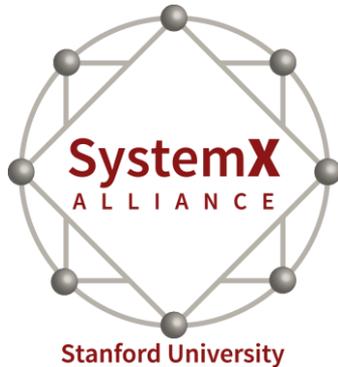


Mixed-Signal Techniques for Embedded Machine Learning Systems

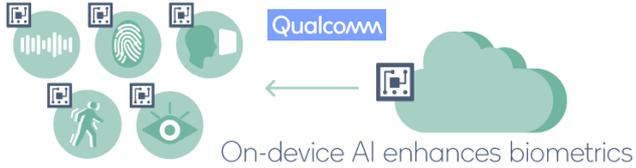


Boris Murmann
June 16, 2019

Edge ←

Applications

→ Cloud



Conversational Interfaces
 ...natural?
 ...seamless?
 ...real-time?



Speed of response

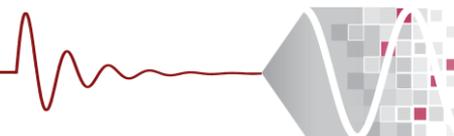
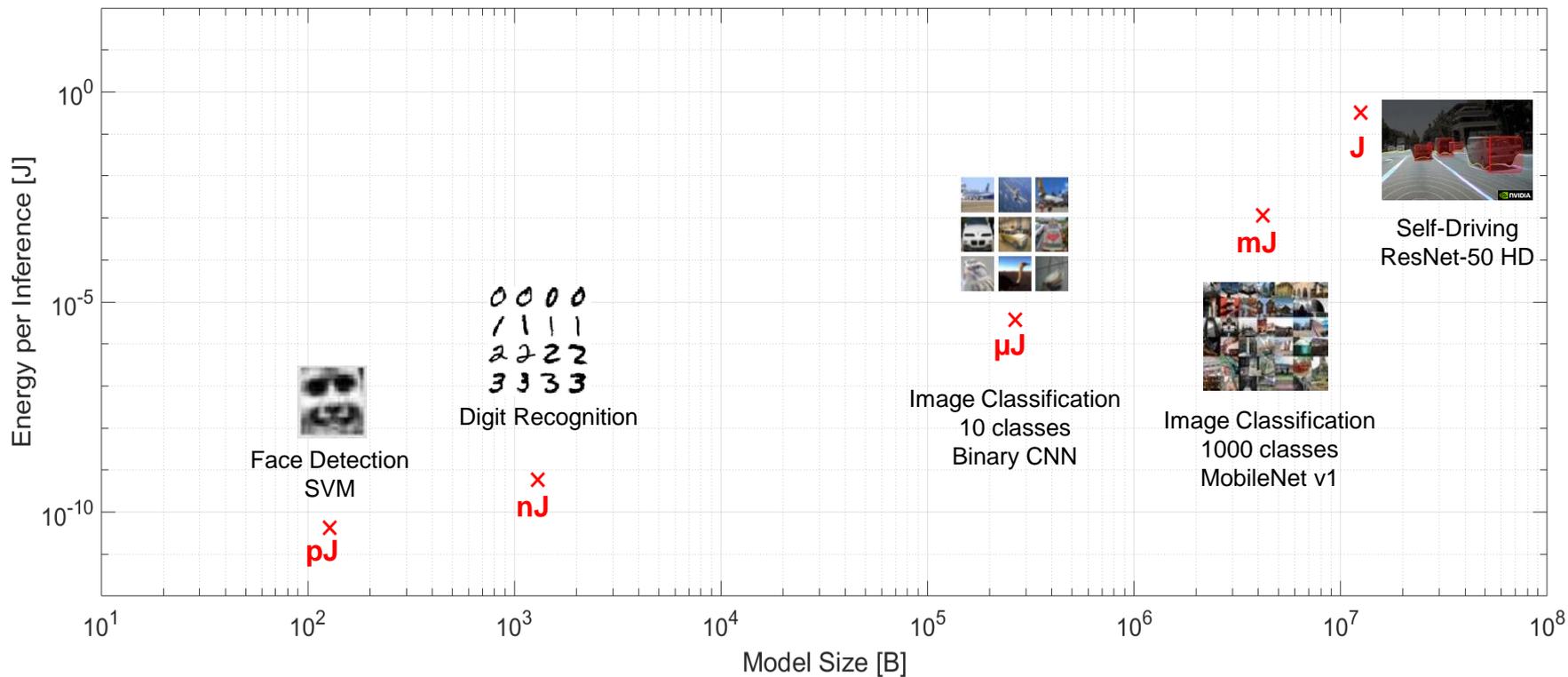
Bandwidth utilized

Privacy

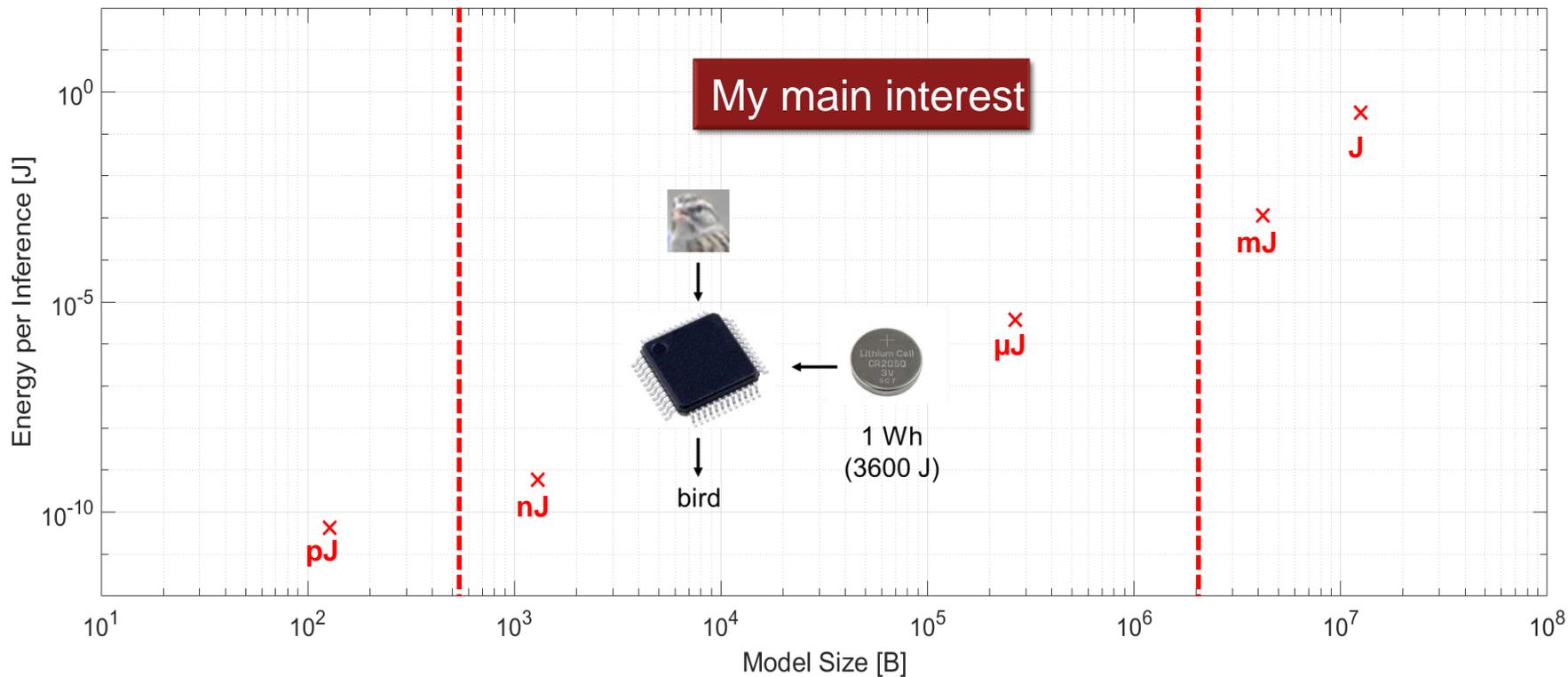
Power consumed



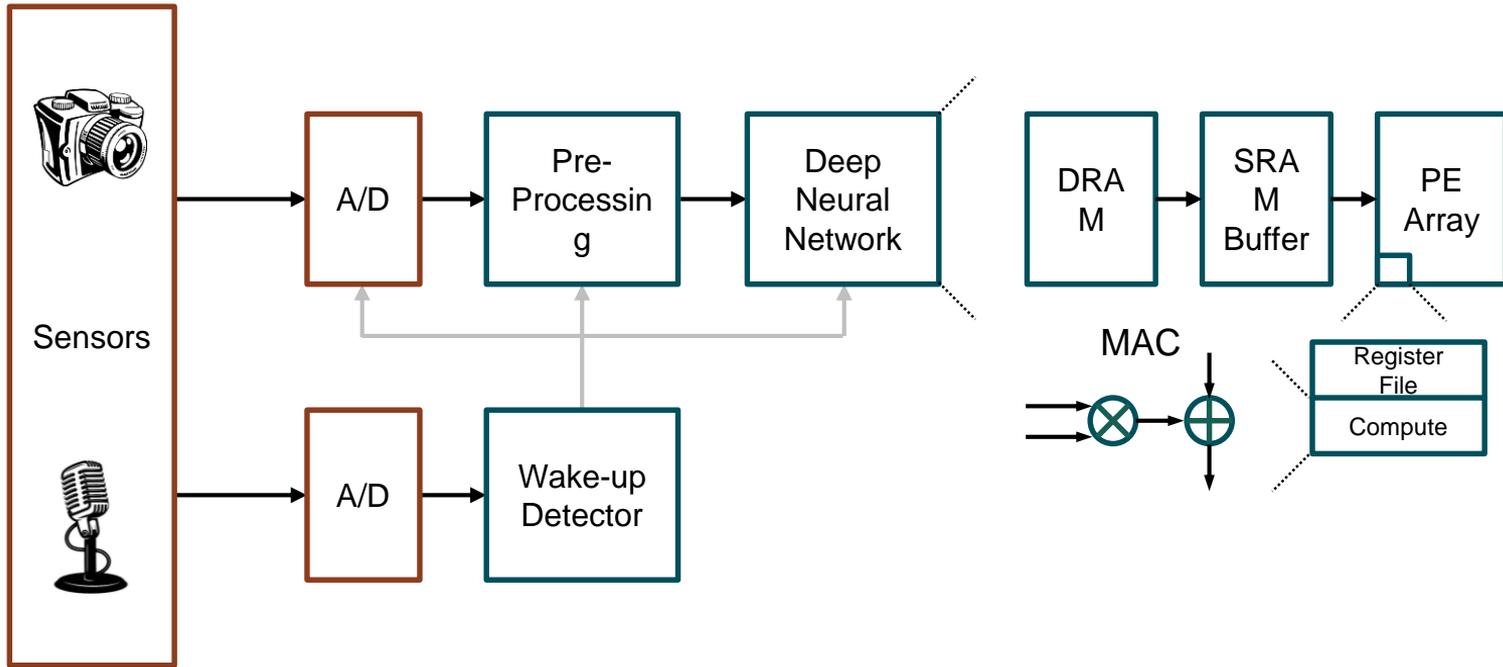
Task Complexity, Memory and Classification Energy



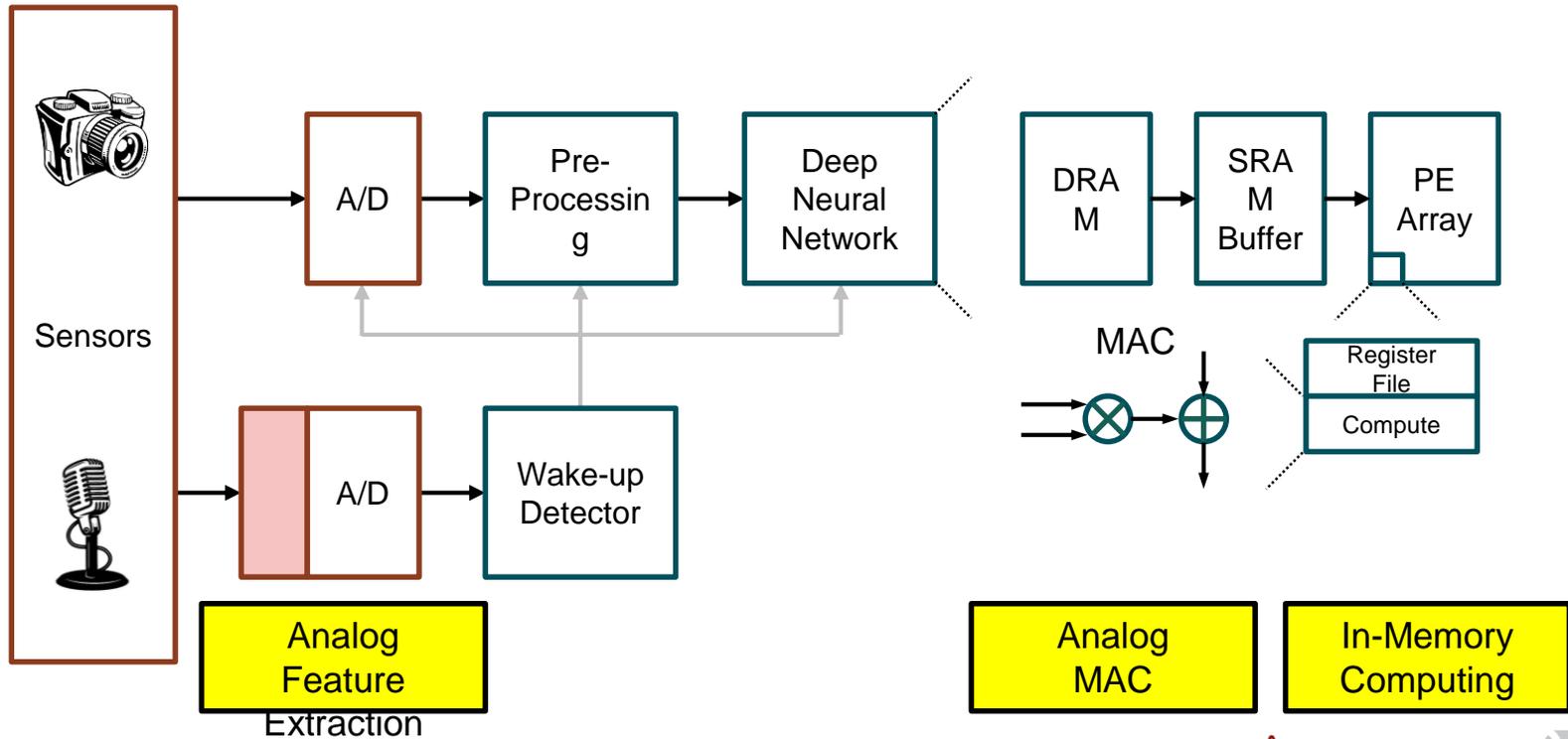
Task Complexity, Memory and Classification Energy



Edge Inference System

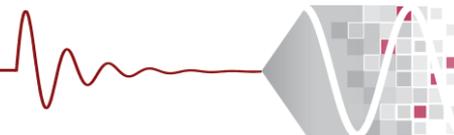


Opportunities for Analog/Mixed-Signal Design

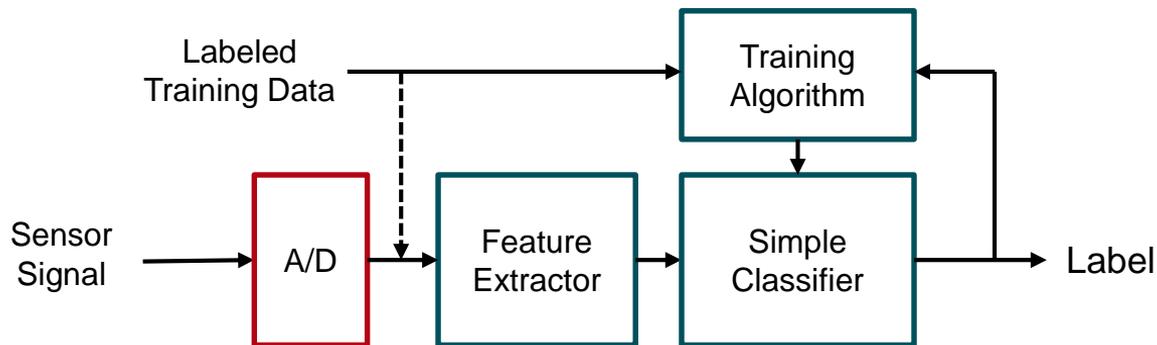


Outline

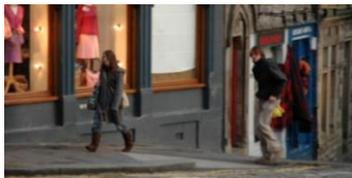
- **Data-Compressive Imager for Object Detection**
 - › Omid-Zohoor & Young, TCSVT 2018 & ISSCC 2019
- **Mixed-Signal ConvNet**
 - › Bankman, ISSCC 2018 & JSSC 2019
- **RRAM-based ConvNet with In-Memory Compute**
 - › Ongoing work



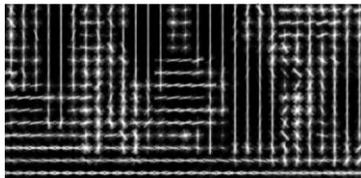
Wake-Up Detector with Hand-Crafted Features



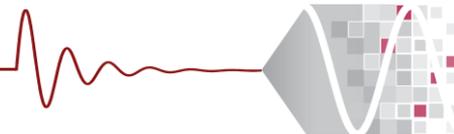
Data Deluge



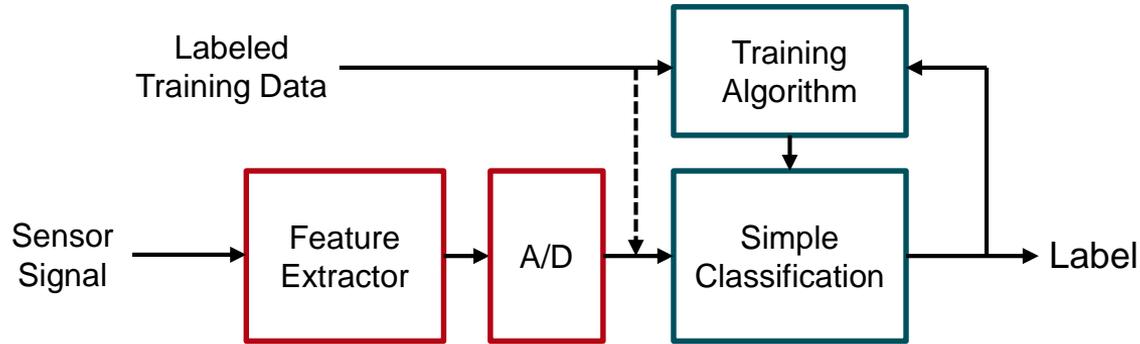
High-dimensional data



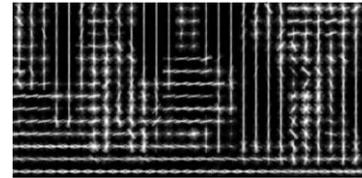
Low-dimensional representation



Analog Feature Extractor



- Low-rate and/or low-resolution ADC
- Low data rate digital I/O
- Reduced memory requirements



Low-dimensional representation

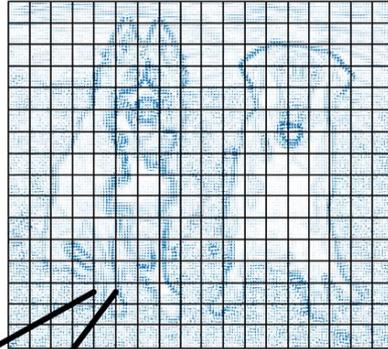


Histogram of Oriented Gradients

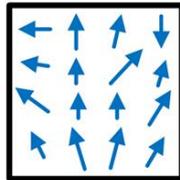
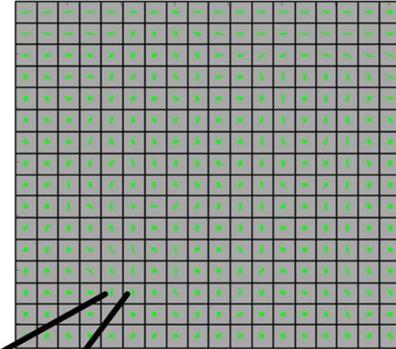
Standard Image



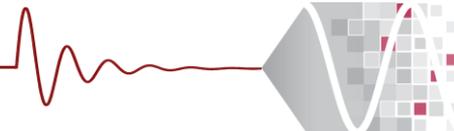
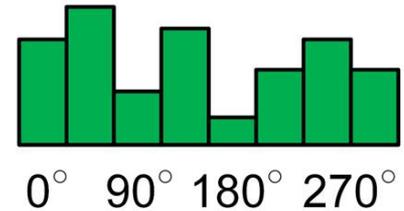
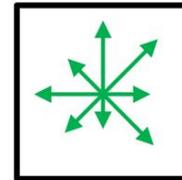
Gradient Image



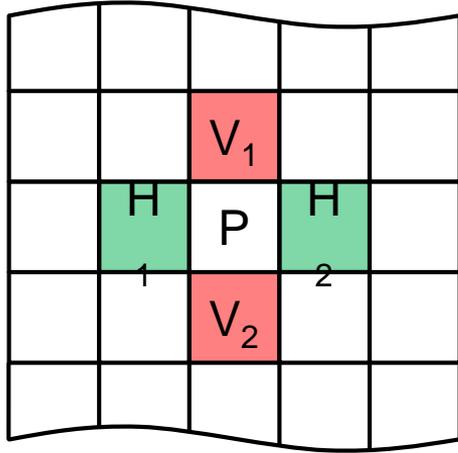
Gradient Histograms



Local histogram in
8x8 pixel patches.



Analog Gradient Computation



$$G_H = \boxed{H_1} - \boxed{H_2} \quad \text{horizontal}$$

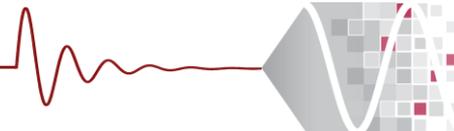
$$G_V = \boxed{V_1} - \boxed{V_2} \quad \text{vertical}$$

Bright patch

$$G_H = 400mV - 100mV = \mathbf{300mV}$$

Dark patch

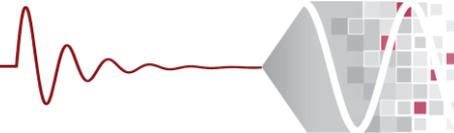
$$G_H = \left(\frac{1}{4}\right) 400mV - \left(\frac{1}{4}\right) 100mV = \mathbf{75mV}$$



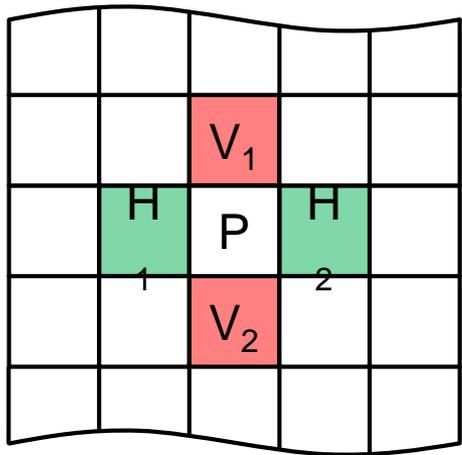
High Dynamic Range Images



- Gradient magnitude varies significantly across image
- Would require high-resolution ADCs ($\geq 9b$) to digitize computed gradients
 - › But, we want to produce as little data as possible



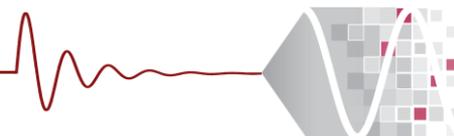
Ratio-Based (“Log”) Gradients



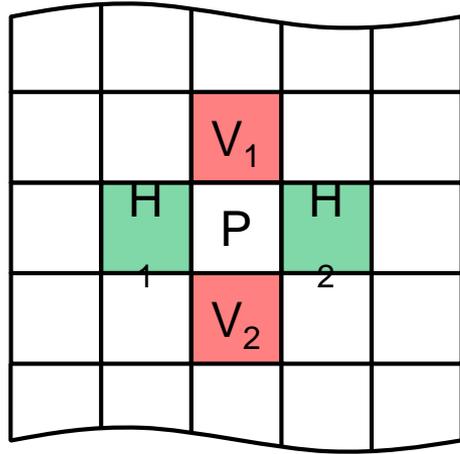
$$G_H = \frac{H_1}{H_2} \quad G_V = \frac{V_1}{V_2}$$

$$G_H = \frac{\alpha \times H_1}{\alpha \times H_2} = \frac{H_1}{H_2}$$

Illumination Invariant

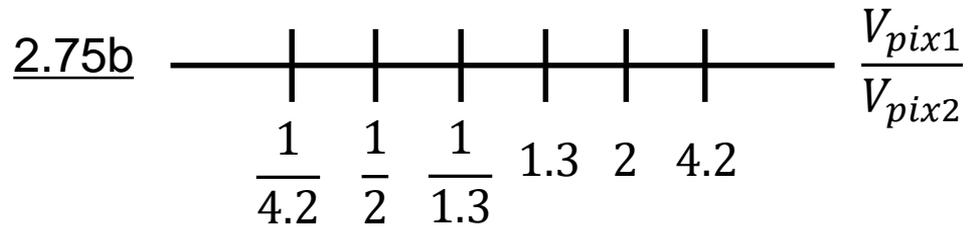
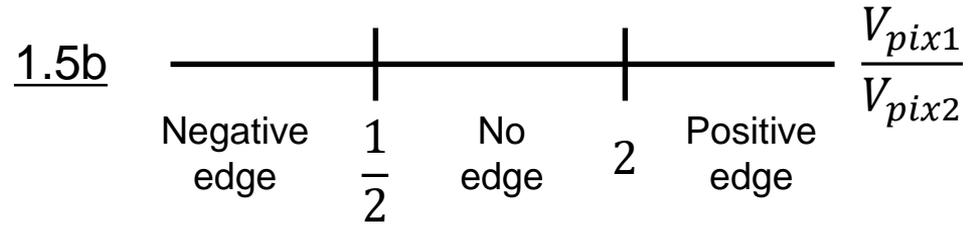


Ratio Quantization

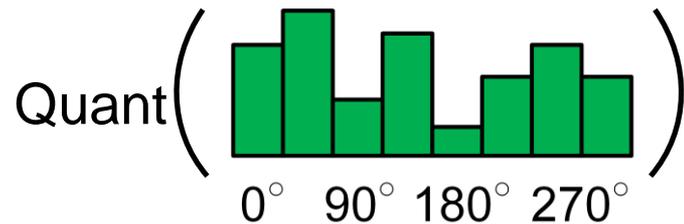
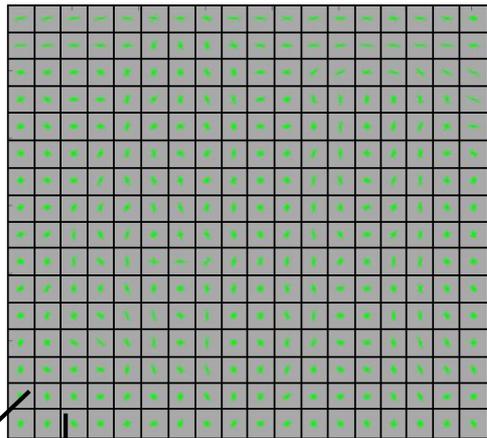
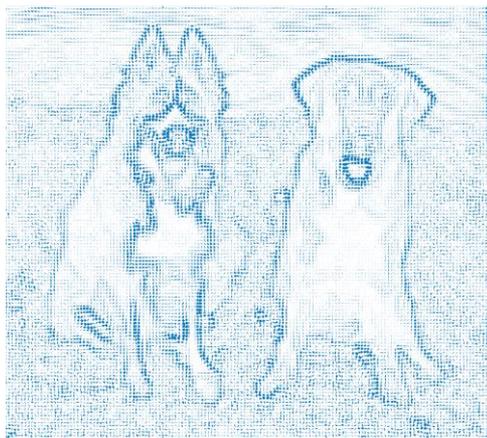


*Empirically determined thresholds

$$G_H = Q \left(\frac{H_1}{H_2} \right) \quad G_V = Q \left(\frac{V_1}{V_2} \right)$$

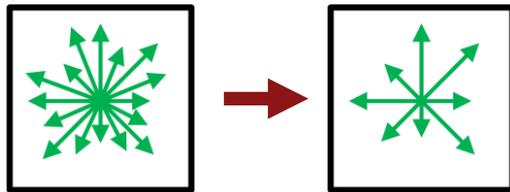


HOG Feature Compression with 1.5b Gradients



Quantizing histogram magnitudes

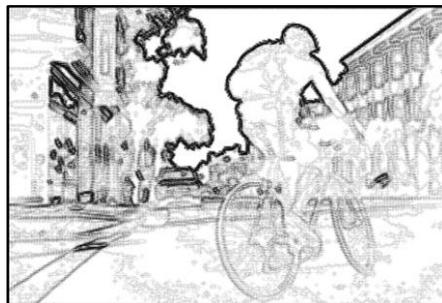
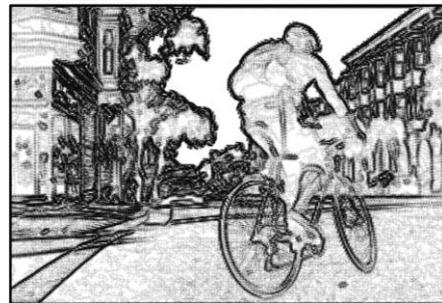
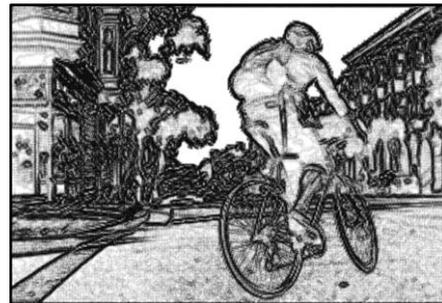
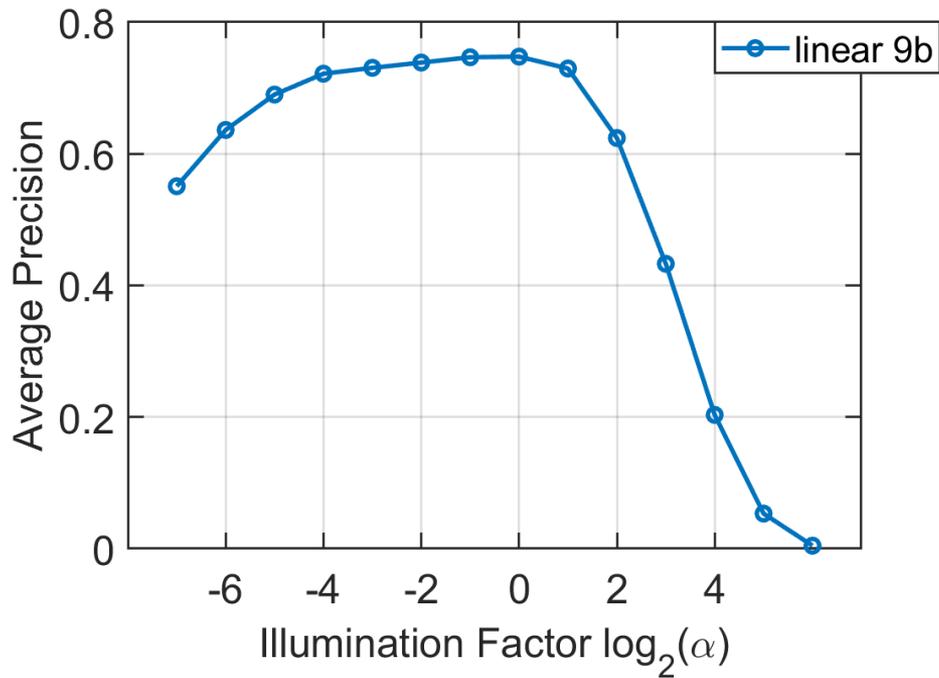
Fewer angle bins



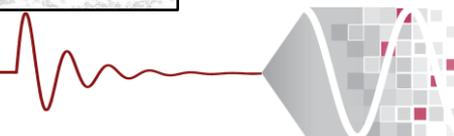
25 × less data in HOG features compared to 8-bit image



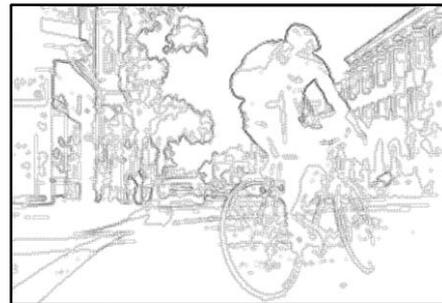
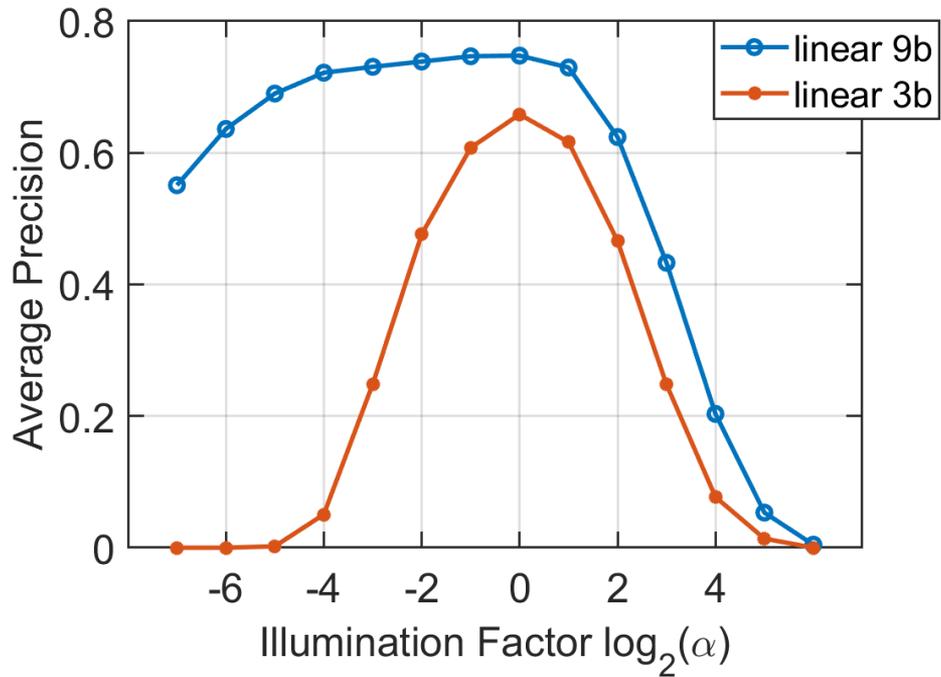
Log vs. Linear Gradients



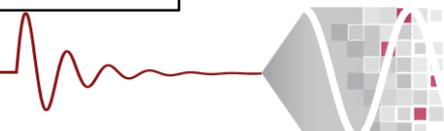
Less Illumination



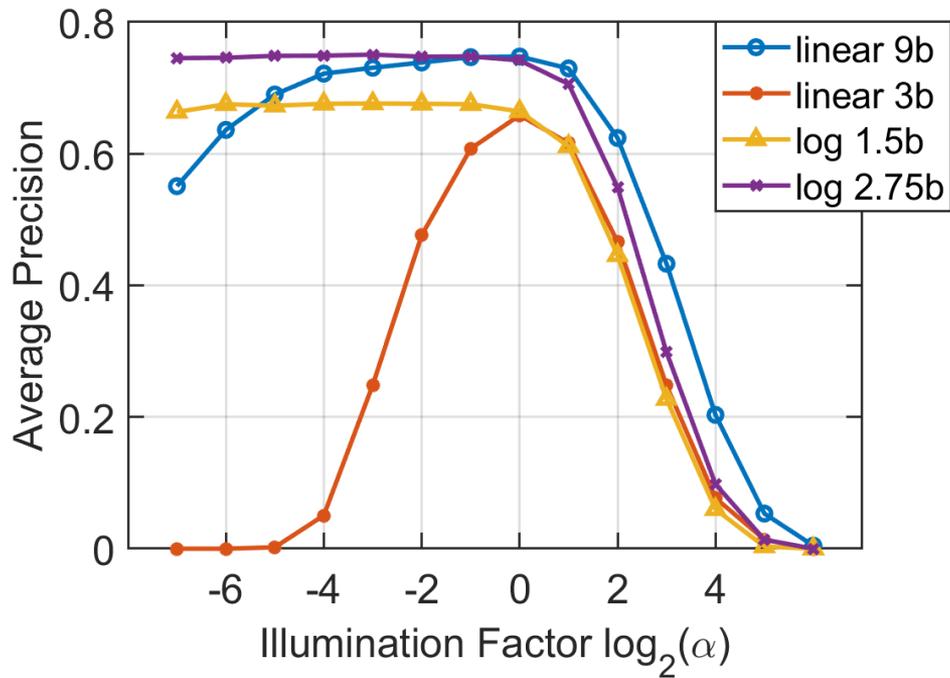
Log vs. Linear Gradients



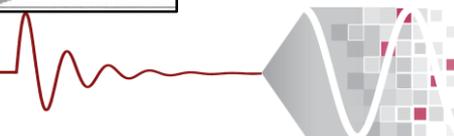
Less Illumination



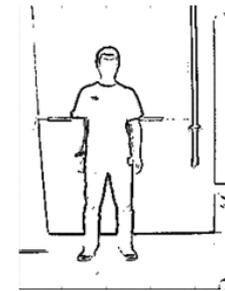
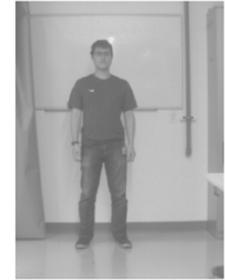
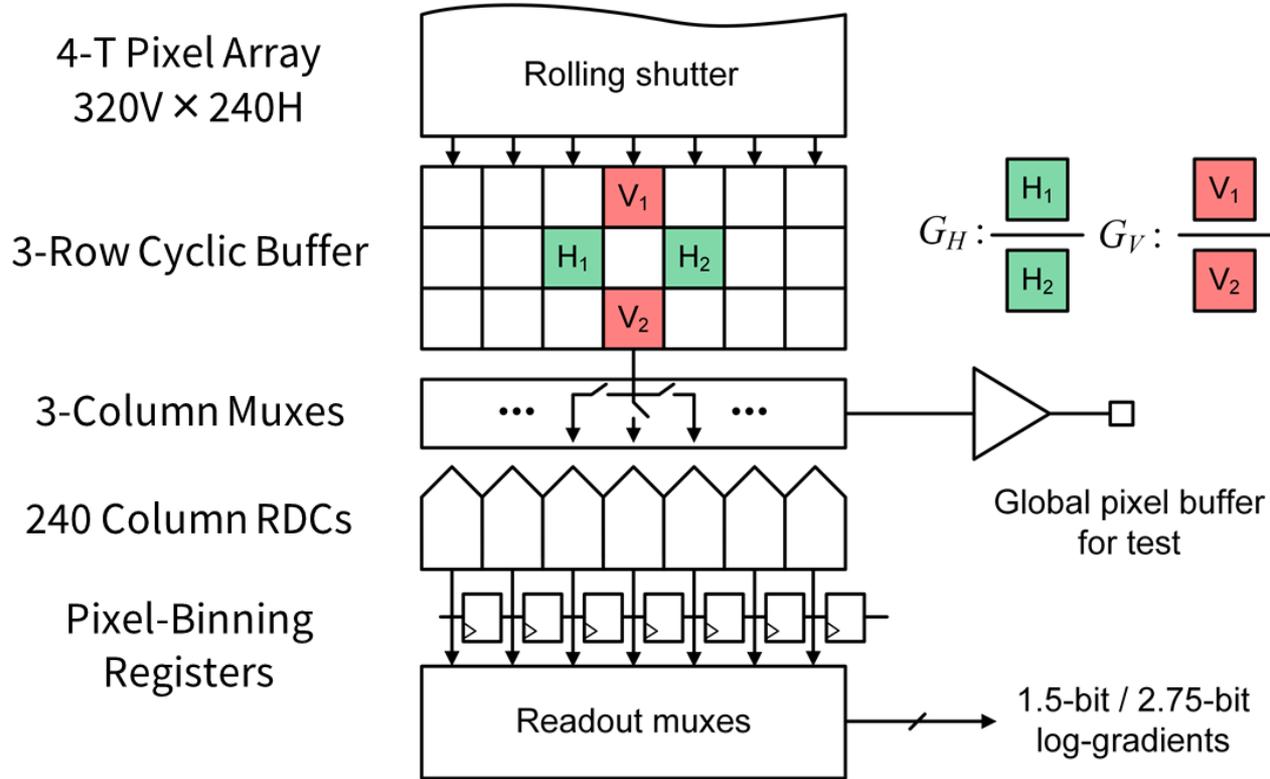
Log vs. Linear Gradients



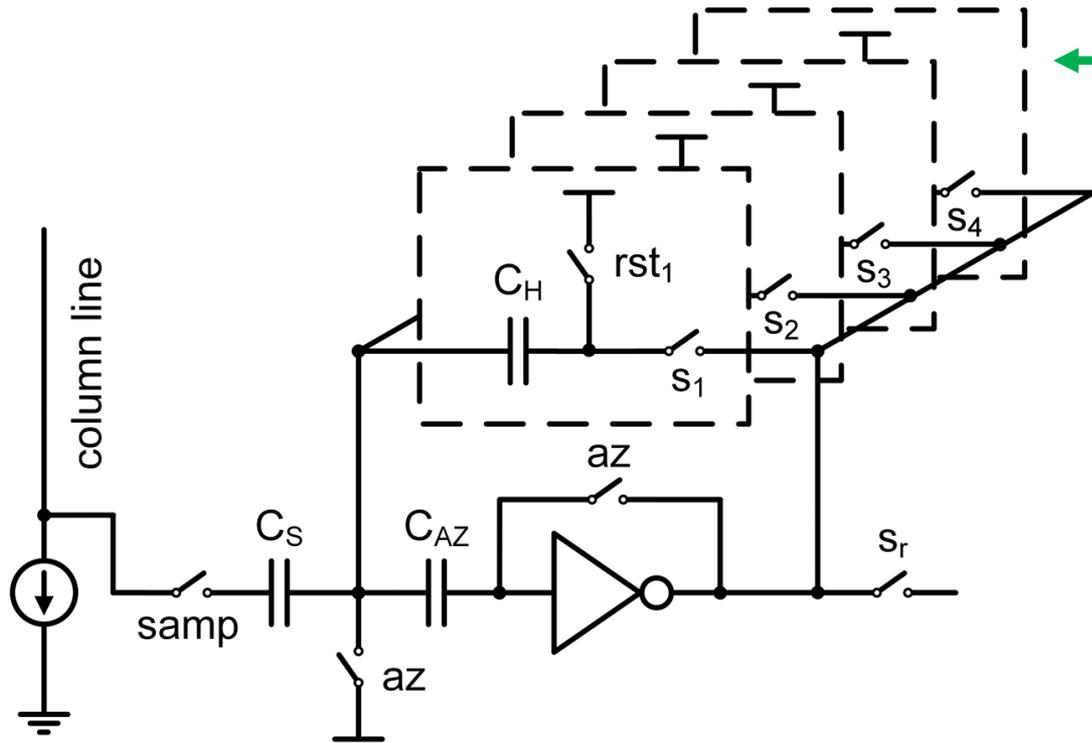
Less Illumination



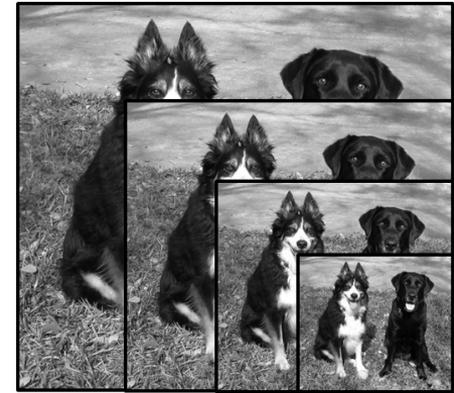
Prototype Chip



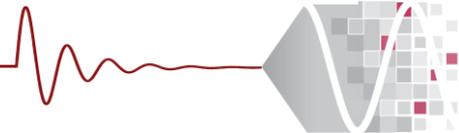
Row Buffers with Pixel Binning (Image Pyramid)



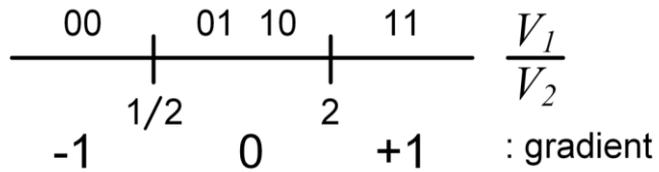
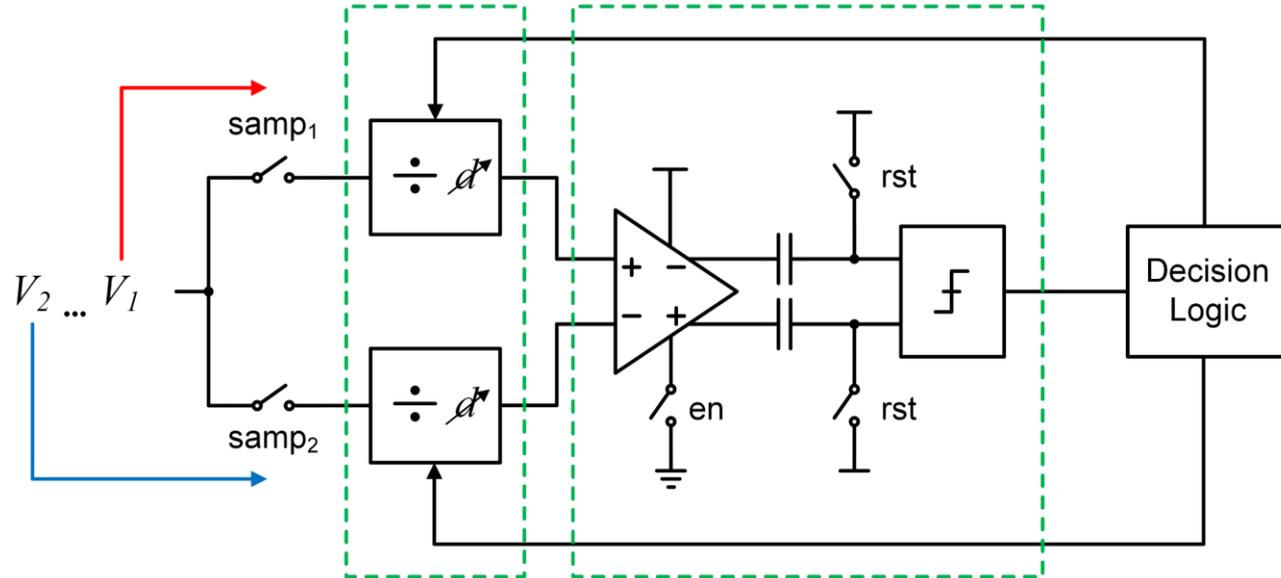
← extra row



Pixel-binning and multi-scale object detection



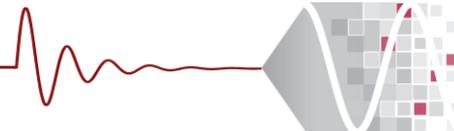
Ratio-to-Digital Converter (RDC)



Ratio Cell

Class-A preamp and
double-tail latch

Decision
Logic



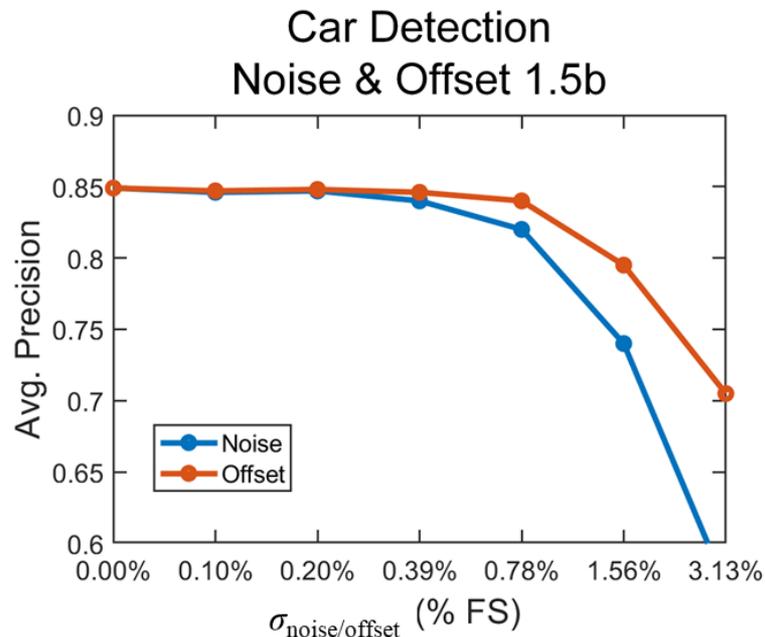
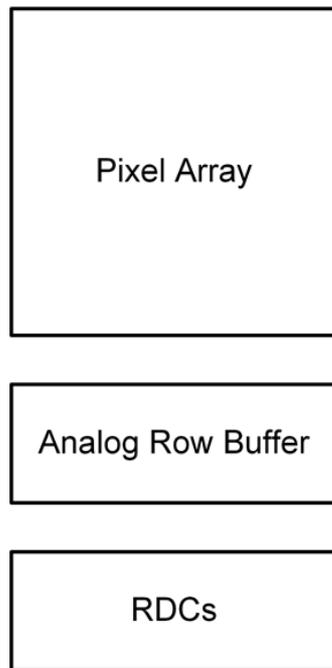
Data-Driven Spec Derivation

$$\frac{H_1 + \epsilon_1}{H_2 + \epsilon_2}$$

referred thermal noise & offset 

Sources

- Pixels
- Row buffers
- Pre-Amps
- Comparators




 ≈ 9 ENOB



Chip Summary

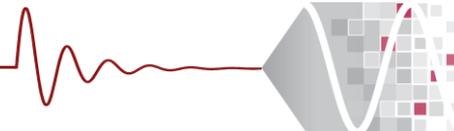
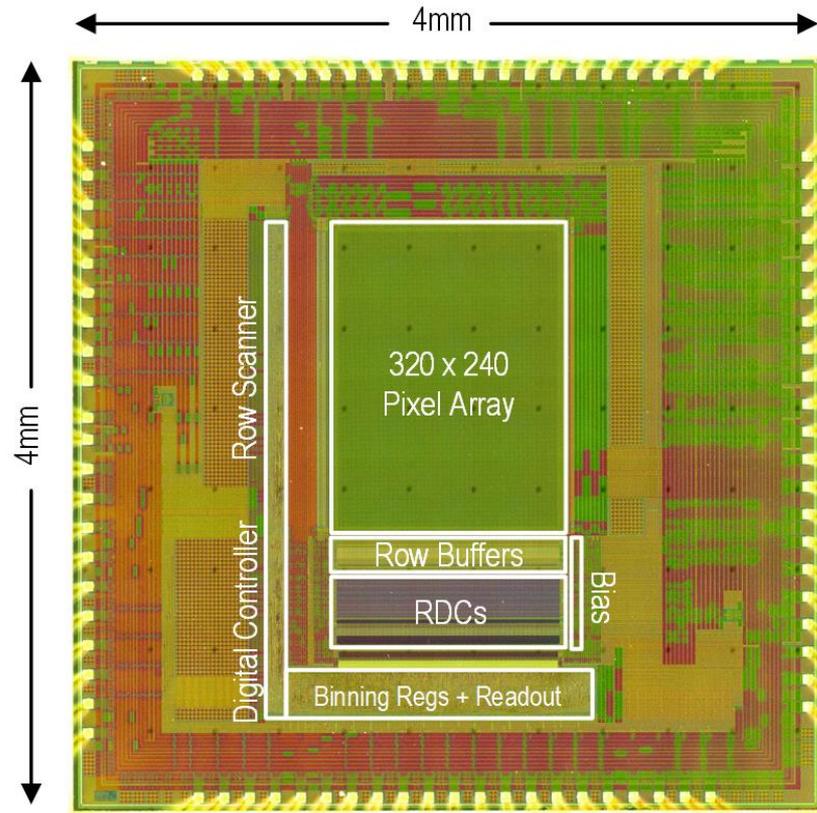
- 0.13 μm CIS 1P4M
- 5 μm 4T pixels
- QVGA 320(V) x 240(H)
- 229 μW @ 30 FPS

Supply Voltages

Pixel: 2.5V

Analog: 1.5V, 2.5V

Digital: 0.9V



Sample Images

Raw, 8b Pixels

9b Linear Gradients

2.75b Log Gradients

1.5b Log Gradients

1.5b Log Gradients
*Truncated

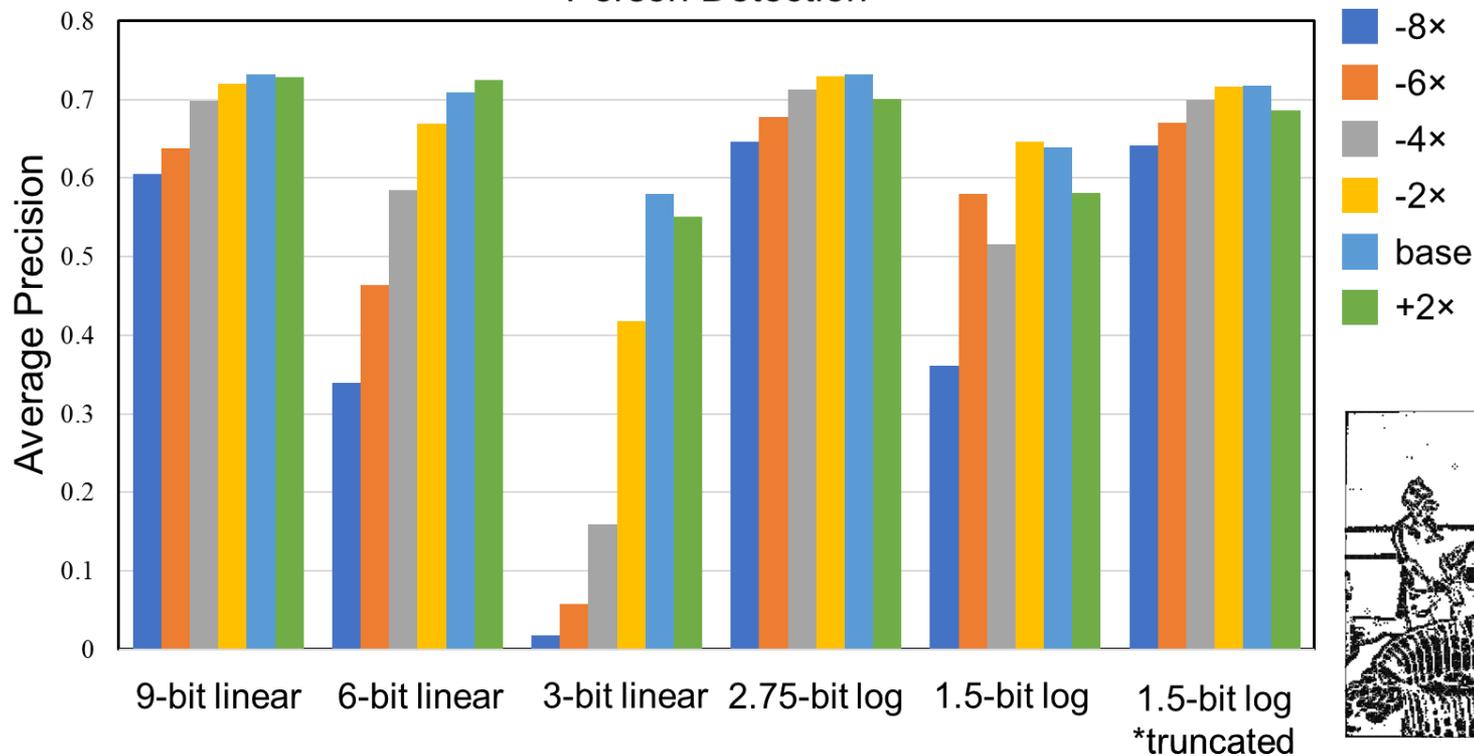
Nominal Exposure



4x less light



Person Detection



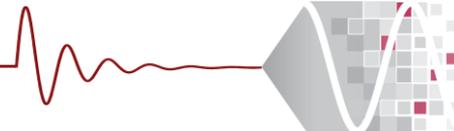
Results using Deformable Parts Model detection & custom database (PascalRAW)



Comparison to State of the Art

	This Work	[Choi, ISSCC'13]	[Katic, Sens.J.'15]	[Bong, ISSCC'17]
Technology	0.13 μm 1P4M	0.18 μm 1P4M	0.18 μm	65 nm 1P8M
Resolution	320x240	256x256	32x32	320x240
Pixel Size	5 μm x 5 μm	5.9 μm x 5.9 μm	31 μm x 26 μm	7 μm x 7 μm
Fill Factor	60.4%	30%	24%	-
Feature Type	log-gradients	linear HOGs	relative ratios between pixels	linear Haar-like w/ face-detector
Frame Rate	30 fps nom. 207 fps max	15 fps - reported	9756 fps nom. 24000 fps max	1 fps - reported
Dynamic Range	59.3 dB¹	54.8 dB	43 dB ²	-
Power Consumption	229 μW @ 30fps	51 μW @ 15 fps	4 mW @ 9765 fps	24-96 μW @ 1fps
Energy Efficiency	1.5-bit: 99 pJ/pixel 2.75-bit: 114 pJ/pixel	52 pJ/pixel	404 pJ/pixel	312 - 1250 pJ/pixel
Multi-Scale	Yes - arbitrary square bins	No	No	Yes - three scales

1. At output of cyclic row buffer, without RDC 2. Pixel-to-pixel dynamic range

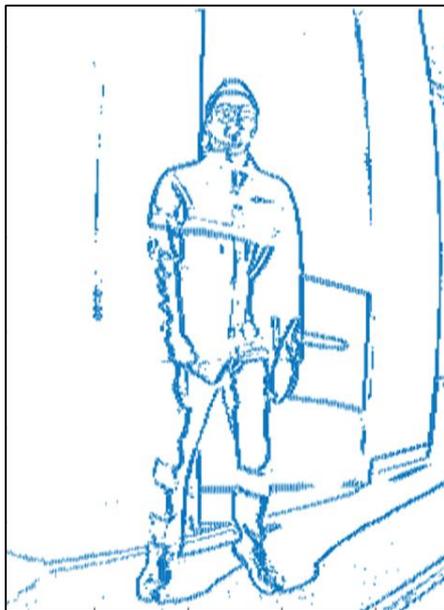


Information Preservation

Raw Pixels



1.5-bit Log Gradients

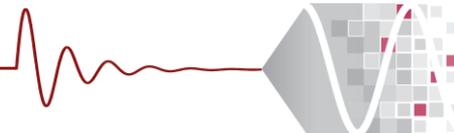


*truncated from 2.75-bit

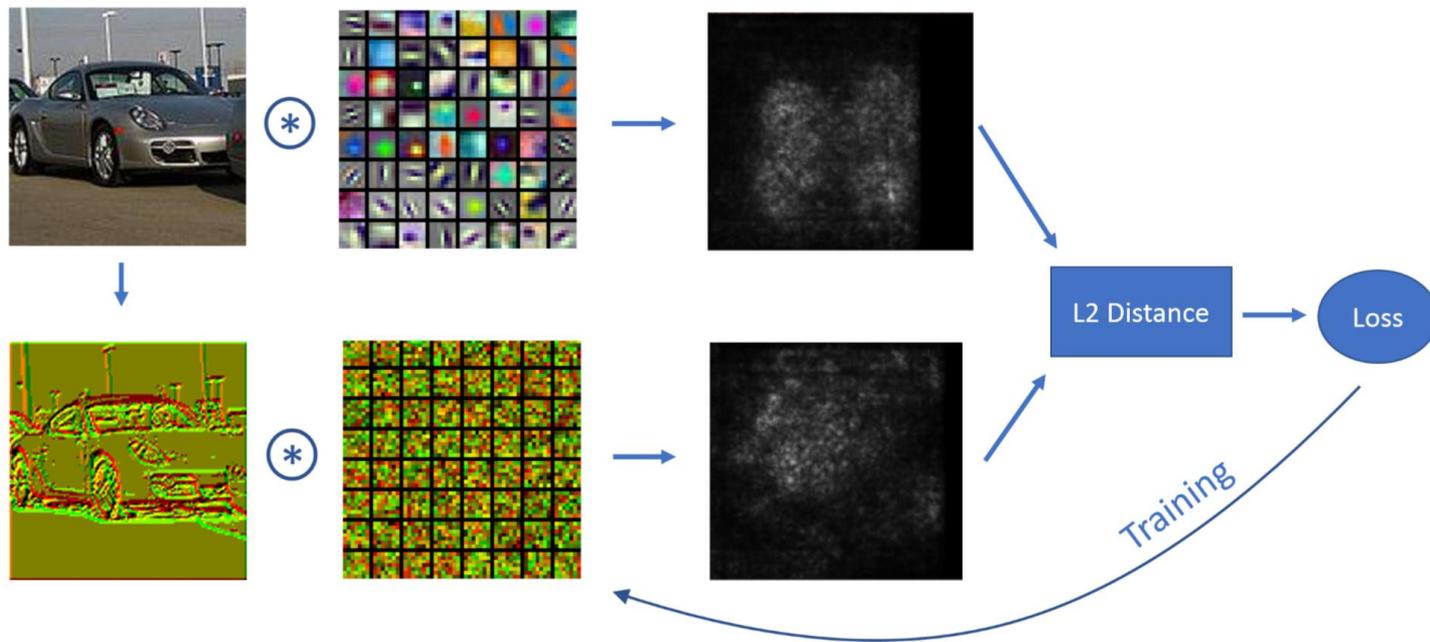
Reconstruction



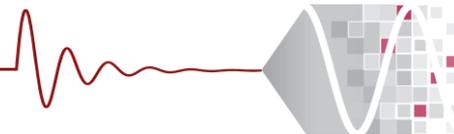
*courtesy Julien Martel



Use Log Gradients as ConvNet Input?



- Ongoing work; comparable performance using ResNet-10 (PascalRaw dataset)

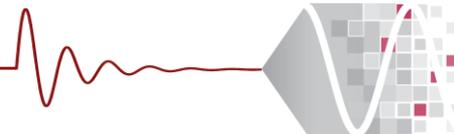


Can We Play Mixed-Signal Tricks in a ConvNet?



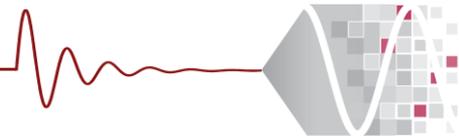
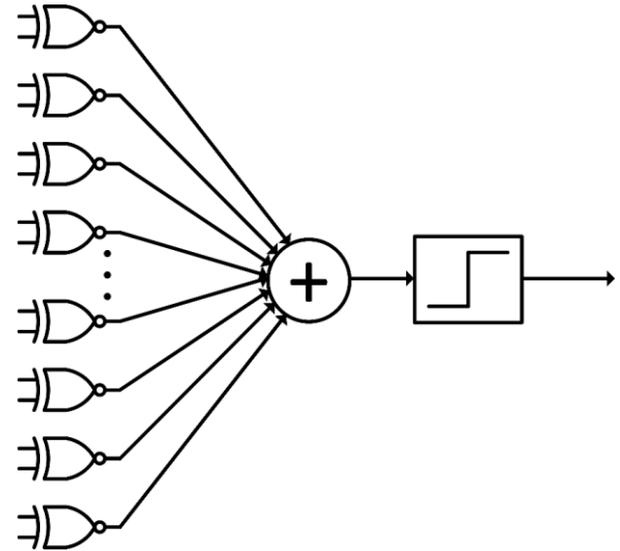
**Analog
Goodness**

Digital ConvNet Fabric



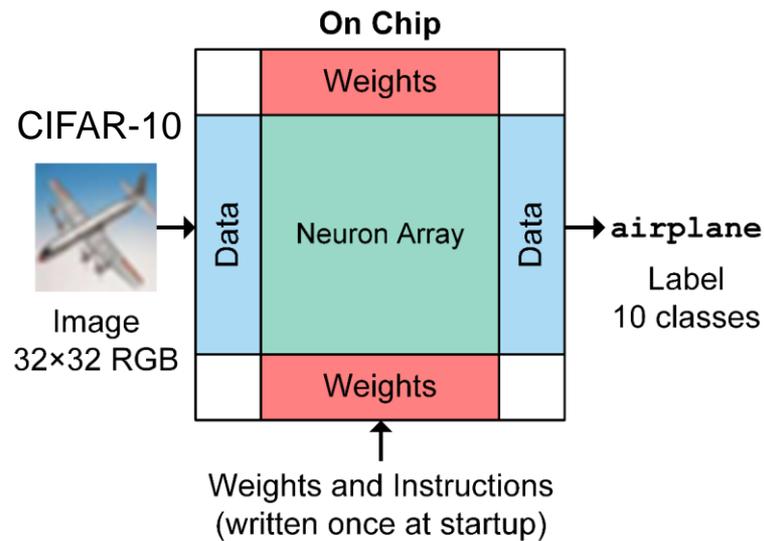
BinaryNet

- Courbariaux et al., NIPS 2016
- Weights and activations constrained to +1 and -1, multiplication becomes XNOR
- Minimizes D/A and A/D overhead
- Nice option for small/medium-size problems and mixed-signal exploration

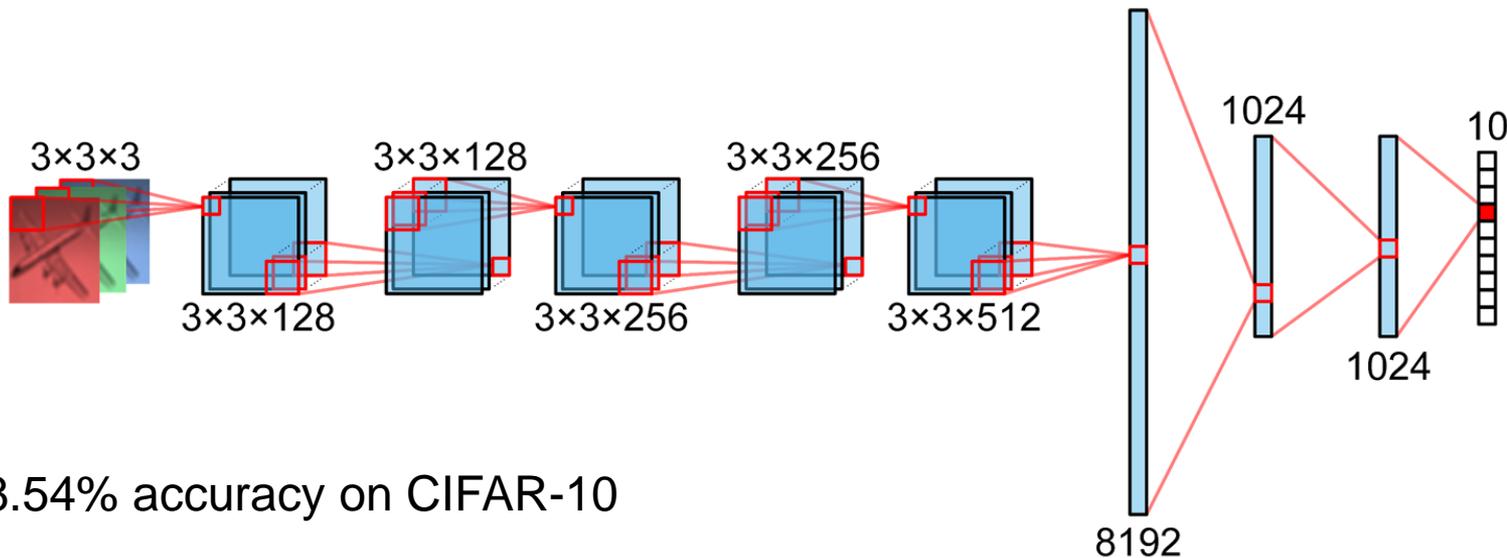


Mixed-Signal Binary CNN Processor

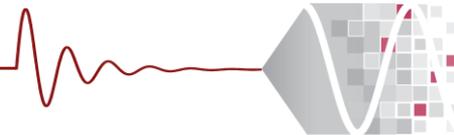
1. Binary CNN with “CMOS-inspired” topology, engineered for minimal circuit-level path loading
2. Hardware architecture amortizes memory access across many computations, with all memory on chip (328 KB)
3. Energy-efficient switched-capacitor neuron for wide vector summation, replacing digital adder tree



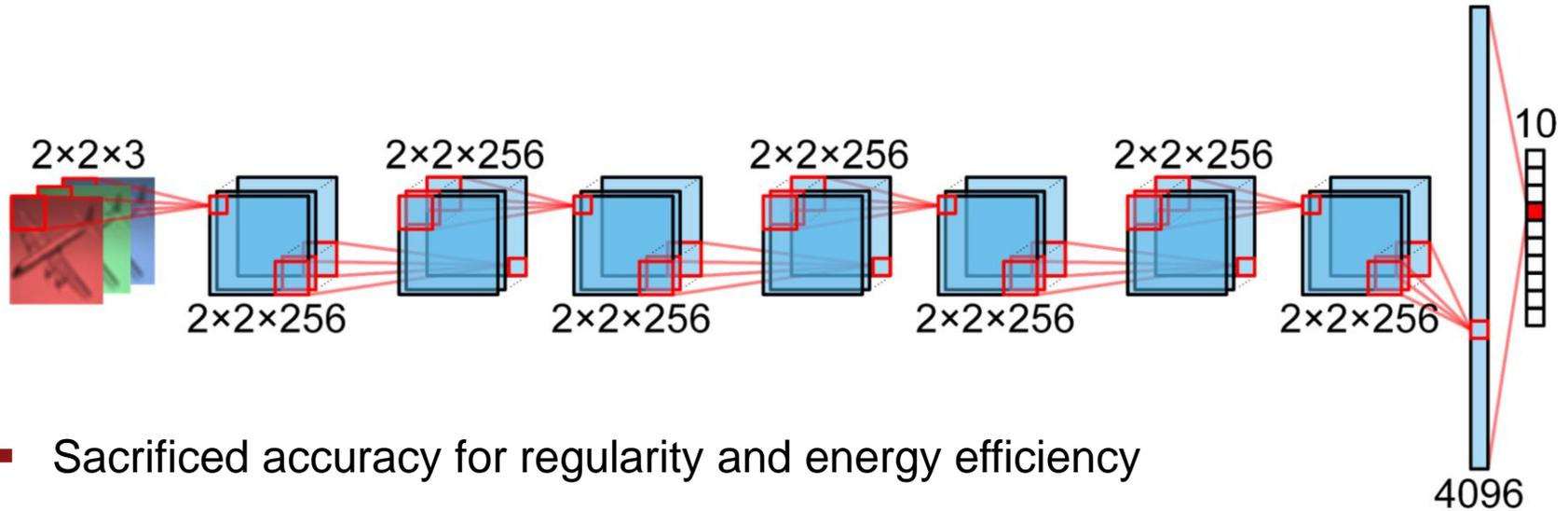
Original BinaryNet Topology



- 88.54% accuracy on CIFAR-10
- 1.67 MB weight memory (68% FC layers)
- 27.9 mJ/classification with FPGA



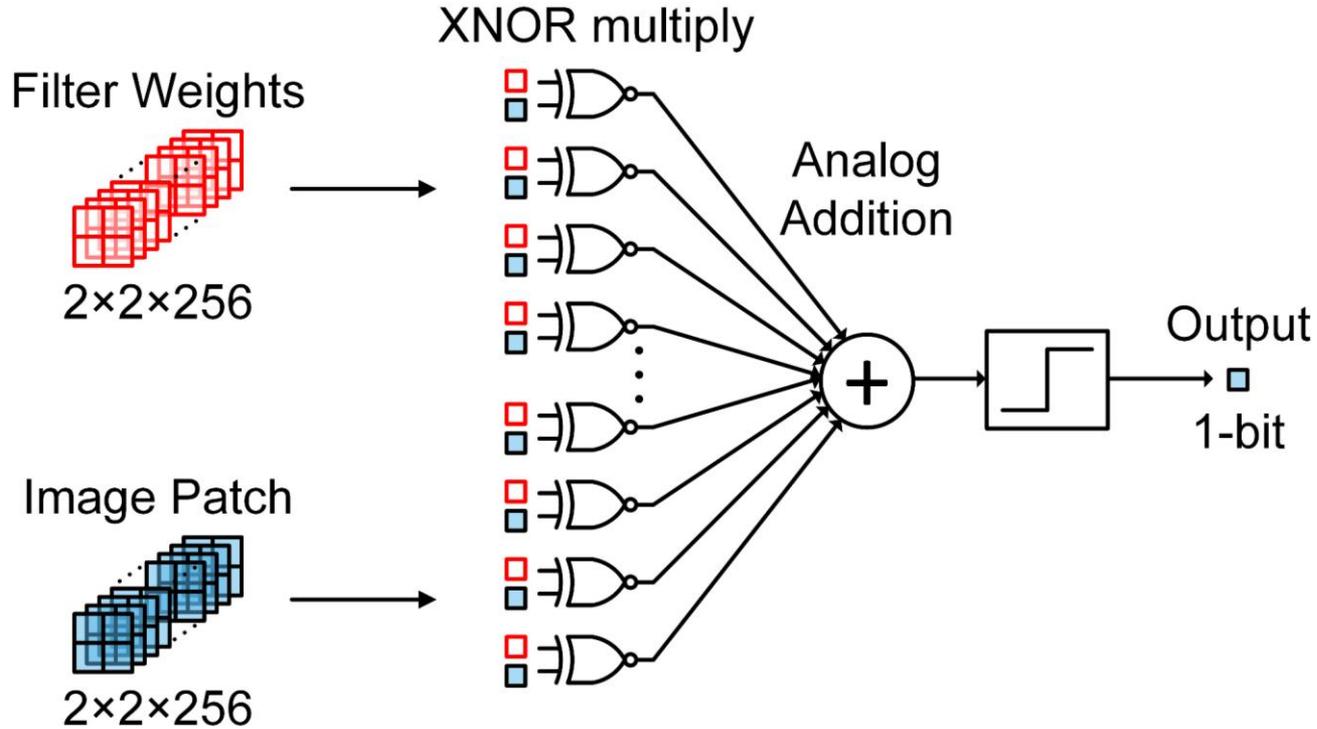
Mixed-Signal BinaryNet Topology



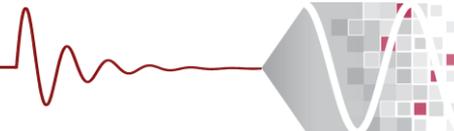
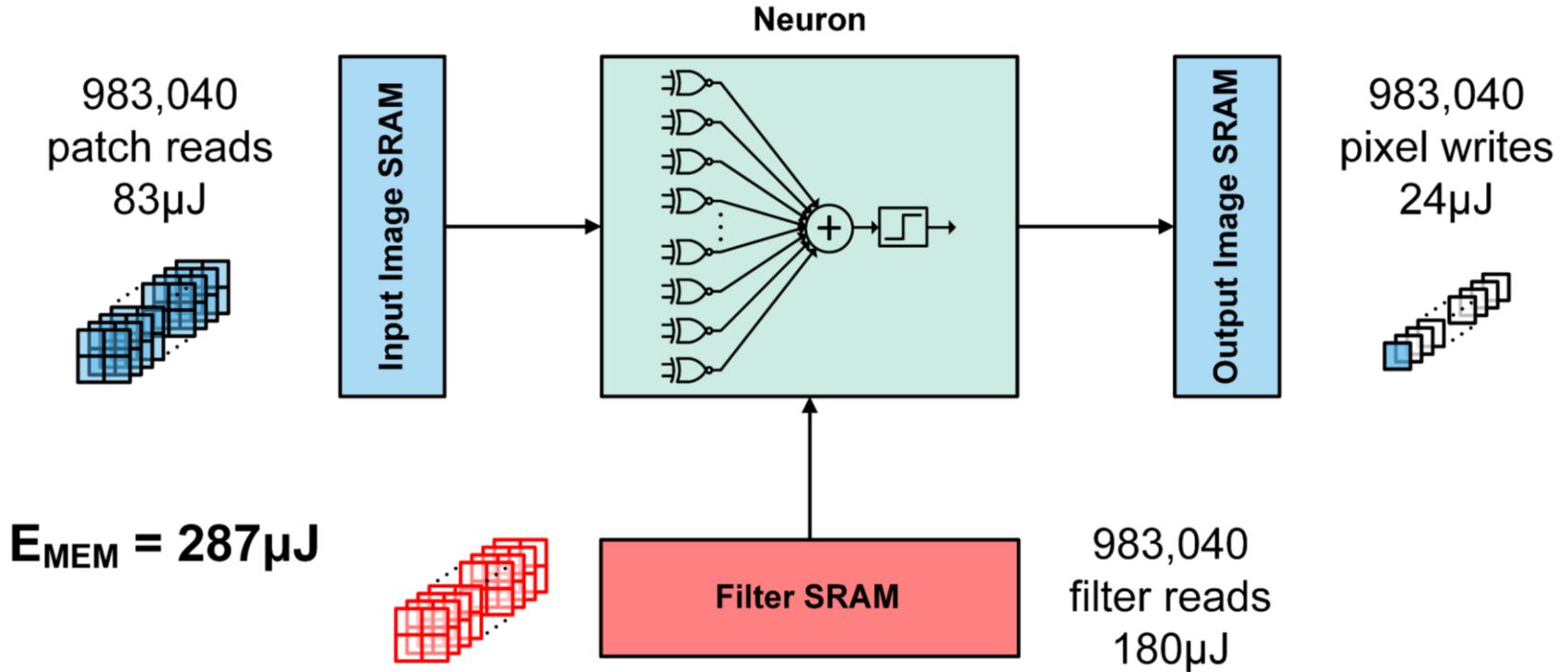
- Sacrificed accuracy for regularity and energy efficiency
- 86.05% accuracy on CIFAR-10
- 328 KB weight memory
- 3.8 μJ per classification



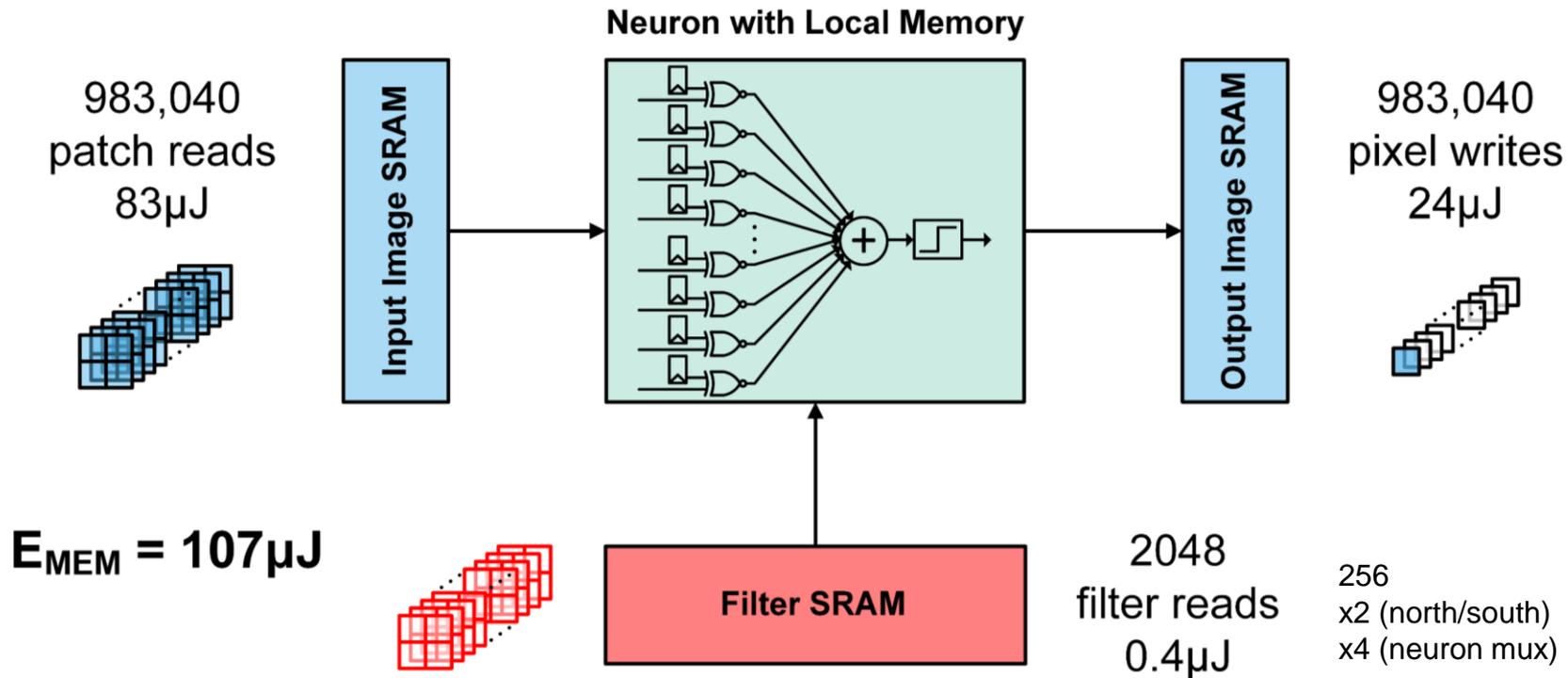
Neuron



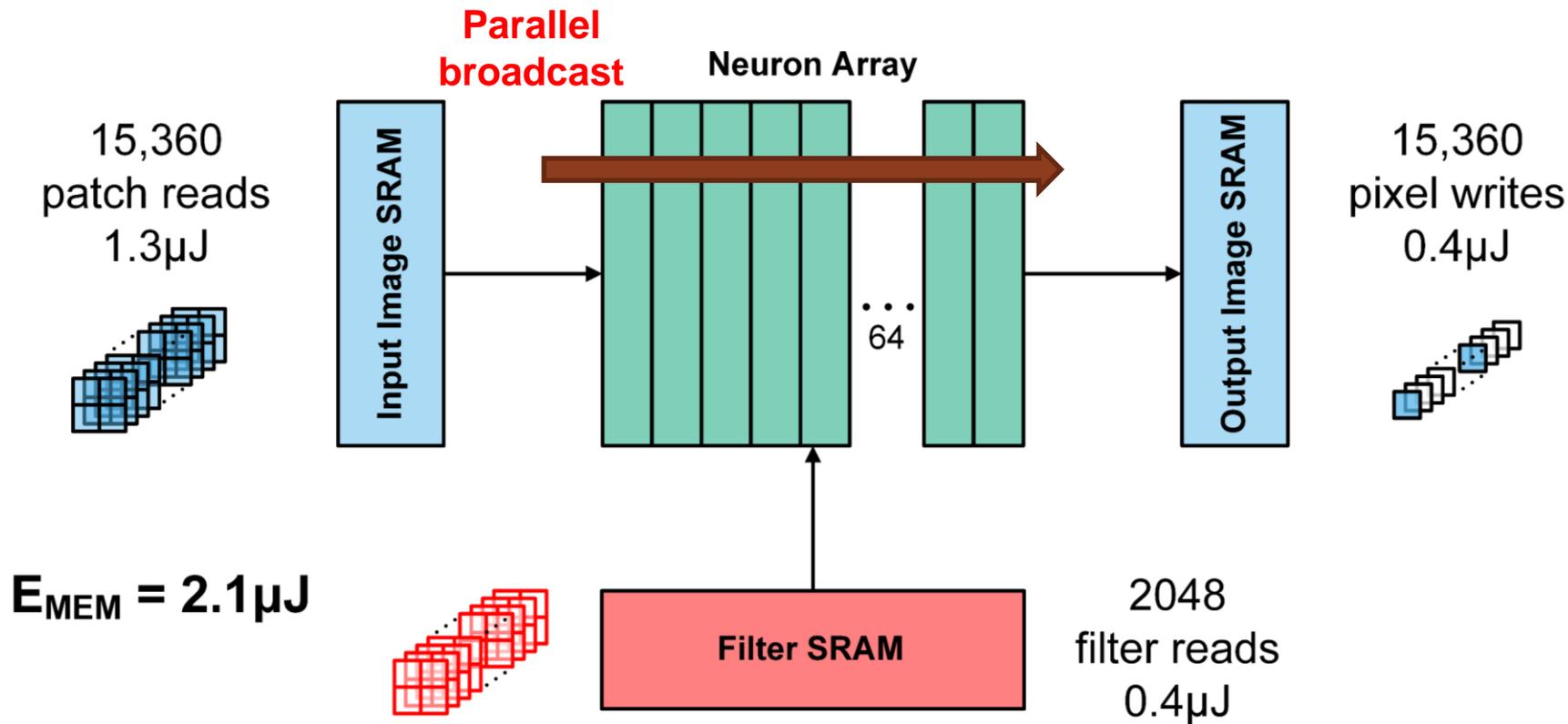
Naïve Sequential Computation



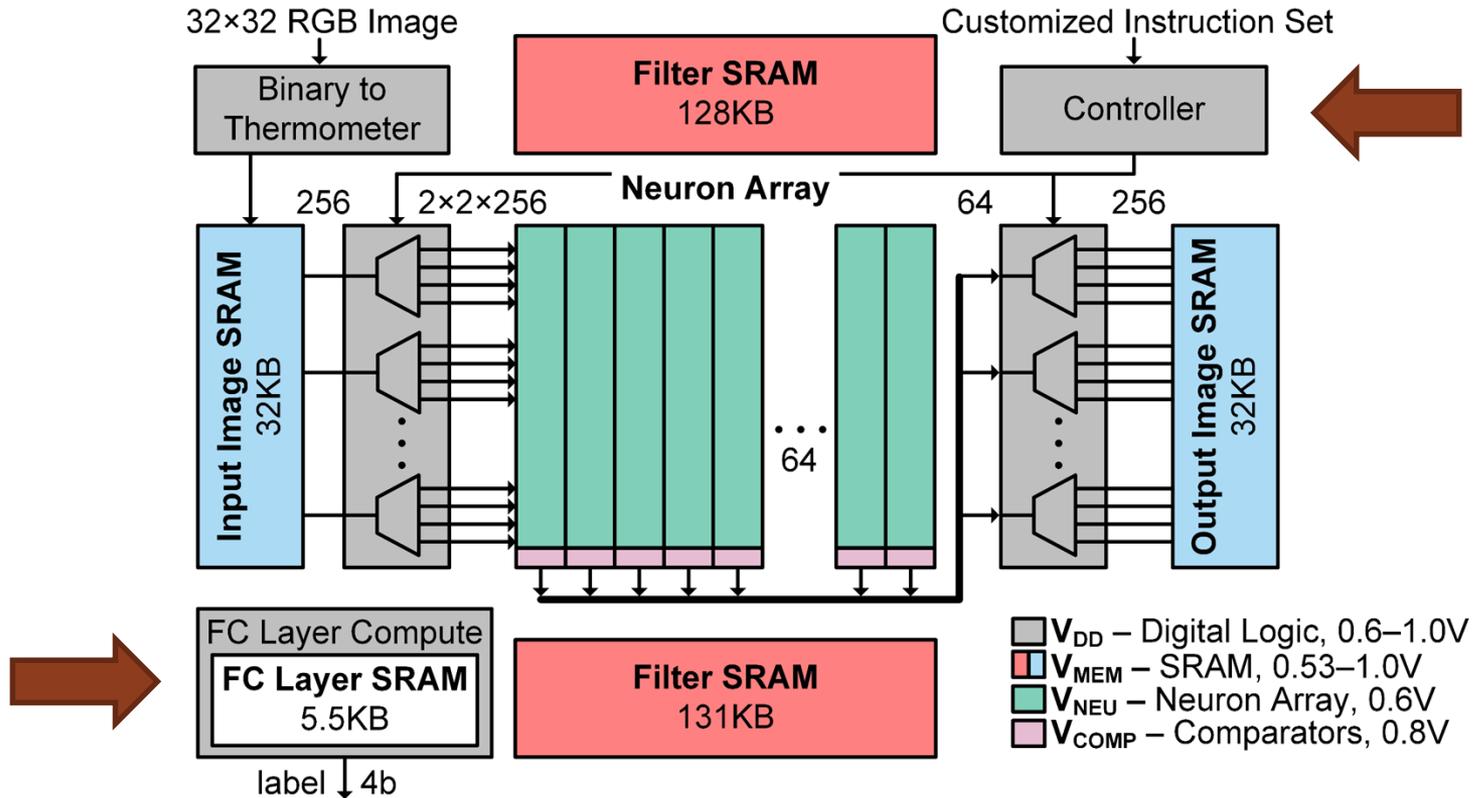
Weight-Stationary



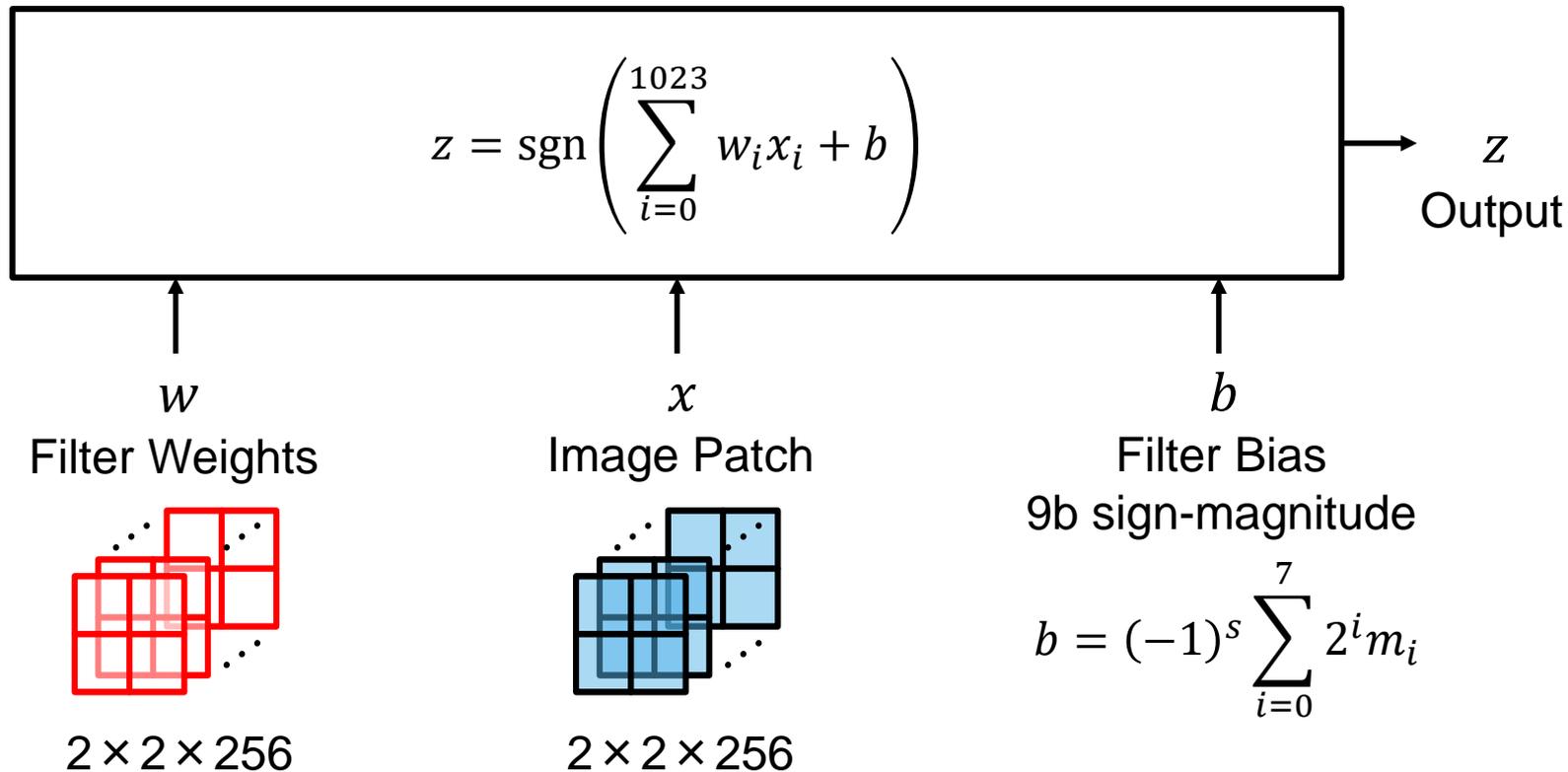
Weight-Stationary and Data-Parallel



Complete Architecture



Neuron Function

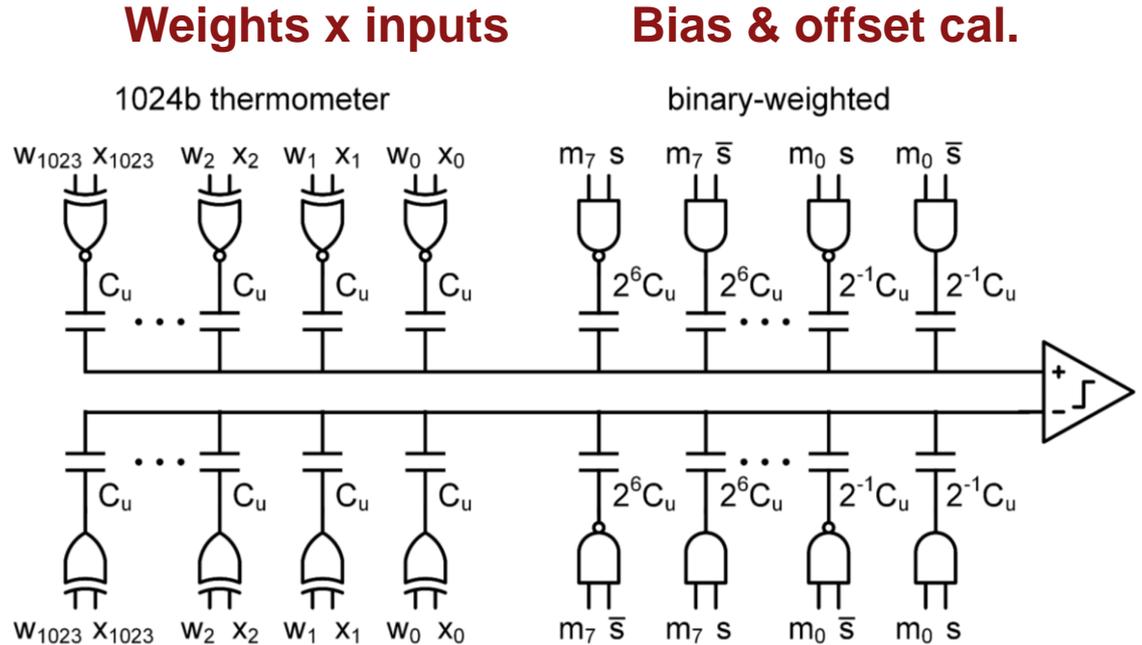


Switched-Capacitor Implementation

$$\frac{v_{\text{diff}}}{V_{DD}} = \frac{C_u}{C_{\text{tot}}} \left(\sum_{i=0}^{1023} w_i x_i + b \right)$$

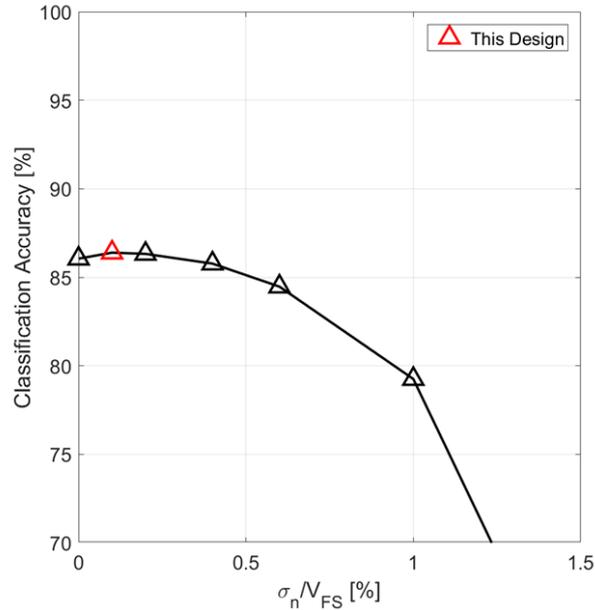
$$b = (-1)^s \sum_{i=0}^7 2^i m_i$$

- Batch normalization folded into weight signs and bias

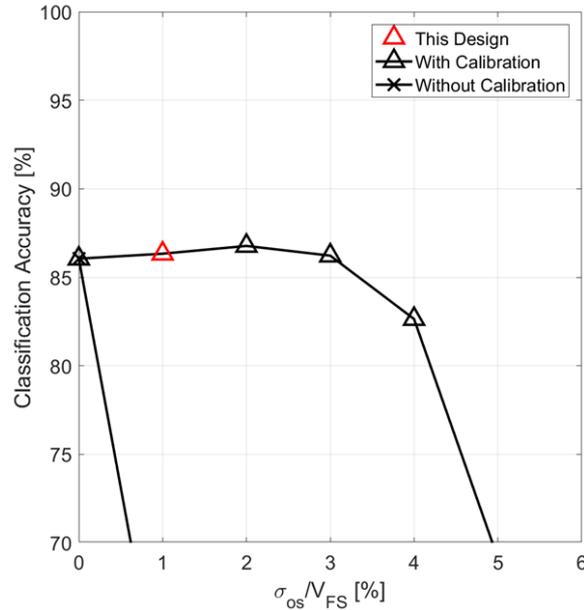


Behavioral Simulations

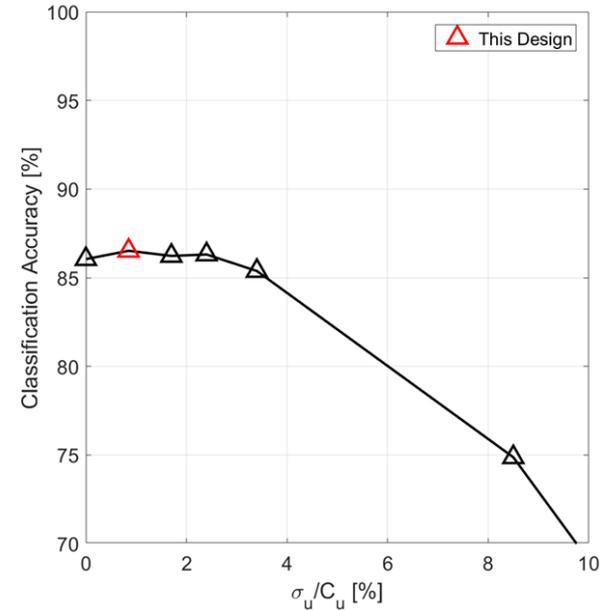
$v_n = 460 \mu\text{V RMS}$



$v_{os} = 4.6 \text{ mV RMS}$



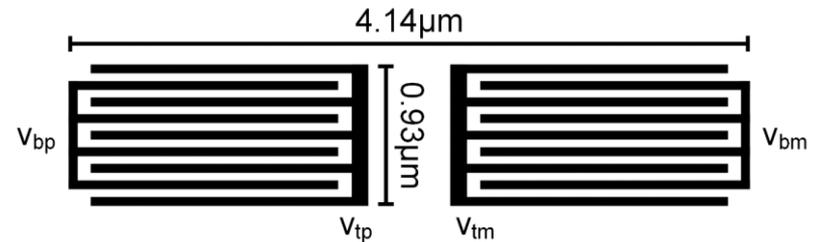
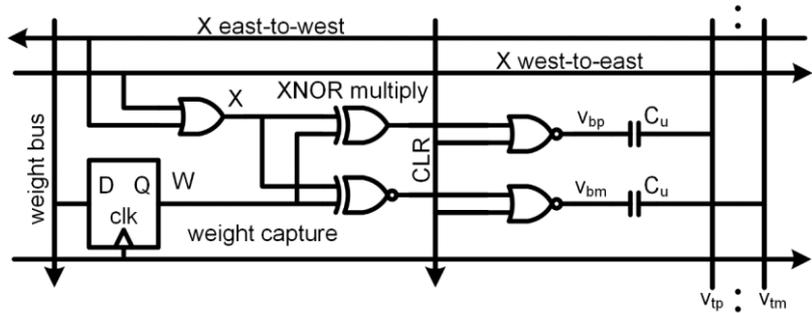
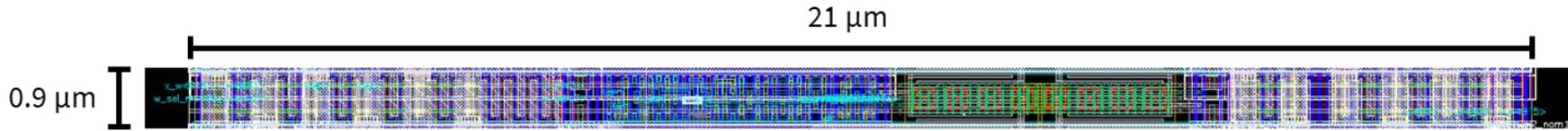
$C_u = 1 \text{ fF}$



- Significant margin in noise, offset, and mismatch ($V_{FS} = 460 \text{ mV}$)

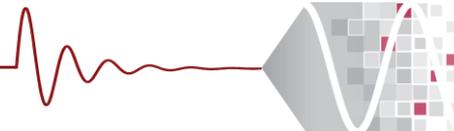


“Memory-Cell-Like” Processing Element



Standard-cell-based
42 transistors
24107 F²

1 fF metal-oxide-metal fringe capacitor



Die Photo

- TSMC 28nm HPL 1P8M
- 6 mm² area
- 328 KB SRAM
- 10 MHz clock

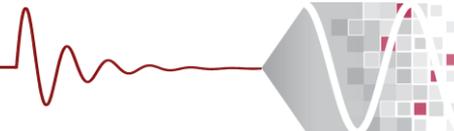
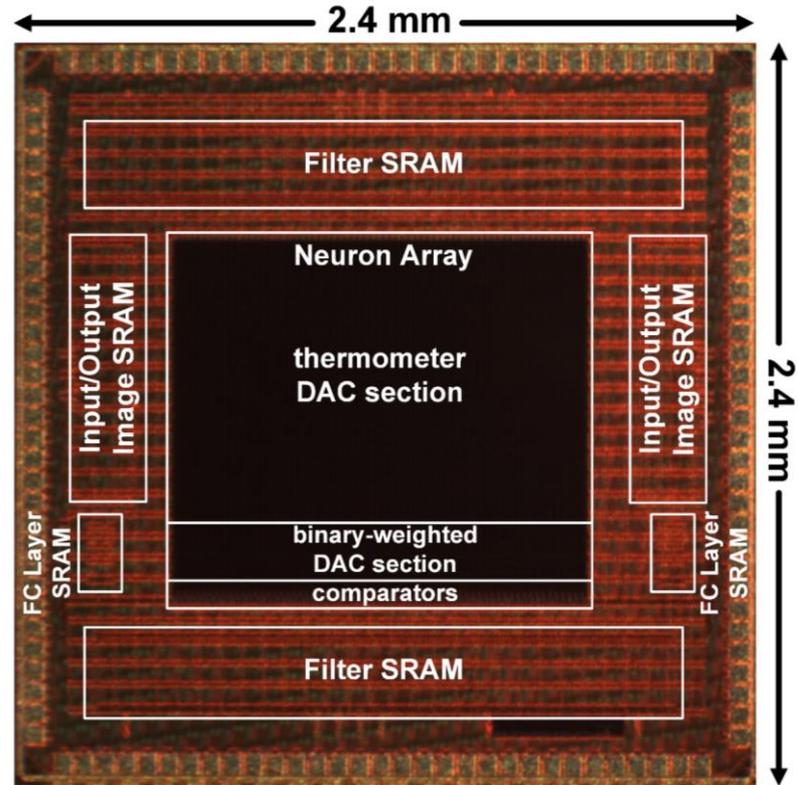
Supply Voltages

V_{DD} – Digital Logic, 0.6V – 1.0V

V_{MEM} – SRAM, 0.53V – 1.0V

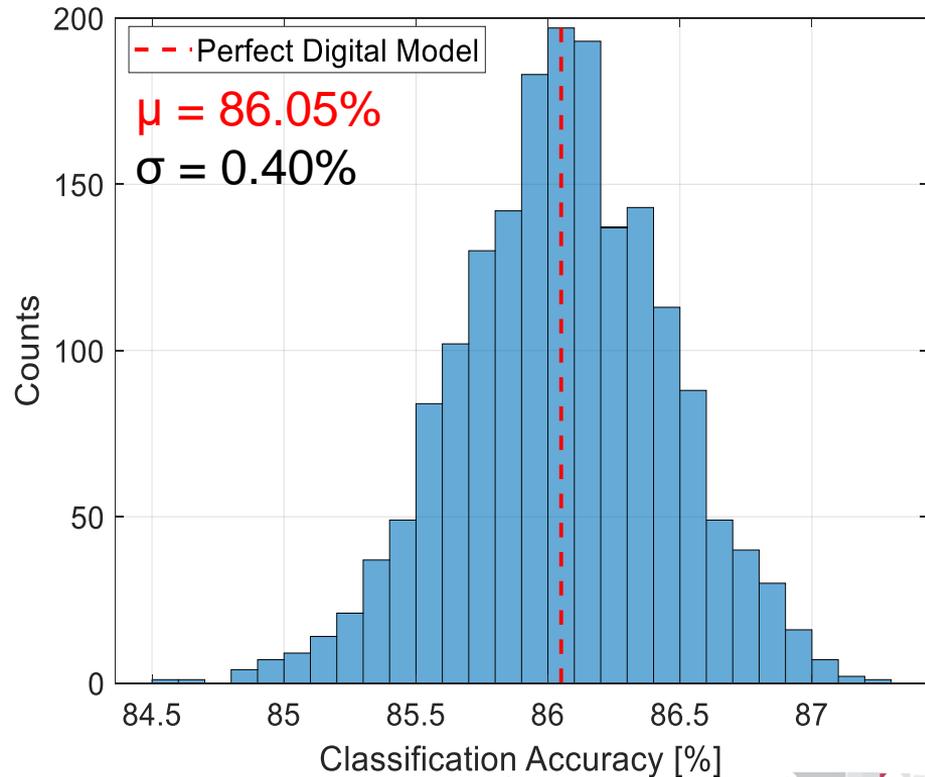
V_{NEU} – Neuron Array, 0.6V

V_{COMP} – Comparators, 0.8V



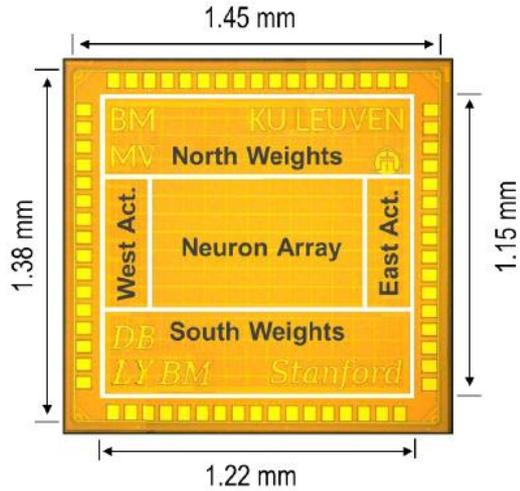
Measured Classification Accuracy

- 10 chips, 180 runs each through 10,000 CIFAR-10 test images
- $V_{DD} = 0.8V$, $V_{MEM} = 0.8V$
- 3.8 μJ /classification
- 237 FPS, 899 μW
- 0.43 μJ in 1.8V I/O
- Mean accuracy $\mu = 86.05\%$ same as perfect digital model



Comparison to Synthesized Digital

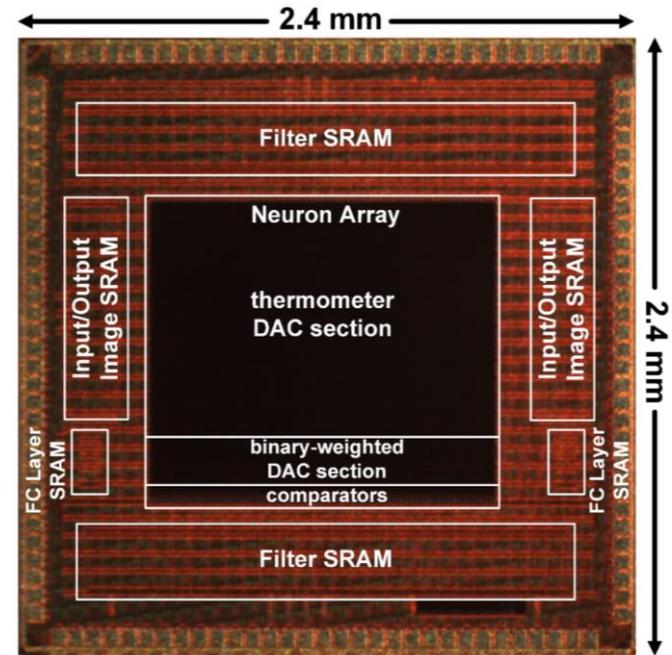
Synthesized Digital



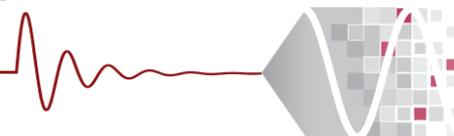
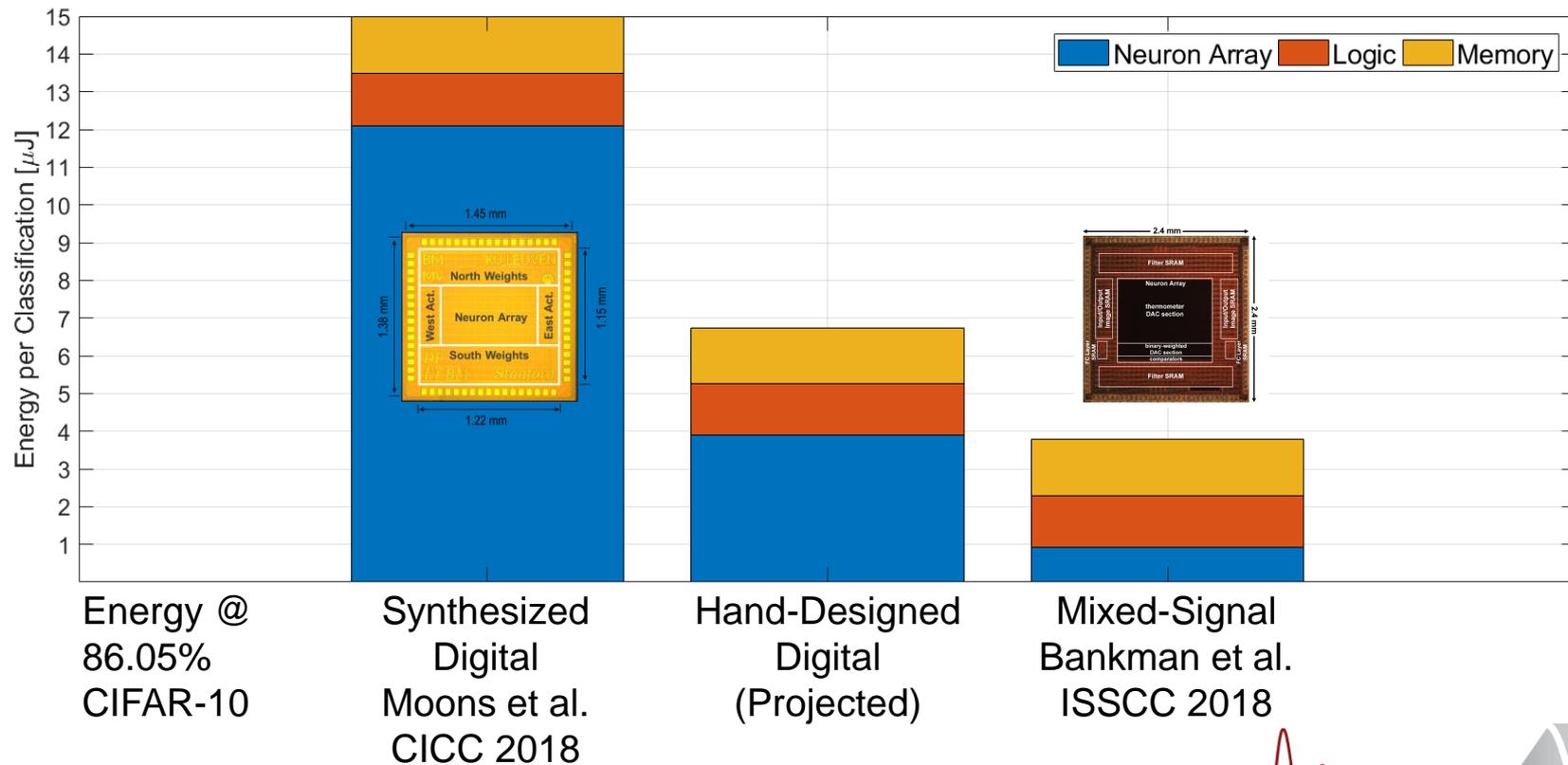
BinarEye

(Moons et al., CICC 2018)

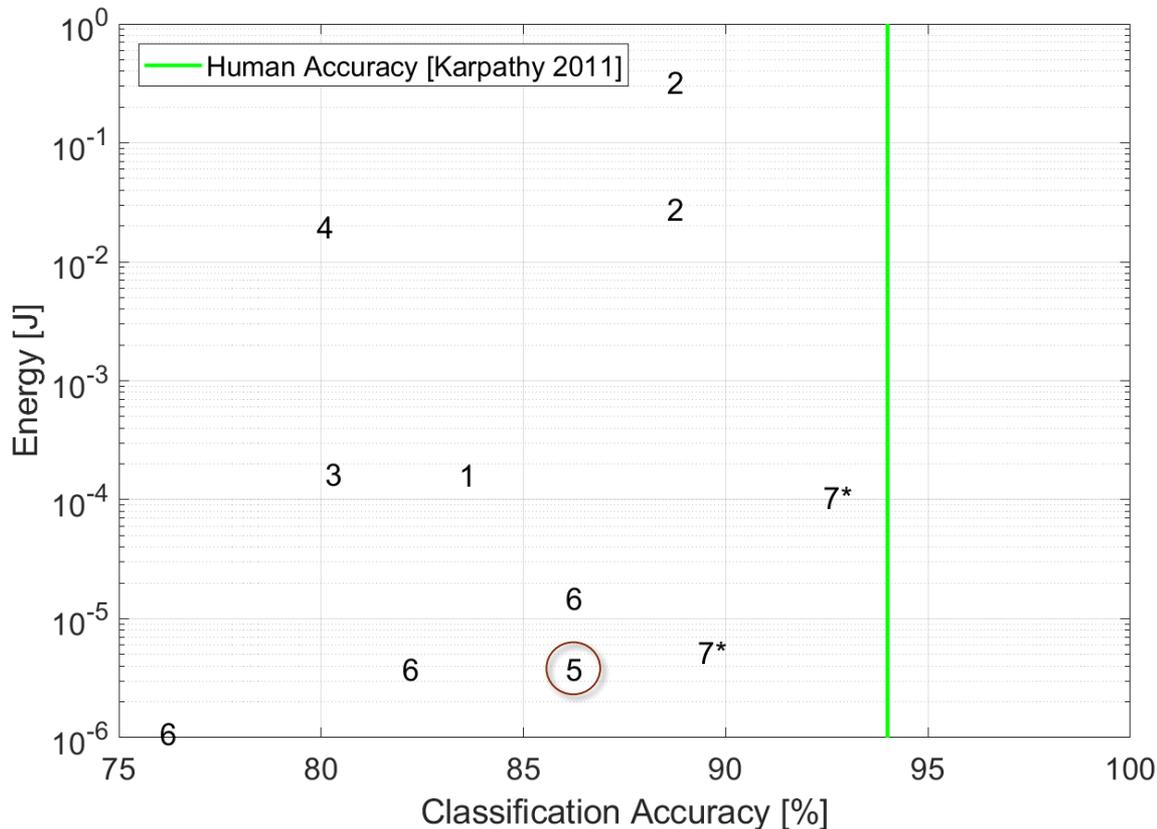
Mixed-Signal



Digital vs. Mixed-Signal Binary CNN Processor



CIFAR-10 Energy vs. Accuracy



- Neuromorphic
 - › [1] TrueNorth, Esser PNAS 2016
 - GPU
 - › [2] Zhao FPGA 2017
 - FPGA
 - › [2] Zhao FPGA 2017
 - › [3] Umuroglu FPGA 2017
 - MCU
 - › [4] CMSIS-NN, Lai arXiv 2018
 - **Memory-like, mixed-signal**
 - › [5] Bankman ISSCC 2018
 - BinarEye, digital
 - › [6] Moons CICC 2018
 - In-memory, mixed-signal
 - › [7] Jia arXiv 2018
- *energy excludes off-chip DRAM

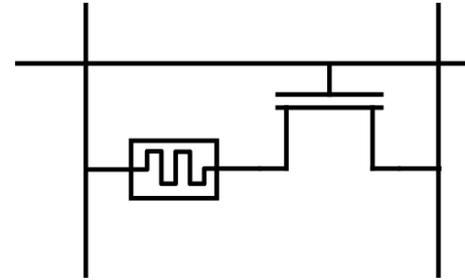
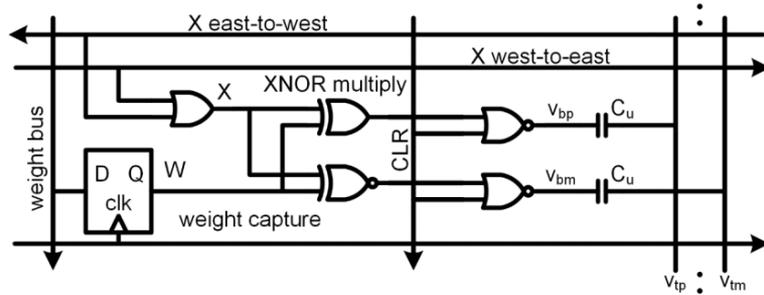


Limitations of Mixed-Signal BinaryNet

- Poor programmability
- Relatively limited accuracy (even on CIFAR-10) due to 1b arithmetic
- Energy advantage over customized digital is not revolutionary
 - › Same SRAM, essentially same dataflow
- **Need a more “analog” memory system to unleash larger gains**
 - › **In-memory computing**

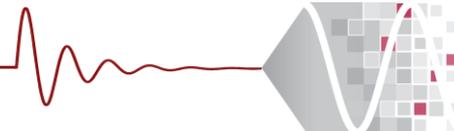


BinaryNet Synapse versus Resistive RAM

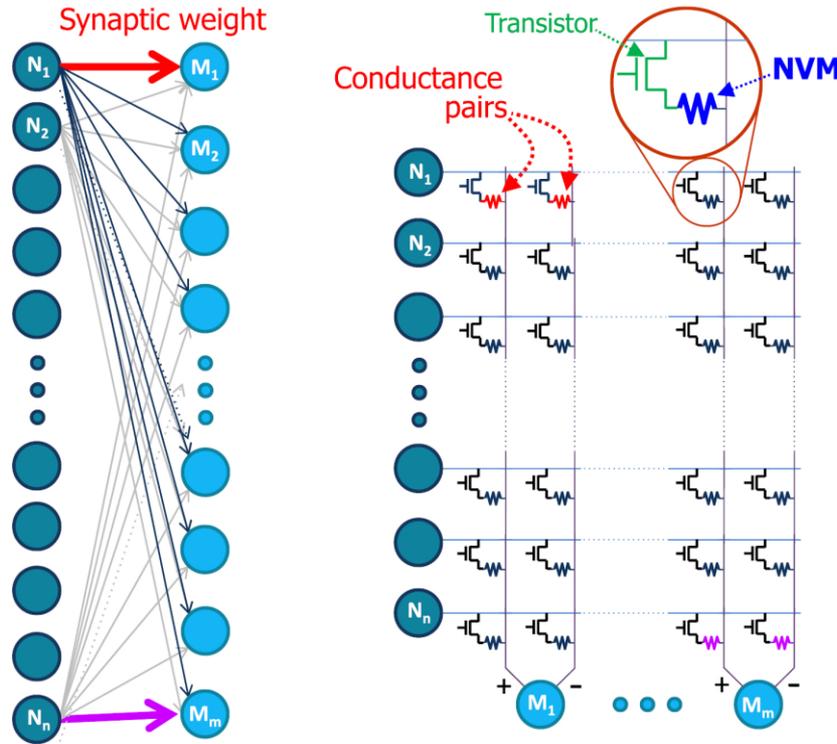


- 0.93 fJ per 1b-MAC in 28 nm
- **24107 F²**
- Single-bit

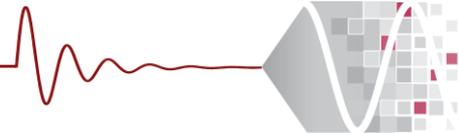
- TBD
- **25 F²**
- Multi-bit (?)



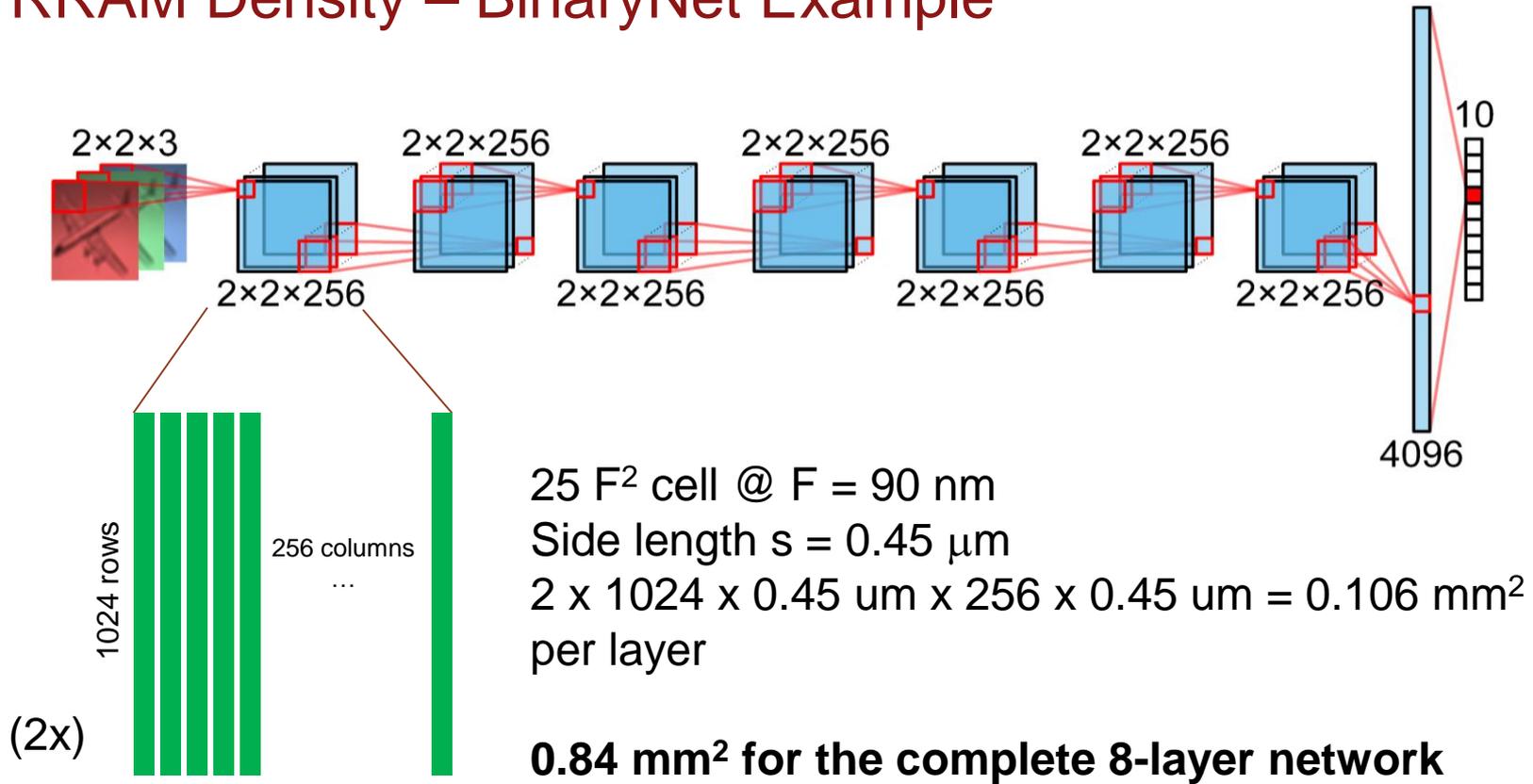
Matrix-Vector Multiplication with Resistive Cells



Typically use two cells to achieve pos/neg weights (other schemes possible)

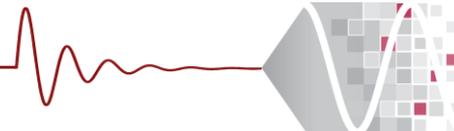
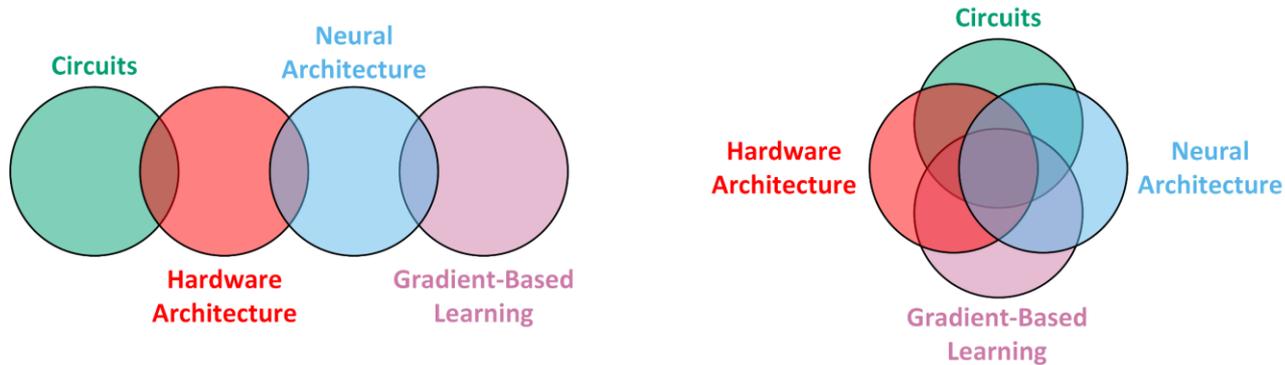


RRAM Density – BinaryNet Example



Ongoing Research

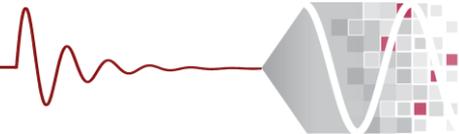
- What is the best architecture?
- How many levels can be stored in each cell?
- What is the most efficient readout?
- Can we cope with nonidealities using training techniques?



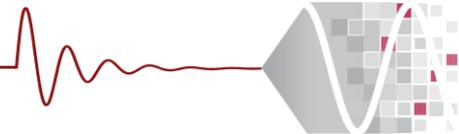
(Content deleted for online distribution)



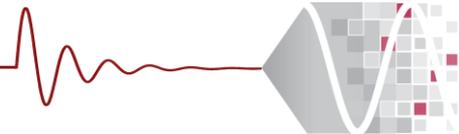
(Content deleted for online distribution)



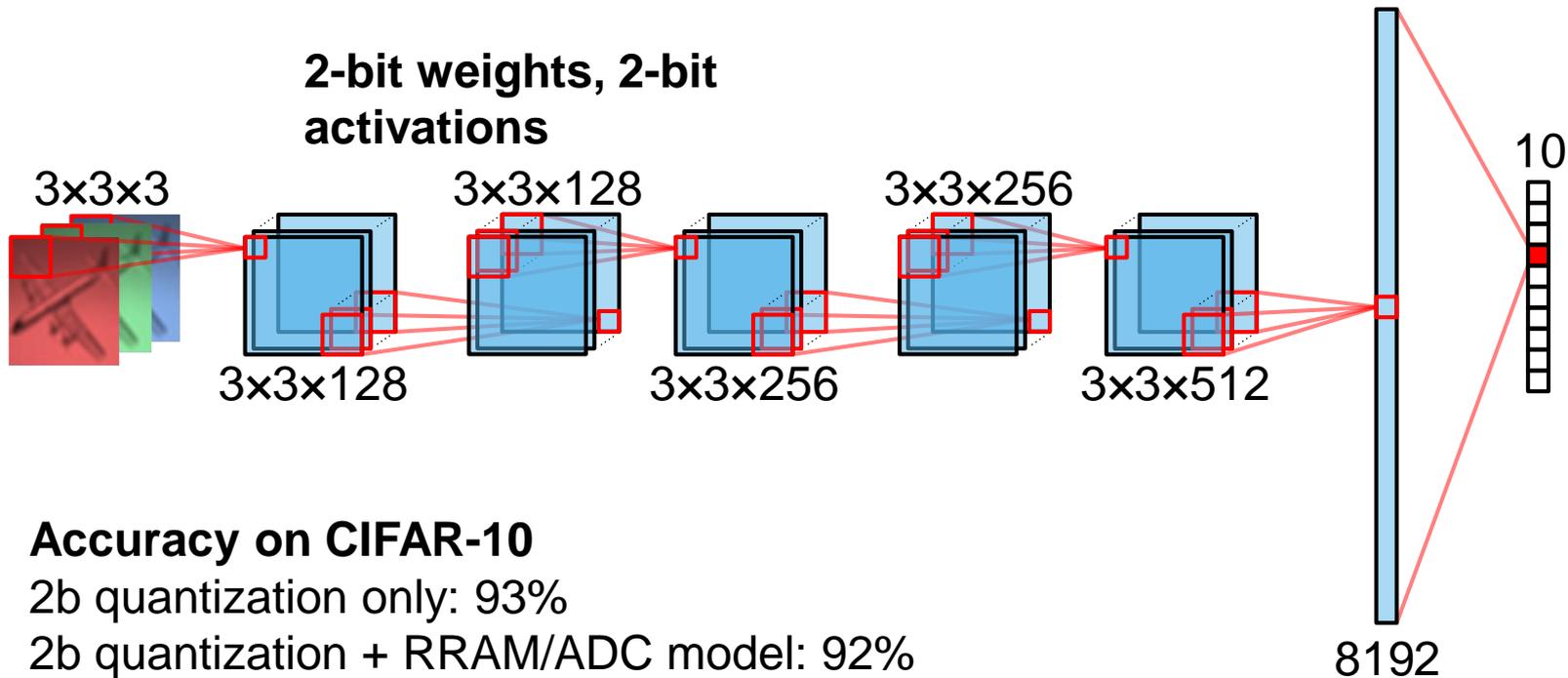
(Content deleted for online distribution)



(Content deleted for online distribution)



VGG-7 Experiment (4.8 Million Parameters)

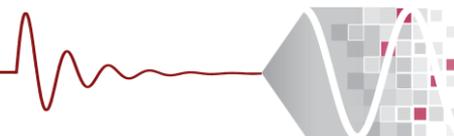


Accuracy on CIFAR-10

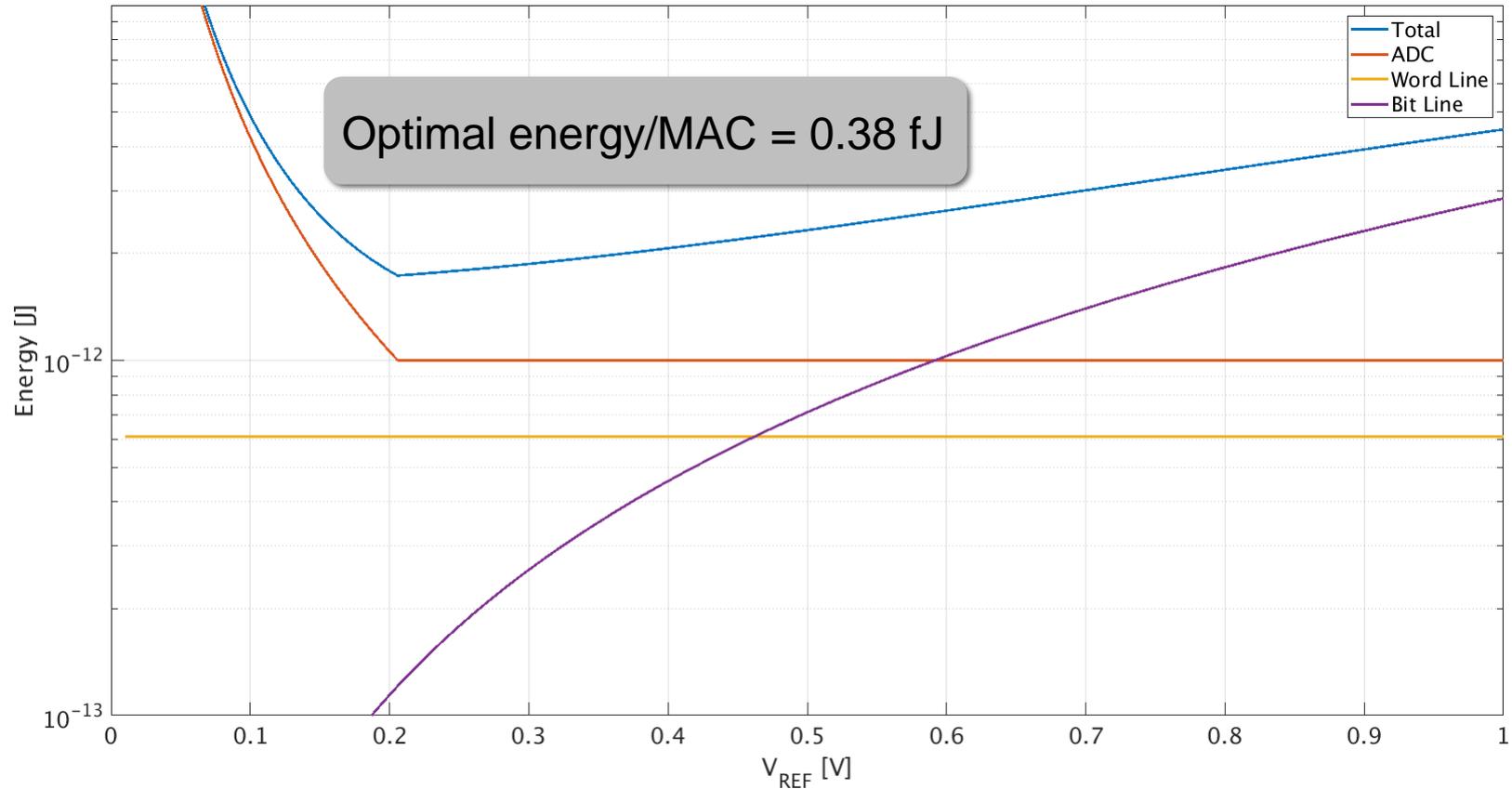
2b quantization only: 93%

2b quantization + RRAM/ADC model: 92%

Work in progress!



Energy Model for Column in Conv6 Layer



Optimal energy/MAC = 0.38 fJ



Summary

- Analog feature extraction is attractive for wake-up detectors
- Adding analog compute in ConvNets can be beneficial when it simultaneously lets us reduce data movement
 - › In-memory analog compute looks most promising
 - › Can consider SRAM or emerging memories (e.g. RRAM)
- Expect significant progress as more application drivers for “machine learning at the edge” emerge

