

High Performance and Energy Efficient Circuit Technologies for Sub-7nm AI Accelerators and In-Memory/Near-Memory Computing

2021 NSF Workshop on Processing In Memory Technology, February 2021

Ram K. Krishnamurthy

Senior Principal Engineer, 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

Motivation: ML in IoT Platforms

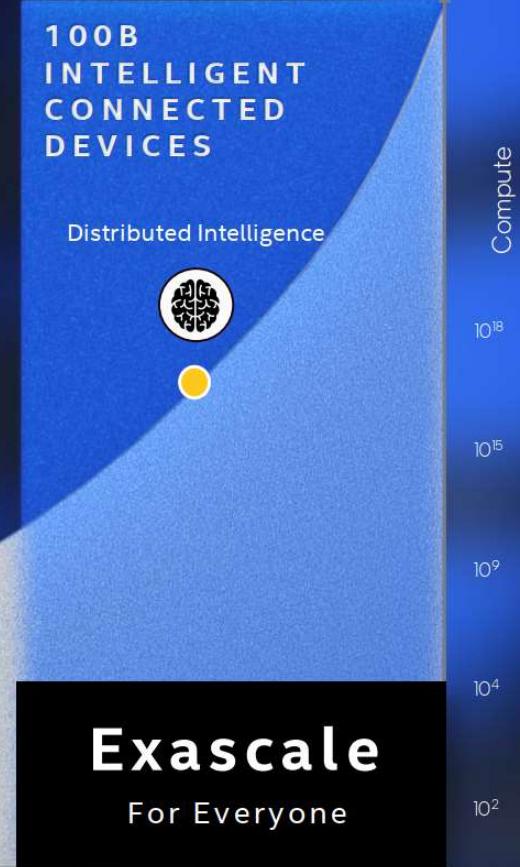
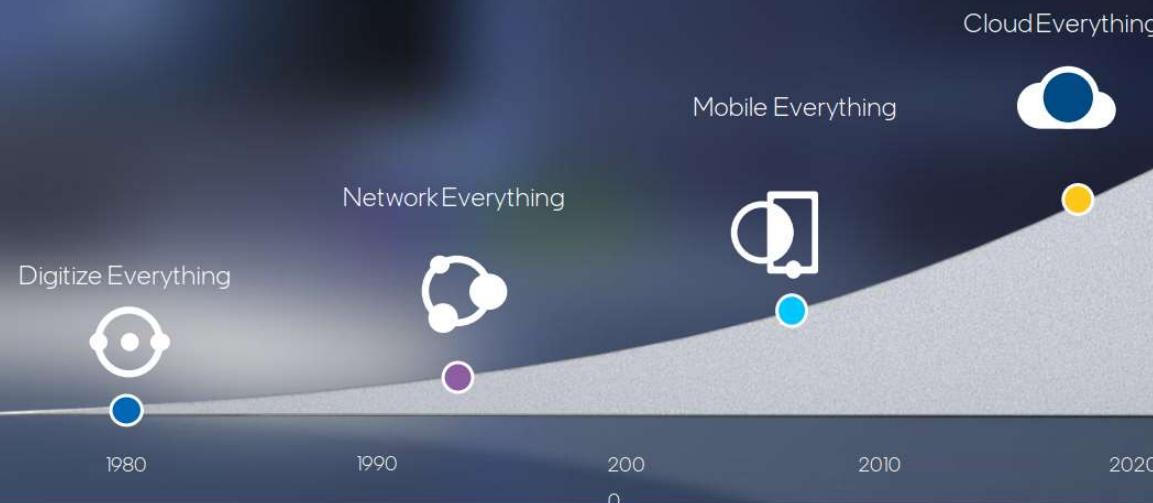


Other names and brands may be claimed as the property of others.

NETWORK + EDGE COMPUTING ACCELERATED BY 5G

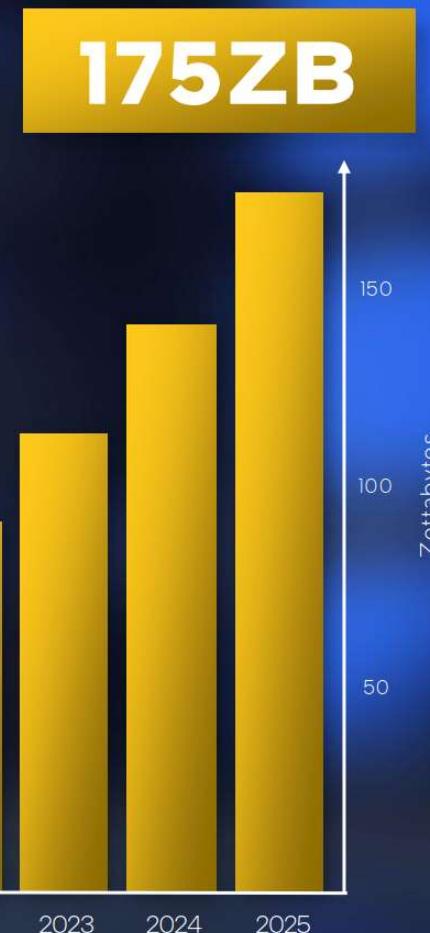


Performance Democratization



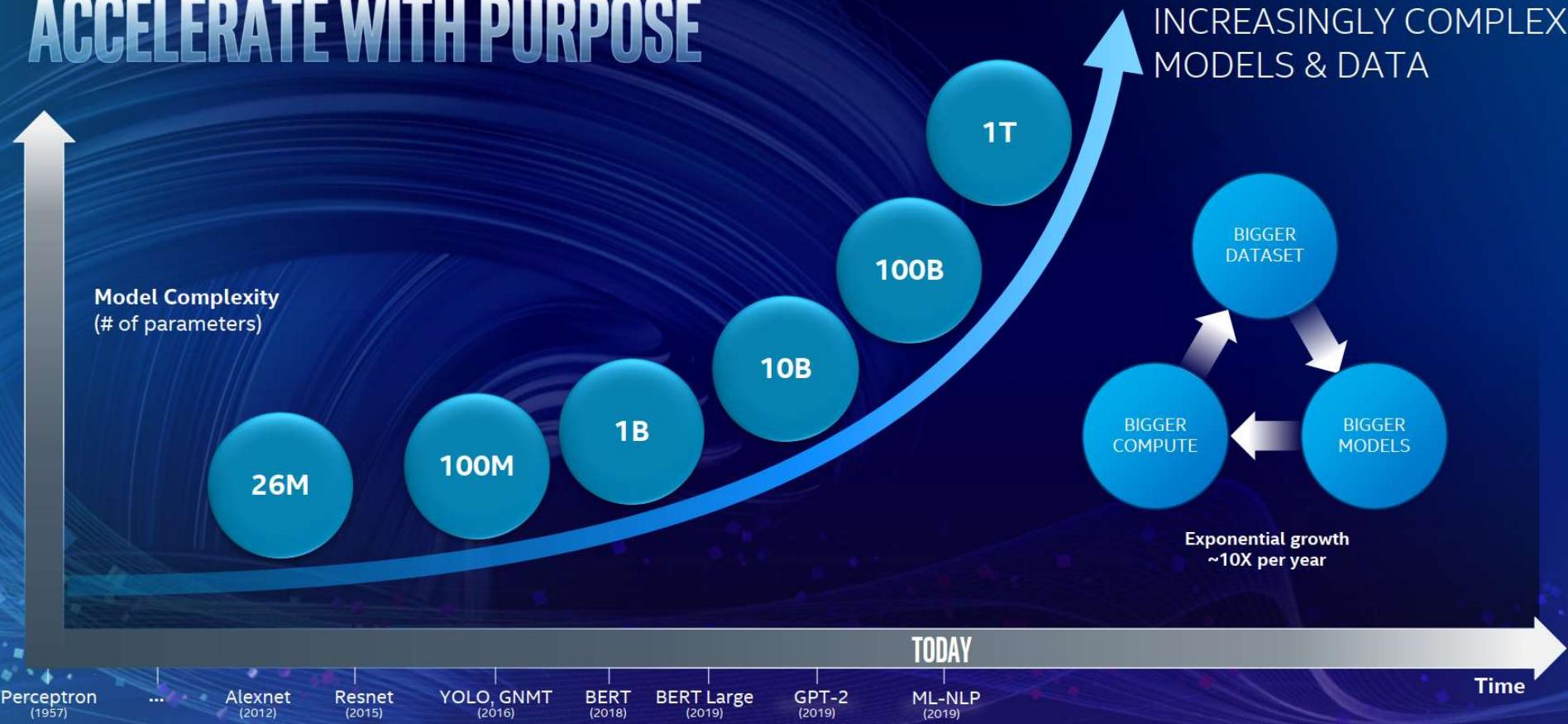
The Data Problem

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



AS MODEL COMPLEXITY GROWS, ACCELERATE WITH PURPOSE

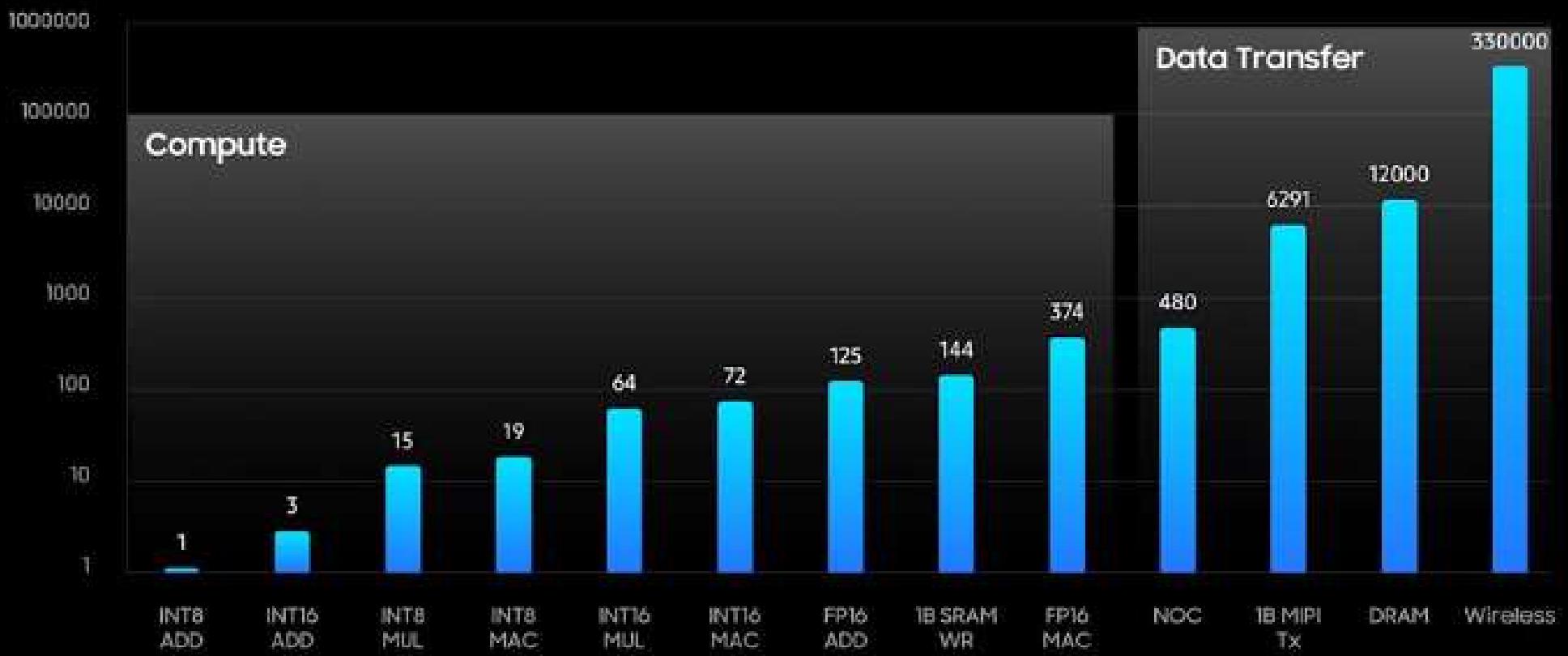
NEXT-GEN DL
INCREASINGLY COMPLEX
MODELS & DATA



AI SUMMIT 2019

Other names and brands may be claimed as the property of others.

Normalized Energy



SOURCE: [facebook](#)

INTRODUCING SECOND GENERATION INTEL® XEON® SCALABLE PROCESSORS



INTEL® XEON®
PLATINUM 9200
PROCESSORS



A NEW CLASS OF
ADVANCED
PERFORMANCE

INTEL® XEON®
PLATINUM 8200
PROCESSORS



INTEL® XEON®
GOLD 6200
PROCESSORS



INTEL® XEON®
GOLD 5200
PROCESSORS



INTEL® XEON®
SILVER 4200
PROCESSORS



INTEL® XEON®
BRONZE 3200
PROCESSORS



BUILT-IN
VALUE

UNINTERRUPTED
LEADERSHIP WORKLOAD
PERFORMANCE

GROUNDBREAKING
MEMORY INNOVATION

EMBEDDED
ARTIFICIAL INTELLIGENCE
ACCELERATION

HW ENHANCED
SECURITY

ENHANCED
AGILITY & UTILIZATION

INTEL.COM/XEONSCALABLE



INTELLIGENCE FOUNDATION



INTEL® XEON® SCALABLE PROCESSORS

2017

1ST GEN
AVX-512

FIRST BUILT-IN AI
ACCELERATION

2019

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*



Cloud Computing

Thermal Budget			Power Source
Memory Capacity			Computing Resource



SOURCE: L. LOH, ISSCC 2020

Diversified Workload & Increasing Demands

0.1 TOPS
1 TOPS/W

Vision Perception

1 TOPS
3 TOPS/W

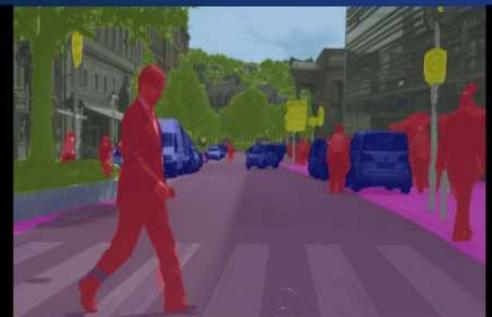
Vision Construction

10 TOPS
10 TOPS/W

Visual Quality

100 TOPS
30 TOPS/W

Multi-Streaming



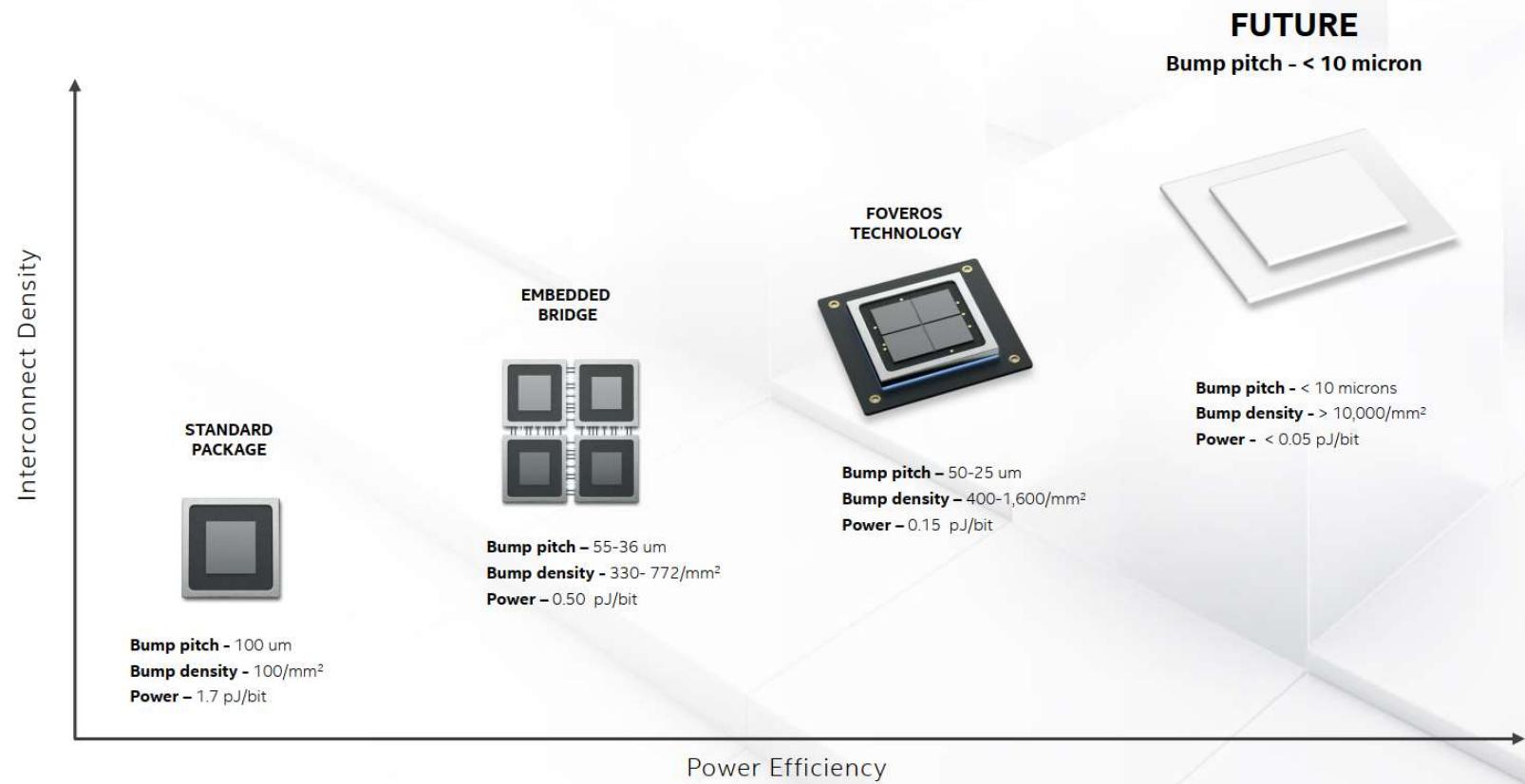
SOURCE: L. LOH, ISSEC 2020

RELENTLESS INNOVATION CONTINUES

Transistor efficiency
(Perf / W)



Packaging Technology Roadmap



TECHNOLOGY
PILARS

Architecture Day **2020**

How can we improve Compute Efficiency by 100x - 1000x?

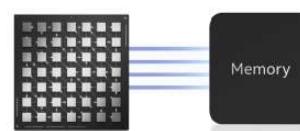
... by thinking differently

... by thinking differently

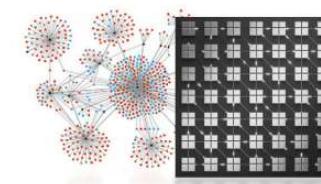
von Neumann Computing



Deep Learning (DNNs)



Neuromorphic Computing



PROGRAMMING BY
ENCODING ALGORITHMS

SYNCHRONOUS
CLOCKING

SEQUENTIAL THREADS
OF CONTROL

OFFLINE TRAINING USING
LABELED DATASETS

SYNCHRONOUS
CLOCKING

PARALLEL
DENSE COMPUTE

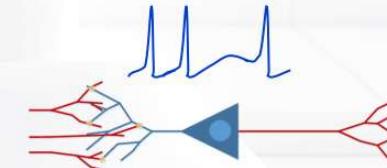
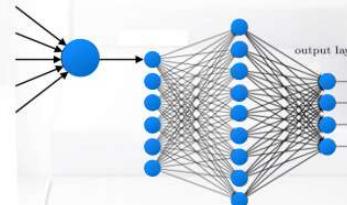
LEARN ON THE FLY THROUGH
NEURON FIRING RULES

ASYNCHRONOUS
EVENT-BASED SPIKES

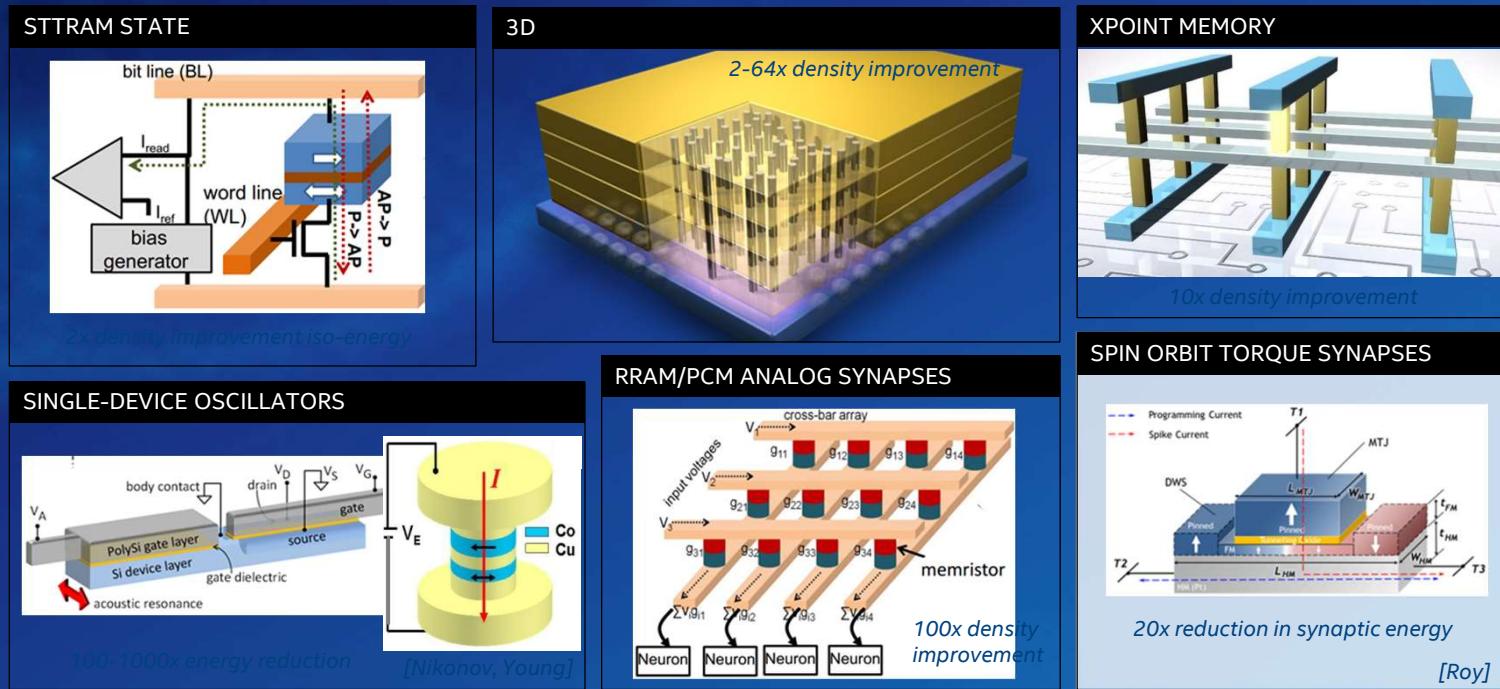
PARALLEL
SPARSE COMPUTE

```
if X then
  ...
else
  ...
...
```

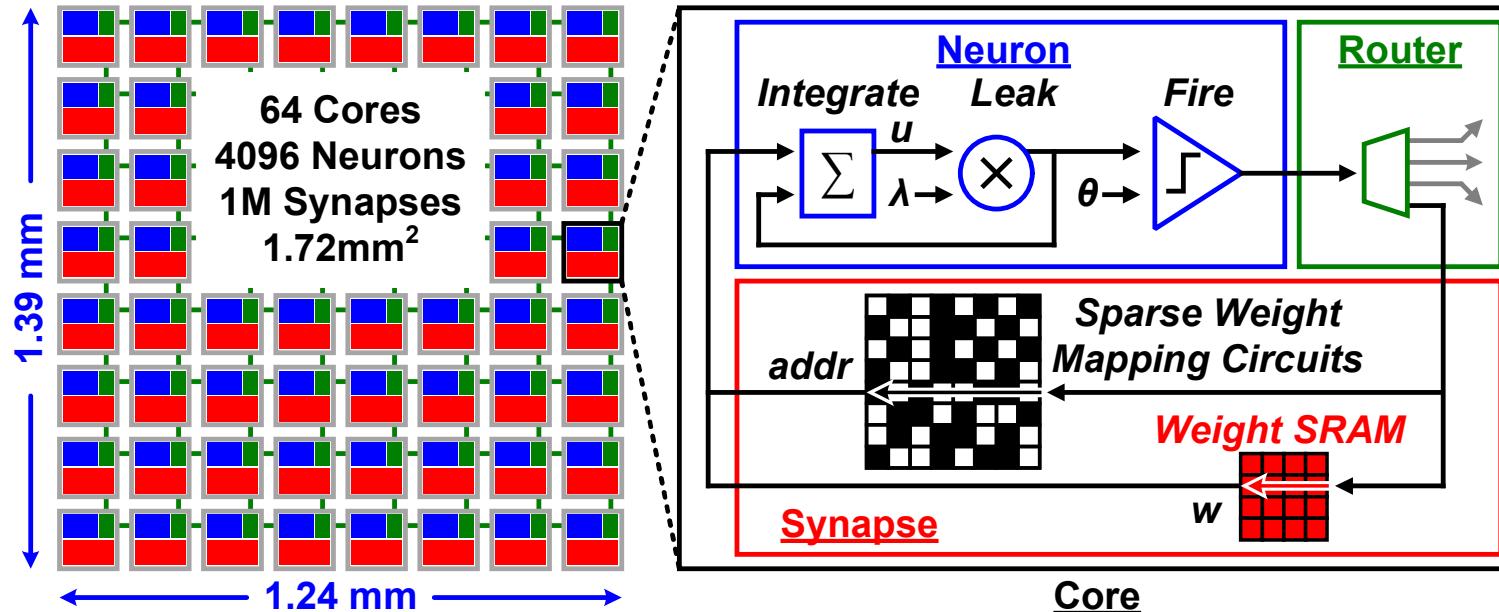
01100
11010
00100



OPPORTUNITIES FOR IN-MEMORY/NEAR-MEMORY PROCESS AND CIRCUIT INNOVATION (BOTH DIGITAL AND ANALOG/MIXED-SIGNAL)



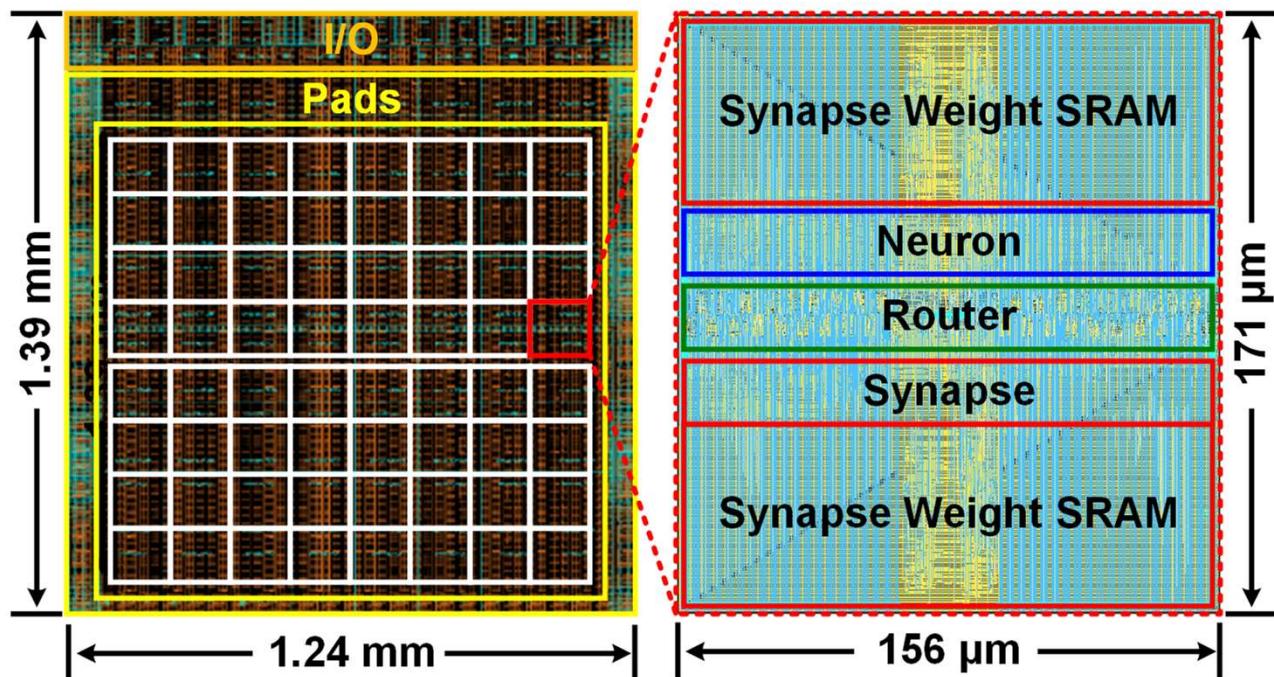
10nm 4096 Neuron Spiking Neural Network Overview



G. Chen, R. Krishnamurthy et al, IEEE Journal of Solid-State Circuits, 2019

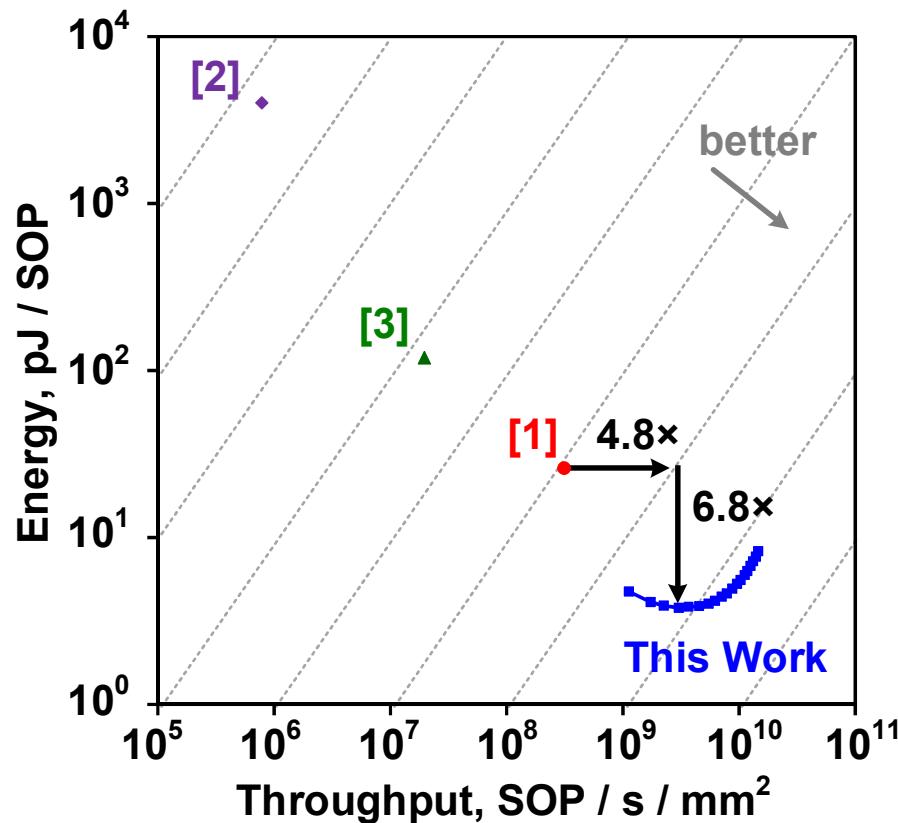
- Neuron stores membrane potentials and performs Leaky Integrate & Fire calculations
- Router multicasts neuron spikes to fan-out synapses in core across the chip
- Synapse stores weights and computes sparse structured connectivity to target neurons

10nm 4096 Neuron Spiking Neural Network Overview



Process	10nm FinFET
Area	1.72 mm ²
Neurons	4096
Synapses	1M × 7b
Throughput	25.2 GSOP/s @ 0.9V
Energy Efficiency	3.8 pJ/SOP @ 525mV
Power / Neuron	2.3 μW @ 450mV

Comparison to Previously Published SNNs

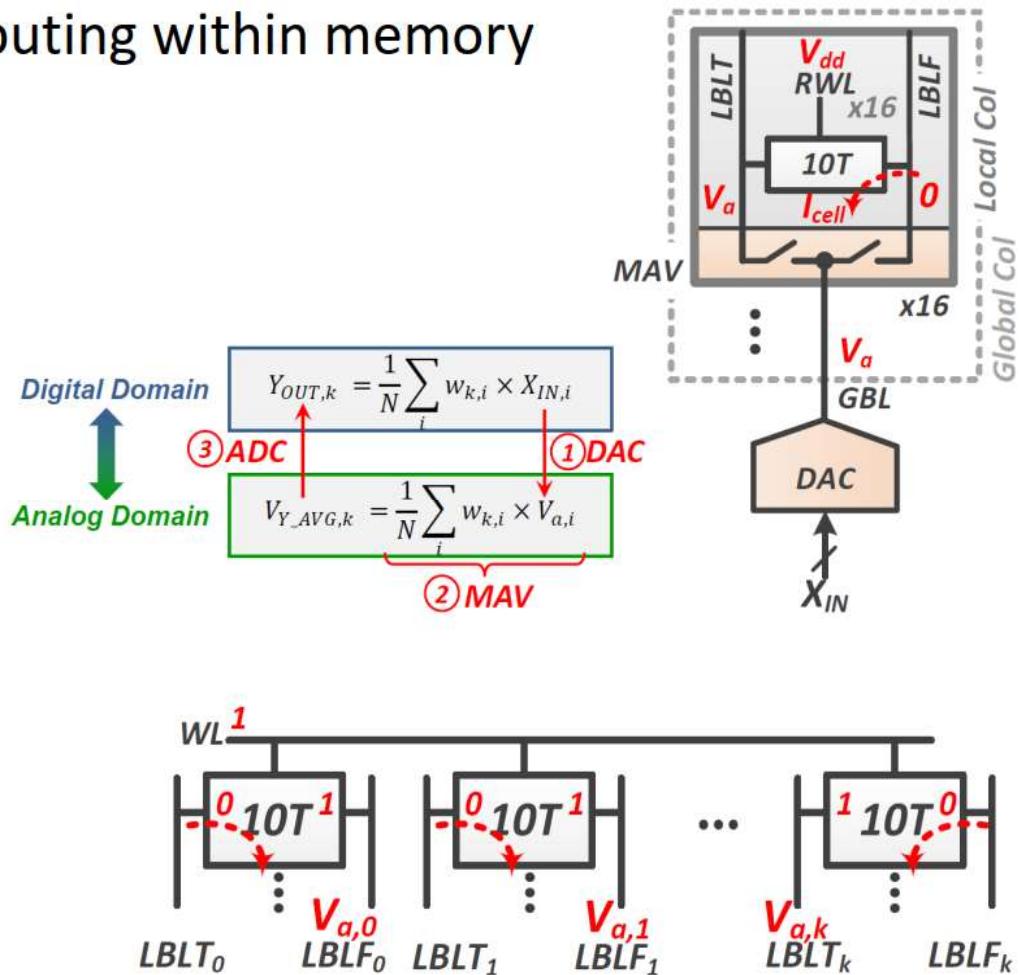
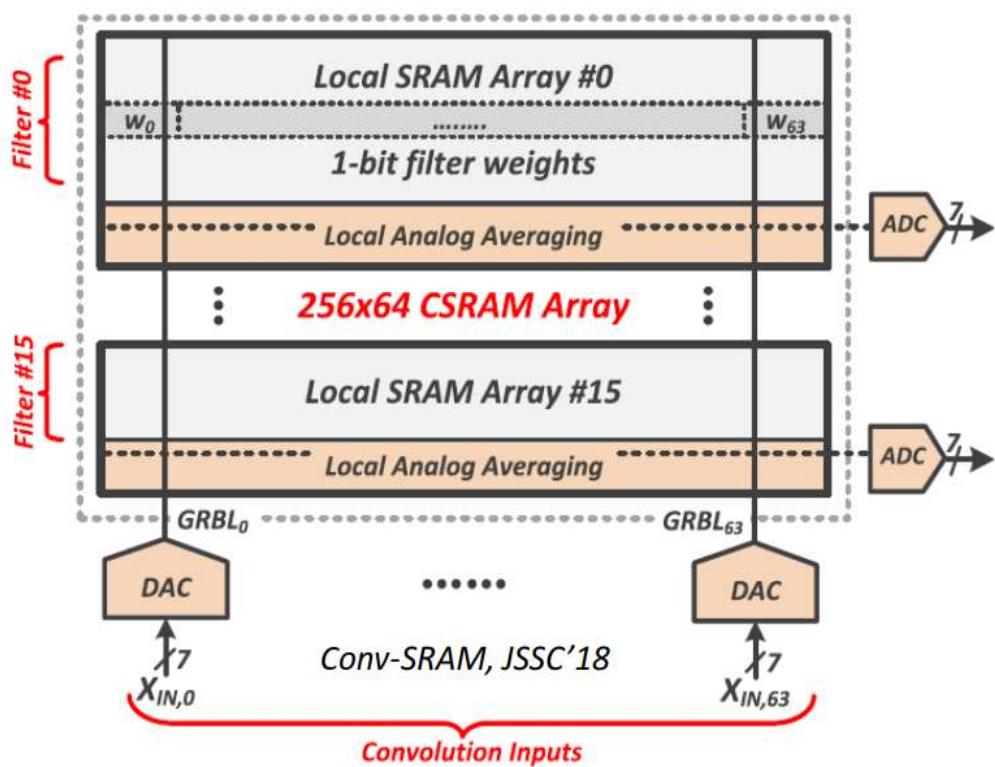


Project	This Work	[1]	[2]	[3]
Process	10nm	28nm	130nm	180nm
Area, mm^2	1.72	430	102	168
Synapse Bits	7b	1b	-	13b
Learning	Yes	No	Yes	No
Sparsity	Yes	No	Yes	No
Voltage	525mV	0.9V	0.775V	1.2V
Freq., MHz	105	506	-	180
SOP/s/ mm^2	3.0G	14.6G	624M	784k
pJ/SOP	3.8	8.3	26	4000
	119			

G. Chen, R. Krishnamurthy et al, IEEE Journal of Solid-State Circuits, 2019

In Memory Computing

- Eliminate weight movements by computing within memory
- Exploit Kirchhoff's law for summation



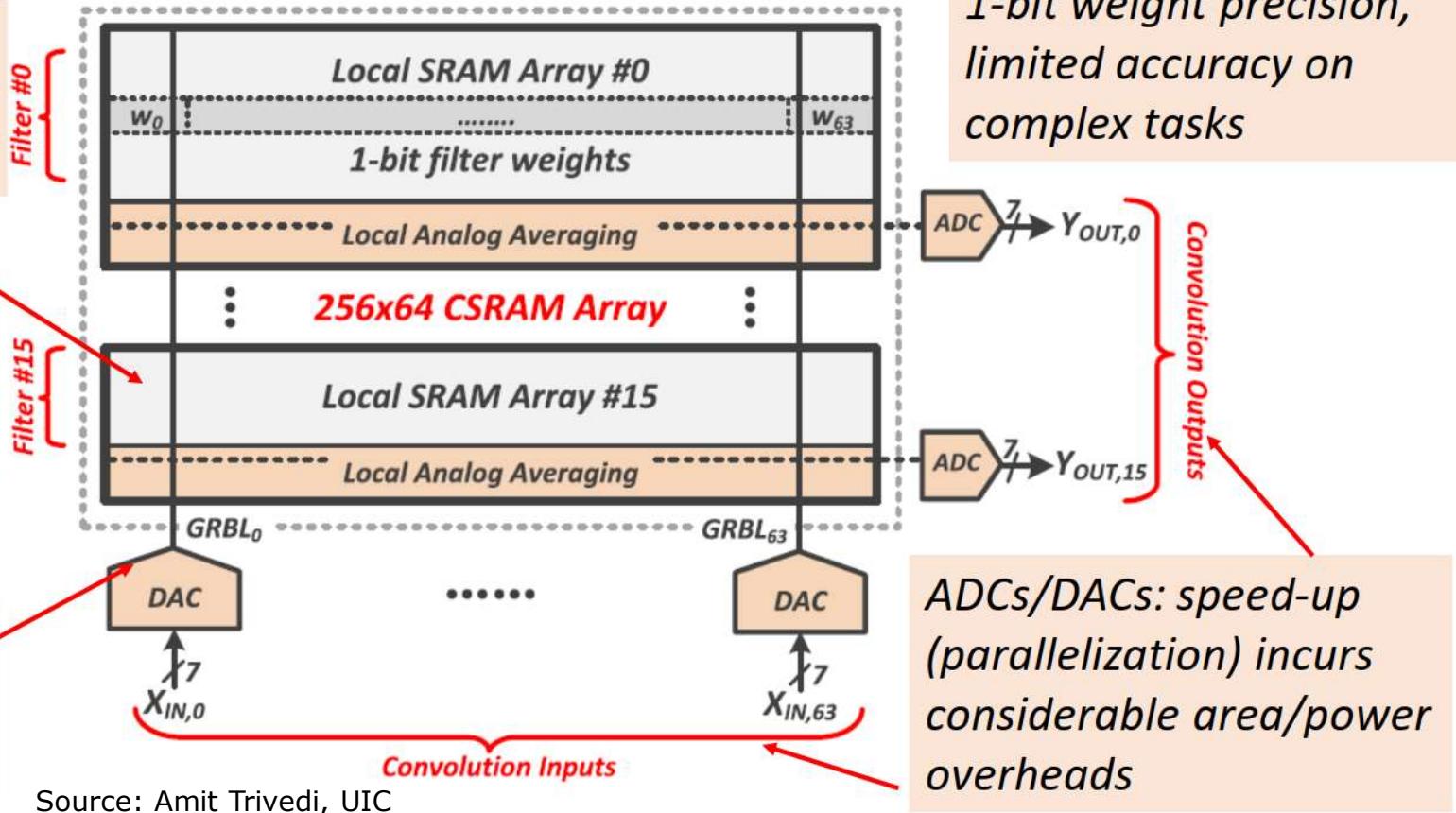
Source: Amit Trivedi, UIC

In Memory Computing Challenges

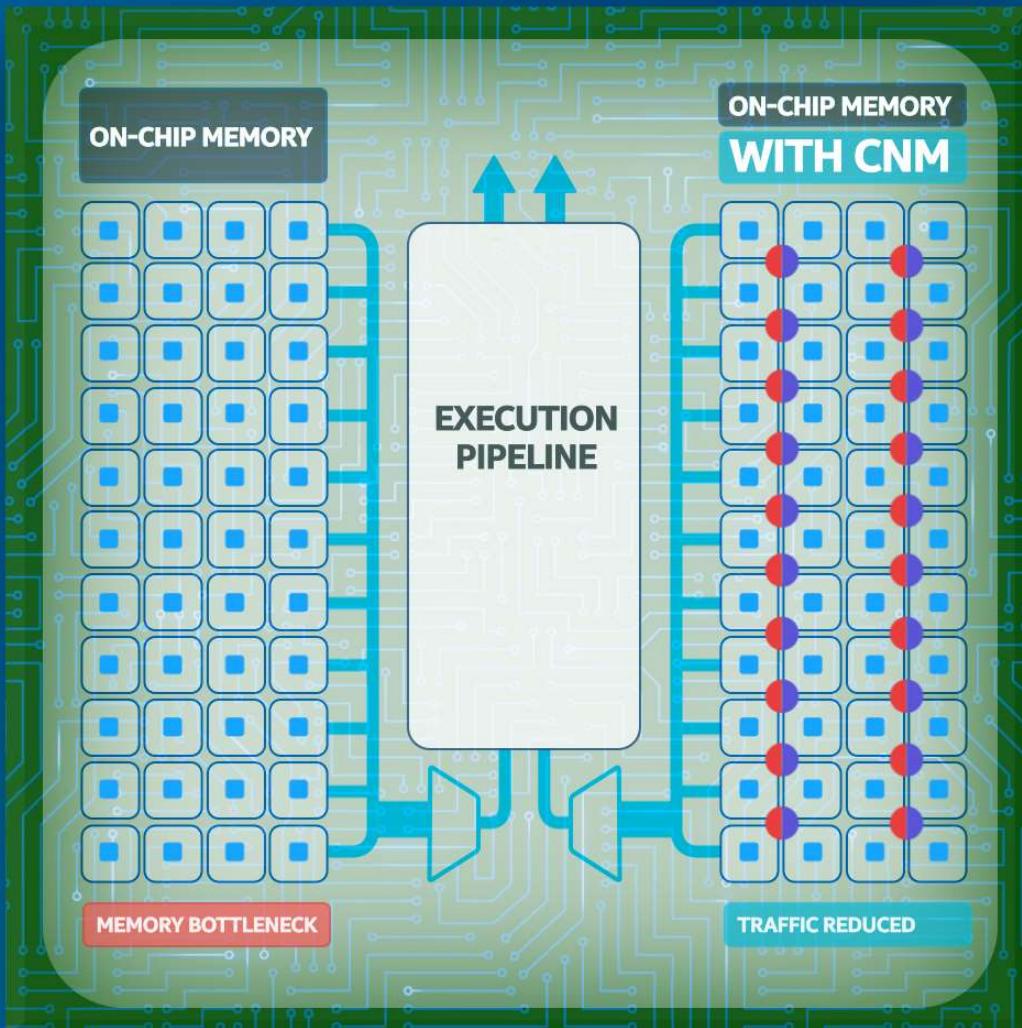
*Inherent inflexibility
in mapping wider
variety of DNNs*

*More area-
expensive
memory cells*

*High susceptibility
to process
variability*



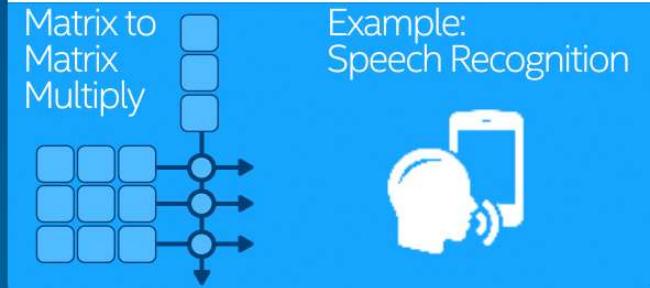
Source: Amit Trivedi, UIC



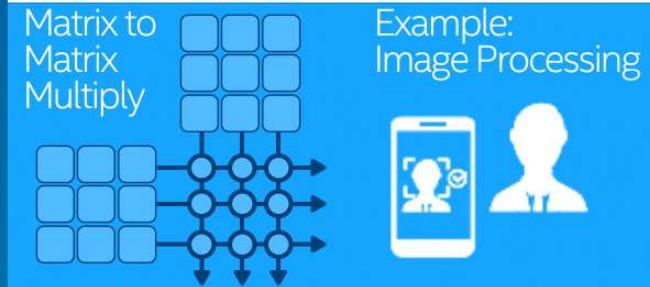
COMPUTE NEAR MEMORY (CNM)

- REDUCE MEMORY BOTTLENECK
- INCREASE PERFORMANCE / WATT
- OPTIMIZE FOR AI OPERATIONS

AI WORKLOAD CONFIGURATION 2

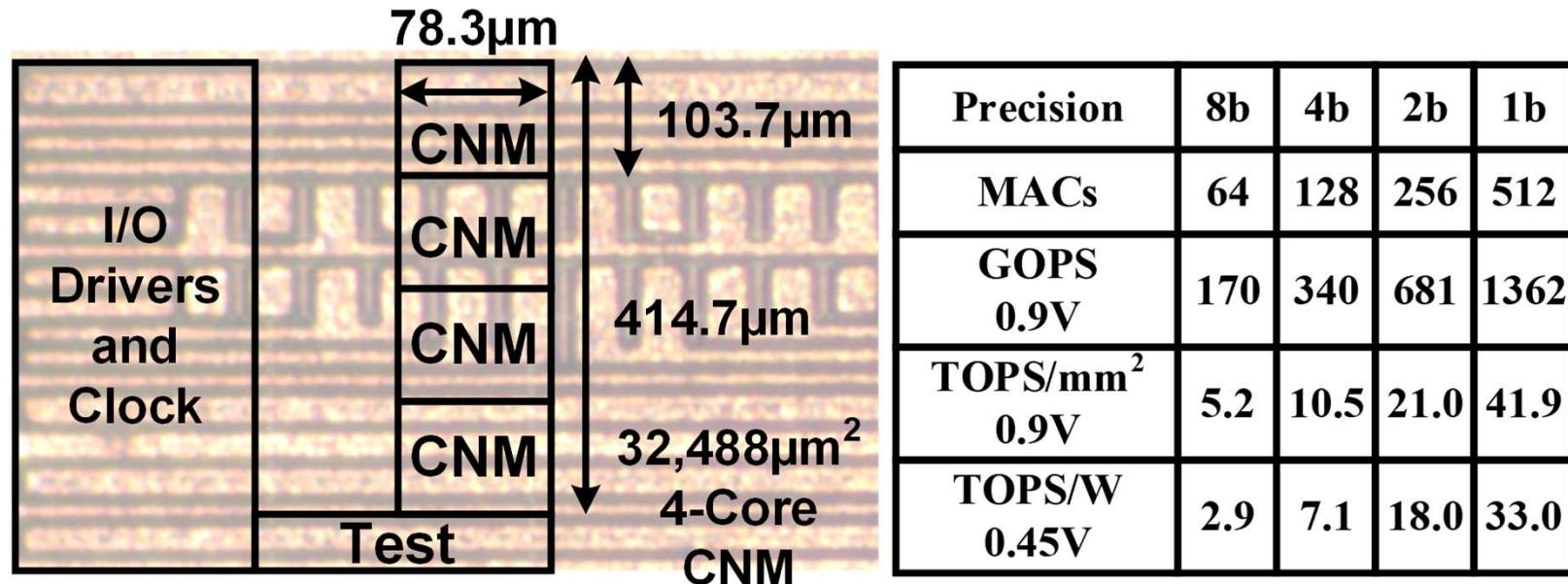


AI WORKLOAD CONFIGURATION 2



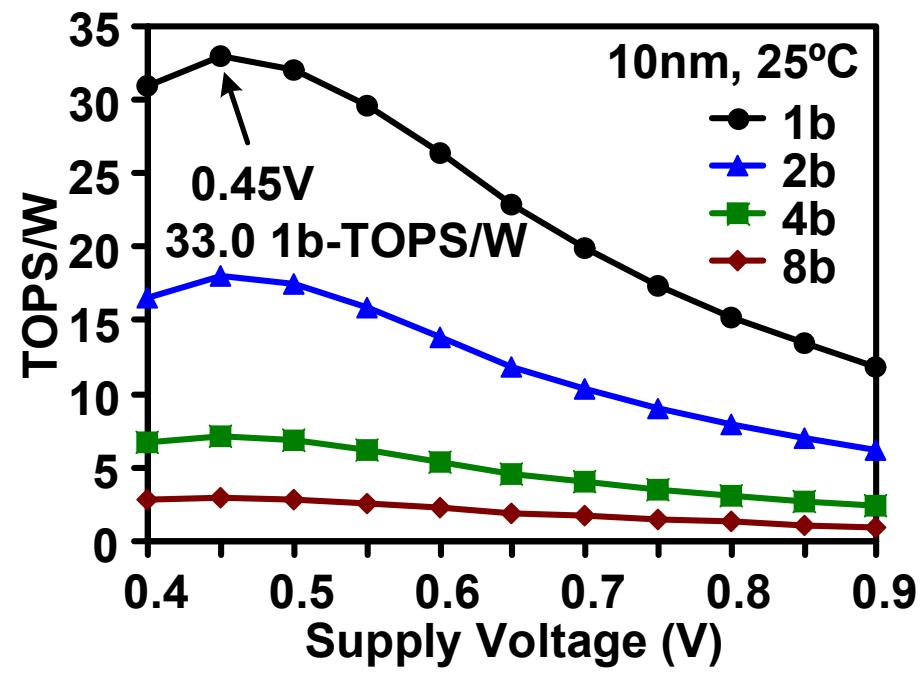
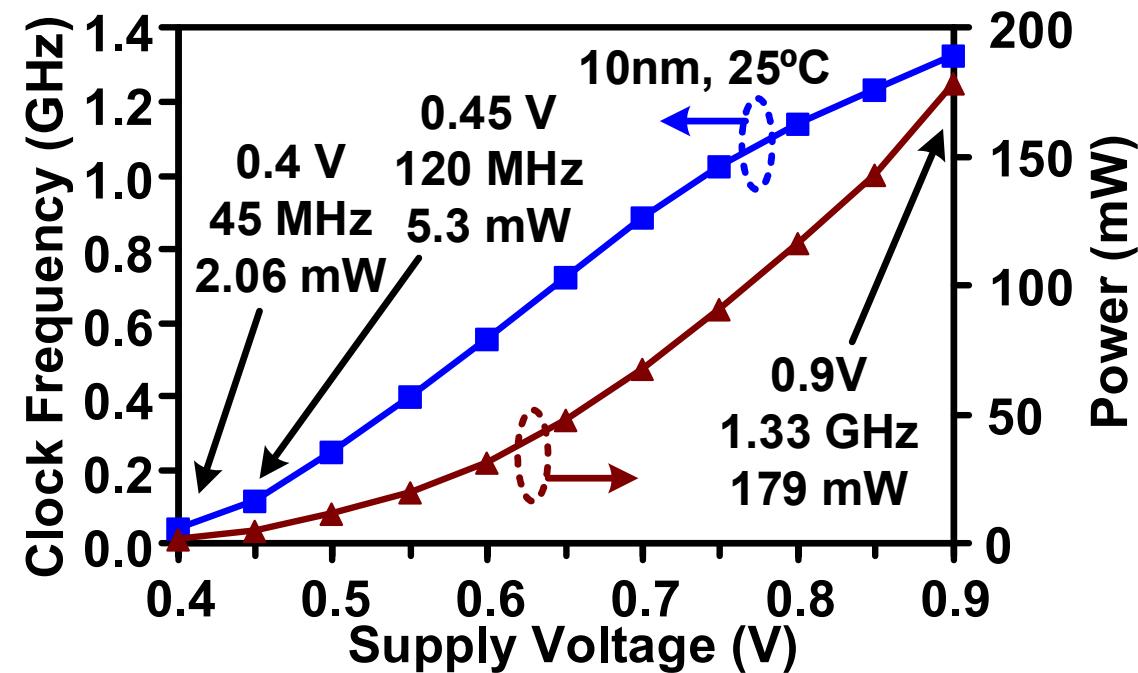
E. Sumbul, R. Krishnamurthy et al, IEEE Solid-State Circuits Letters, 2020

10nm Near Memory Computing AI Accelerator



- 4 CNM cores with 8KB of weight memory and 64 8b multipliers
- Supports memory-intensive batch-1, large-batch, and in-place convolution

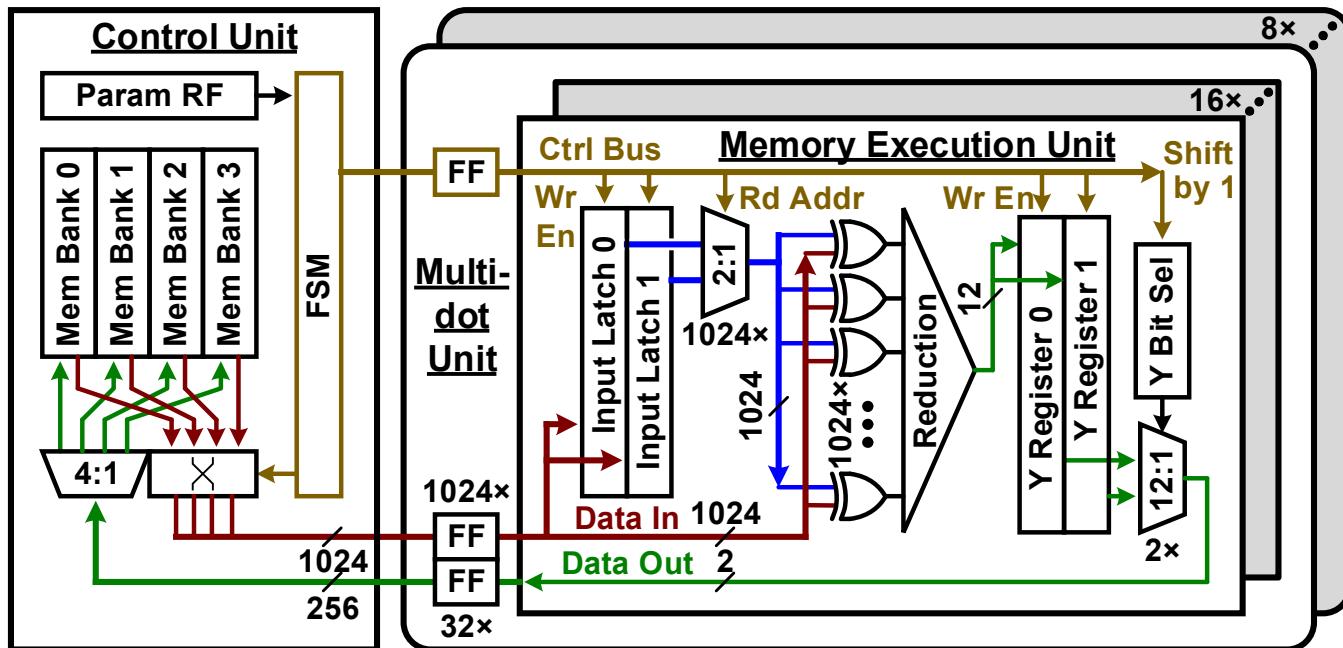
10nm Near-Memory Computing Measurement Results



- Peak throughput of 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

E. Sumbul, R. Krishnamurthy et al, IEEE Solid-State Circuits Letters, 2020

10nm Binary Neural Network Accelerator Overview

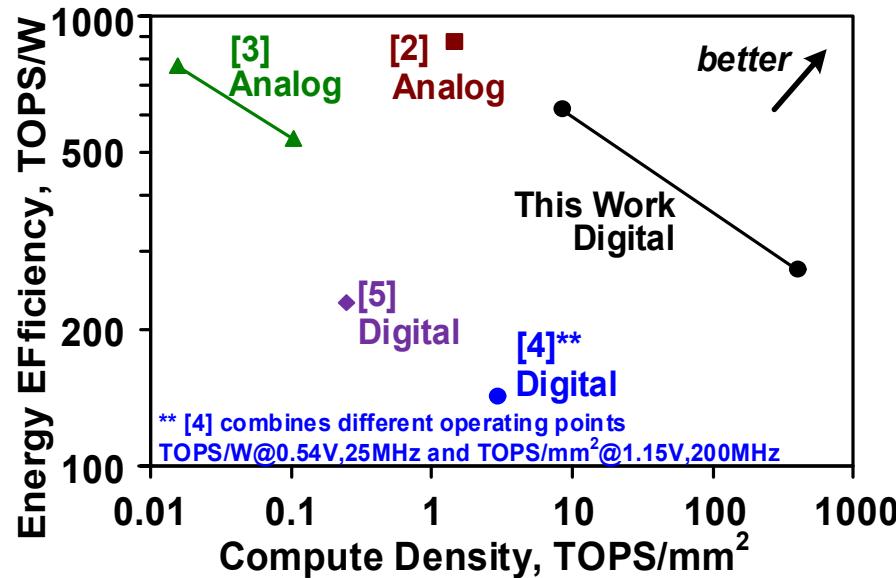


- Array of 128 Memory Execution Units (MEU) combine latch base memory and inner product compute in fine grain manner to minimize interconnect energy
- Central controller manages data flow from four 256b memory banks to MEUs
- 2 latch words per MEU enables data reuse reducing input bandwidth by 2x

P. Knag, R. Krishnamurthy et al, IEEE VLSI Circuits Symposium, 2020

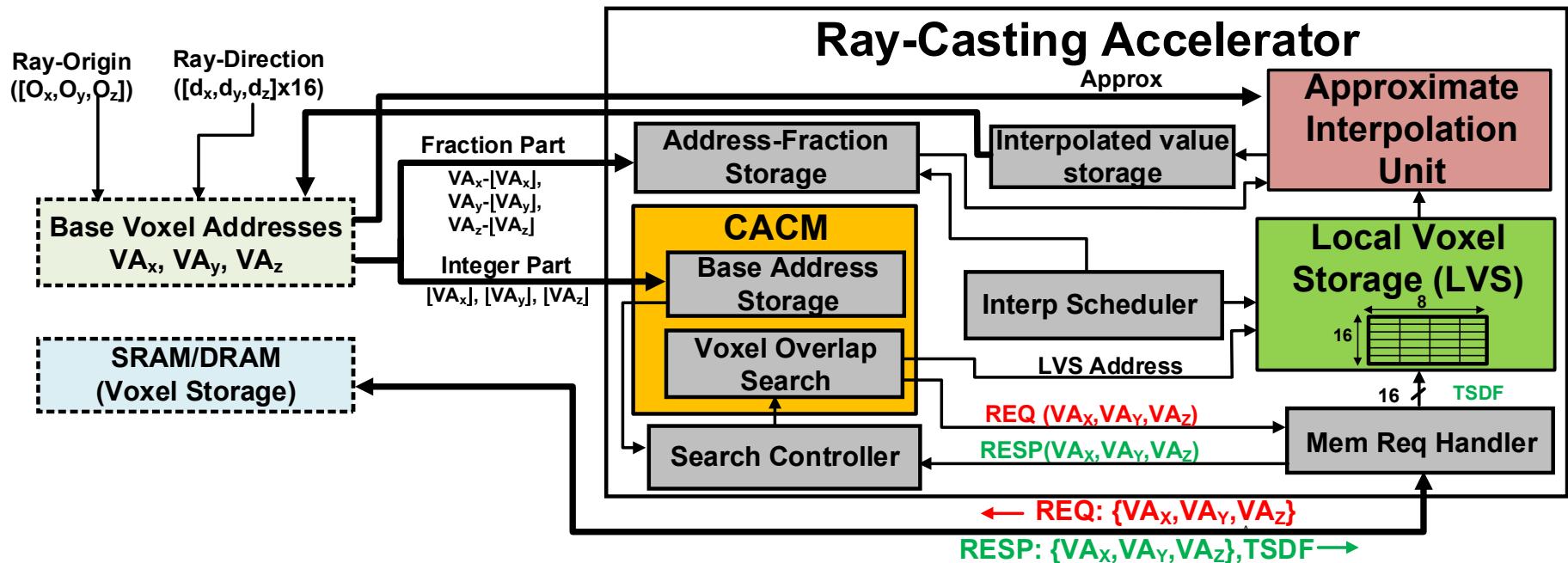
Comparison to Previously Published BNNs

	[2]	[3]	[4]	[5]	This Work	
Digital or Analog	Analog	Analog	Digital	Digital	Digital	
Technology (nm)	65	28	65	28	10	
Area (mm ²)	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/mm ²	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	-
Frequency (MHz)	100	1.5	10	25	200	6.0
Power (mW)	22	0.094	0.899	-	-	1.5
TOPS	19	0.072	0.478	-	14.9	0.35
TOPS/W	866	772	532	140	-	230
TOPS/mm ²	1.5	0.02	0.10	-	3.1	0.25
						418



- Industry-leading peak energy efficiency of 617TOPS/W
- 418 TOPS/mm² area efficiency

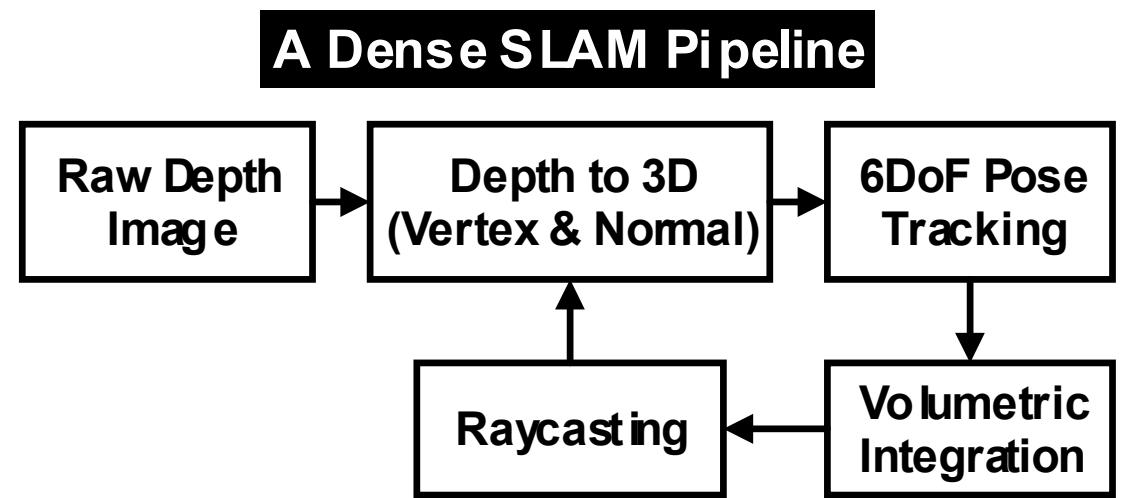
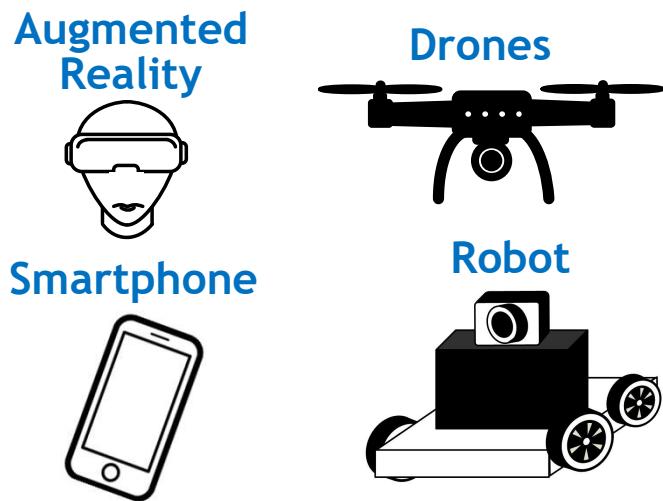
SLAM Ray-Casting Accelerator for Edge AI



- Content addressable compute memory (CACM) searches spatially overlapping voxels for rays of different pixels
- Interpolation unit with exact and approximate computation modes

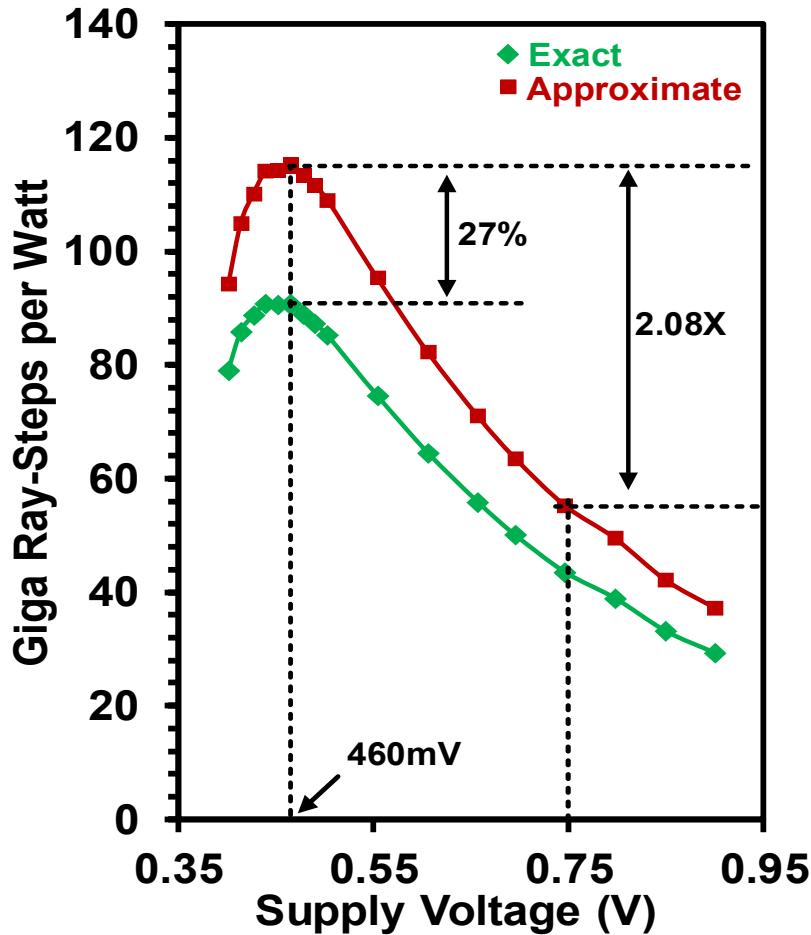
S. Hsu, R. Krishnamurthy et al, IEEE VLSI Circuits Symposium, 2020

Applications with 3D Scene Reconstruction



- Wide range of edge applications demand accurate volumetric reconstruction of surrounding scenes
- Dense Simultaneous Location and Mapping (SLAM) required for accurate surface estimation and 3D scene reconstruction

10nm Ray Casting Energy-Efficiency Measurements

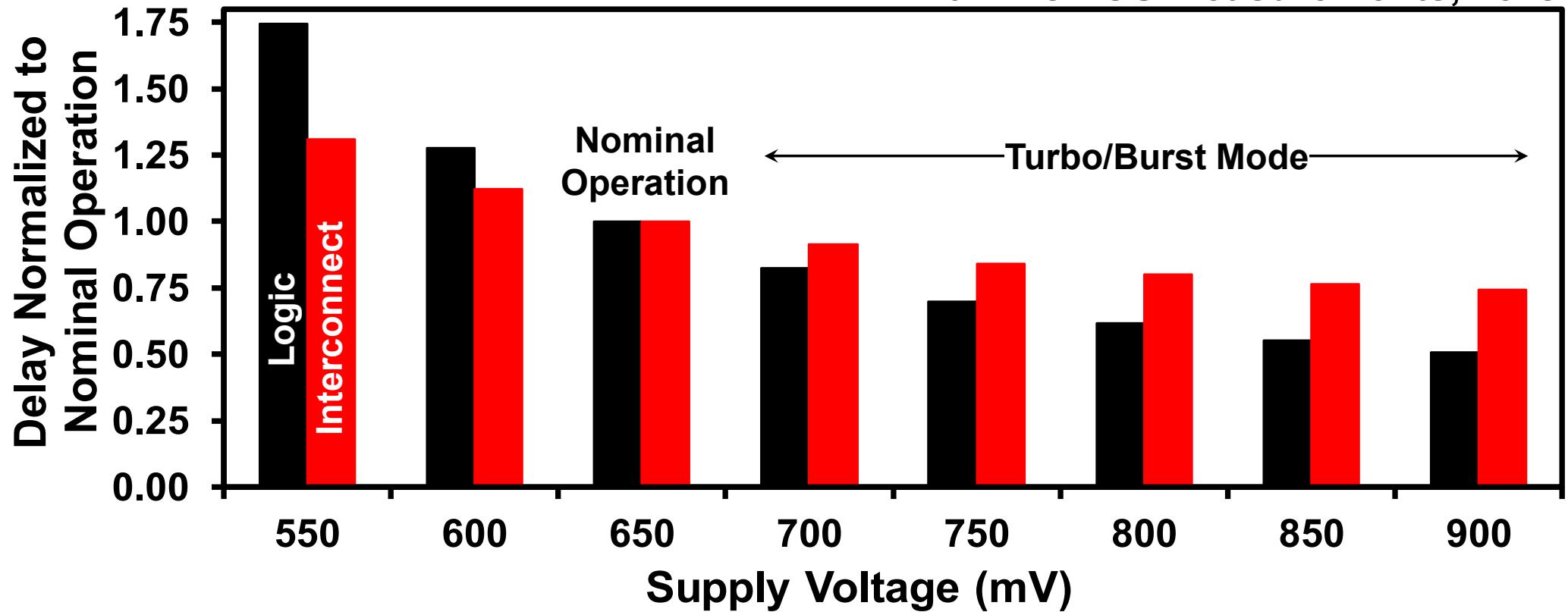


- A peak 115.3 Giga Ray-Steps/Watt energy-efficiency at 460mV, 25°C
- 27% improvement in peak energy-efficiency with Approximate Computing mode at 460mV
- Low-voltage operation provides 2.08x improvement in energy-efficiency

S. Hsu, R. Krishnamurthy et al, IEEE VLSI Circuits Symposium, 2020

On-Die Interconnects Data Movement Challenges

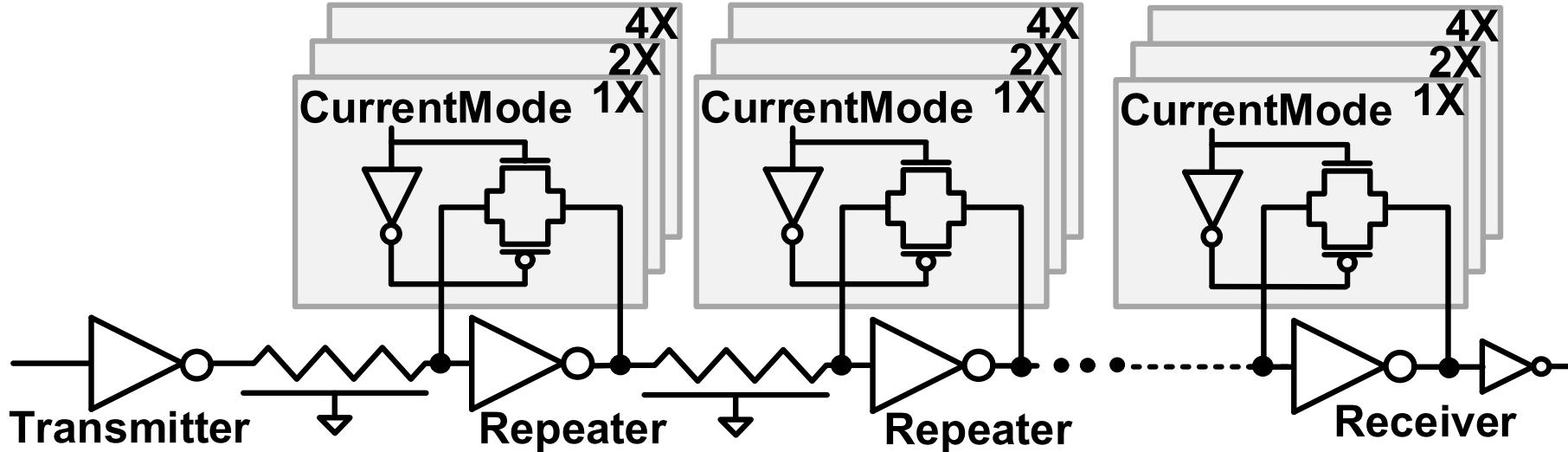
10nm CMOS Measurements, 25°C



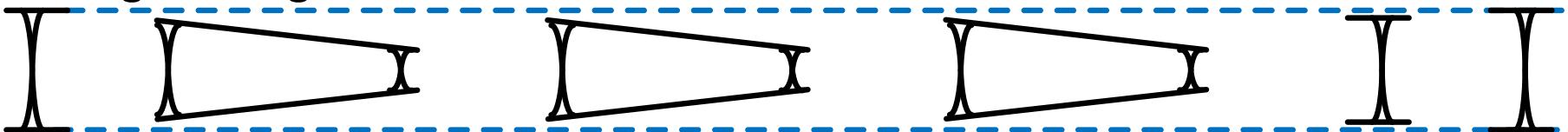
- Clock frequencies limited by interconnect RC delay at high voltages
- Reduces power/performance tunability of high-performance processors

M. Anders, R. Krishnamurthy et al, ISSCC 2020

Reconfigurable Voltage/Current-Mode Interconnects

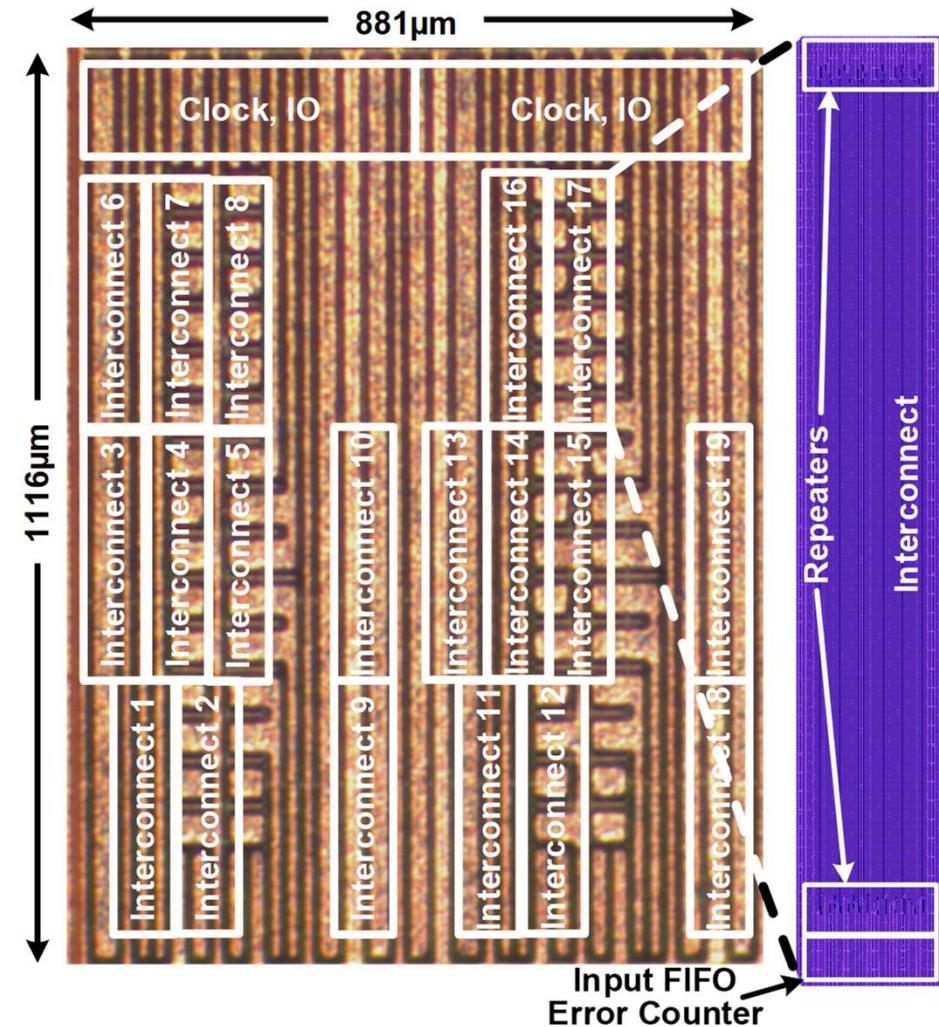


Voltage Swing in Current Mode



- Programmable transmission gates provide feedback resistance
 - Less than 10% additional transistor width and 5% voltage-mode delay overhead
- Voltage swing centered around inverter switching threshold
 - Swing reduces along interconnect and is amplified at each repeater

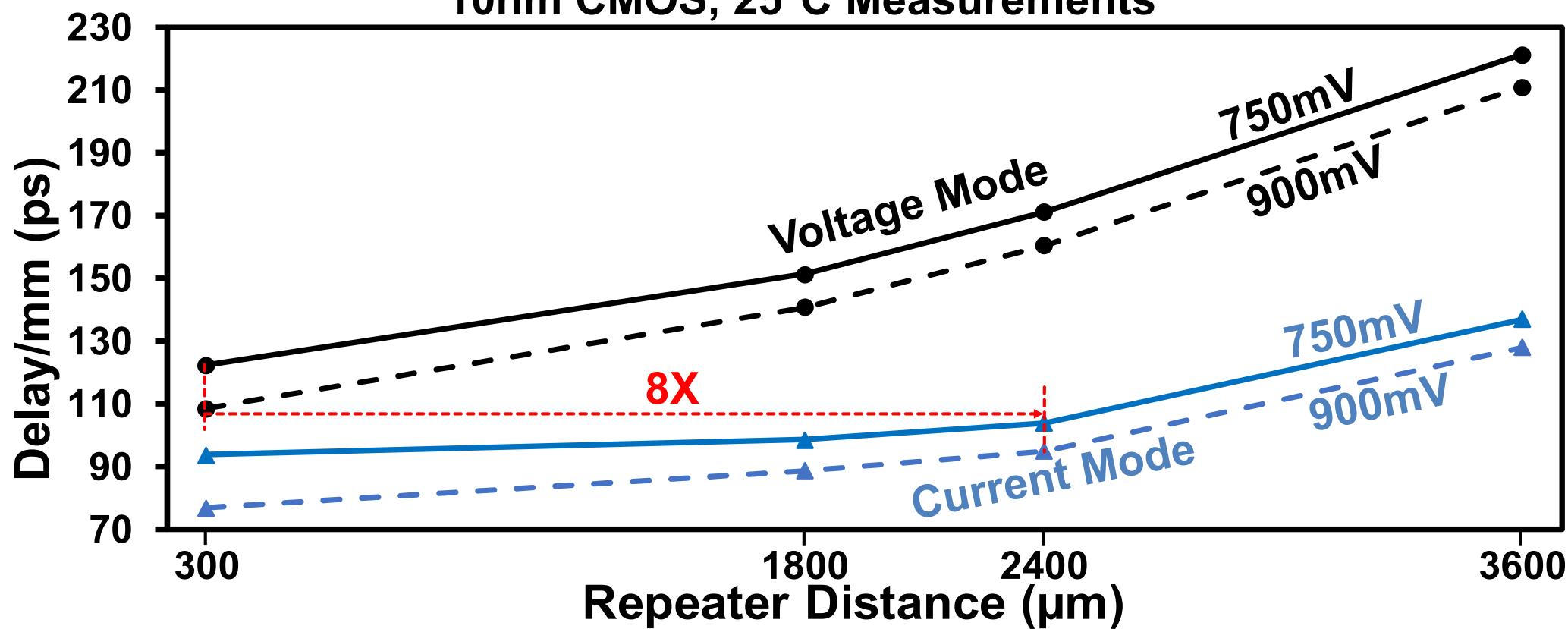
10nm Current Mode Interconnects Test Chip



Experiments	Conventional Voltage-mode Reconfigurable Current-mode Transient Current-mode
Interconnect	Metal 13, 1X - 3X minimum pitch
Repeater Distance	0.3mm – 3.6mm
Measurements	Delay, Power, Noise, Voltage Swing
Configuration Parameters	Supply Voltage Feedback Resistance Pulse Width/Delay Data Pattern

Current Mode vs. Voltage Mode Interconnects

10nm CMOS, 25°C Measurements



- Improved current-mode delay translated to longer inter-repeater distance
- 8X longer distance at same delay target as voltage-mode

M. Anders, R. Krishnamurthy et al, ISSCC 2020

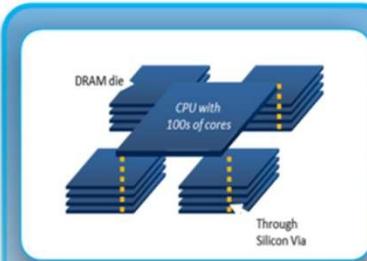
“Extreme” efficiency research



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