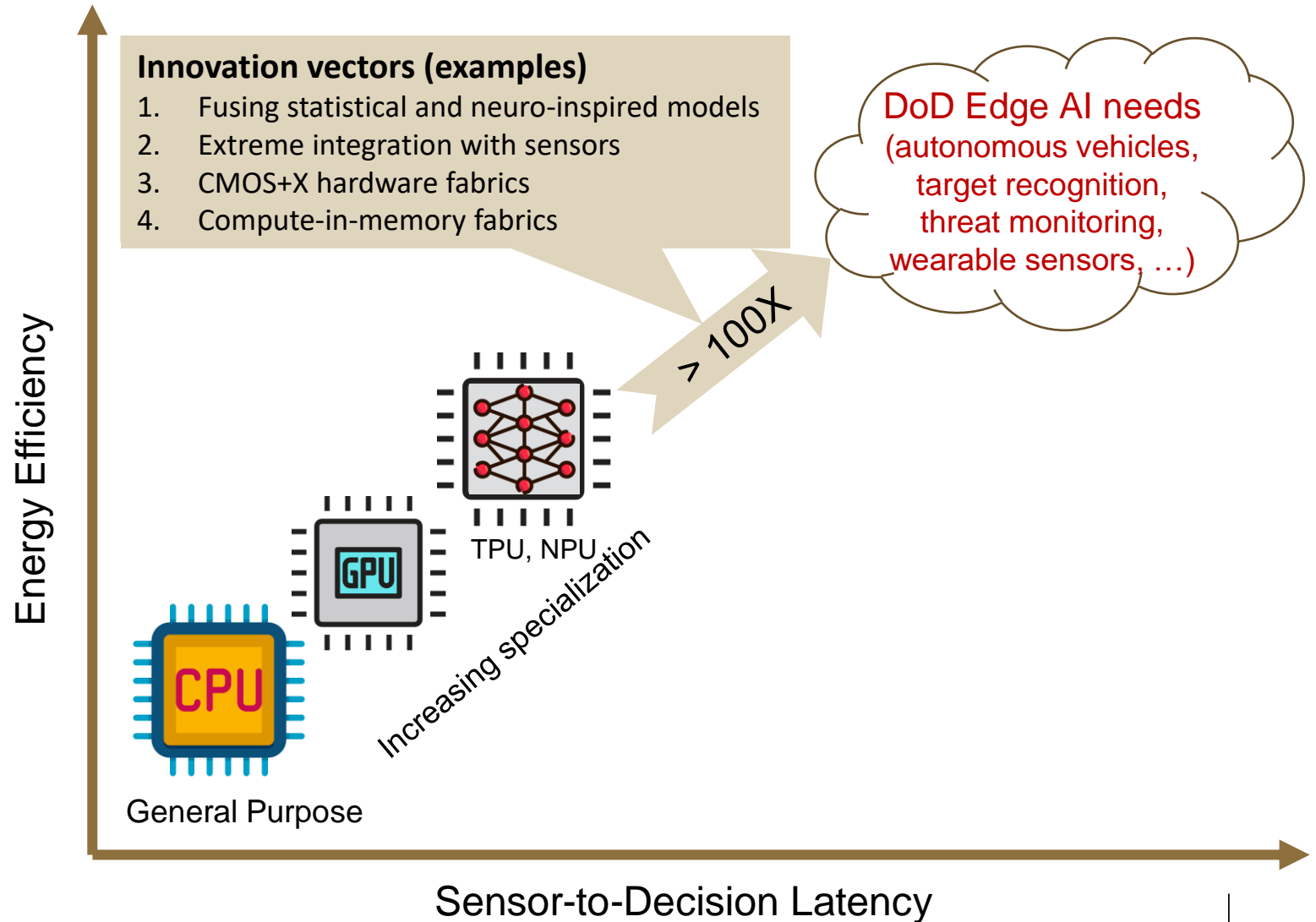
**AI Hardware Landscape**

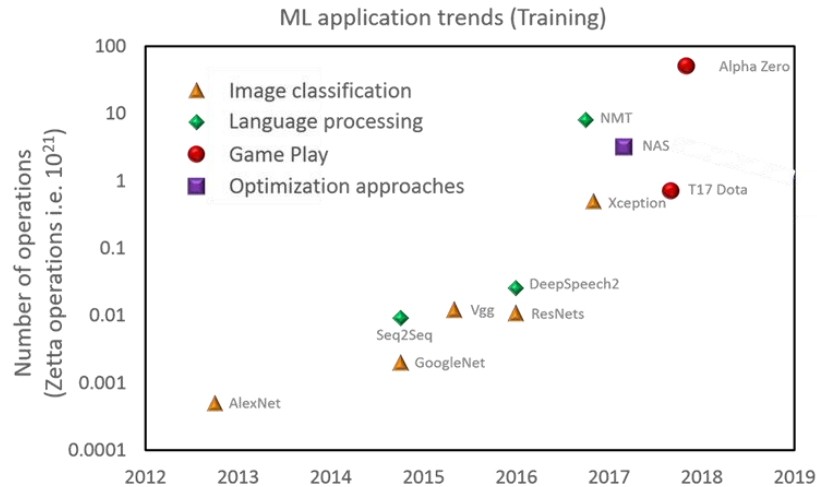# AI Hardware: Goals, Current Landscape

## How's it done today?

➤ **AI is powered by GPUs, TPUs/NPUs**

➤ **Most of the improvements are from data-parallel architectures, specialization and high-bandwidth memory**

## Need for radically different approach: beyond von-Neuman based hardware

**Innovation vectors (examples)**

1. Fusing statistical and neuro-inspired models
2. Extreme integration with sensors
3. CMOS+X hardware fabrics
4. Compute-in-memory fabrics

> 100X

DoD Edge AI needs (autonomous vehicles, target recognition, threat monitoring, wearable sensors, …)

Energy Efficiency

TPU, NPU

GPU

CPU

Increasing specialization

General Purpose

Sensor-to-Decision Latency

# *Why co-design is pivotal for AI hardware?*

Data-intensive AI workloads



ML application trends (Training)

- ▲ Image classification
- ◆ Language processing
- ● Game Play
- ■ Optimization approaches

Enormous energy/performance overheads in traditional von-Neumann architectures

Solution: Computing-in-Memory (CiM) of Matrix Vector Multiplications

Additional design conflicts/complexities due to CiM
- Larger range of CiM output currents (than simple read) → More non-idealities
- Need for complex peripheral circuits

Aggravation of Technology Non-Idealities
- Wire resistivity/resistance
- Process variations

Technology scaling to support high data storage and energy/performance demands of AI hardware
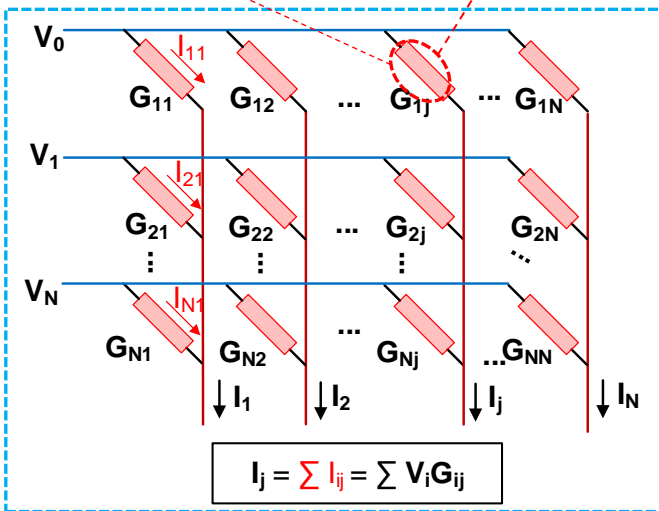
**Computational Errors, Energy Inefficiencies, Latency Increase !!**

**Need a technology-circuit-architecture-algorithm co-design solutions to tackle the increase in challenges**
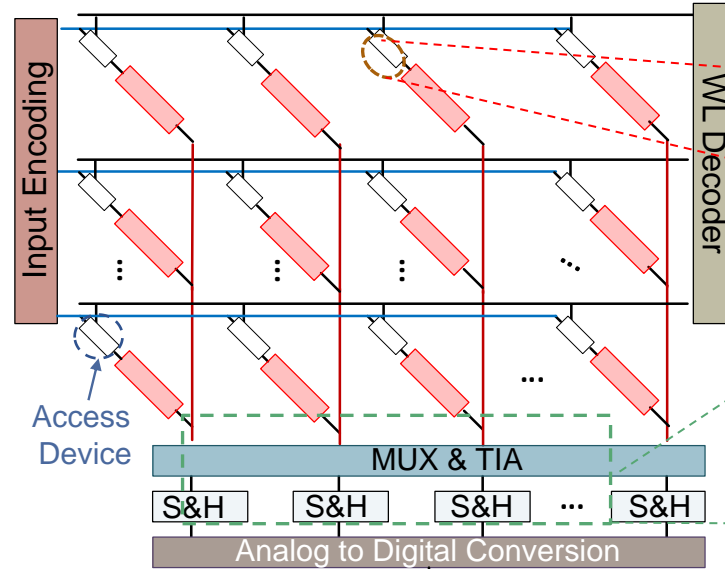
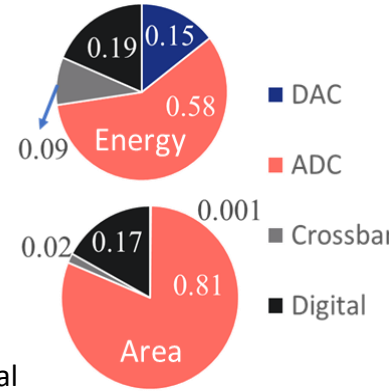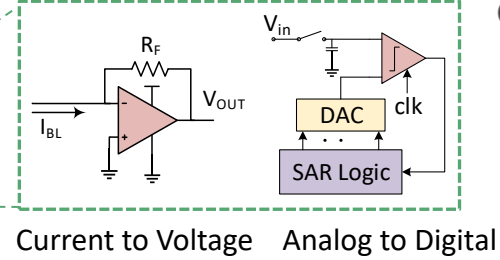# *Major Compute Requirement: MMMs*



**In-Memory Computing Devices**

Two-terminal Devices

PCM · RRAM · STT-MTJ · CMOS

**Efficient MVM Unit**

Input Encoding · WL Decoder

Access Device

MUX & TIA · S&H · Analog to Digital Conversion

2-Terminal Selector or Transistor

Chakraborty et. al, Resistive Crossbars…: Opportunities and Challenges, Proc. IEEE, 2020

**The Peripherals**

$R_F$ · $V_{OUT}$ · $I_{BL}$ · $V_{in}$ · DAC · clk · SAR Logic

Current to Voltage     Analog to Digital

Energy — 0.15 DAC, 0.58 ADC, 0.09, 0.19, 0.15

Area — 0.001, 0.02, 0.17, 0.81

DAC · ADC · Crossbar · Digital

**Efficient MVM**

$V_0$, $I_{11}$, $G_{11}$, $G_{12}$, $G_{1j}$, $G_{1N}$

$V_1$, $I_{21}$, $G_{21}$, $G_{22}$, $G_{2j}$, $G_{2N}$

$V_N$, $I_{N1}$, $G_{N1}$, $G_{N2}$, $G_{Nj}$, $G_{NN}$

$I_1$, $I_2$, $I_j$, $I_N$

$$I_j = \sum I_{ij} = \sum V_i G_{ij}$$

**NoC Architecture**

Network-on-chip · System Bus · Compute Core
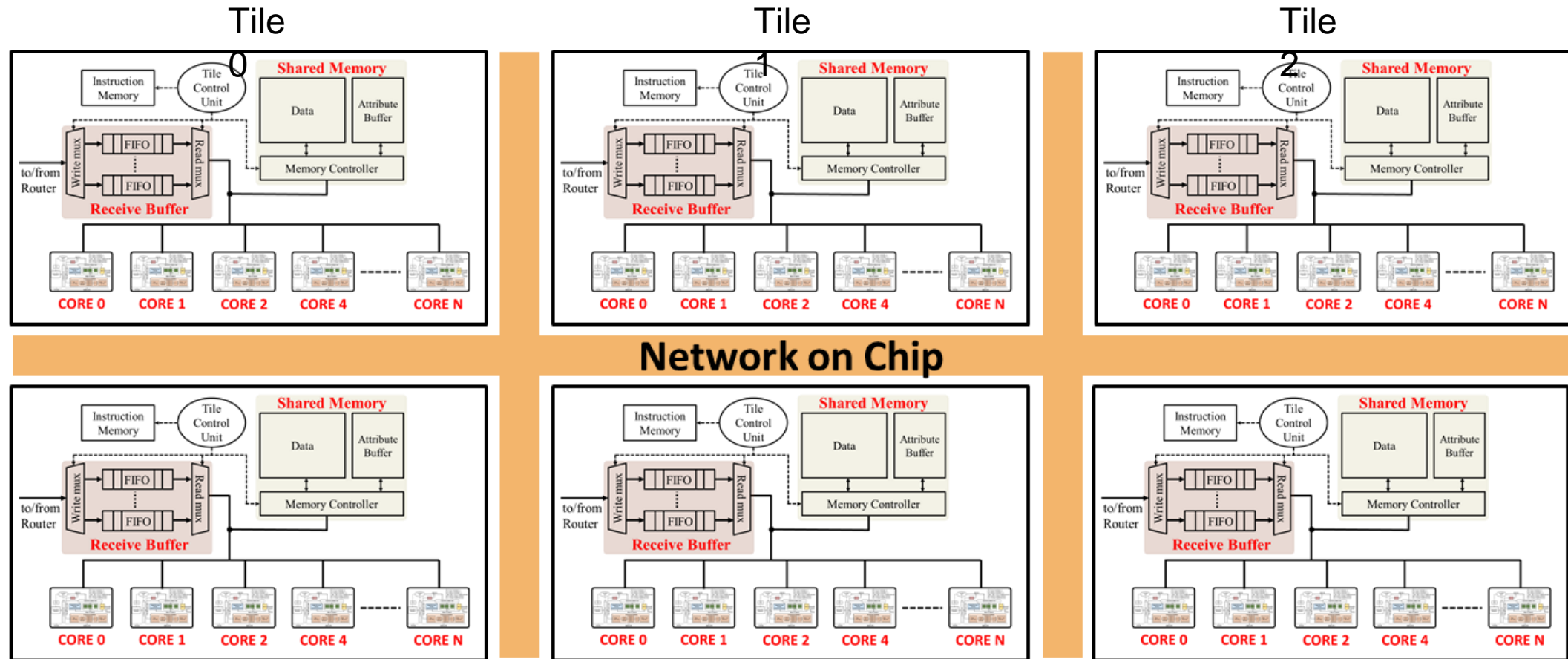
**High Density off-chip memory**

**PUMA: A Programmable Ultra-efficient Memristor-based Accelerator for Machine Learning Inference**



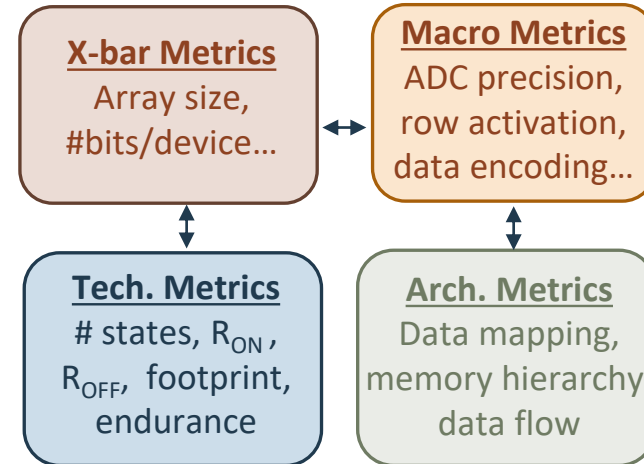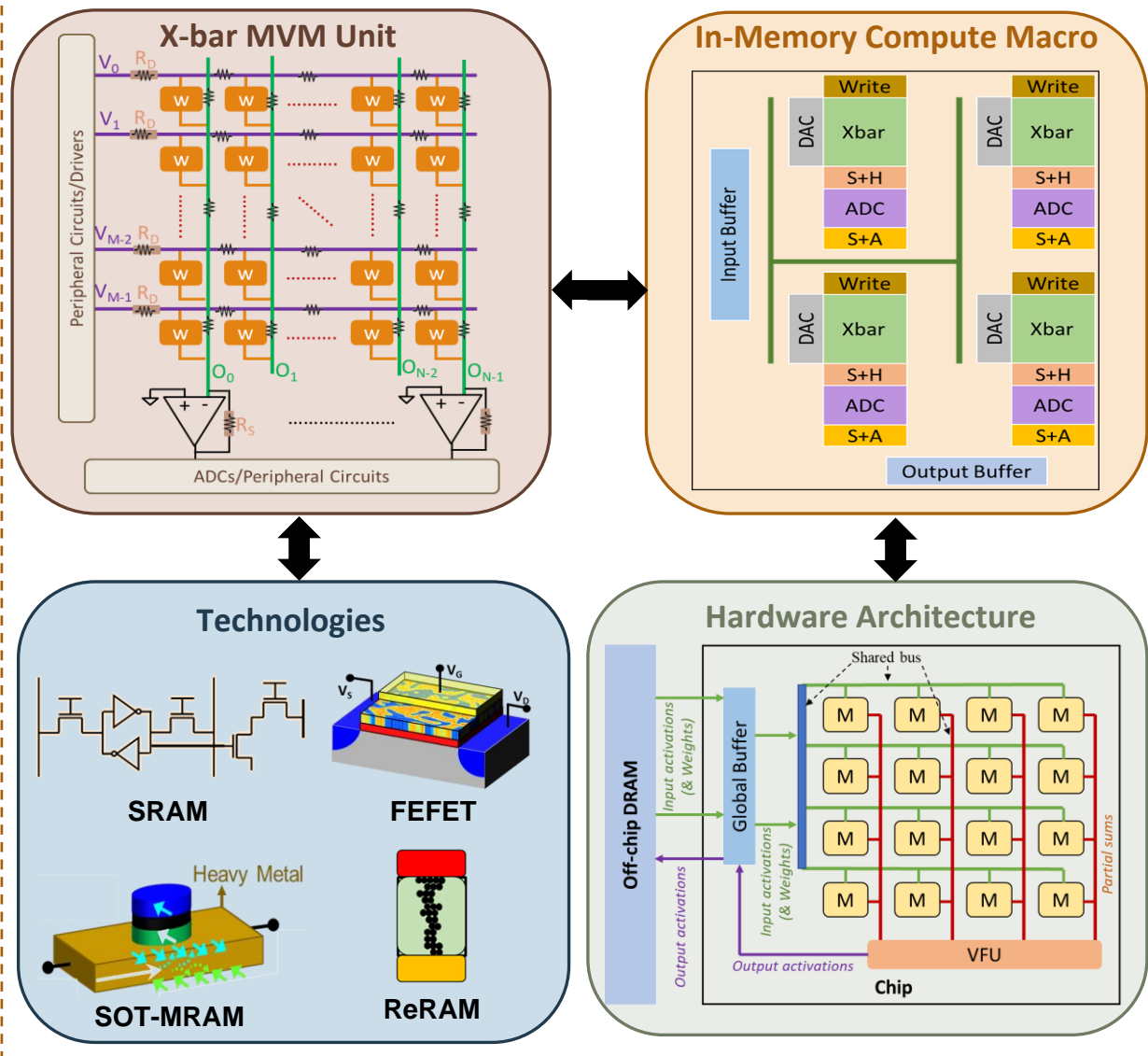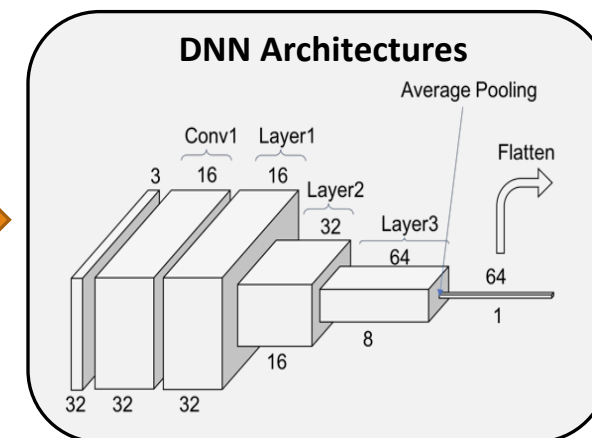Massively parallel accelerator –> Amenable to Data-Level Parallelism -> Highly efficient ML inference

Ankit et. al, PUMA:…, ASPLOS 2019

# AI Hardware Design: Intricate Cross-Layer Interactions

**X-bar MVM Unit**

**In-Memory Compute Macro**

**Technologies**
SRAM · FEFET · SOT-MRAM · ReRAM

**Hardware Architecture**

**X-bar Metrics**
Array size, #bits/device…

**Macro Metrics**
ADC precision, row activation, data encoding…

**Tech. Metrics**
# states, $R_{ON}$, $R_{OFF}$, footprint, endurance

**Arch. Metrics**
Data mapping, memory hierarchy, data flow

Need to consider the cross-layer optimization techniques to capture the effects of such interactions

**DNN Architectures**

$$NF = \frac{I_{Ideal} - I_{Non-Ideal}}{I_{Ideal}}$$

Low distinguishability
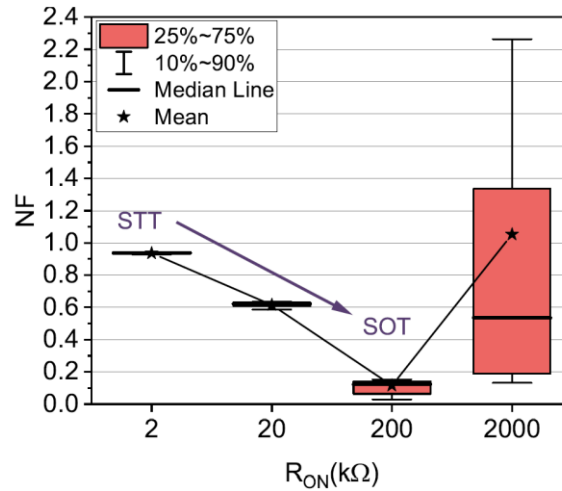
Algorithm-hardware co-design to enhance system accuracy
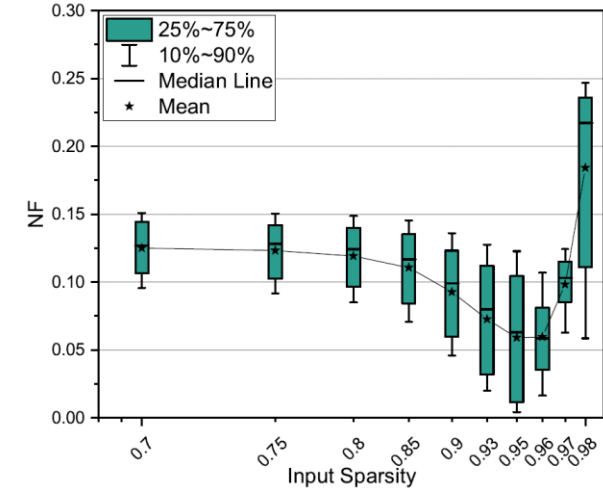
Algorithmic Sparsity

Lower number of activated row → Less impact of 'false' ON states

Synergize to reduce computational errors

Device Design
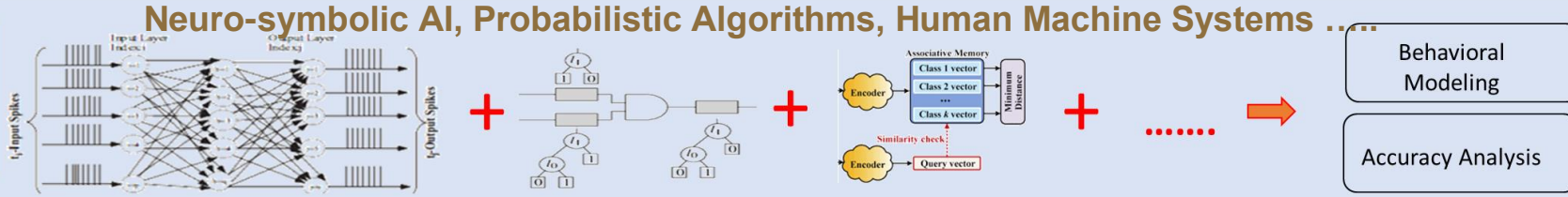
• Read-Write Path Separation
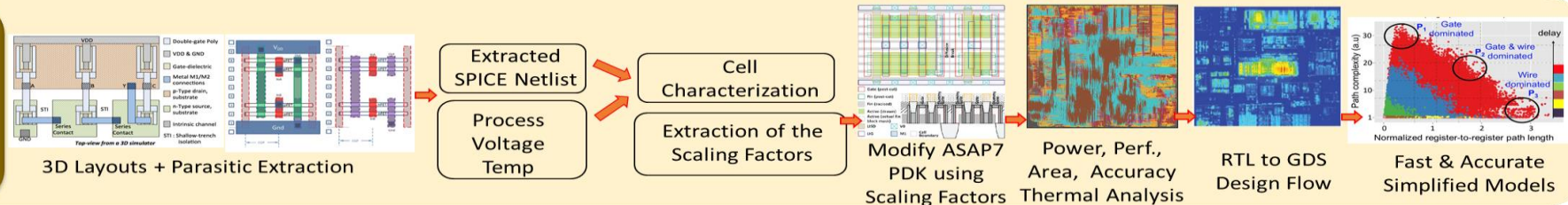• $t_{MgO}$ Optimization

$R_{ON}$ optimization to reduce NF

Algorithm design to reduce NF

T. Sharma, K. Roy *et al*, ISLPED 2021

# Cross-layer flow for System-Technology Co-Design



**Next Generation AI Architectures**

Neuro-symbolic AI, Probabilistic Algorithms, Human Machine Systems .....

Behavioral Modeling

Accuracy Analysis

**Automated Design Flow**

3D Layouts + Parasitic Extraction

Extracted SPICE Netlist

Process Voltage Temp

Cell Characterization

Extraction of the Scaling Factors

Modify ASAP7 PDK using Scaling Factors

Power, Perf., Area, Accuracy Thermal Analysis

RTL to GDS Design Flow

Fast & Accurate Simplified Models

**Library of Primitive Components**

Neural Primitives

Reconfigurable synaptic sub-systems, Adaptive spiking neurons, 3D integrable neural fabrics

Symbolic Logic Primitives

Real-valued Logic CiM

CiM of A+B, A-B, AB, Comparison (A,B)

Peripherals

Compute-enabled memories supporting real-valued logic operations

Probabilistic Primitives

Probabilistic circuits leveraging inherent stochasticity in emerging devices

HD Computing Primitives

Content-Addressable Memories and other primitives for HD computing

**Modeling and Simulation Frameworks (Calibrated with experiments)**

Physical Models

FinFET

Stacked Nanowire FETs

TCAD
CFET (IMEC)

Magnetoelectric Content Addressable Memory

In-House Models

Circuit-Compatible Models

Variation and Stochasticity Models

**Technologies**

CMOS

Ferroelectric

2D TMD Devices

Spin/Magnetoelectric

BEOL FETs

Interconnects

Technology Re-Design

Area required for fully mapping the workload

He, Roy et al, ICCAD 2022
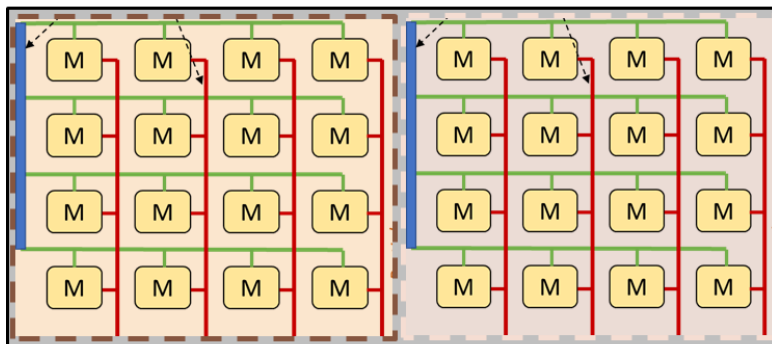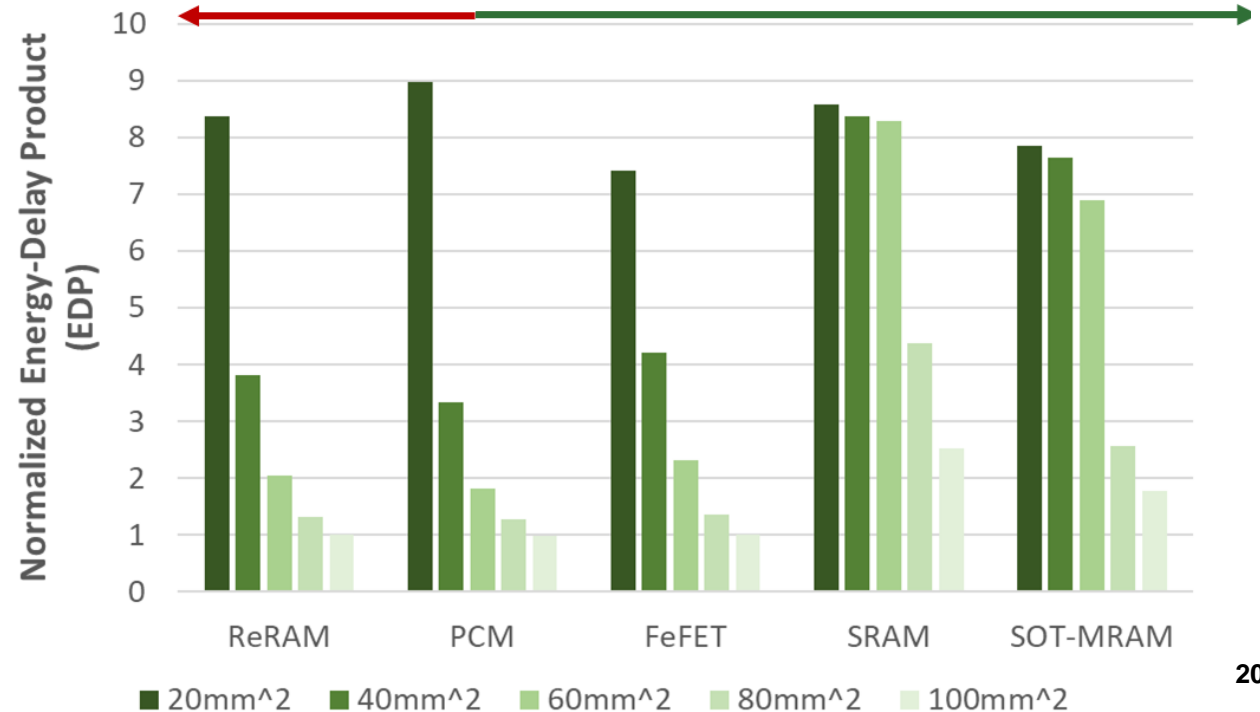
Small area budget

Off-Chip DRAM

Weight reloading (Needs DRAM access)

Large area budget

Weight replication (Increases parallelism)

Weight reloading from DRAM     Weight replication



Normalized Energy-Delay Product (EDP)

ReRAM     PCM     FeFET     SRAM     SOT-MRAM

■ 20mm^2  ■ 40mm^2  ■ 60mm^2  ■ 80mm^2  ■ 100mm^2

20

# HW-SW Codesign for future AI compute systems

Application and algorithm optimization
Low-precision neural network exploration

LLM,
neuromorphic exploratory algorithm

Use/drive emerging technologies

SRAM, STT/SOT-MRAM, Flash, FeFET, IGZO, PCM

ALGORITHM

DEVICE

imec-PU AI system exploration

ARCHITECTURE

CIRCUITS

System architecture exploration
Digital hardware acceleration

CiM accelerator component, middleware for accelerator integration in the system

Compute-in-Memory (CiM)
macro using emerging memories

1T1R, 2T1R, 2T1C, 2T0C cells for CiM macro

imec

21

# THANK YOU!

PURDUE
UNIVERSITY®