# A survey of Emerging Beyond CMOS Devices and Architectures
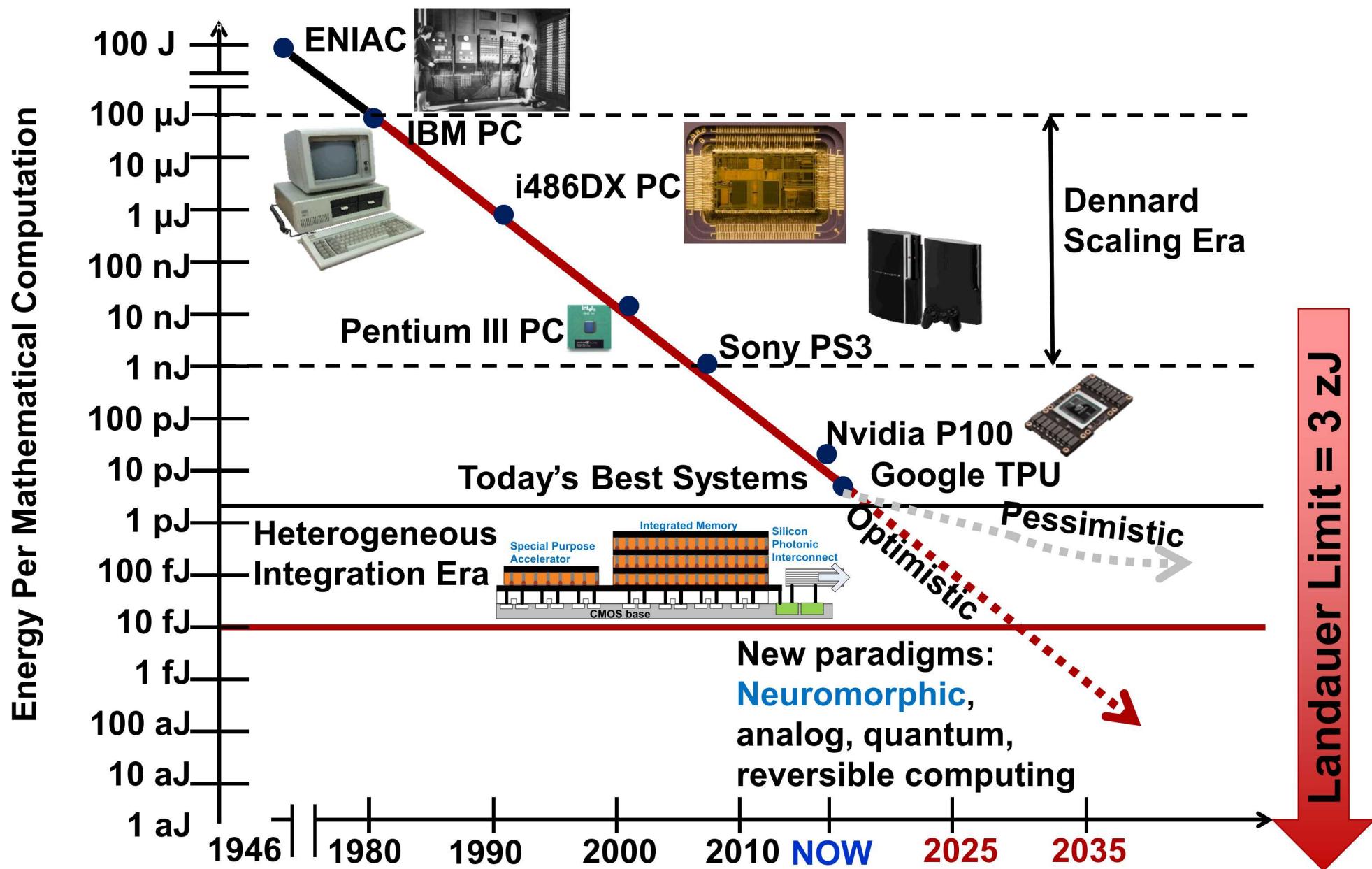
Sapan Agarwal

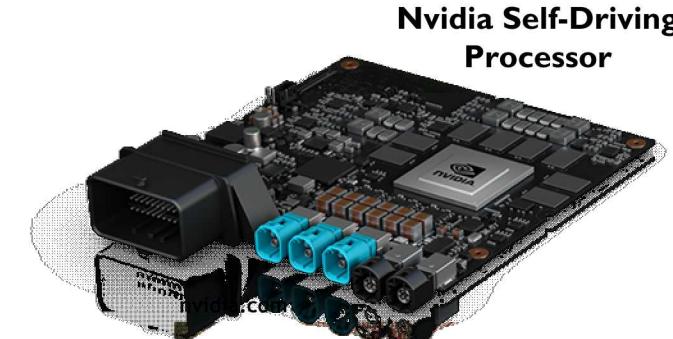Sandia National Laboratories

# Evolution of Computing Machinery

# Computing Across Power Envelopes

**IoT, Edge, and Mobile Computing**

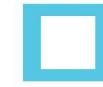**Self Driving Cars, Unmanned Arial Vehicles, and Satellite Computing**

**Datacenters, HPC**

**Nvidia Self-Driving Processor**

nest.com

wikimedia.org

**ASCI Red Supercomputer**

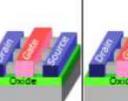$1W \qquad 10W \qquad 10^2W \qquad 10^3W \qquad 10^4W \qquad 10^5W \qquad 10^6W$
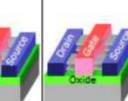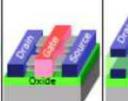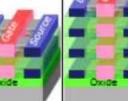
## Successor to ITRS (International Technology Roadmap for Semiconductors)

### Objective of the Beyond CMOS (BC) Chapter

Still road mapping near term semiconductors (Moore Moore)

**Novel computing paradigms and application pulls**

↑ Efficiency and performance

* Big data
* IoT and trillions of edge sensors
* Deep learning and artificial intelligence
* Exascale supercomputing
* Robotics and autonomous systems

**Beyond CMOS**

☞ *Emerging Architectures*
☞ *Emerging Devices / Processes*
☞ *Emerging Materials*

**More Moore trend**

→ Size

100 nm    10 nm    1 nm

| YEAR OF PRODUCTION | 2020 | 2022 | 2025 | 2028 | 2031 | 2034 |
|---|---|---|---|---|---|---|
| | G48M36 | G45M24 | G42M20 | G40M16 | G38M16T2 | G38M16T4 |
| Logic industry "Node Range" Labeling (nm) | "5" | "3" | "2.1" | "1.5" | "1.0 eq" | "0.7 eq" |
| IDM-Foundry node labeling | i7-f5 | i5-f3 | i3-f2.1 | i2.1-f1.5 | i1.5e-f1.0e | i1.0e-f0.7e |
| Logic device structure options | FinFET | finFET LGAA | LGAA | LGAA | LGAA-3D | LGAA-3D |
| Mainstream device for logic | finFET | finFET | LGAA | LGAA | LGAA-3D | LGAA-3D |
| **LOGIC DEVICE GROUND RULES** | | | | | | |
| Mx pitch (nm) | 36 | 32 | 24 | 20 | 16 | 16 |
| M1 pitch (nm) | 32 | 30 | 21 | 20 | 19 | 19 |
| M0 pitch (nm) | 30 | 24 | 20 | 16 | 16 | 16 |
| Gate pitch (nm) | 48 | 45 | 42 | 40 | 38 | 38 |
| $L_g$: Gate Length - HP (nm) | 18 | 16 | 14 | 12 | 12 | 12 |
| Lg: Gate Length - HD (nm) | 20 | 18 | 14 | 12 | 12 | 12 |
| Channel overlap ratio - two-sided | 0.20 | 0.20 | 0.20 | 0.20 | 0.20 | 0.20 |
| Spacer width (nm) | 7 | 6 | 5 | 4 | 4 | 4 |
| Contact CD (nm) - finFET, LGAA | 16 | 17 | 18 | 20 | 18 | 18 |
| Contact CD (nm) - VGAA | | | | | | |
| Device architecture key ground rules | | | | | | |
| FinFET pitch (nm) | 28.0 | 24.0 | | | | |
| FinFET Fin width (nm) | 7.0 | 6.0 | | | | |
| FinFET Fin height (nm) | 50 | 60 | | | | |
| Footprint drive efficiency - finFET | 3.82 | 5.25 | | | | |
| Lateral GAA lateral pitch (nm) | | | 22.0 | 20.0 | 20.0 | 20.0 |
| Lateral GAA vertical pitch (nm) | | | 18.0 | 16.0 | 14.0 | 14.0 |
| Lateral GAA (nanosheet) thickness (nm) | | | 7.0 | 6.0 | 5.0 | 5.0 |
| Number of vertically stacked nanosheets | | | 3 | 3 | 4 | 4 |
| LGAA width (nm) - HP | | | 30 | 20 | 15 | 10 |
| LGAA width (nm) - HD | | | 20 | 11 | 6 | 6 |
| LGAA width (nm) - SRAM | | | 7 | 6 | 6 | 6 |
| LGAA total height (nm) | | | 53 | 48 | 57 | 57 |
| Footprint drive efficiency - lateral GAA - HP | | | 4.80 | 4.59 | 5.52 | 5.00 |
| Device effective width (nm) - HP | 107.0 | 126.0 | 192.0 | 156.0 | 160.0 | 120.0 |
| Device effective width (nm) - HD | 107.0 | 126.0 | 132.0 | 102.0 | 88.0 | 88.0 |
| Device lateral pitch (nm) | 28 | 24 | 22 | 20 | 20 | 20 |
| Device height (nm) | 50.0 | 60.0 | 53.0 | 48.0 | 57.0 | 57.0 |
| Device width (nm) - HP | 7 | 6 | 25 | 20 | 15 | 10 |
| Device width (nm) - HD | 7 | 6 | 15 | 11 | 6 | 6 |
| Device width (nm) - SRAM | 7 | 6 | 7 | 6 | 6 | 6 |

# 2019 Beyond CMOS Team Members and Contributors

| | | | |
|---|---|---|---|
| **Sapan Agarwal** | Michael Fuhrer | Tsu-Jae King Liu | Takahiro Shinada |
| Brad Aimone | Mike Garner | **Matthew Marinella** | Urmita Sikder |
| **Hiro Akinaga** | Chakku Goplan | Bicky A. Marque | Greg Snider |
| Otitoaleke Akinola | Bogdan Govoreanu | Rivu Midya | John-Paul Strachan |
| Mustafa Badaroglu | Cat Graves | Yoshiyuki Miyamoto | Dimitri Strukov |
| Gennadi Bersuker | Kohei Hamaya | Johannes Muller | Naoyuki Sugiyama |
| Christian Binek | Masami Hane | Azad Naeemi | Tarek Taha |
| **Geoffrey Burr** | Jennifer Hasler | Mitchell A. Nahmias | Alexander N. Tait |
| Leonid Butov | Yoshihiro Hayashi | Emre Neftci | Shinichi Takagi |
| Kerem Camsari | Toshiro Hiramoto | Mike Niemier | Norikatsu Takaura |
| Gert Cauwenberghs | **D. Scott Holmes** | Dmitri Nikonov | Tsutomu Teduka |
| **An Chen** | Sharon Hu | Yutaka Ohno | Yasuhide Tomioka |
| Winston Chern | **Francesca Iacopi** | Chenyun Pan | **Wilman Tsai** |
| Supriyo Datta | Danielle Ilmeni | Ferdinand Peper | Tohru Tsuruoka |
| John Dallesasse | Jean Anne Incorvia | Shriram Ramanathan | Zhongrui Wang |
| **Shamik Das** | Engin Ipek | Mingyi Rao | R. Stanley Williams |
| **Erik DeBenedictis** | Satoshi Kamiyama | **Shashi Paul** | Justin Wong |
| Peter Dowben | Kiyoshi Kawabata | Paul R. Prucnal | Dirk Wouters |
| Tetsuo Endoh | Asif Khan | Titash Rakshit | Patrick Xiao |
| Ben Feinberg | Hajime Kobayashi | Arijit Raychowdhury | Kojiro Yagami |
| Thomas Ferreira de Lima | Suhas Kumar | Sayeef Salahuddin | J. Joshua Yang |
| Akira Fujiwara | Ilya Krivorotov | Shintaro Sato | Noboyuki Yoshikawa |
| Elliot Fuller | Xiuling Li | Michael Schneider | Victor Zhirnov |
| **Michael Frank** | Xiang (Shaun) Li | Bhavin J. Shastri | |
| Paul Franzon | Shy-Jay Lin | Xia Sheng | |

# What's the Minimum Energy to Operate a Transistor?

Consider a signal Energy $E_{signal}$
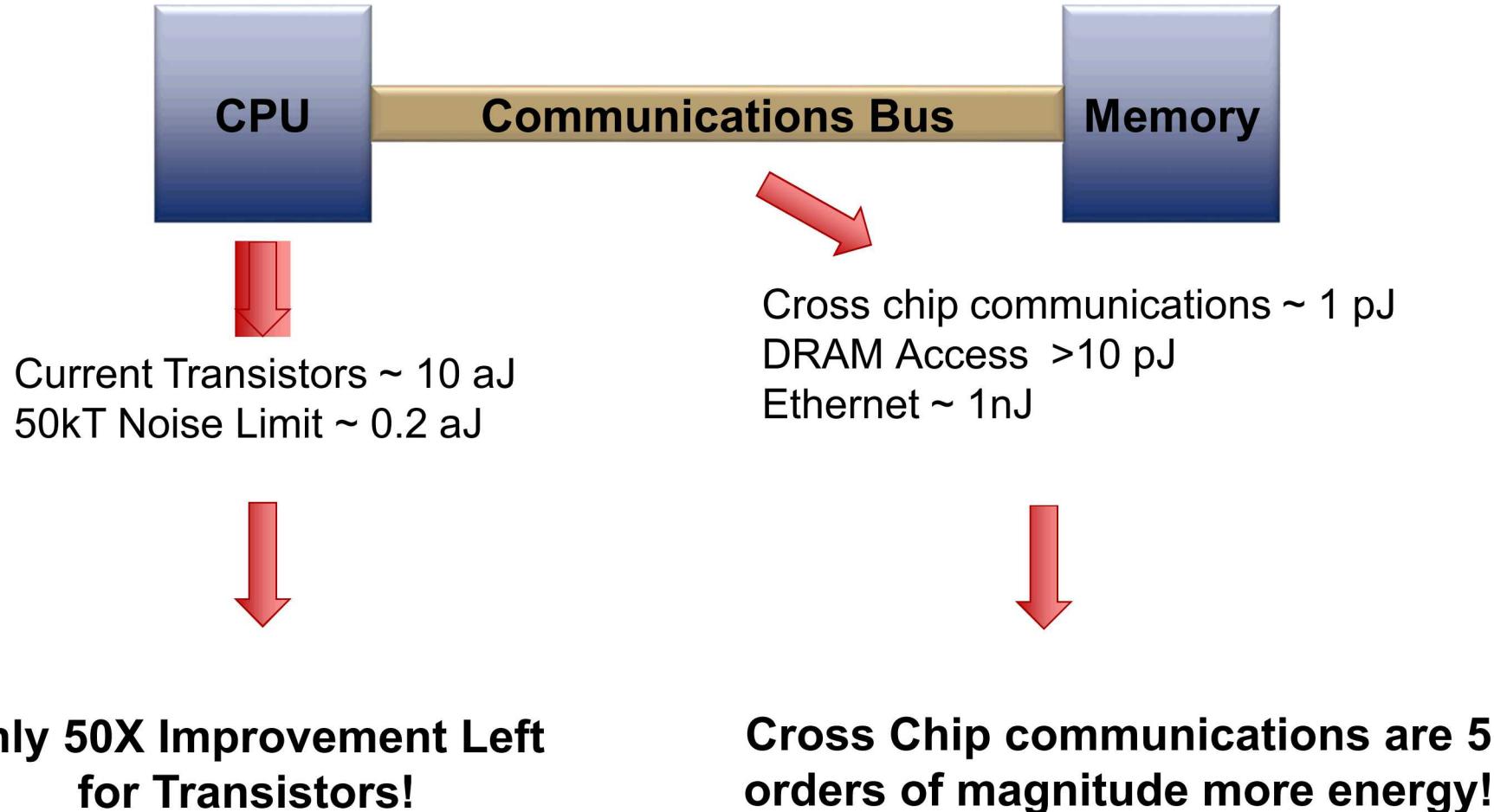
The probability of an error due to thermal noise is:

$$P(Error) = e^{-E_{signal}/kT}$$

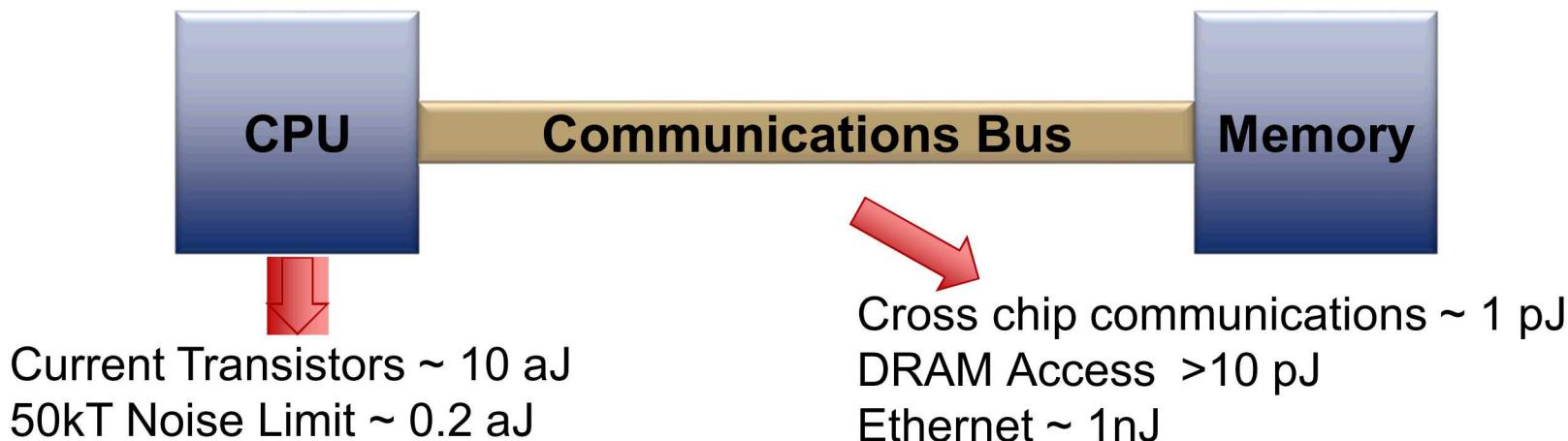In order to ensure a full system with billons of transistors is reliable, we need:

$$E_{signal} \sim 50 \text{ kT}$$

Landauer – Shannon Limit

# How Efficient are Current Systems?

**CPU** | **Communications Bus** | **Memory**

Current Transistors ~ 10 aJ
50kT Noise Limit ~ 0.2 aJ

Cross chip communications ~ 1 pJ
DRAM Access  >10 pJ
Ethernet ~ 1nJ

**Only 50X Improvement Left for Transistors!**

**Cross Chip communications are 5 orders of magnitude more energy!**

# Beyond Moore Technologies

**CPU** ——— **Communications Bus** ——— **Memory**

Current Transistors ~ 10 aJ
50kT Noise Limit ~ 0.2 aJ

Cross chip communications ~ 1 pJ
DRAM Access  >10 pJ
Ethernet ~ 1nJ

**Extending Von Neumann**

- Low Voltage or Novel Transistors
- Optical Communications
- Reduced Data Movement
  - New On-Chip Memory
  - Processing near Memory

**Alternate Computing Paradigms**

- Neuromorphic
- Analog
  - Computing with memory devices
- Quantum
- Stochastic
- Approximate

**Going Below 50 kT**

- Error Correction
- Reversible Computing
  - Adiabatic Computing / Energy Recycling
- Superconducting

# Beyond Moore Technologies

**CPU** — **Communications Bus** — **Memory**

Current Transistors ~ 10 aJ
50kT Noise Limit ~ 0.2 aJ

Cross chip communications ~ 1 pJ
DRAM Access  >10 pJ
Ethernet ~ 1nJ

**Extending Von Neumann**

- Low Voltage or Novel Transistors
- Optical Communications
- Reduced Data Movement
  - New On-Chip Memory
  - Processing near Memory

**Alternate Computing Paradigms**

- Neuromorphic
- Analog
  - Computing with memory devices
- Quantum
- Stochastic
- Approximate

**Going Below 50 kT**

- Error Correction
- Reversible Computing
  - Adiabatic Computing / Energy Recycling
- Superconducting

# Extending Von Neumann – New Transistors

*Lowering voltage lowers CV² energy of communications*

## Logic and Information Processing Devices



**Domain Wall Logic**



J.A. Incorvia et al, Nature Comm 7, 2016

**Negative Capacitance FET**



Wong and Salahuddin, TED 2019

# A New Switch has to Satisfy Three Specifications

## Low Active Power ($CV^2$ energy)

- Steepness (or sensitivity)
  - switches with only a few milli-volts
  - 60mV/decade $\Rightarrow$ **1mV/decade**

## Low Leakage Power

- On/Off ratio:  **$10^5 : 1$**

## High Speed (RC delay)

- High Conductance Density
  - **1 milli-Siemen/micron**

log{I}

Current

steeper
sub-threshold
swing

Gate Voltage

$V_g$

# Need Restoring Logic for Novel State Variables

- Need to be able to drive multiple output transistors across multiple stages of logic

Not Restoring

V(out)

V(in)



Figure from Elad Elon, 2010 E3S Retreat

Restoring

Inputs slightly below $V_{dd}$ restored to $V_{dd}$ and inputs slightly above 0 are driven to zero

https://www.allaboutcircuits.com/technical-articles/restoring-digital-signals-in-pass-transistor-logic/

# The Memory Hierarchy



10⁷ difference in speed from HDD to Registers

10¹⁰ difference in density from HDD to Registers

Source: http://www.ts.avnet.com/uk/products_and_solutions/storage/hierarchy.html

# Storage Class Memory - Intel/Micron 3D XPoint



Also Samsung Z-NAND:
- Re-optimize flash for speed rather than density (single level per cell)
- Comparable to Intel Optane (3D XPoint) products
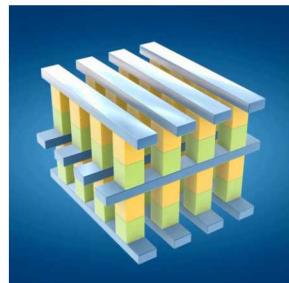
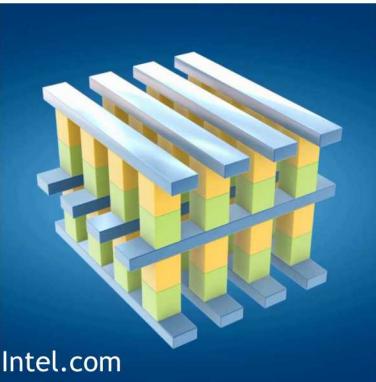# Extending Von Neumann – Reduced Data Movement

## Optical Interconnects



Processor Layer

Photonic Layer

http://www.bu.edu/ipl/research.html

100 fJ to 1 pJ

## 2.5D & 3D Integration



Richard Goering, "Three Die Stack -- A Big Step "Up" for 3D-ICs with TSVs" Cadence blog
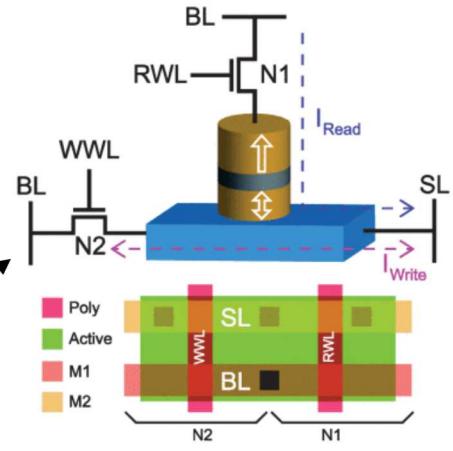
## Embedded Memory



- Embed new denser nonvolatile memories in the processor
- Add simple processing to DRAM or disk drive controllers "Processing in Memory"

# Emerging Memory Devices

**SOT Memory**



Z. Wang EDL 39, 2018

**Memory**

- **Volatile**
  - SRAM
  - DRAM
    - Stand-alone
    - Embedded
- **Nonvolatile**
  - **Baseline**
    - Flash
      - NOR
      - NAND
  - **Prototypical**
    - FeRAM
    - MRAM
    - PCM
    - STT-RAM
    - ReRAM
  - **Emerging**
    - Novel Magnetic Memory
    - Ferroelectric Memory
      - FeFET
      - FTJ
    - OxRAM
    - CBRAM
    - Macromolecular Memory
    - Mott Memory
    - Massive Storage Devices

**Storage Class Memory**



Intel.com

**Successfully tracked and transferred to More Moore**

**STT-MRAM**



Free Layer
Tunnel Barrier
Fixed Layer

**DNA Memory**



STORAGE LIMITS
Estimates based on bacterial genetics suggest that digital DNA could one day rival or exceed today's storage technology.

WEIGHT OF DNA NEEDED TO STORE WORLD'S DATA

|  | Hard disk | Flash memory | Bacterial DNA |
|---|---|---|---|
| Read–write speed (μs per bit) | ~3,000–5,000 | ~100 | <100 |
| Data retention (years) | >10 | >10 | >100 |
| Power usage (watts per gigabyte) | ~0.04 | ~0.01–0.04 | <$10^{-10}$ |
| Data density (bits per cm³) | ~$10^{13}$ | ~$10^{16}$ | ~$10^{19}$ |

~1 kg

©nature

# Advantages of Emerging Memories
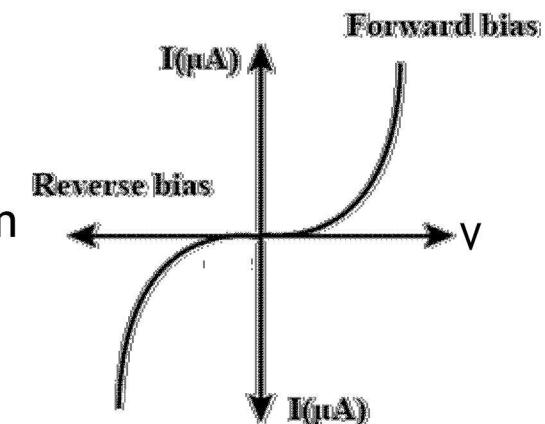
- Back end of line integration
  ◦ Can integrate in the metal layers directly on top of logic

- 3D stackable

- Higher endurance ($10^9$ to $10^{12}$) relative to flash ($10^4$), but not as good as DRAM ($>10^{16}$)

- Nonvolatile: >10 year retention, no standby leakage

- O(10ns) read and write times

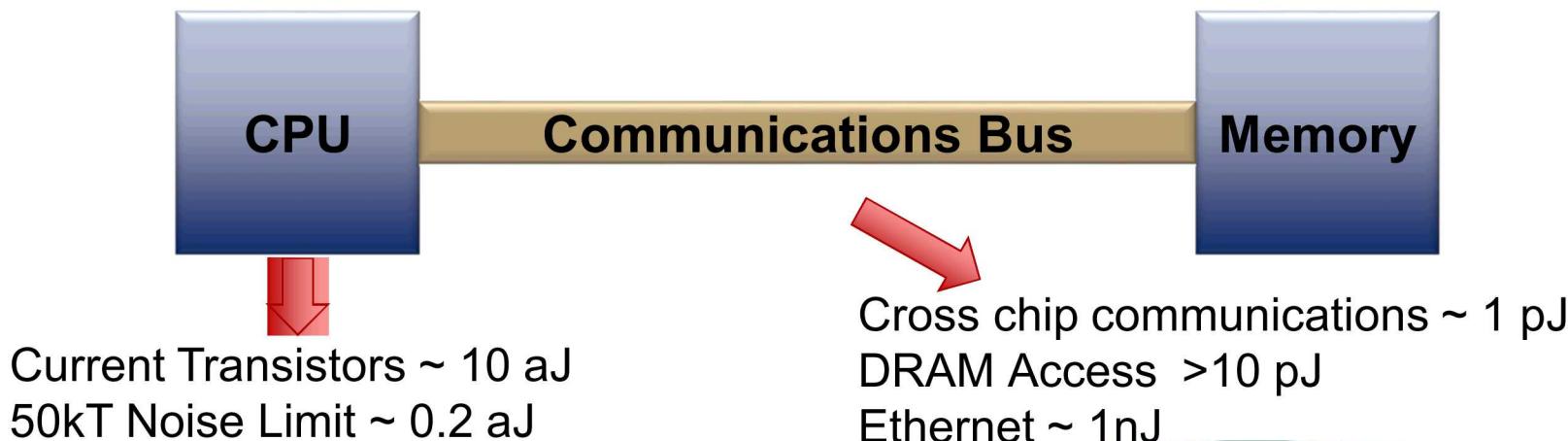# The Need for a Select Device for Resistive Memories

## V/3 write scheme

Target device sees $W_{write}$

$V_W/3$

$-V_W/3$

$-V_W/3$

$-V_W/3$

0        $-2V_W/3$        0        0

Half selected device sees $V_{write}/3$

**Memory Selector Devices**

| Transistor | Diodes | Volatile switch | Nonlinear devices |
|---|---|---|---|
| Planar | Si diodes | Threshold switch | Tunneling-based nonlinear selector |
| Vertical | Oxide/oxide heterojunctions | Mott switch | MIEC |
| | Metal/oxide Schottky junctions | | Complementary structures |
| | Reverse-conduction diodes | | Intrinsic nonlinearity |
| | Self-rectification | | |

Forward bias

$I(\mu A)$

Reverse bias

V

$I(\mu A)$

- Need a nonlinear 2 terminal back end of line compatible device to block current from half selected and unselected devices
- The larger the on/off ratio, the larger the possible array and therefore density
- Typically needs to be bi-directional

# Beyond Moore Technologies

**CPU** — **Communications Bus** — **Memory**

Current Transistors ~ 10 aJ
50kT Noise Limit ~ 0.2 aJ

Cross chip communications ~ 1 pJ
DRAM Access >10 pJ
Ethernet ~ 1nJ

**Extending Von Neumann**

- Low Voltage or Novel Transistors
- Optical Communications
- Reduced Data Movement
  - New On-Chip Memory
  - Processing near Memory

**Alternate Computing Paradigms**

- Neuromorphic
- Analog
  - Computing with memory devices
- Quantum
- Stochastic
- Approximate

**Going Below 50 kT**

- Error Correction
- Reversible Computing
  - Adiabatic Computing / Energy Recycling
- Superconducting

# Alternate Computing Paradigms

**Crossbar Based Computing Architectures**

- Vector Matrix Multiplication

- Outer Product Update

- Crossbar Based Matrix Solvers

- Ternary Content Addressable Memory

**Neuro-inspired Computing**

- Hyperdimensional Computing

- Local Learning Rules

- Spiking Neural Networks

**Probabilistic and Stochastic Circuits**

**Computing With Dynamical Systems**

- Simulated Annealing

- Coupled Oscillator based energy minimization

# Alternate Computing Paradigms

**Crossbar Based Computing Architectures**

- Vector Matrix Multiplication
- Outer Product Update

Neural Network
Training Accelerators

- Crossbar Based Matrix Solvers

- Ternary Content Addressable Memory

**Neuro-inspired Computing**

- Hyperdimensional Computing

- Local Learning Rules

- Spiking Neural Networks

**Probabilistic and Stochastic Circuits**

**Computing With Dynamical Systems**

- Simulated Annealing

- Coupled Oscillator based energy minimization

# Use Resistive Memories for Local Computation

$$V = I \times R$$

$$I = G \times V \longleftarrow$$

multiplication

- A resistive memory or ReRAM is a programmable resistor
  - Apply small voltages allows the conductance to be read: $I = G \times V$
  - Apply large voltages to change the resistance

$I_1$

$I_2$

Addition: $I = I_1 + I_2$

Current

Read Window

$V_{RESET}$

SET

Voltage

$V_{READ}$   $V_{SET}$

Write

RESET

OFF

Pt

TaO$_x$

Ta

ON

Pt

Ta

# Vector Matrix Multiply: Directly Process in the Memory Itself

## Mathematical

$$V^T W = I$$

$$\begin{bmatrix} V_1 & V_2 & V_3 \end{bmatrix} \begin{bmatrix} W_{1,1} & W_{1,2} & W_{1,3} \\ W_{2,1} & W_{2,2} & W_{2,3} \\ W_{3,1} & W_{3,2} & W_{3,3} \end{bmatrix} =$$

$$\begin{bmatrix} I_1 = \Sigma V_{i,1} W_{i,1} & I_2 = \Sigma V_{i,2} W_{i,2} & I_3 = \Sigma V_{i,3} W_{i,3} \end{bmatrix}$$

## Electrical



$I_1 = \Sigma V_{i,1} G_{i,1} \quad I_2 = \Sigma V_{i,2} G_{i,2} \quad I_3 = \Sigma V_{i,3} G_{i,3}$

**<10 fJ MAC**
**<10 fJ Update**
**>100 TOPS/W**

Analog is efficiently and naturally able to combine computation and data access

Large-scale processing in memory with a multiplier and adder at each real-valued memory location

- Energy to charge the crossbar is $CV^2$
- $E \propto C \propto$ number of RRAMs $\propto N \times M$

# SRAM Arrays Require Charging Columns Multiple Times

N rows

WL[0]

WL[1]

WL[2]

BL[0]    BL[1]    BL[2]

M columns

SRAMs must be read one row at a time, charging M columns
Each column wire length is O(N).

Energy = N Rows × M Columns × O(N) wire length
Energy ~ $O(N^2 \times M)$
O(N) times worse than a crossbar!

# The Noise Limited Energy to Read a Crossbar Column is Independent of Crossbar Size

$$I_o = G_o V$$

$$I_o = G_o V$$

$$I_o = G_o V$$

Measure N resistors and determine the total output current with some signal to noise ratio (SNR)[*]

What is the minimum energy?

$$Energy = V^2 G_O \times N \times \frac{1}{\Delta f}$$

Power in each resistor × number of resistors

Determined by noise and SNR

$$\text{Thermal Noise} = \langle \Delta I^2 \rangle$$

$$= N \times \left( 4 k_b T \times G_o \times \Delta f \right)$$

$$SNR^2 = \frac{(N I_o)^2}{\langle \Delta I^2 \rangle}$$

$$\frac{1}{\Delta f} = 4 k_b T \times SNR^2 \times \frac{1}{V^2 G_o \times N}$$

If we double the number of resistors, we can double the speed to get the same energy and SNR.

This is because the noise scales as sqrt(N) while the signal scales as N

$$Energy = 4 k_b T \times SNR^2$$

*we are assuming we need some fixed precision on the output, and don't need full floating point accuracy

# Need to Use Analog to Efficiently Discard Precision

$V_1 = x_1$  - +

$w_{11}$

$V_2 = x_2$  - +

$w_{21}$

$V_3 = x_3$  - +

$w_{31}$

$V_4 = x_4$  - +

$w_{41}$

Sum 1024 8 bit weights X 8 bit inputs:
- Result has 26 bits of information!
- A 26 bit ADC would eliminate any analog advantage!

The sum can be done at full precision in analog, but a lower precision approximation is needed when digitizing
- i.e. digitize only 8 bits or fewer

To get the highest 8 bits of information, digital would need to keep a 26 bit intermediate result

Can design an ADC to choose non uniform values to digitize

Neuron Function

Analog Sum

# Outer Product Update: Parallel Write



Energy to charge the crossbar is $CV^2$
$E \propto C \propto$ number of RRAMs $\propto N{\times}M$

$$E \sim O(N{\times}M)$$

# Example: Design an Neural Network Training Accelerator

## Vector Matrix Multiply



## Matrix Vector Multiply



## Outer product Update



## Chip Architecture



## Neuron Circuitry

# Energy Area and Latency Advantages of an Analog Accelerator

1024 x1024 = 1M array operations, sum over 1 training cycle, 3 operations:
- Vector Matrix Multiply
- Matrix Vector Multiply
- Outer Product Update



Energy
430 – 6,900X over SRAM

Latency
35 – 800X over SRAM

Area
11 – 20X over SRAM

Legend:
- 8 bit in/out 8 bit weights
- 4 bit in/out 8 bit weights
- 2 bit in/out 8 bit weights

Used a commercial 14/16 nm PDK        ***Requires 100 MΩ on state devices

M. J. Marinella, S. Agarwal, A. Hsia, I. Richter, R. Jacobs-Gedrim, J. Niroula, S. J. Plimpton, E. Ipek, and C. D. James, "Multiscale Co-Design Analysis of Energy, Latency, Area, and Accuracy of a ReRAM Analog Neural Training Accelerator," *IEEE Journal on Emerging and Selected Topics in Circuits and Systems,* 2018.

# Go from Measurement to Accuracy

**Fabricate Device**

TiN

TaO$_x$ – 10 nm

Ta– 50 nm

TiN

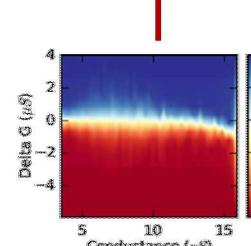**Measured Pulsing**

**ΔG Scatterplot**

**Cumulative Probability of ΔG**

positive weights

negative weights

D/A

A/D

R  Bus  R  Bus  R

Digital Core

Digital Core

Neural Core(s)

Neural Core(s)

R  Bus  R  Bus  R

Digital Core

Digital Core

Neural Core(s)

Neural Core(s)

R  Bus  R  Bus  R

Router

**Large Digits**

Exp. Derived

Ideal Numeric

# Multiscale Model of a Neural Training Accelerator



**Small Digits** — **File Types** — **Large Digits**

**Target Algorithms**
- Deep Learning
- Sparse Coding
- Liquid State Machines

**Algorithms**

**Architecture**

*Modified McPAT/CACTI:* Model performance and energy requirements

**Sandia Cross-Sim:** Translates device measurements and crossbar circuits to algorithm-level performance

**Circuits**

**Sandia's Xyce Circuit Sim:** Simulate crossbar circuits based on our devices

*Memristor fabrication and measurements in MESAFab*

*Drift-diffusion model of ReRAM band diagram & transport (REOS, Charon)*

**Devices**

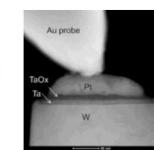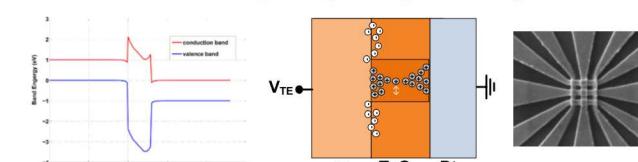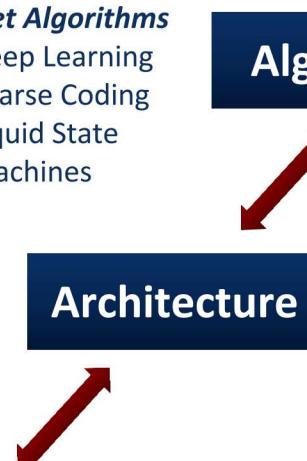**Materials**

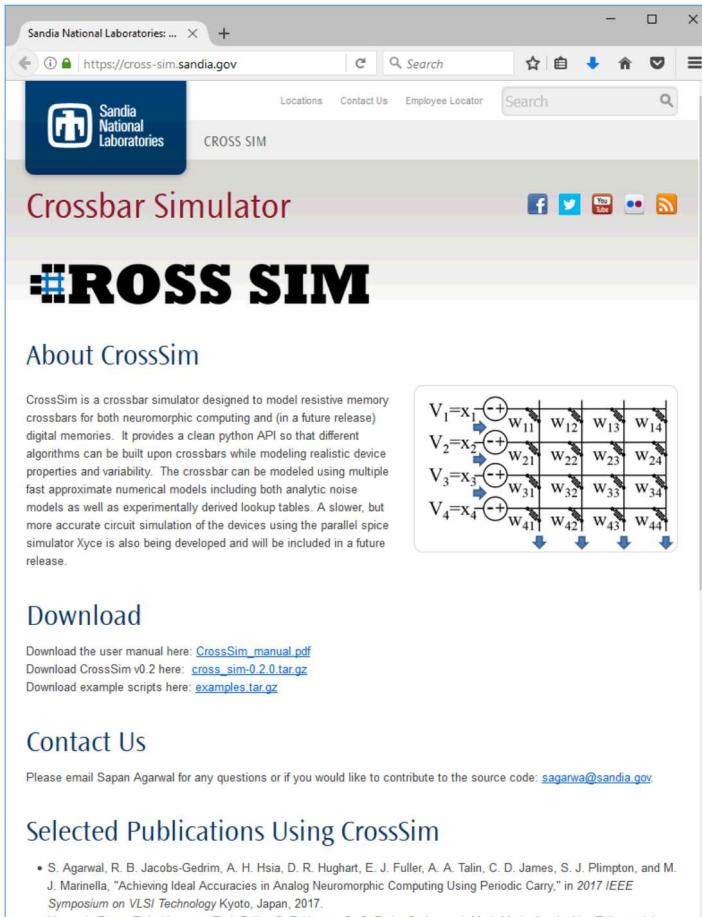*In situ TEM of filament switching:* Use DFT model to interpret EELS signature

*DFT of model of oxide physics, bands*

# ✇ROSS SIM

## https://cross-sim.sandia.gov



**Simple Python API:**

*# Do a matrix vector multiplication*
result = neural_core.run_xbar_mvm(vector)



Learning Algorithm

Neural Core Simulator

$$\begin{pmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{pmatrix}$$

Xyce Crossbar Circuit Model

Physical Hardware Crossbar

Numeric Crossbar Simulator

**Detailed but slow**

**Fast but approximate**

**Measured Devices**
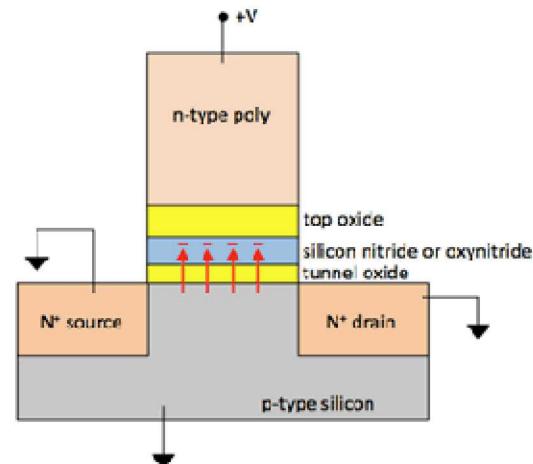
**Algorithmic Performance**
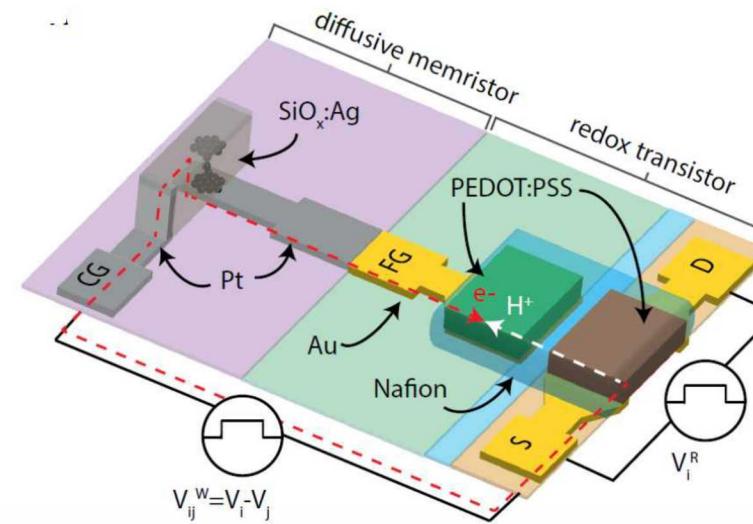
# Compare Analog Devices

**ReRAM**

**SONOS
Silicon-Oxygen-
Nitrogen-Oxygen-Silicon**
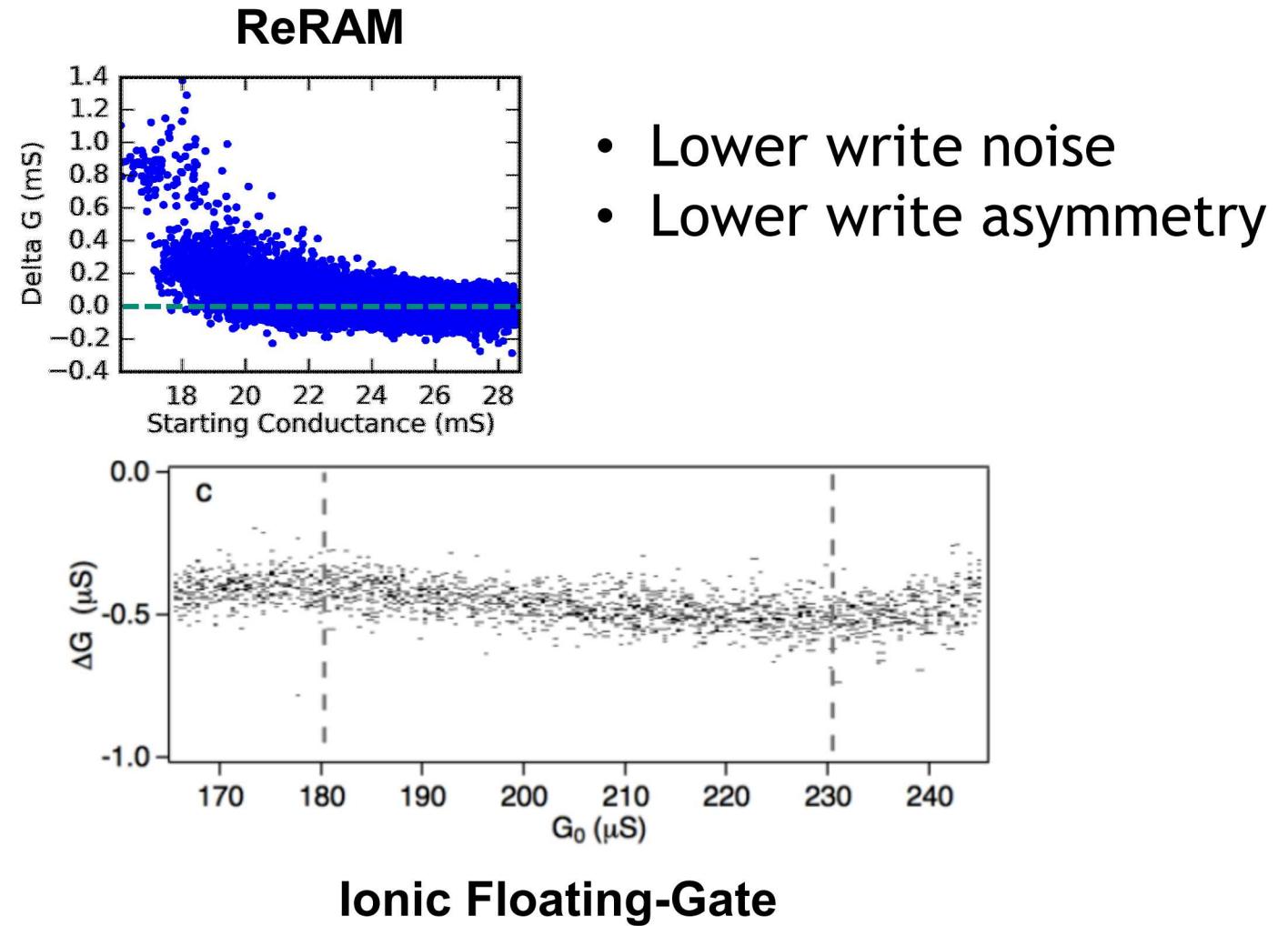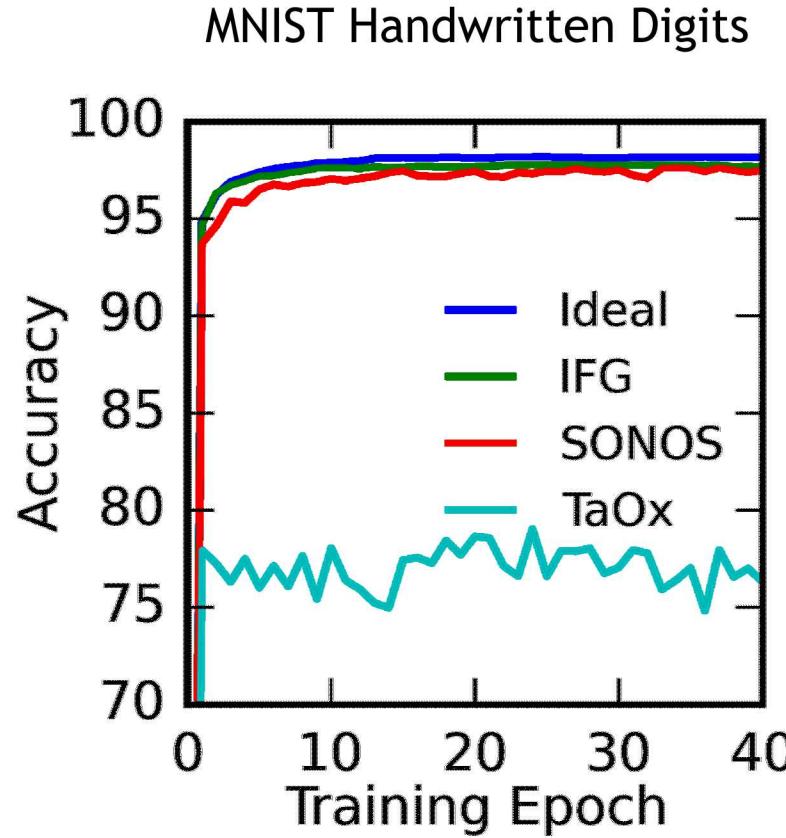
**Ionic Floating-Gate Memory**

R. B. Jacobs-Gedrim *et al.*, "Impact of Linearity and Write Noise of Analog Resistive Memory Devices in a Neural Algorithm Accelerator," IEEE International Conference on Rebooting Computing (ICRC) Washington, DC, November 2017.

S. Agarwal *et al.*, "Using Floating Gate Memory to Train Ideal Accuracy Neural Networks," *IEEE Journal of Exploratory Solid-State Computational Devices and Circuits,* 2019
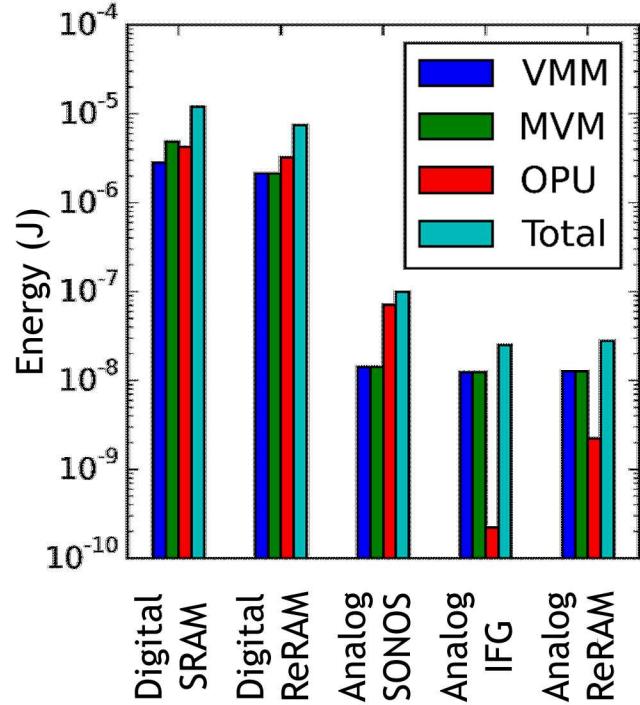
E. J. Fuller *et al.*, "Parallel programming of an ionic floating-gate memory array for scalable neuromorphic computing," *Science,* vol. 364, no. 6440, pp. 570-574, 2019.

# Three Terminal Devices Tend to Have Higher Accuracy

MNIST Handwritten Digits



**ReRAM**



- Lower write noise
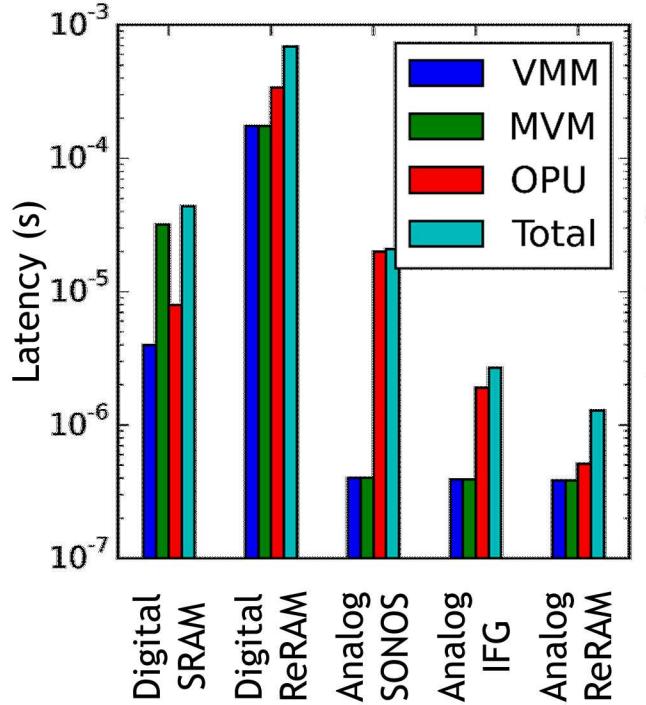- Lower write asymmetry



**Ionic Floating-Gate**

# Compare Architectural Advantages

### 120-430X Energy Advantage



### 2-34X Latency Advantage



### 5-11X Area Advantage



1024 x1024 = 1M array operations, sum over 1 training cycle, 3 operations:
- Vector Matrix Multiply    - Matrix Vector Multiply    - Outer Product Update

Used a commercial 14/16 nm PDK

***Requires 100 MΩ on state devices

# Compare Architectural Advantages: Vector Matrix Multiply

120-430X Energy Advantage   2-34X Latency Advantage   5-11X Area Advantage

All Analog Vector Matrix Multiply and Matrix Vector Multiply have same energy and latency
- Entirely dominated by ADC, device properties irrelevant

**120-430X Energy Advantage**     **2-34X Latency Advantage**     **5-11X Area Advantage**

Outer Product Update is device dependent
- SONOS has slow write (~1 ms) and high write voltage (11V)
- IFG and ReRAM write energy negligible compared to VMM
- IFG has extra delay over ReRAM for access device to turn off

37

# Compare Architectural Advantages: Area

120-430X Energy Advantage     2-34X Latency Advantage     5-11X Area Advantage



SONOS area cost reasonable, roughly doubles area

IFG and ReRAM go over transistors, area dominated by ADC and DAC

## ReRAM



- Large Energy/Area/Latency advantage over digital
- Accuracy not good enough
- Back end of line compatible
- Under commercial development

## SONOS
### Silicon-Oxygen-Nitrogen-Oxygen-Silicon



- Moderate Energy/Area/Latency advantages over digital
- High Accuracy
- Commercially available
- Need to prove endurance and device to device variability

## Ionic
### Floating-Gate Memory



- Large Energy/Area/Latency advantages over digital
- High Accuracy
- Not clear how to integrate
- Has retention challenges

# Alternate Computing Paradigms

**Crossbar Based Computing Architectures**

- Vector Matrix Multiplication
- Outer Product Update
- Crossbar Based Matrix Solvers
- Ternary Content Addressable Memory

**Neuro-inspired Computing**

- Hyperdimensional Computing
- Local Learning Rules
- Spiking Neural Networks

**Probabilistic and Stochastic Circuits**

**Computing With Dynamical Systems**

- Simulated Annealing
- Coupled Oscillator based energy minimization

# Analog Matrix Inversion

## Analog matrix inversion can perform a dense approximate matrix solve

$$b = Ax \implies I_{row} = \sum_{col} G_{row,col}(V_{col} - V_{row}) \to 0$$

$V_1 = x_1 \qquad V_2 = x_2 \qquad V_3 = x_3$

$I_1 = b_1$

$A_{11} \qquad A_{12} \qquad A_{13}$

$I_2 = b_2$

$A_{21} \qquad A_{22} \qquad A_{23}$

$I_3 = b_3$

$A_{31} \qquad A_{32} \qquad A_{33}$

### Challenges:

- Matrix inversion is non-linear, limiting how the computation can be split for large matrices
- Analog non-idealities can cause significant errors

# Ternary Content Addressable Memory (TCAM)

- Can do very efficient fast pattern matching to search stored data
  - Data analytics, k-nearest neighbors machine learning
  - Sparse matrix multiplication
  - Associative Computing
- Crossbars can implement extremely efficient TCAMs

| Input | Stored State | $M_1$ State | $M_1$ Output | $M_2$ State | $M_2$ Output | Total Output |
|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| 1 | 0 | 0 | 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 0 | 0 | 0 | 0 |
| 1 | 1 | 1 | 1 | 0 | 0 | 1 |
| 0 | x | 1 | 0 | 1 | 1 | 1 |
| 1 | x | 1 | 1 | 1 | 0 | 1 |

$M_1$

Input

$M_2$

Output

$x_1$

$x_2$

$x_3$

# Alternate Computing Paradigms

**Crossbar Based Computing Architectures**

- Vector Matrix Multiplication
- Outer Product Update
- Crossbar Based Matrix Solvers
- Ternary Content Addressable Memory

**Neuro-inspired Computing**

- Hyperdimensional Computing
- Local Learning Rules
- Spiking Neural Networks

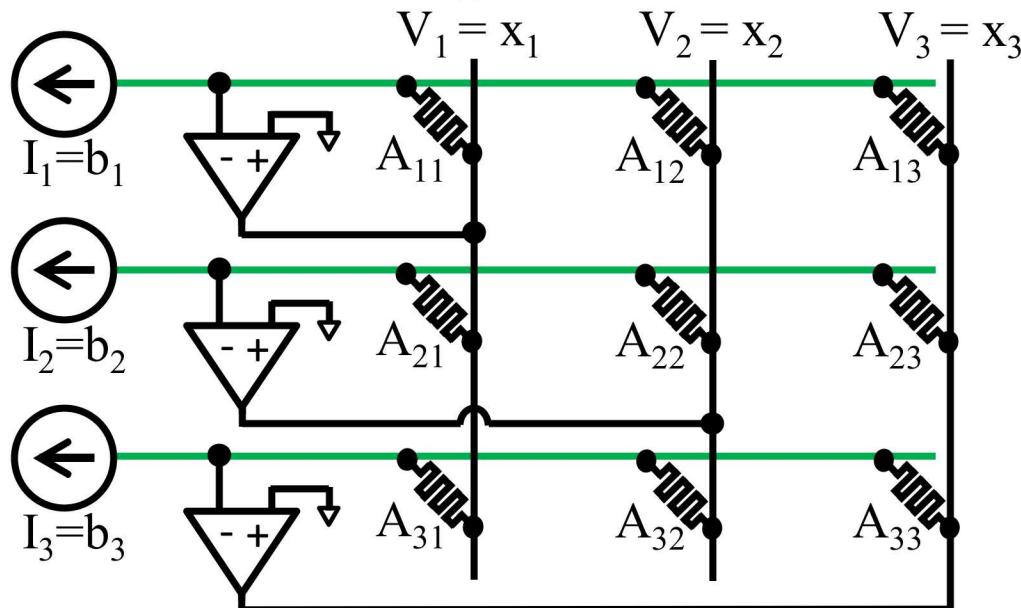**Probabilistic and Stochastic Circuits**

**Computing With Dynamical Systems**

- Simulated Annealing
- Coupled Oscillator based energy minimization

# Hyper-Dimensional Computing

Store data in redundant high dimensional vectors
`1 0 0 1 0 1 1 0 1 1 0 1 1 0 0 1 0 1 0 1 1 1 0 0 0 1 1 0 0 1 0`

As dimensionality increases, two random vectors are nearly orthogonal



Justin C. Wong. "Negative capacitance and hyperdimensional computing for unconventional low-power computing." PhD thesis, EECS Department, University of California, Berkeley, Dec 2018.

Encode data by combining vectors such that the more similar the data is, the smaller the angle between the vectors

Use ferroelectric content addressable memory to enable pattern matching



New hardware allows for processing large scale vectors and therefore new algorithms that would otherwise be computationally inefficient

Figures from Justin Wong and Sayeef Salahuddin

## Minimizing Data Movement Directly Minimizes Energy and Latency

Training neural networks requires backpropagating information across all layers resulting in long range communication and storage of intermediate values

Want learning rules that can train using only local information that is present at a given moment in time



Each neural network layer has its own local classifier

Kaiser, Jacques, Hesham Mostafa, and Emre Neftci. "Synaptic Plasticity Dynamics for Deep Continuous Local Learning (DECOLLE)." *Frontiers in Neuroscience* 14 (2020): 424.

# Spiking Neural Networks

## Minimizing Data Movement Directly Minimizes Energy and Latency

- For sparse data, communicating only non-zero values is more energy efficient than communicating all data
  - Need to account for overhead of including an address in flexible routing based networks

- Analog system energy is limited by analog to digital conversion
  - Binary outputs from an analog system are far more efficient

- Key challenge is developing high accuracy algorithms with binary inputs and outputs

Analog Neural Network Training Accelerator



Legend:
- 8 bit in/out, 8 bit weights (blue)
- 4 bit in/out, 8 bit weights (green)
- 2 bit in/out, 8 bit weights (red)

# Alternate Computing Paradigms

**Crossbar Based Computing Architectures**

- Vector Matrix Multiplication

- Outer Product Update

- Crossbar Based Matrix Solvers

- Ternary Content Addressable Memory

**Neuro-inspired Computing**

- Hyperdimensional Computing

- Local Learning Rules

- Spiking Neural Networks

**Probabilistic and Stochastic Circuits**

**Computing With Dynamical Systems**

- Simulated Annealing

- Coupled Oscillator based energy minimization

# Probabilistic and Stochastic Circuits

- Generating good random numbers is very computationally intensive

- Compact devices that provide true randomness with tunable probabilities enable new stochastic computing paradigms

**Magnetic Tunnel Junction**



- Single Electron Bipolar Avalanche Transistor
  - Avalanche breakdown is stochastic
- ReRAM
  - The intrinsic variability of memristive switching provides a source of randomness
- Contact-Resistive RAM
- CMOS -  ring oscillator jitter
- Stochastic Josephson Junctions

Borders, William A., et al. "Integer factorization using stochastic magnetic tunnel junctions." *Nature* 573.7774 (2019): 390-393

# Can map optimization problems to a set of connected stochastic bits: The Ising problem

minimize:

$$H = -\sum_{i,j} J_{ij} S_i S_j$$

coupling

spins
$\{-1, +1\}$

A. Lucas, "**Ising formulations of many NP problems**," 2014

Slide from Tianyao P. Xiao and Eli Yablonovitch

# Landscape of the Ising problem



Ising energy $H$ vs possible solutions, with label "global minimum"

Slide from Tianyao P. Xiao and Eli Yablonovitch

# Method 1: **Simulated annealing**
## (and other digital heuristics)



Slide from Tianyao P. Xiao and Eli Yablonovitch

# Method 2: **Adiabatic quantum optimization**

Prepare the ground state
of a simple problem

Transform the system into the
desired problem



If done *slowly enough*, the system is guaranteed to remain in the
ground state during the full evolution

Slide from Tianyao P. Xiao and Eli Yablonovitch

# Method 3: **First to threshold**

Map the rate of power loss
in each mode to $H$

Ising energy $H$

possible solutions = **physical modes**

# Method 3: **First to threshold**



**power loss**

Ising energy $H$

**power gain**

possible solutions = **physical modes**

No mode is stable – noise dominates in circuit

Slide from Tianyao P. Xiao and Eli Yablonovitch

# Method 3: **First to threshold**



A stable mode emerges, representing the ground state!

Slide from Tianyao P. Xiao and Eli Yablonovitch

# Analog electronic Ising machine
## (high-level view)



bistable LC oscillators

resistive coupling elements

Slide from Tianyao P. Xiao and Eli Yablonovitch

T. P. Xiao, "Optoelectronics for refrigeration and analog circuits for combinatorial optimization," Ph.D. dissertation, Chapter 6, University of California, Berkeley, 2019

# Beyond Moore Technologies

CPU — Communications Bus — Memory

Current Transistors ~ 10 aJ
50kT Noise Limit ~ 0.2 aJ

Cross chip communications ~ 1 pJ
DRAM Access >10 pJ
Ethernet ~ 1nJ

## Extending Von Neumann

- Low Voltage or Novel Transistors
- Optical Communications
- Reduced Data Movement
  - New On-Chip Memory
  - Processing in Memory

## Alternate Computing Paradigms

- Neuromorphic
- Analog
  - Computing with memory devices
- Quantum
- Stochastic
- Approximate

## Going Below 50 kT

- Error Correction
- Reversible Computing
  - Adiabatic Computing / Energy Recycling
- Superconducting

# Temporal Error Correction
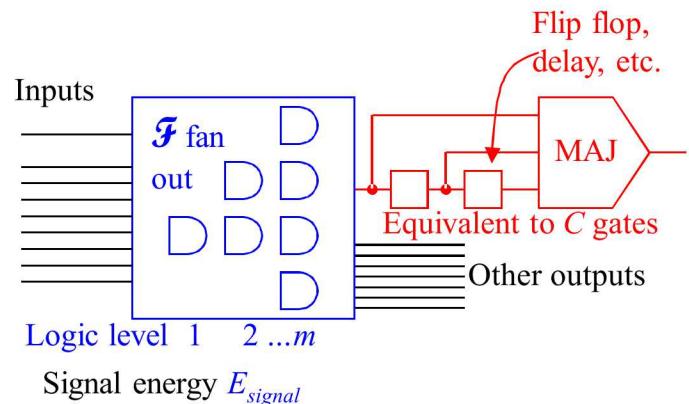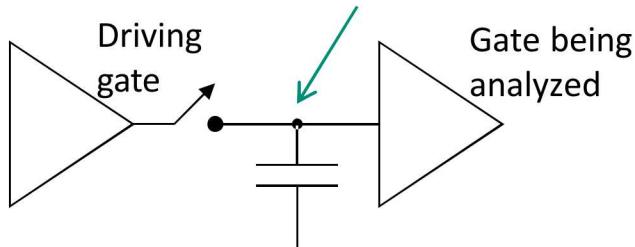
$$E_{signal} = CV^2 \longrightarrow P(Error) = e^{-E_{signal}/kT}$$



Driving gate

Gate being analyzed

Flip flop, delay, etc.

Inputs

$\mathcal{F}$ fan out

MAJ

Equivalent to $C$ gates

Other outputs

Logic level  1    2 ...*m*

Signal energy $E_{signal}$

Use multiple samples and take the majority of the result

Large Capacitance Buffer (requires extra energy to switch / guarantee stability)

Inputs

$\mathcal{F}$ fan out

Other outputs

Logic level  1    2 ...*m*

Signal energy $E_{signal}$

Use high capacitance stage / low pass filter to integrate out transient thermal errors

# How Much Energy is Saved by a Majority Gate?

**Take the majority of α samples**

### Signal Energy Needed to Maintain a Given Accuracy



Minimum Energy ~ 10 kT

### Error correction circuitry starts to dominate



Consider a 16x16 multiplier
- 48 levels of logic depth
- 32 inputs/outputs (any input can affect any output)

**Can get around a 2X reduction in energy**

**Reversible Computing**

Consider a signal Energy $E_{signal}$

The probability of an error due to thermal noise is:

$$P(Error) = e^{-E_{signal}/kT}$$

In order to ensure a full system with billons of transistors is reliable, we need:

$$E_{signal} \sim 50 \text{ kT}$$

Landauer – Shannon Limit

What is $E_{signal}$?

It could be the energy on a single irreversible gate

It could also be the energy in a reversible system that computes a complex logic function the comprises many logical functions

In both cases the signal energy is the same!

# Adiabatic Computing

To switch, need Q= C x V$_{dd}$

Conventional
- supply charge through a resistor, R, with voltage V$_{dd}$ across it.
- The time it takes is RC

$$E = \frac{V_{dd}^2}{R} \times RC = CV_{dd}^2 = V_{dd} \times Q$$

Adiabatic
- Reduce the power burned in the resistor, by minimizing the voltage across it
  - Charge the circuit with a lower current, I$_{low}$
- This takes a longer time, $\tau$, to get the required charge: $Q = C \times V_{dd} = I_{low} \times \tau$

$$E = I_{low}^2 \times R \times \tau$$
$$= I_{low} \times R \times (I_{low} \times \tau)$$
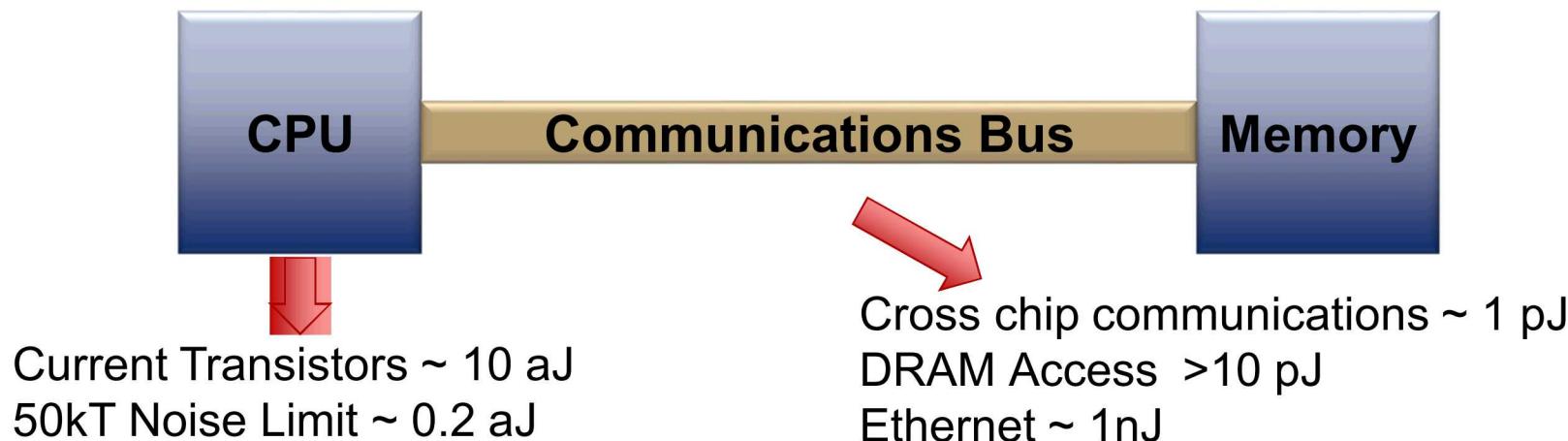$$= I_{low} \times R \times Q$$
$$= V_{low} \times Q$$

Energy Reduced by:

$$\frac{E_{adiabatic}}{E_{classical}} = \frac{V_{low}}{Vdd} = \frac{I_{low}R}{V_{dd}} = \frac{\frac{CV_{dd}}{\tau}R}{V_{dd}} = \frac{RC}{\tau}$$

Delay Increased By: $\dfrac{\tau}{RC}$

(ignoring factors of 2 for simplicity)

# Summary: There are many ways to extend Moore's Law!

**CPU** — **Communications Bus** — **Memory**

Current Transistors ~ 10 aJ
50kT Noise Limit ~ 0.2 aJ

Cross chip communications ~ 1 pJ
DRAM Access >10 pJ
Ethernet ~ 1nJ

**Extending Von Neumann**

- Low Voltage or Novel Transistors
- Optical Communications
- Reduced Data Movement
  - New On-Chip Memory
  - Processing near Memory

**Alternate Computing Paradigms**

- Neuromorphic
- Analog
  - Computing with memory devices
- Quantum
- Stochastic
- Approximate

**Going Below 50 kT**

- Error Correction
- Reversible Computing
  - Adiabatic Computing / Energy Recycling
- Superconducting

# Sandia is Hiring! sandia.gov/careers

~12,500 people in Albuquerque, NM

~1,800 people in Livermore, CA

## Fulfilling Our National Security Mission

Nuclear Deterrence

Defense Nuclear Nonproliferation

National Security Programs

Energy & Homeland Security

Advanced Science & Technology