# Compute-in-Memory and AI Accelerator Technologies for the Sub-18Å Era

**SRC Industry Talk, August 2023**
**Ram K. Krishnamurthy**
**High Performance and Low Voltage Circuits Research**
**Circuits Research Lab, Intel Labs**
**Intel Corporation, Hillsboro, OR 97124, USA**
**ram.krishnamurthy@intel.com**

# Internet of Everything (IoE)



*Need end-to-end energy efficiency, ML everywhere*

# DATA DEFINES THE FUTURE



**The Economist**

- Obama the warrior
- Misgoverning Argentina
- The economic shift from West to East
- Genetically modified crops blossom
- The right to eat cats and dogs

## The data deluge

AND HOW TO HANDLE IT: A 14-PAGE SPECIAL REPORT

**The Economist**

- Crunch time in France
- Ten years on: banking after the crisis
- South Korea's unfinished revolution
- Biology, but without the cells

## The world's most valuable resource

Data and the new rules of competition

## POPULAR SCIENCE

THE FUTURE NOW

### THE CONTROL CENTERS
Using Data to Feed the World, Solve Cold Cases, Battle Malware, Predict Our Fate

### OFFICER ALGORITHM
Can a Crime Be Prevented Before It Begins?

### NEW WAYS OF SEEING
A Gallery of Extraordinary Infographics

**PLUS**
- Juan Enriquez Reprograms Life
- James Gleick Unsplits the Bit
- AND Lawrence Weschler Questions the Cloud

SPECIAL ISSUE

## DATA IS POWER

HOW INFORMATION IS DRIVING THE FUTURE

## COSMOS

8-PAGE SPECIAL POSTGRADUATE SURVIVAL GUIDE · 8TH BIRTHDAY ISSUE!

THE SCIENCE OF EVERYTHING

### THE END OF VIOLENCE
Steven Pinker on the new peace

### DEFEATING POLIO
Will politics jeapordise a cure?

### FRAUDS AND FAKES
Science's biggest scams

### GENIUS OF DOGS
Inside the canine brain

## IS data THE NEW GOD?
How tracking your digital trail could predetermine your future – and why you'll benefit from today's data deluge.

GALAXIES AND NEBULAE · CANCER VACCINES · WHALES · FICTION · REVIEWS

(intel)

# The Data Problem

We are **generating data** at a **faster** rate than our ability to **analyze, understand, transmit, secure** and **reconstruct** in real-time

**175ZB**



Zettabytes

150

100

50

2010 2011 2012 2013 2014 2015 2016 2017 2018 2019 2020 2021 2022 2023 2024 2025

intel. labs

# Performance Democratization

**100B INTELLIGENT CONNECTED DEVICES**

Distributed Intelligence

Cloud Everything

Mobile Everything

Network Everything

Digitize Everything

**Exascale**

For Everyone

Compute

$10^{18}$

$10^{15}$

$10^{9}$

$10^{4}$

$10^{2}$

1980    1990    2000    2010    2020

The Future Begins Here

intel. labs

# Compute and Memory Challenges for AI



- Compute demand growth rate: Doubling every 3-4 months
- Memory capacity growth rate: 10X per year

# INTELLIGENCE FOUNDATION

## INTEL® XEON® SCALABLE PROCESSORS

**2017**

**1ST GEN**
### AVX-512

**FIRST BUILT-IN AI ACCELERATION**

**2019**

**2ND GEN**
### INTEL® DEEP LEARNING BOOST

UP TO **30X** IMPROVEMENT IN AI INFERENCE PERFORMANCE

**2020**

**3RD GEN**
### INTEL® DL BOOST EXTENSIONS

UP TO **60%** INCREASE IN AI TRAINING PERFORMANCE

Up to 14X AI performance improvement with Intel® Deep Learning Boost (Intel DL Boost) compared to Intel Xeon Platinum processor (April 2019). See configuration disclosure for details. Up to 60% performance improvement with Intel® Deep Learning Boost (Intel DL Boost) is a projection based on Intel internal measurements using pre-production hardware/software as of December 2019. All products, computer systems, dates, and figures are preliminary based on current expectations, and are subject to change without notice. No product or component can be absolutely secure. For more complete information about performance and benchmark results, visit www.intel.com/benchmarks.

# AI Has Moved to the Edge

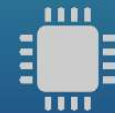**Edge Devices**

High Privacy

Low Latency

High Availability

Energy Efficient

*Algorithm & Hardware Advancement*

**AI**

**Cloud Computing**

Thermal Budget

Power Source

Memory Capacity

Computing Resource

Source: L. LOH, isscc 2020

# THE EDGE OPPORTUNITY...AND CHALLENGE



Deep Learning Chipset Revenue, Training vs. Inference, World Markets: 2018-2025

Legend: Training, Inference

"Cameras grow at highest CAGR."
MARKETSANDMARKETS™

"75% of AI hardware will be at the edge."
Tractica

"The success of AI on edge needs clever optimization techniques on limited power."
Gartner®

# Diversified Workload & Increasing Demands

| 0.1 TOPS | 1 TOPS | 10 TOPS | 100 TOPS |
|---|---|---|---|
| 1 TOPS/W | 3 TOPS/W | 10 TOPS/W | 30 TOPS/W |

**Vision Perception** | **Vision Construction** | **Visual Quality** | **Multi-Streaming**



Source: L. LOH, isscc 2020

# Intel Process Technology

# Go Wider: Within Package Interconnect Scaling

**Foveros Direct**

**Embedded Bridge (EMIB)**

**Foveros**

Interconnect Density →

**Standard**

| Bump Pitch | 100 um |
|---|---|
| Bump Density | 100/mm² |
| Power | 1.7 pJ/bit |

| Bump Pitch | 55-36 um |
|---|---|
| Bump Density | 330-772/mm² |
| Power | 1.7 pJ/bit |

| Bump Pitch | 50-25 um |
|---|---|
| Bump Density | 400-1600/mm² |
| Power | 0.15 pJ/bit |

| Bump Pitch | < 10 um |
|---|---|
| Bump Density | >10,000/mm² |
| Power | < 0.05 pJ/bit |

*Graphics are for illustrative purposes only and not to scale

Power Efficiency →

# Normalized Energy

**Compute**

**Data Transfer**

| | INT8 ADD | INT16 ADD | INT8 MUL | INT8 MAC | INT16 MUL | INT16 MAC | FP16 ADD | 1B SRAM WR | FP16 MAC | NOC | 1B MIPI Tx | DRAM | Wireless |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Value | 1 | 3 | 15 | 19 | 64 | 72 | 125 | 144 | 374 | 480 | 6291 | 12000 | 330000 |

# Memory Bottleneck

- Performance gap between processor and memory

   **Von Neumann Bottleneck**

- Huge energy consumption of memory

- Reduce data movement between processor and memory

- Computation-in-Memory (CiM)



Source: "Taming the Power Hungry Data Center" by Fusion-IO.

| Operation | Energy [pJ] | Relative Cost |
|---|---|---|
| 32 bit int ADD | 0.1 | 1 |
| 32 bit float ADD | 0.9 | 9 |
| 32 bit Register File | 1 | 10 |
| 32 bit int MULT | 3.1 | 31 |
| 32 bit float MULT | 3.7 | 37 |
| 32 bit SRAM Cache | 5 | 50 |
| **32 bit DRAM Memory** | **640** | **6400** |

Source: Song Han et al. "EIE: efficient inference engine on compressed deep neural network," ISCA 2016.

SC2-6: Alternate Technologies for SRAM, Hai Li, Duke University, *IEDM*, 2020. Source: K. Takeuchi, IRPS 2023
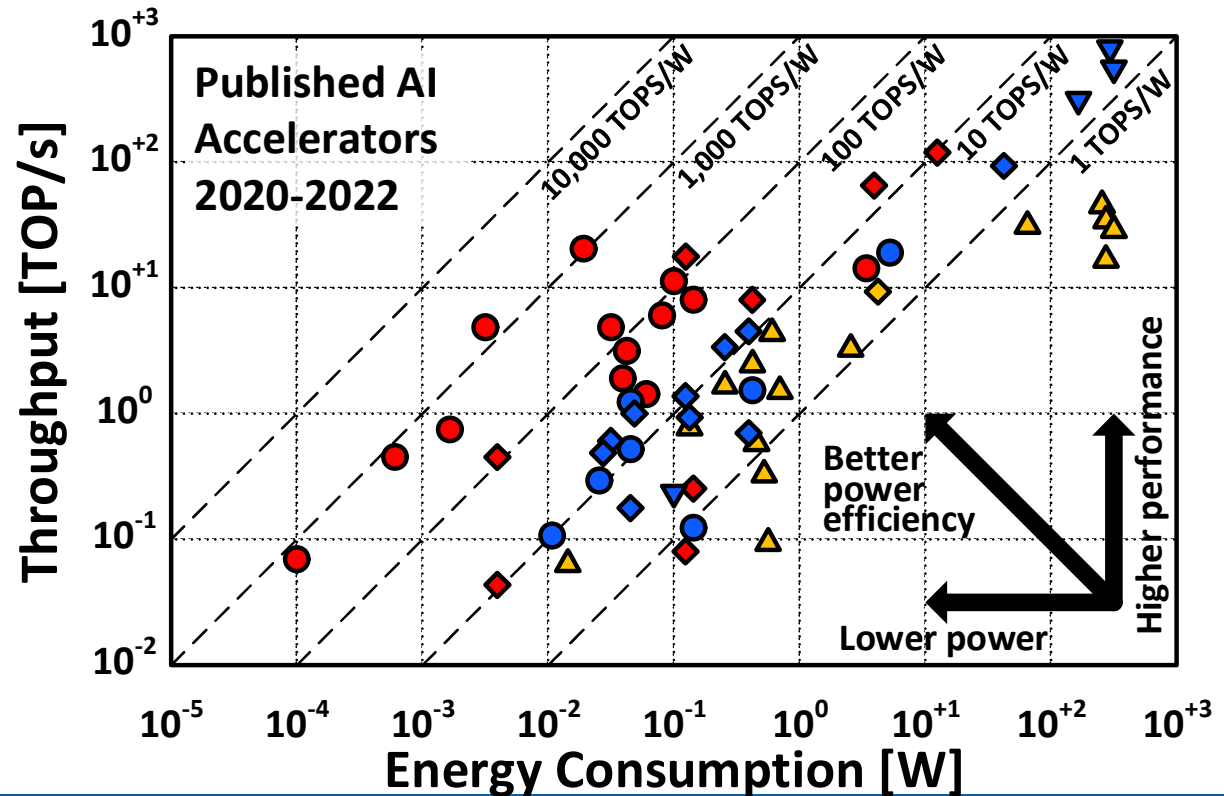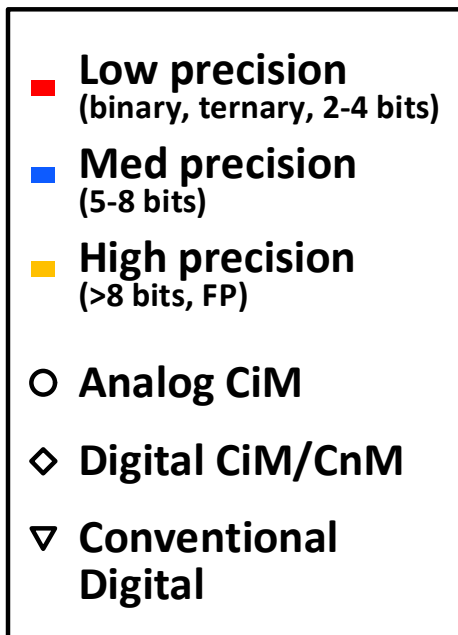
# Compute-in-Memory Motivation

- Data movement is costly
- Multiply–accumulate (MAC) operation
- Massively parallel processing
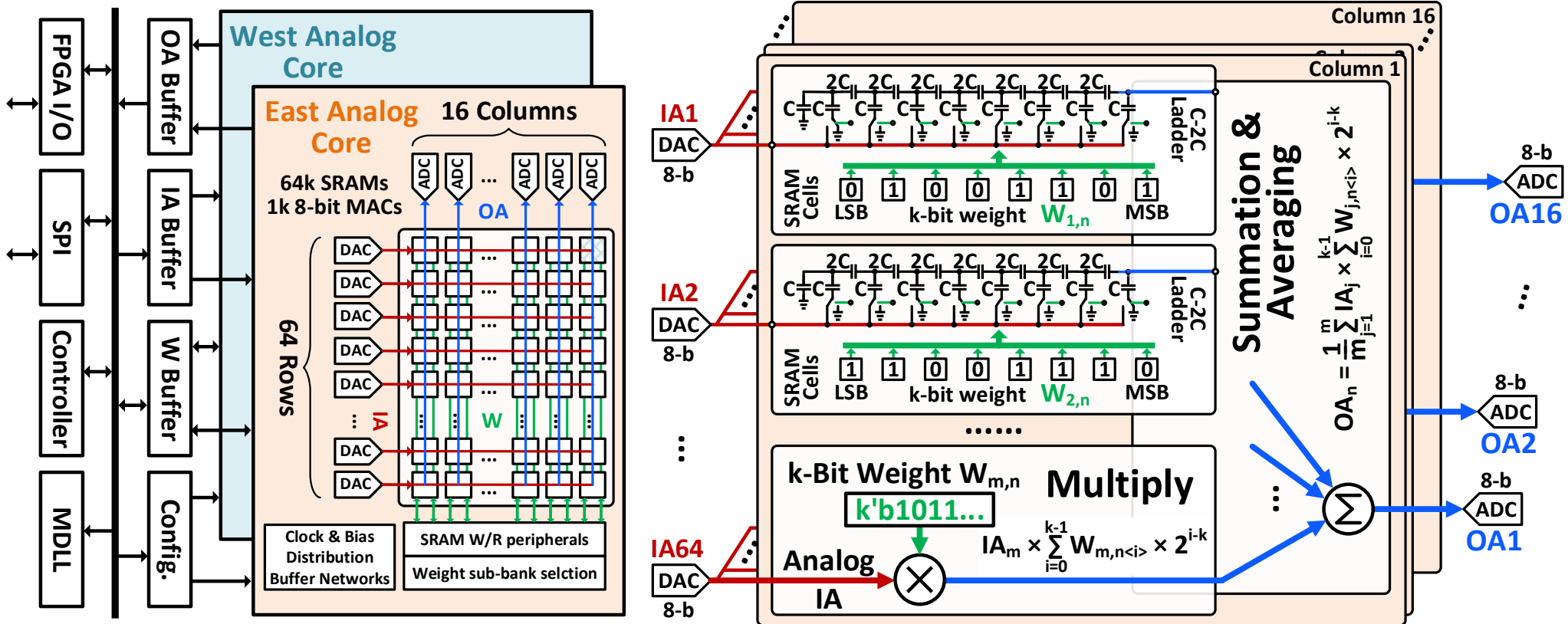- Beyond von Neumann architecture

# Compute-in-Memory Challenges

- TOPS/W versus Precision

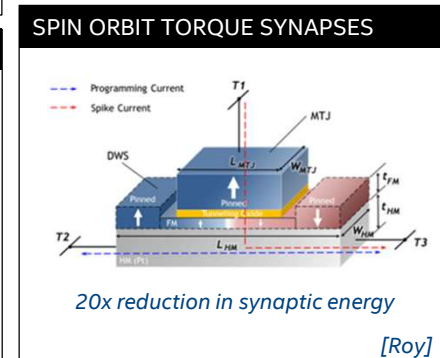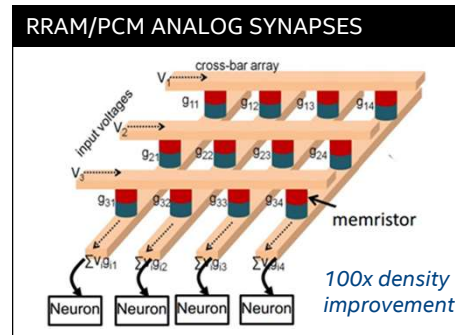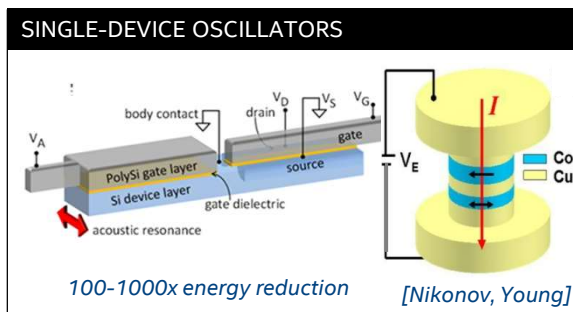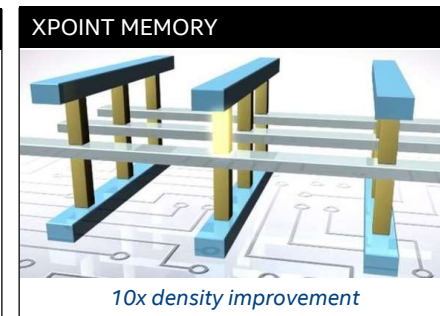# Intel Labs Analog CIM Architecture Overview



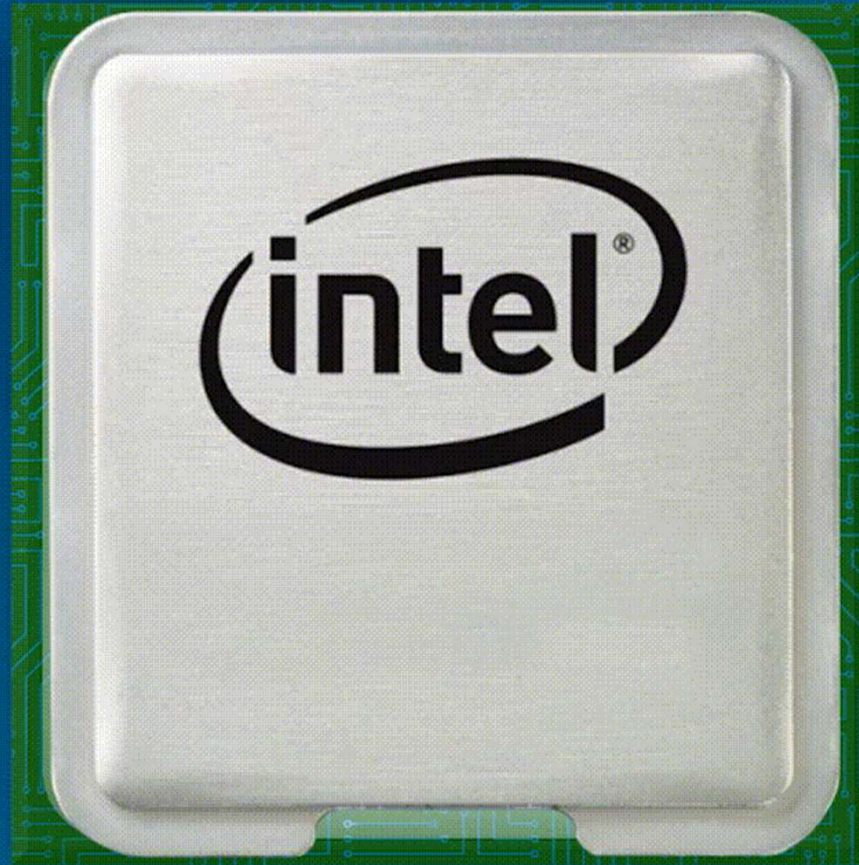H. Wang et al, IEEE VLSI Circuits Symposium 2022

# Intel Labs Analog CIM Measured Performance

- Energy efficiency: 15.5-32.2 TOPS/W
- Area efficiency: 2.4-4.0 TOPS/mm$^2$
- Clock frequency: 145-240 MHz
- Supply Voltage: 0.7-1.1 V

# Opportunities for in-memory/near-memory Process and Circuit Innovation
## (Both Digital and Analog/Mixed-Signal)



STTRAM STATE

*2x density improvement iso-energy*

3D

*2-64x density improvement*

XPOINT MEMORY

*10x density improvement*

SINGLE-DEVICE OSCILLATORS

*100-1000x energy reduction*

*[Nikonov, Young]*

RRAM/PCM ANALOG SYNAPSES

*100x density improvement*

SPIN ORBIT TORQUE SYNAPSES

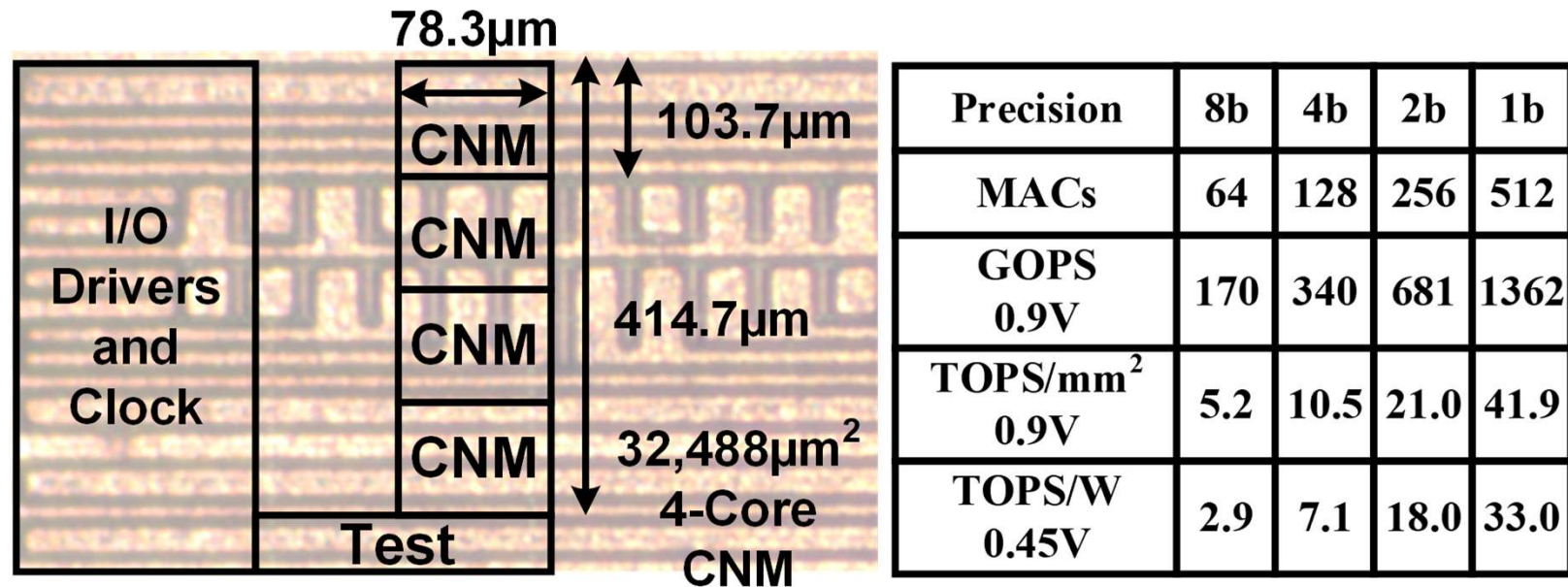*20x reduction in synaptic energy*

*[Roy]*

# COMPUTE NEAR MEMORY CHALLENGES AND OPPORTUNITIES



**E. Sumbul, R. Krishnamurthy et al, IEEE ESSCIRC 2021**

# 10nm Near Memory Computing AI Inference Accelerator



| Precision | 8b | 4b | 2b | 1b |
|---|---|---|---|---|
| MACs | 64 | 128 | 256 | 512 |
| GOPS 0.9V | 170 | 340 | 681 | 1362 |
| TOPS/mm$^2$ 0.9V | 5.2 | 10.5 | 21.0 | 41.9 |
| TOPS/W 0.45V | 2.9 | 7.1 | 18.0 | 33.0 |

- 4 CNM cores with 8KB of weight memory and 64 8b multipliers
- Supports memory-intensive batch-1, large-batch, and in-place convolution
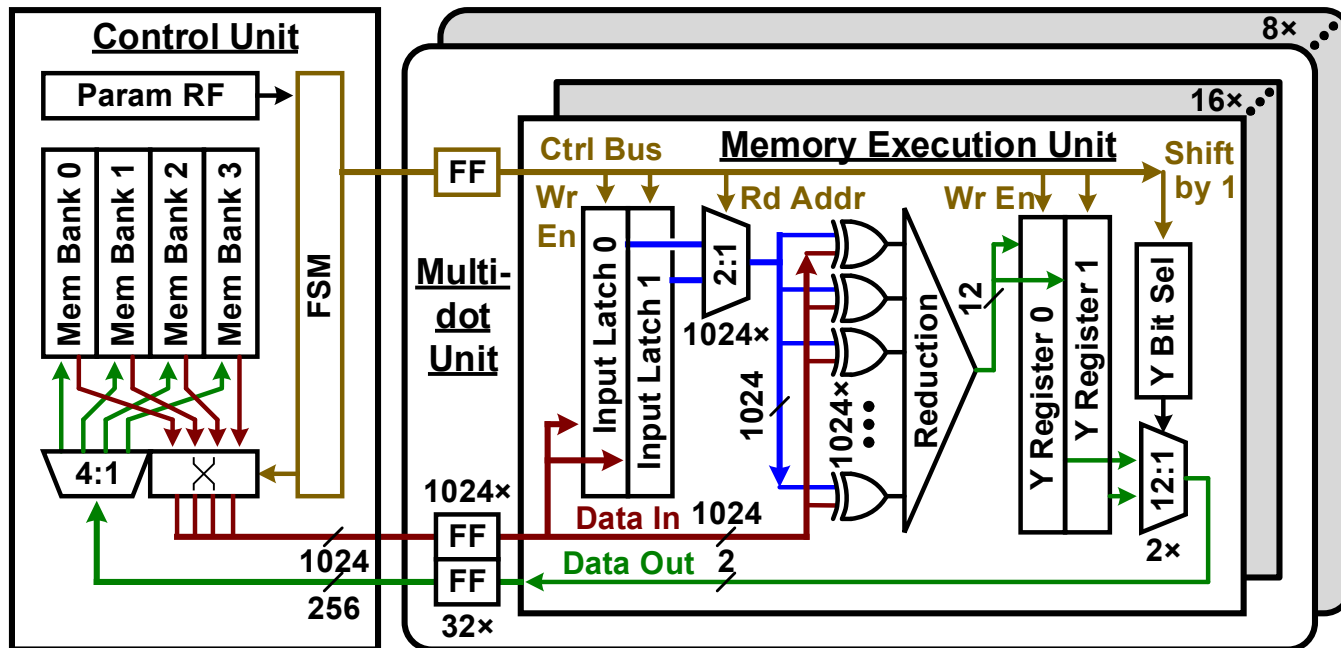
# 10nm Near Memory Computing Measurement Results



- Peak throughput 170 8b TOPS @ 0.9V that scaled up with number of CNM cores
- NTV operation down to 450mV decreases energy by 3.1x to 2.9 8b TOPS/W
- Variable precision improves energy efficiency by 11.4x to 33.0 1b TOPS/W

**G. Chen, R. Krishnamurthy et al, IEEE European Solid-State Circuits Conference 2021**
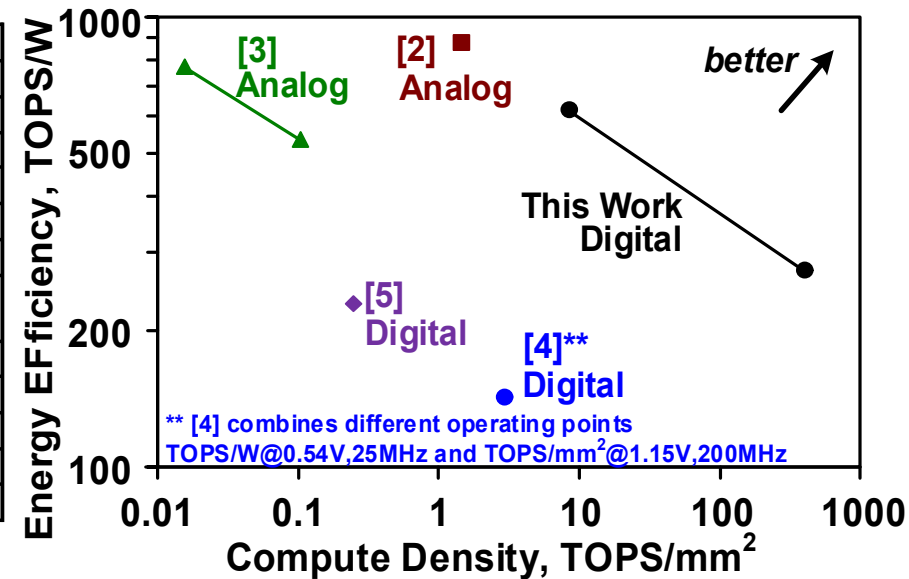
# 10nm Binary Neural Network Inference Accelerator



- Array of 128 Memory Execution Units (MEU) combine latch base memory and inner product compute in fine grain manner to minimize interconnect energy
- Central controller manages data flow from four 256b memory banks to MEUs
- 2 latch words per MEU enables data reuse reducing input bandwidth by 2x
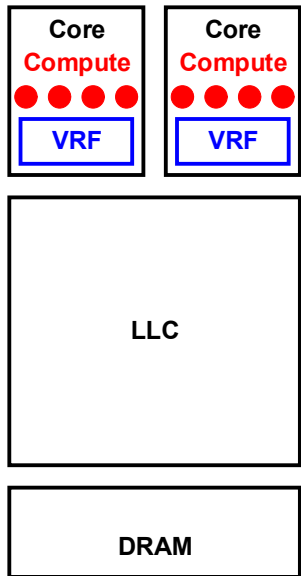
# Comparison to Previously Published BNNs

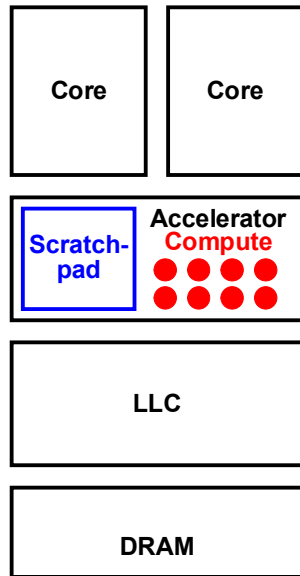| | [2] | [3] | | [4] | | [5] | This Work | |
|---|---|---|---|---|---|---|---|---|
| Digital or Analog | Analog | Analog | | Digital | | Digital | Digital | |
| Technology (nm) | 65 | 28 | | 65 | | 28 | 10 | |
| Area (mm²) | 12.6 | 4.6 | | 4.8 | | 1.4 | 0.39 | |
| Num MACs | - | 65,536 | | - | | 65,536 | 131,072 | |
| Mem Capacity (KB) | 295 | 328 | | 104 | | 328 | 161 | |
| KB/mm² | 23 | 71 | | 22 | | 234 | 414 | |
| Voltage (V) | 0.68, 0.94, 1.2 | 0.6, 0.8, 0.6, 0.53 | 0.6, 0.8, 0.8, 0.8 | 0.54 | 1.15 | - | 0.37 | 0.75 |
| Frequency (MHz) | 100 | 1.5 | 10 | 25 | 200 | 6.0 | 13 | 622 |
| Power (mW) | 22 | 0.094 | 0.899 | - | - | 1.5 | 5.6 | 607 |
| TOPS | 19 | 0.072 | 0.478 | - | 14.9 | 0.35 | 3.4 | 163 |
| TOPS/W | 866 | 772 | 532 | 140 | - | 230 | 617 | 269 |
| TOPS/mm² | 1.5 | 0.02 | 0.10 | - | 3.1 | 0.25 | 8.8 | 418 |



**P. Knag, R. Krishnamurthy et al, IEEE Journal of Solid-State Circuits Invited Paper, April 2021**
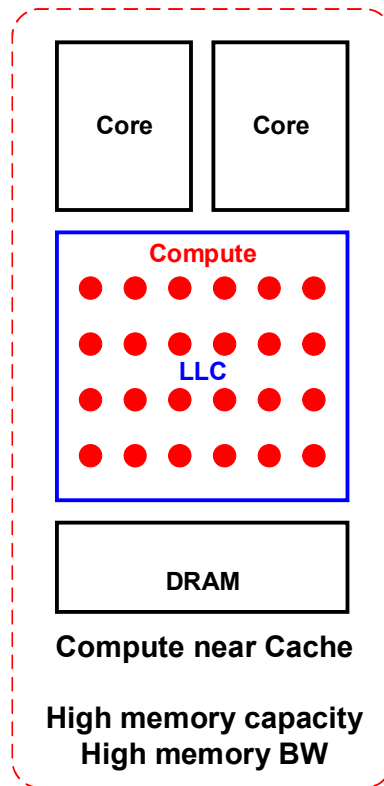
# Compute Near Last Level Cache (CNC)



**Vector Processing**
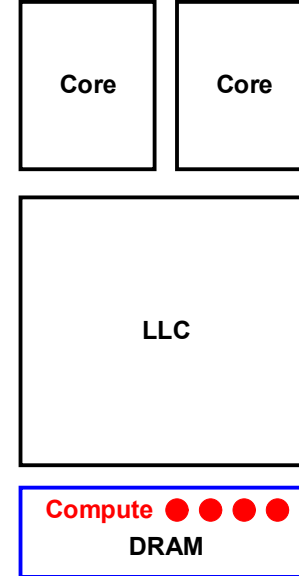
Small Vector RF Capacity
Low BW to LLC

**Accelerator**

Highest Energy Efficiency
Least General Purpose

**Compute near Cache**

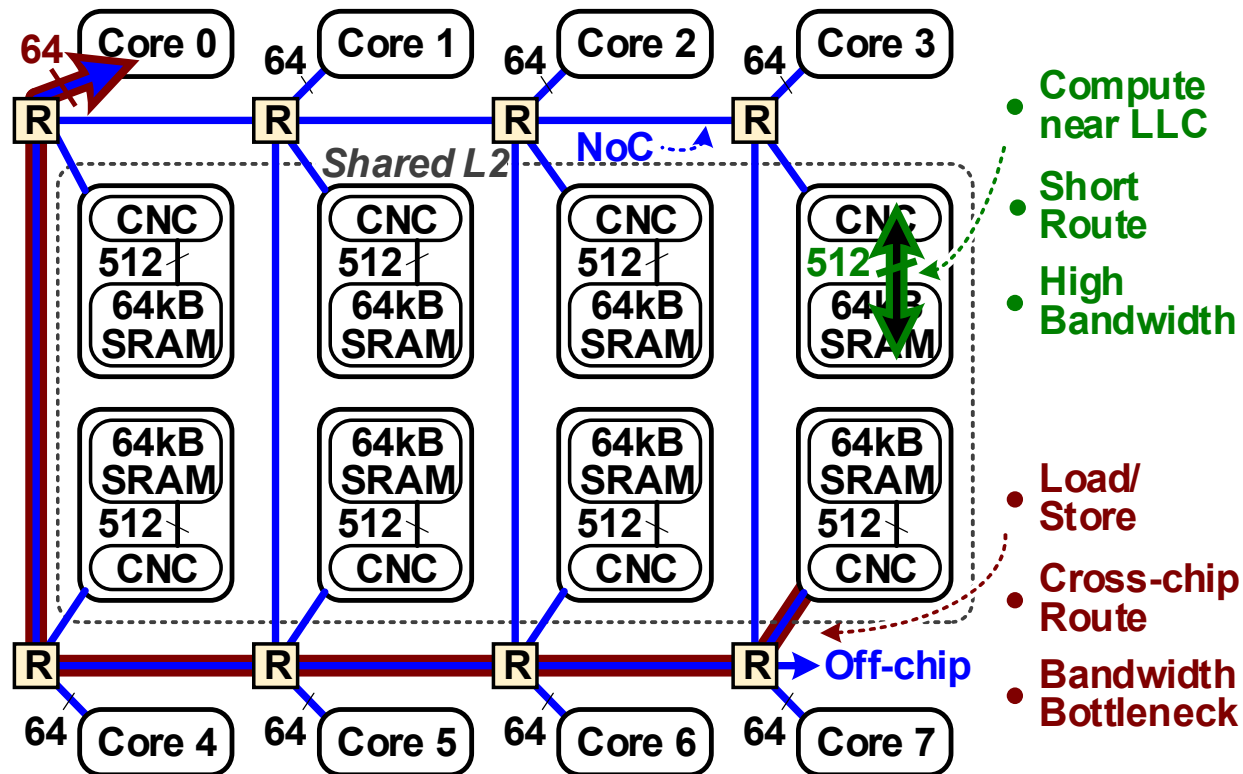High memory capacity
High memory BW

**Processing in DRAM**

Highest memory capacity
DRAM processing

- CNC enables fine grain mixing of near-memory vector and GP scalar computation
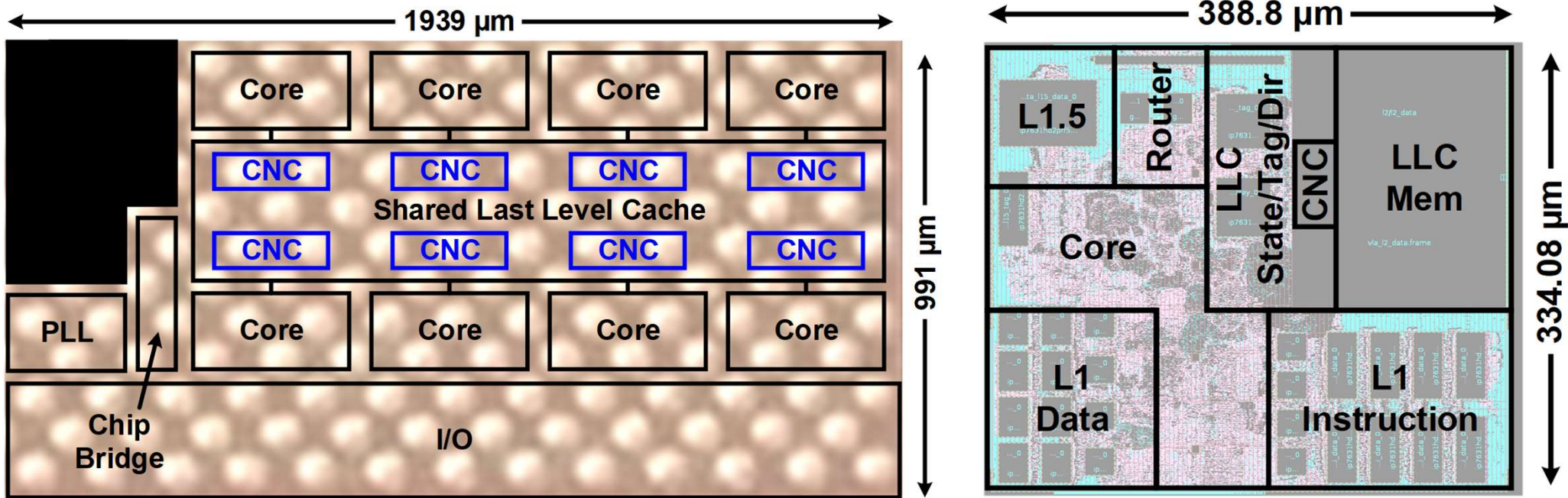- High BW access to highest capacity on-chip memory instead of RF/scratchpad

# Compute Near Last Level Cache of RISC-V Multiprocessor



- 8-core RV64GC processor with 128 INT8 MACs near 512kB shared, distributed LLC

- CNC ISA extension with support for virtual addressing and cache coherence
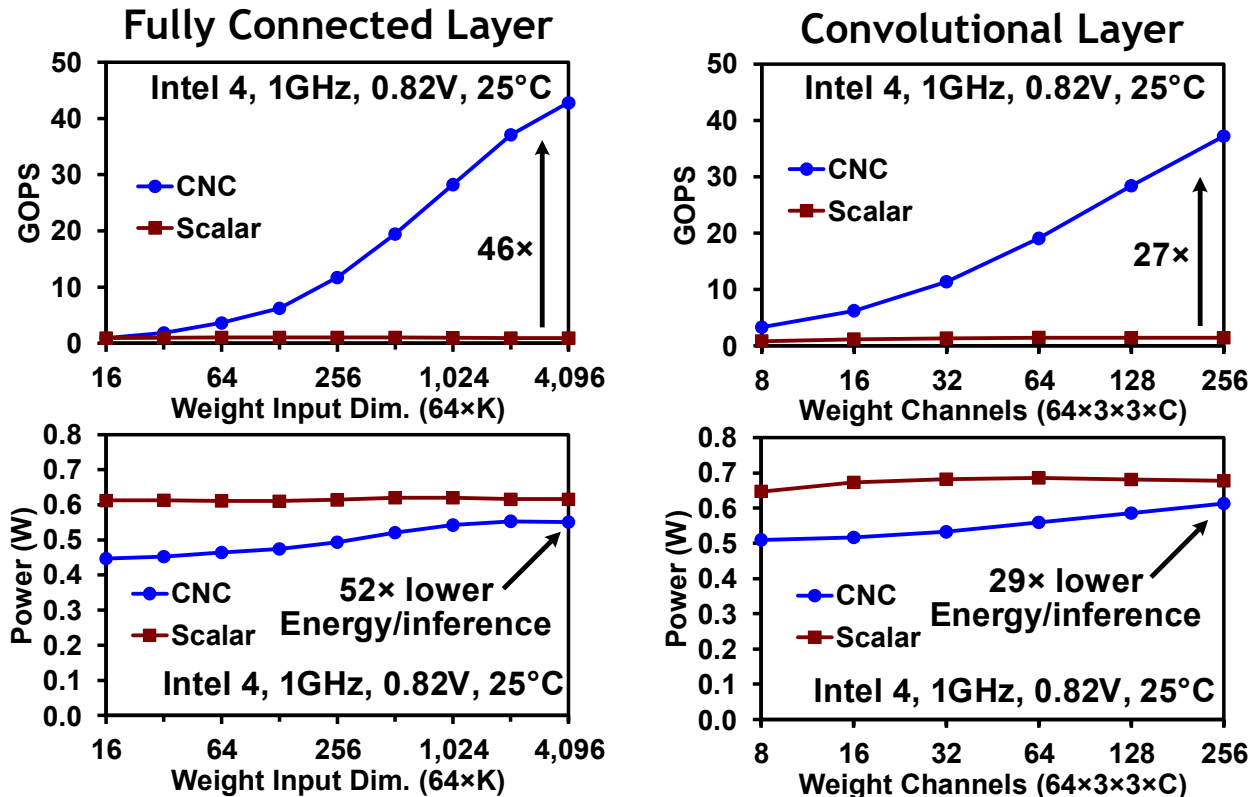
# Intel 4 Silicon Implementation of 8-Core RISC-V



- 1.15GHz Intel 4 test-chip runs programs in C++ with inline CNC and boots Linux

- CNC circuits add 1.4% area overhead over baseline core + LLC design

- Flip-chip packaged with PLL and 32b IO to FPGA chipset

# 8-Core RISC-V DNN Layer Performance in Intel 4



**Fully Connected Layer**

Intel 4, 1GHz, 0.82V, 25°C

- CNC
- Scalar

46×

**Convolutional Layer**

Intel 4, 1GHz, 0.82V, 25°C

- CNC
- Scalar

27×
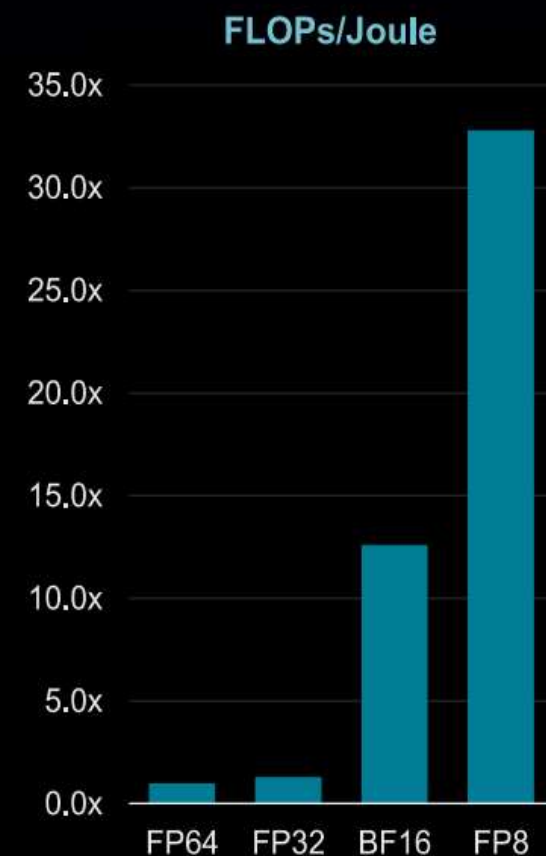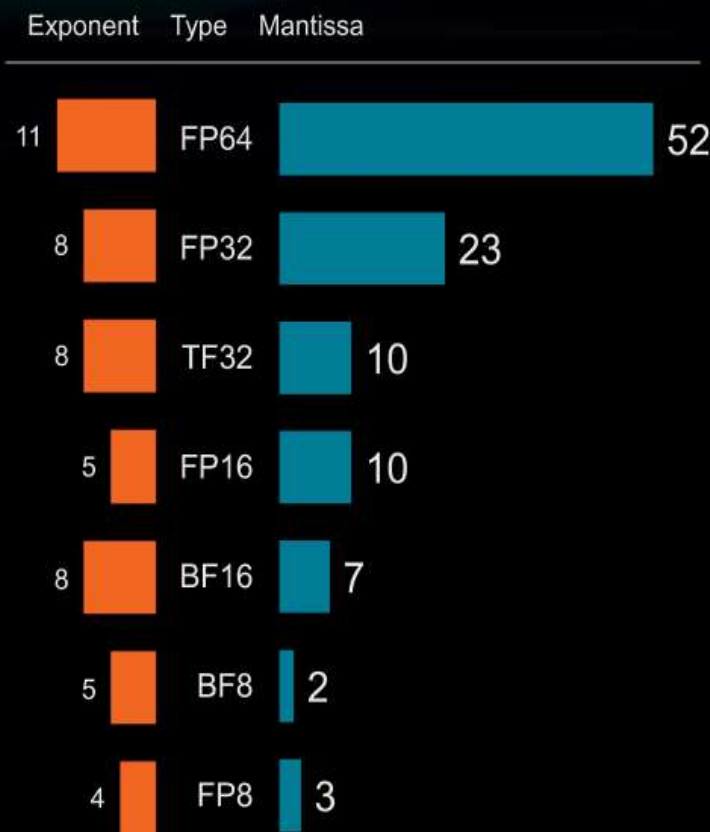
52× lower Energy/inference

29× lower Energy/inference

Measured for dense operation (no sparsity)

- Fully Connected Layers: up to 46× higher performance and 52× lower energy

- Convolutional Layers: up to 27× higher performance and 29× lower energy

# Domain-Specific Computation Enables Workload Optimization which Drives Performance and Efficiency

- Tailor architecture by application

- Adapt algorithms to use lower precision math formats for significant improvements in energy efficiency

| Exponent | Type | Mantissa |
|---|---|---|
| 11 | FP64 | 52 |
| 8 | FP32 | 23 |
| 8 | TF32 | 10 |
| 5 | FP16 | 10 |
| 8 | BF16 | 7 |
| 5 | BF8 | 2 |
| 4 | FP8 | 3 |

**FLOPs/Joule**

Source: L. Su, ISSCC 2023

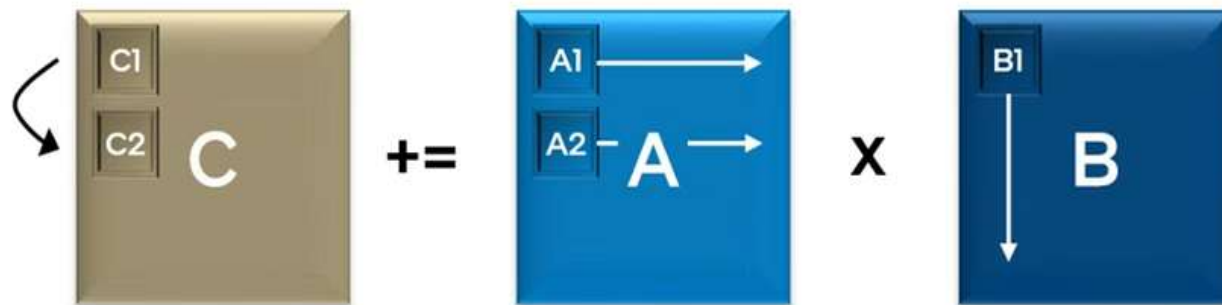# Intel Advanced Matrix Extensions (Intel AMX)

*Tiled Matrix Multiplication Accelerator*

## TILES – Data Structure
- New expandable 2D register file – 8 new registers, 1Kb each
- Supports basic data operators: load/store, clear, set to constant, etc.
- TILES declares state and is OS-managed by XSAVE architecture
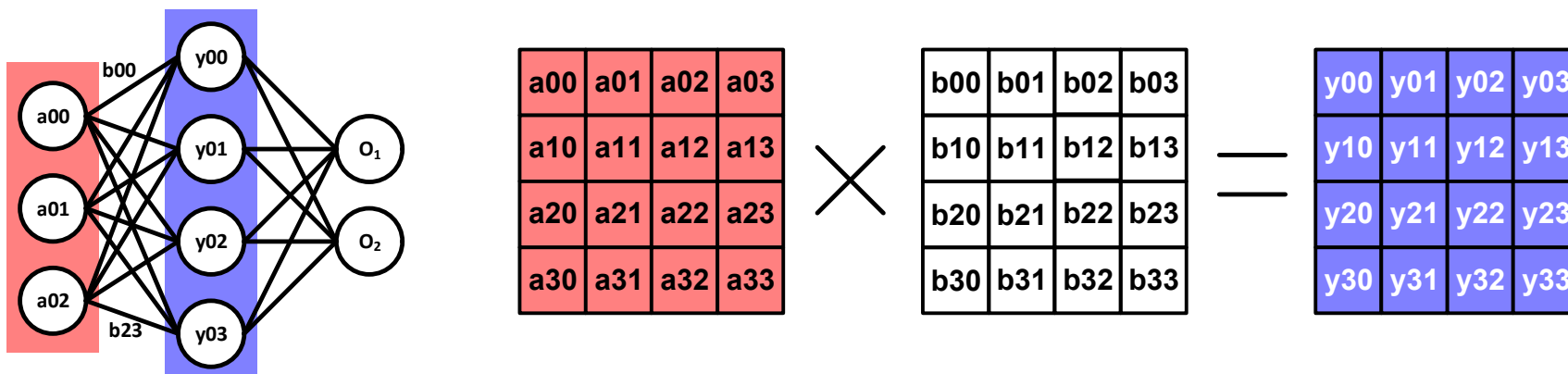
## TMUL – Accelerator Operations
- Set of matrix multiplication instructions, first operators on TILES resgister files
- A MAC computation grid calculates "tiles" of data
- TMUL – performs Matrix ADD-MULTIPLY (C=+A*B) using three Tile register (T2=+T1*T0)
- TMUL requires TILE to be present

C $+=$ A X B

**8x**
operations / cycle / core
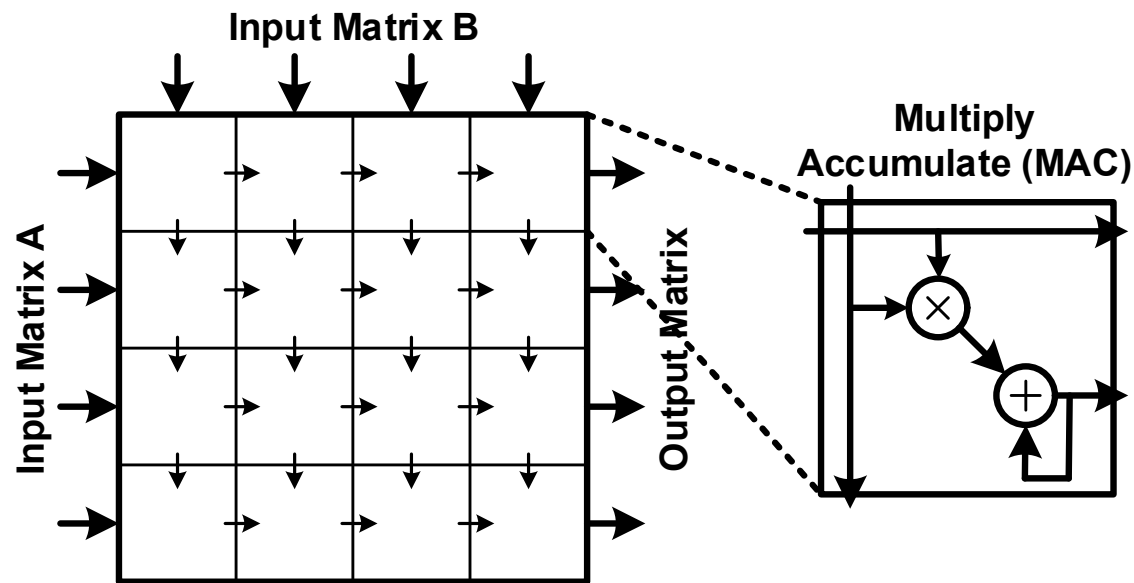relative to VNNI 256 int8

# Multi-precision Neural Networks Matrix Multipliers



**Simple Neural Network**

- Matrix-multiply: power, performance, and area limiter
- Large matrices with many iterations
- Specialized architectures enable higher performance and energy efficiency
- Varying numeric requirements (FP16/INT16/INT8) across applications
  – Require low overhead reconfigurable circuits
- Varying matrix sparsity across applications
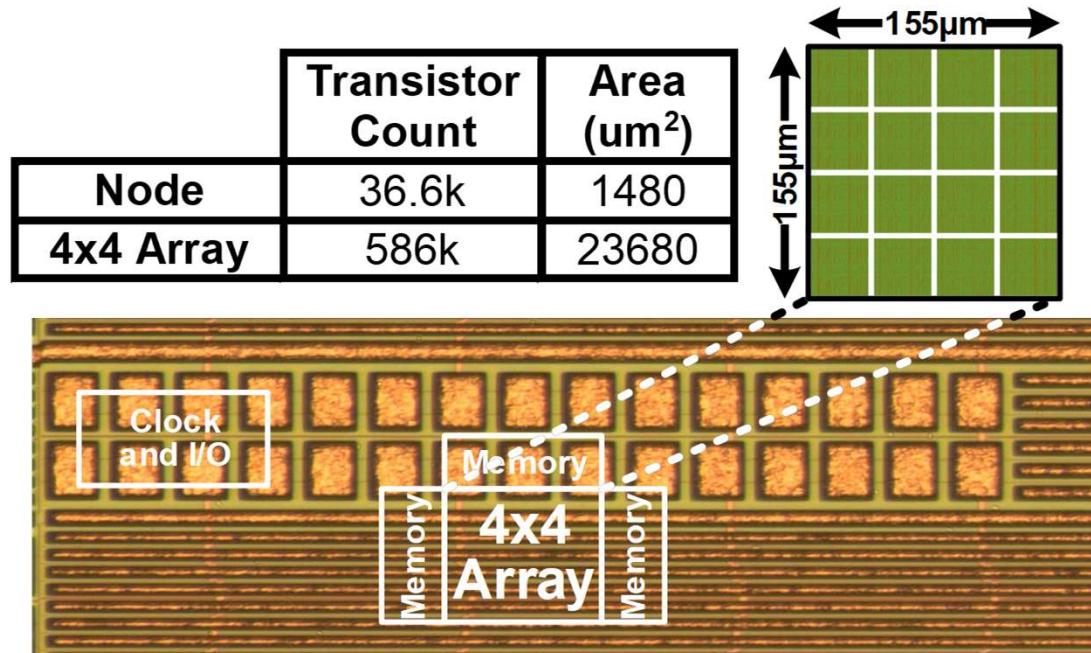  – Optimized circuits can take advantage of sparsity

# Variable Precision Matrix Multiply Accelerator



- 4x4 systolic array
- Fabric reconfigures to optimize data movement in dense/sparse mode
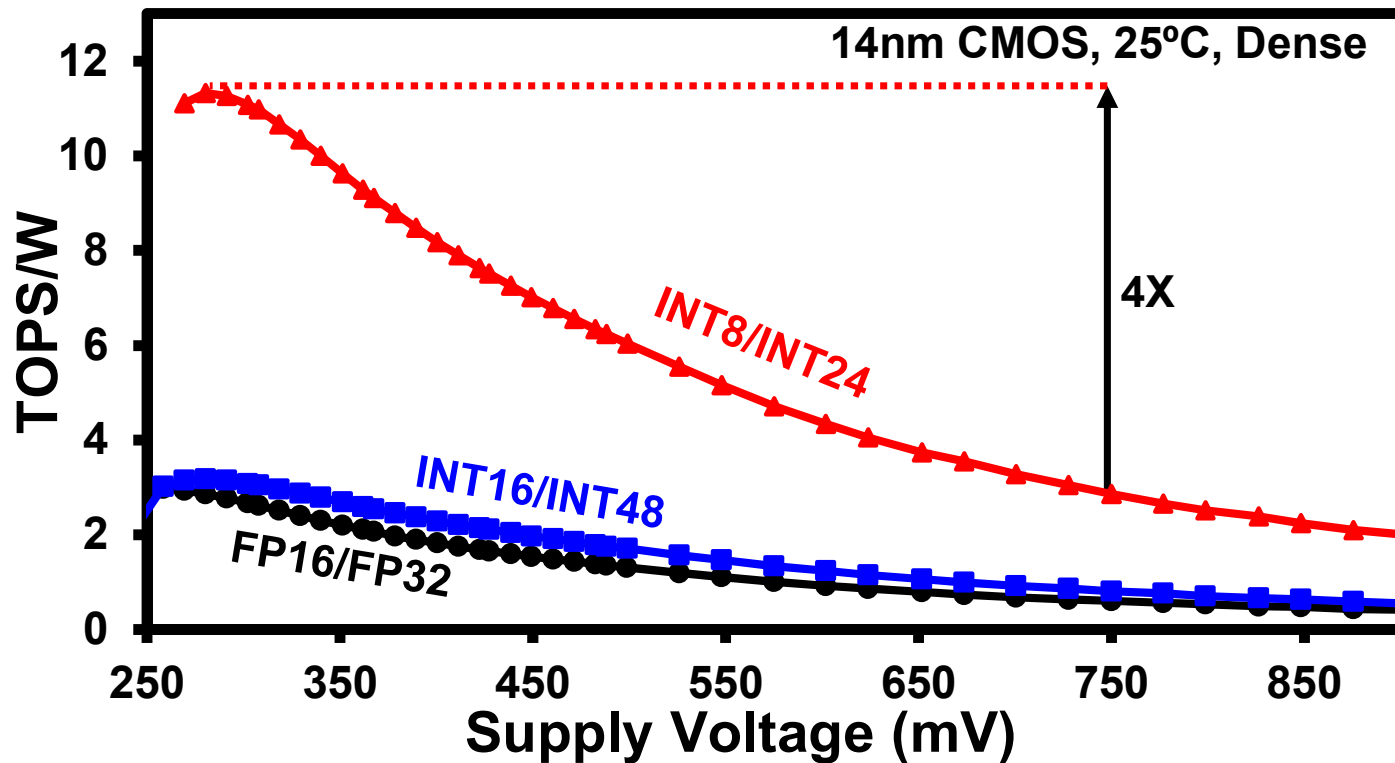- Reconfigurable MAC with signed/unsigned INT16/4xINT8/FP16 support

M. Anders, R. Krishnamurthy et al, VLSI Circuits Symposium 2018
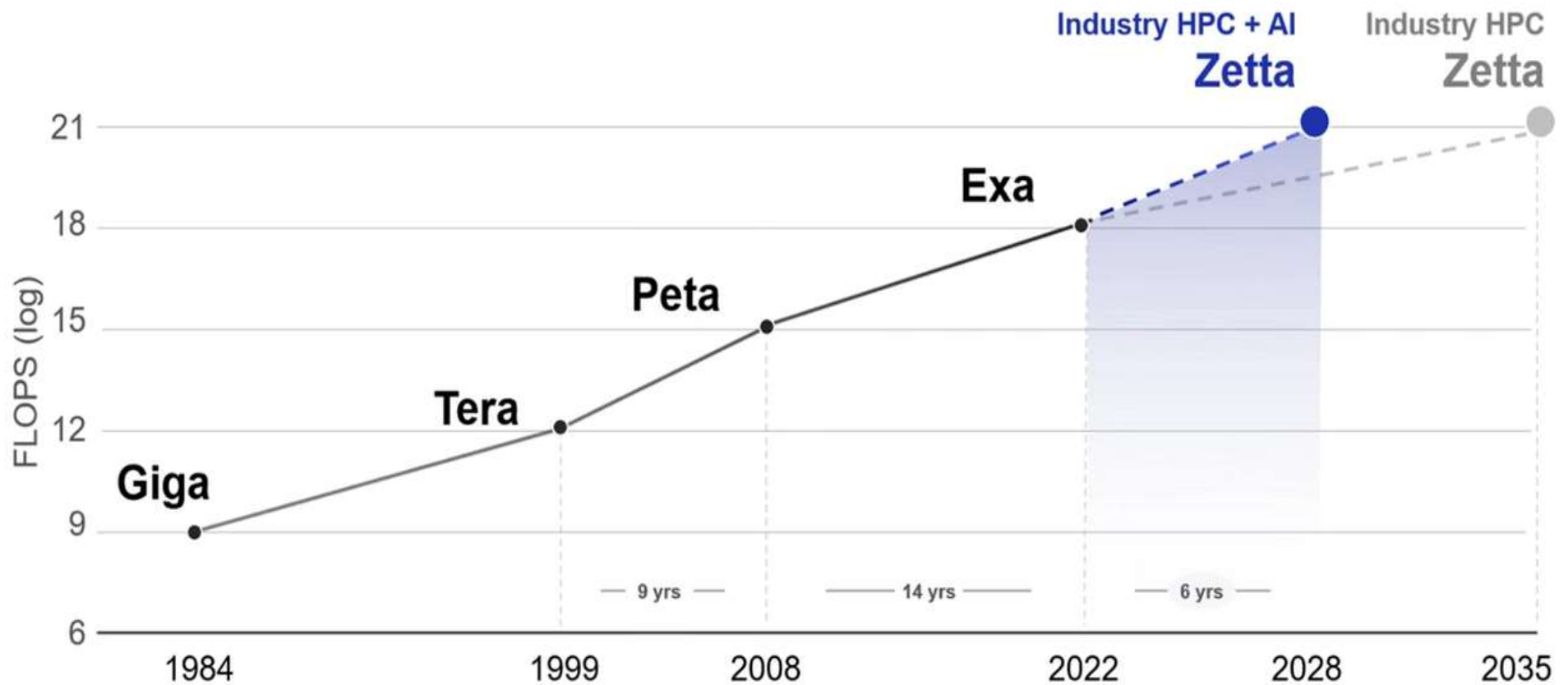
# 14nm Chip Micrograph and Nominal Performance



| | Transistor Count | Area (um$^2$) |
|---|---|---|
| Node | 36.6k | 1480 |
| 4x4 Array | 586k | 23680 |

| Mult/Acc Mode | Nominal (750mV, 25ºC) |
|---|---|
| FP16/FP32 | 800MHz, 42.7mW, 0.6TFLOPS/W |
| INT16/INT48 | 940MHz, 37.6mW, 0.8TOPS/W |
| INT8/INT24 | 1.06GHz, 47.7mW, 2.9TOPS/W |

# Matrix Multiplier Energy Efficiency Measurements
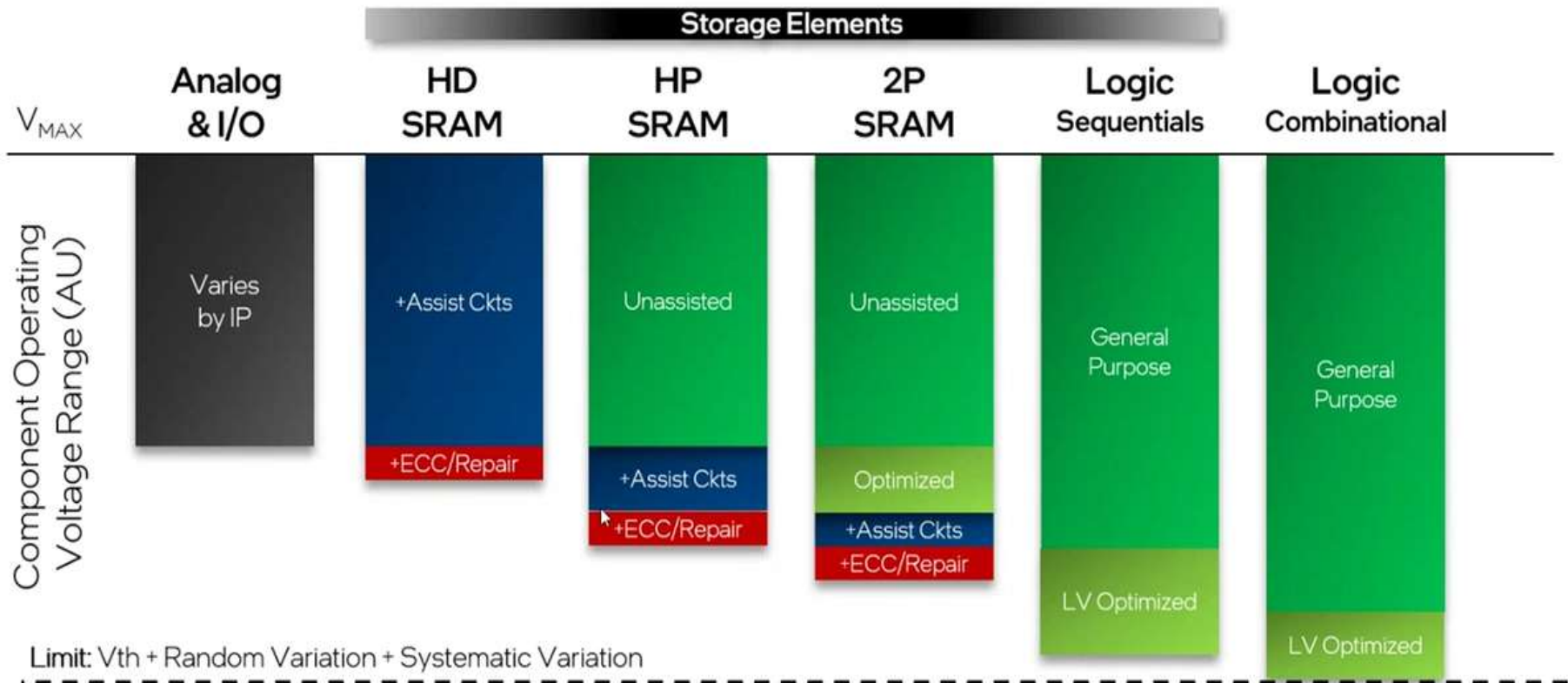


- Efficiency increases 4X from nominal 750mV to near threshold voltage
- Peak energy efficiency range from 2.97TFLOPS/W (FP16) to 11.3TOPS/W (INT8)
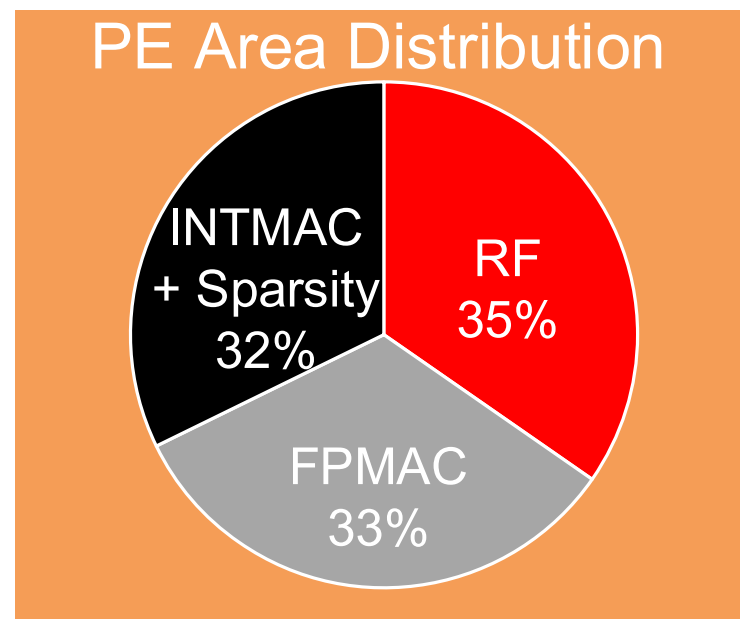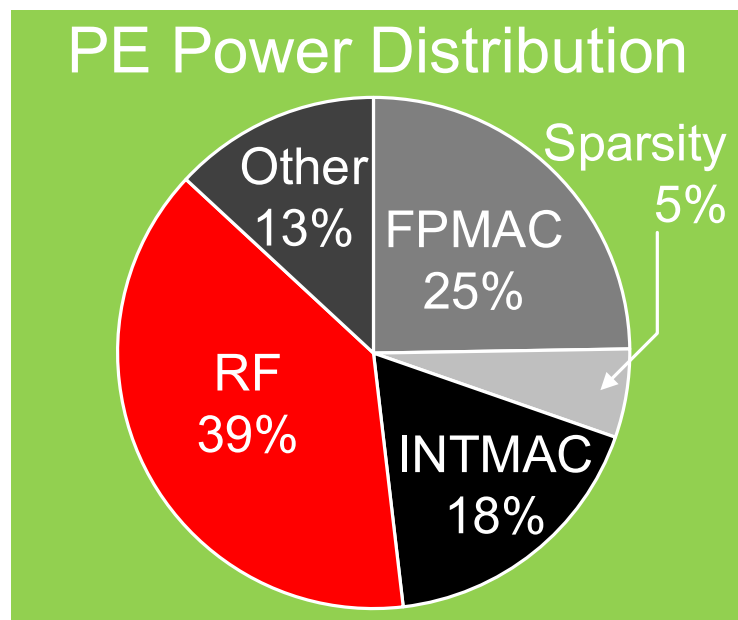
# The Future – Zetta Flop Systems

# Roadmap to Lower Voltage Operation



**Storage Elements**

Component Operating Voltage Range (AU)

$V_{MAX}$

| Analog & I/O | HD SRAM | HP SRAM | 2P SRAM | Logic Sequentials | Logic Combinational |

Analog & I/O: Varies by IP

HD SRAM: +Assist Ckts / +ECC/Repair

HP SRAM: Unassisted / +Assist Ckts / +ECC/Repair

2P SRAM: Unassisted / Optimized / +Assist Ckts / +ECC/Repair

Logic Sequentials: General Purpose / LV Optimized

Logic Combinational: General Purpose / LV Optimized

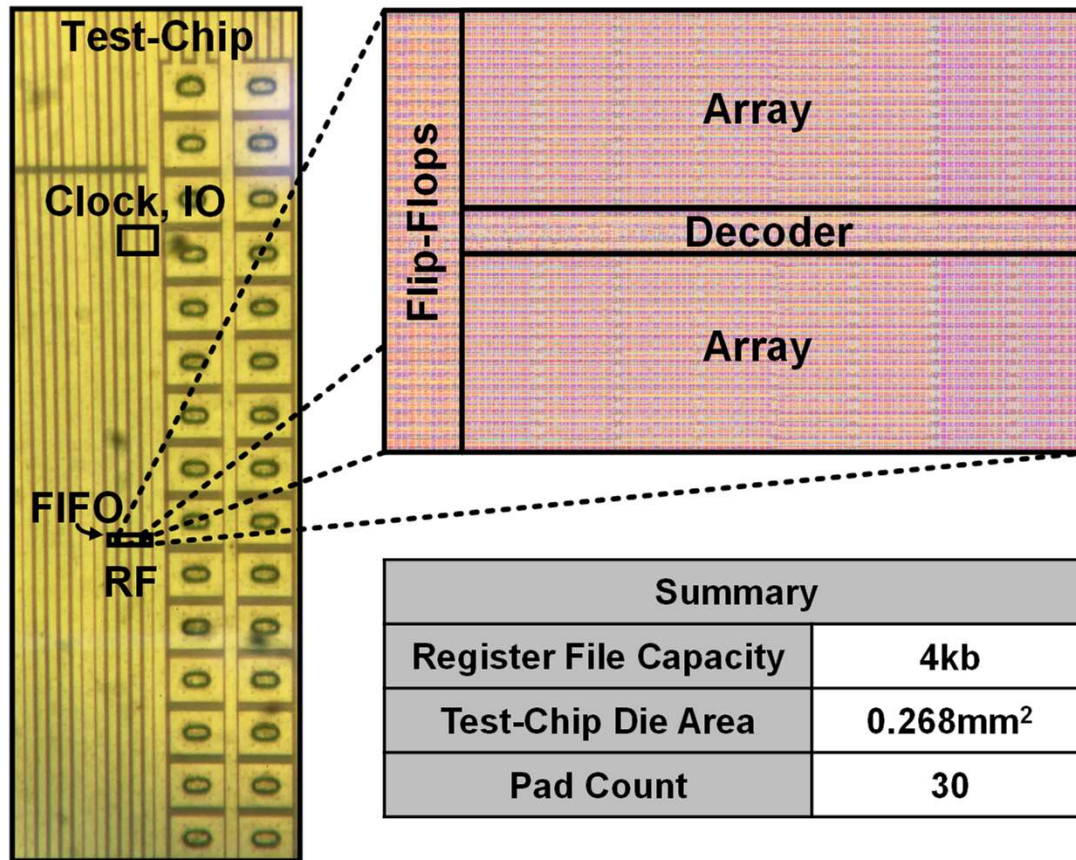Limit: Vth + Random Variation + Systematic Variation

**Low voltage operation requires careful selection and optimization of storage elements**
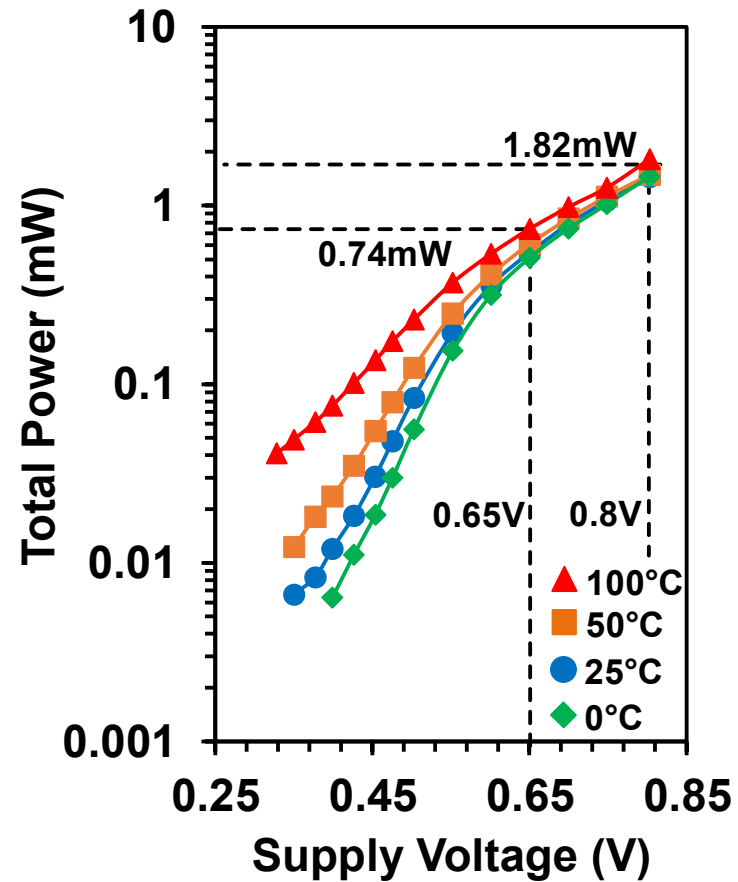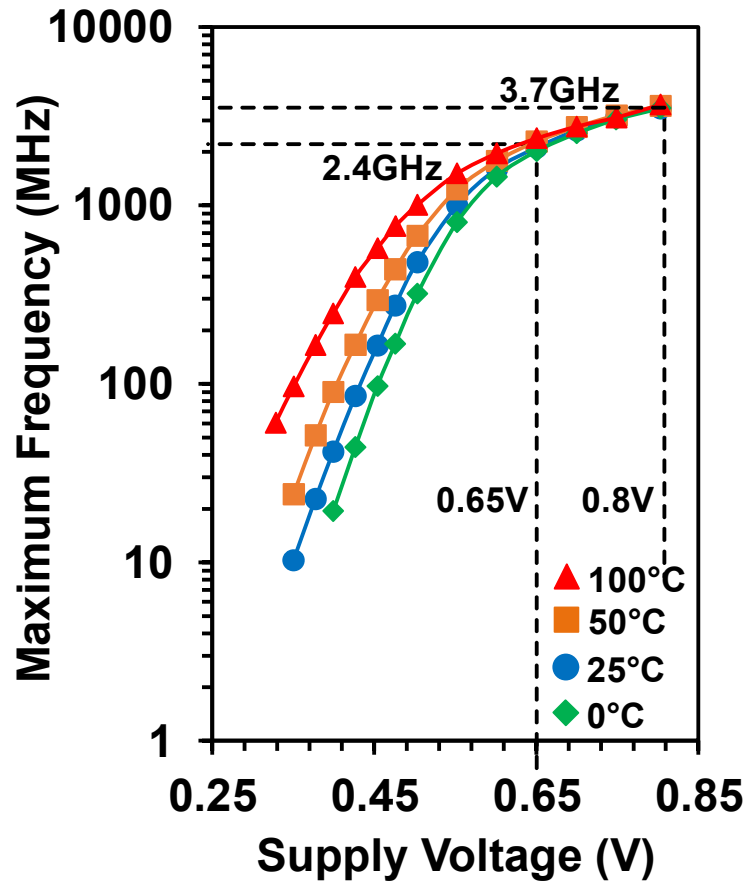
# Memory Challenges in AI Accelerators



- AI accelerators are built using a large array of processing elements (PEs) containing small capacity local register files (RFs)

- Register files contribute a significant amount of power (39%) and area (35%) within the PEs

# Static AI Register File Micrograph



S. Hsu, R. Krishnamurthy et al, IEEE VLSI Circuits Symposium 2022

# Register File Frequency/Power Measurements



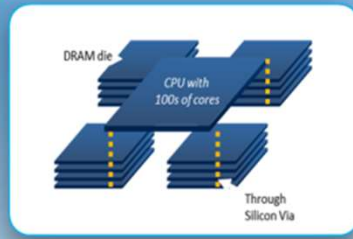- Ultra-low voltage operation at 325mV, 100°C consuming 36.7μW, 60MHz

# "Extreme" efficiency research

Extreme Energy Efficiency

Fine-Grain Power Management

Efficient Memory Subsystem

Self-Aware Computing Operation

Programming for Extreme Parallelism

System-Wide Breakthroughs Needed Across the Board

intel® Look Inside™

intelligence Inside