This paper describes objective technical results and analysis. Any subjective views or opinions that might be expressed in the paper do not necessarily represent the views of the U.S. Department of Energy or the United States Government.

SAND2022-17053C

# Analog Neural Network Inference Accuracy in One-Selector One-Resistor Memory Arrays

**Joshua E. Kim[1]\*, T. Patrick Xiao[2], Christopher H. Bennett[2], Donald Wilson[1], Matthew Spear[1], Maximilian Siath[1], Ben Feinberg[2], Sapan Agarwal[2], Matthew J. Marinella[1]**
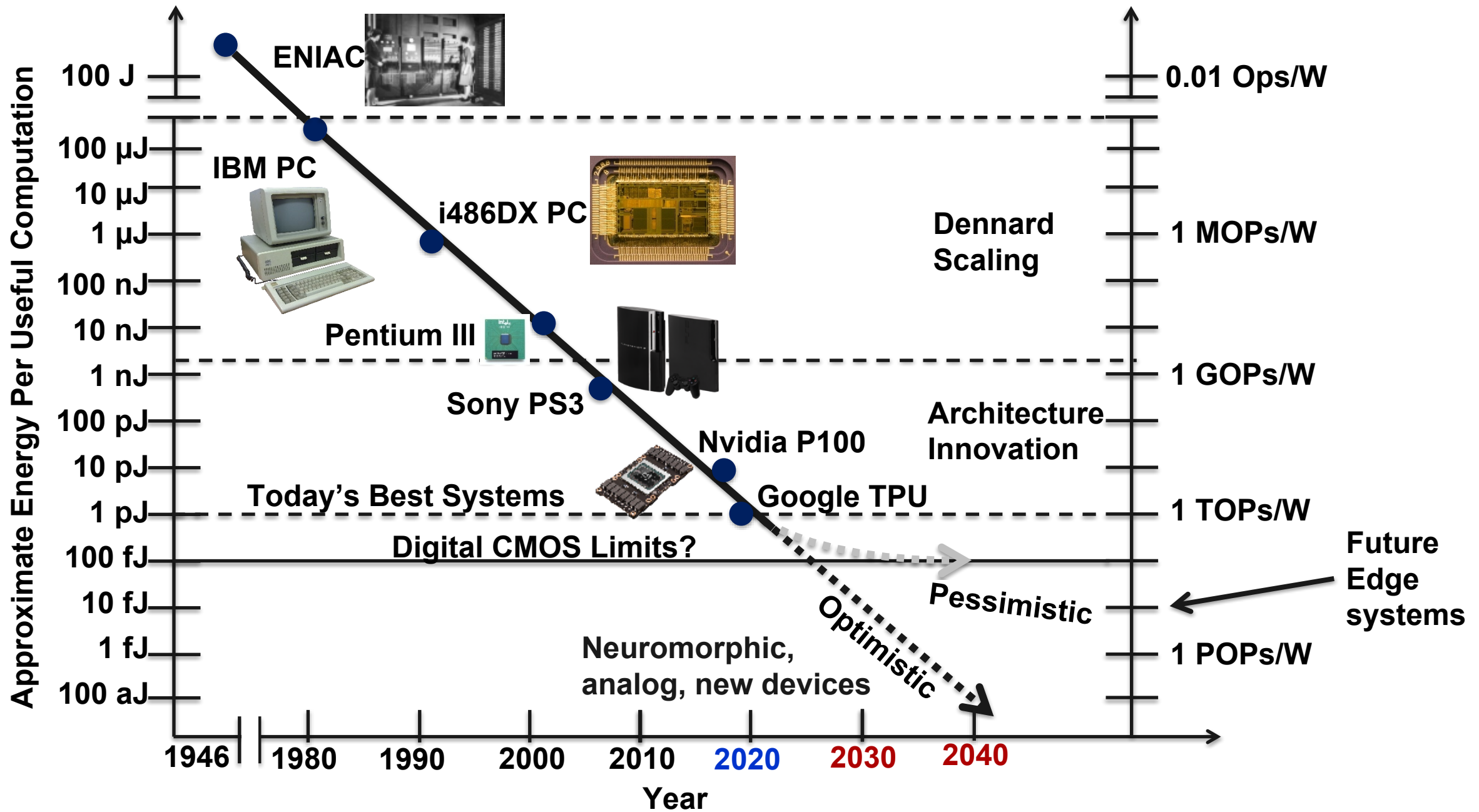
1 – School of Electrical, Computer and Energy Engineering, Arizona State University, Tempe AZ
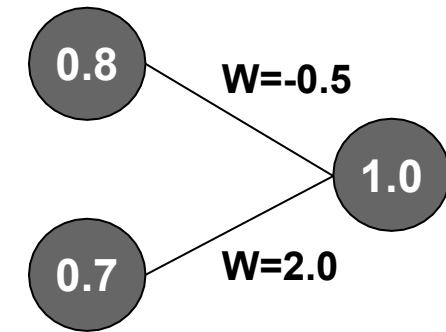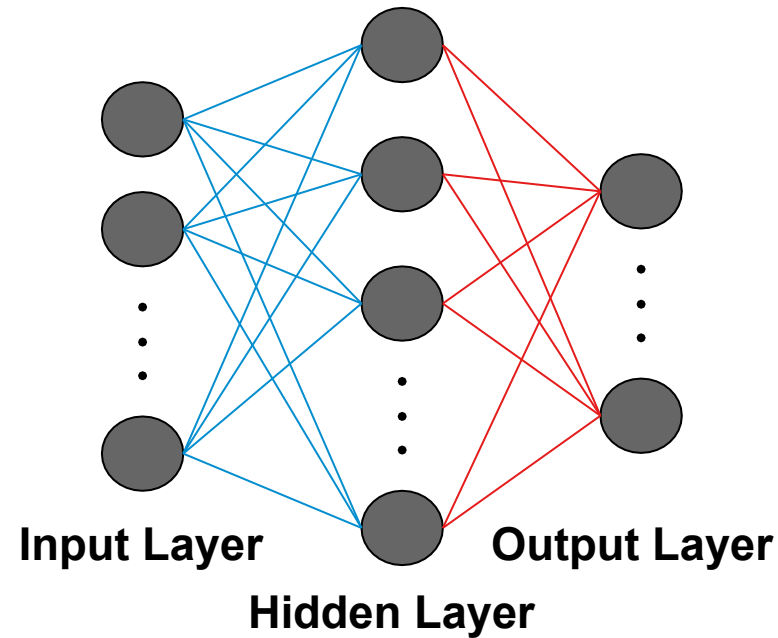2 – Sandia National Laboratories, Albuquerque, NM
jekim7@asu.edu

December 8, 2022

Y-axis (left): Approximate Energy Per Useful Computation
- 100 J
- 100 µJ
- 10 µJ
- 1 µJ
- 100 nJ
- 10 nJ
- 1 nJ
- 100 pJ
- 10 pJ
- 1 pJ
- 100 fJ
- 10 fJ
- 1 fJ
- 100 aJ

Y-axis (right):
- 0.01 Ops/W
- 1 MOPs/W
- 1 GOPs/W
- 1 TOPs/W
- 1 POPs/W

X-axis (Year): 1946, 1980, 1990, 2000, 2010, 2020, 2030, 2040

Data points and labels:
- ENIAC
- IBM PC
- i486DX PC
- Pentium III
- Sony PS3
- Nvidia P100
- Google TPU

Annotations:
- Dennard Scaling
- Architecture Innovation
- Today's Best Systems
- Digital CMOS Limits?
- Pessimistic
- Optimistic
- Neuromorphic, analog, new devices
- Future Edge systems

ASU

Adapted from: **Marinella and Agarwal, Nature Electronics 2, 437, 2019**
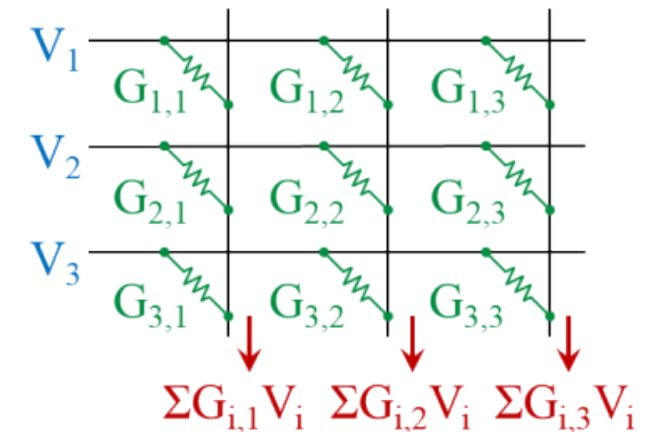
# Neural Networks and Analog Accelerators

**Neural Networks Basics:**
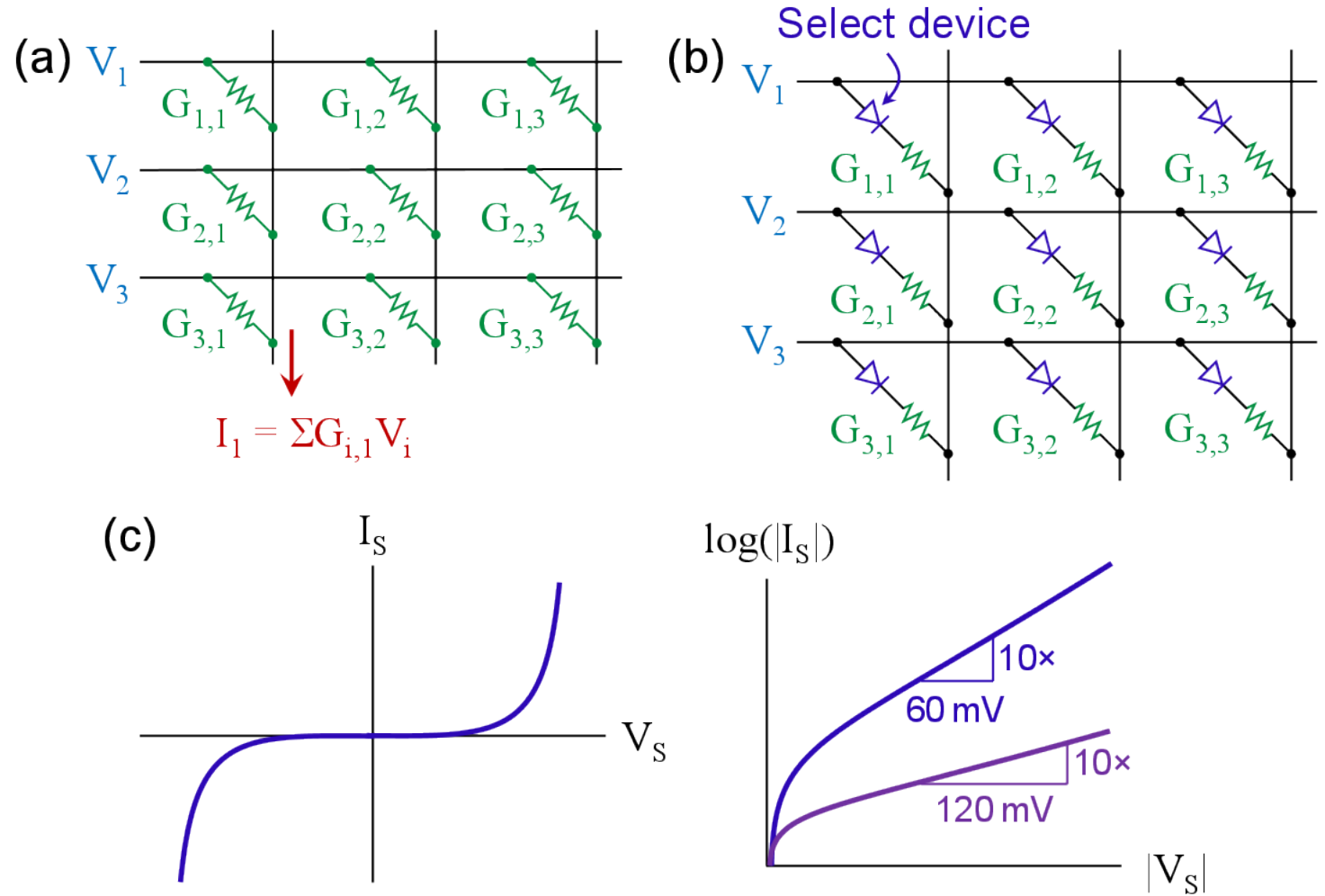


**Matrix-vector multiplication:**

$$\mathbf{Ax}$$

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}^T \begin{bmatrix} A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,1} & A_{3,2} & A_{3,3} \end{bmatrix}$$

$$= \begin{bmatrix} \Sigma A_{i,1}x_i & \Sigma A_{i,2}x_i & \Sigma A_{i,3}x_i \end{bmatrix}$$
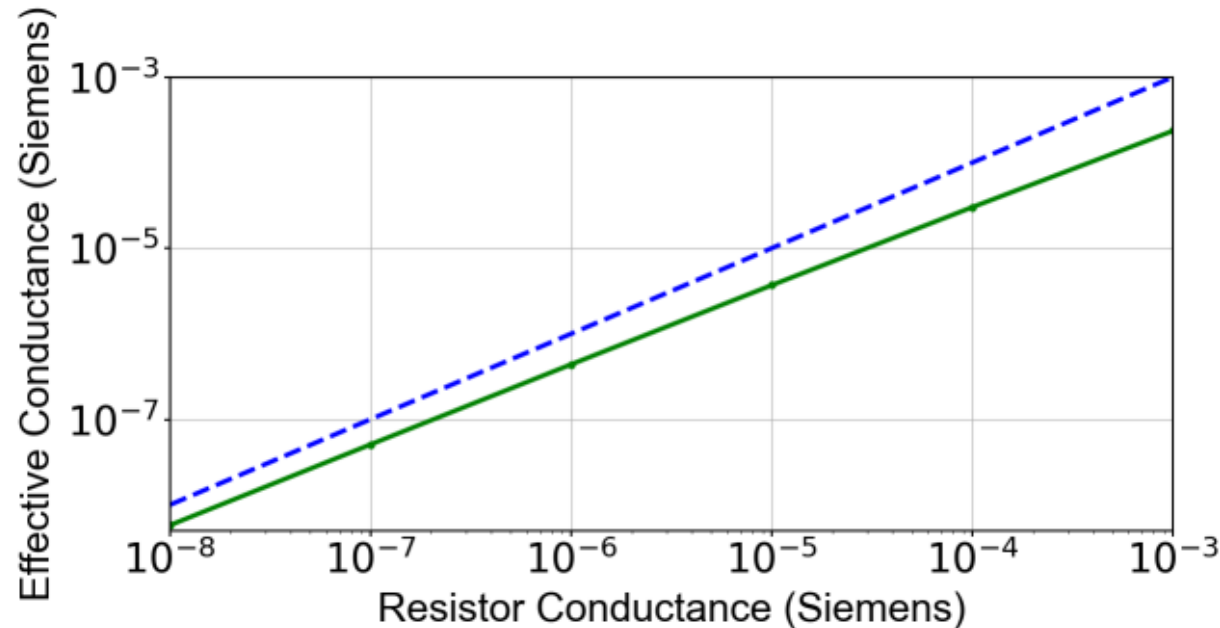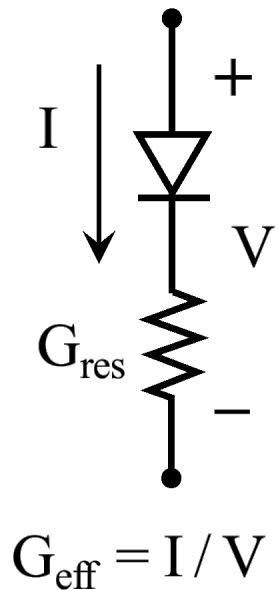


**ASU**

# Select Device Necessary for Write Operation

- Many devices "half-selected" during write operation

- Draws current ruining power efficiency

- Select devices are a solution

- 1T1R memory arrays solve this better than 1S1R arrays do

- However, 1T1R is less compact and may be incompatible with back-end-of-line integration of dense memory arrays

(a) $V_1$ $V_2$ $V_3$

$G_{1,1}$ $G_{1,2}$ $G_{1,3}$

$G_{2,1}$ $G_{2,2}$ $G_{2,3}$

$G_{3,1}$ $G_{3,2}$ $G_{3,3}$

$I_1 = \Sigma G_{i,1} V_i$

(b) Select device

$V_1$ $V_2$ $V_3$

$G_{1,1}$ $G_{1,2}$ $G_{1,3}$

$G_{2,1}$ $G_{2,2}$ $G_{2,3}$

$G_{3,1}$ $G_{3,2}$ $G_{3,3}$

(c) $I_S$ $V_S$

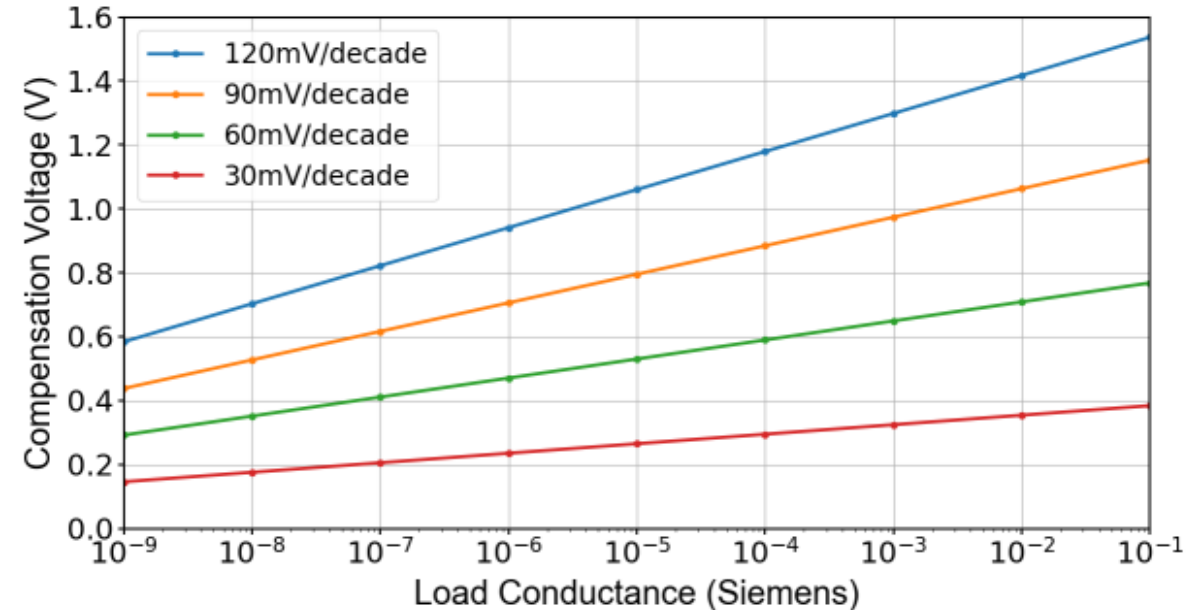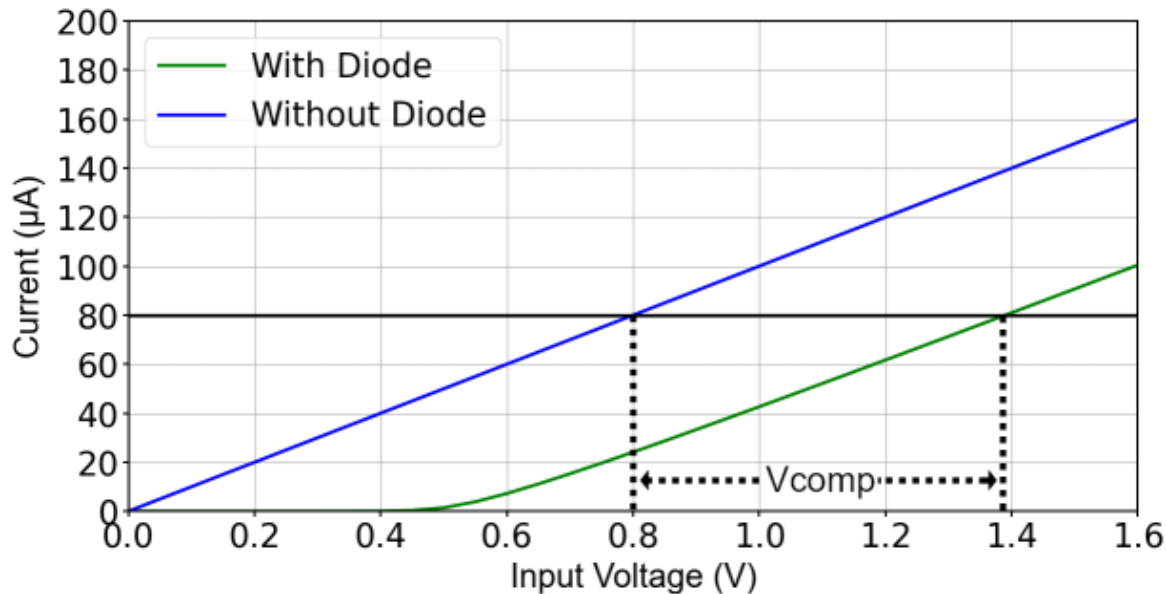$\log(|I_S|)$ $|V_S|$

$10\times$
60 mV

$10\times$
120 mV

# Effective Conductance of 1S1R

- Some voltage is dropped across the select device
- To achieve correct output current, a higher voltage must be applied
- This additional voltage changes as a function of the conductance



The blue dashed line represents y=x or the conductance of the resistor by itself.

# Compensation Voltage



- Compensation voltage can be found perfectly for a single cell
- Function of select device steepness, resistor conductance, and "nominal" voltage
- Nominal voltage is the voltage across the resistor necessary to output the correct current
- Not practical to implement individual compensation for each cell
- Pick one compensation voltage for entire array with goal of minimizing error
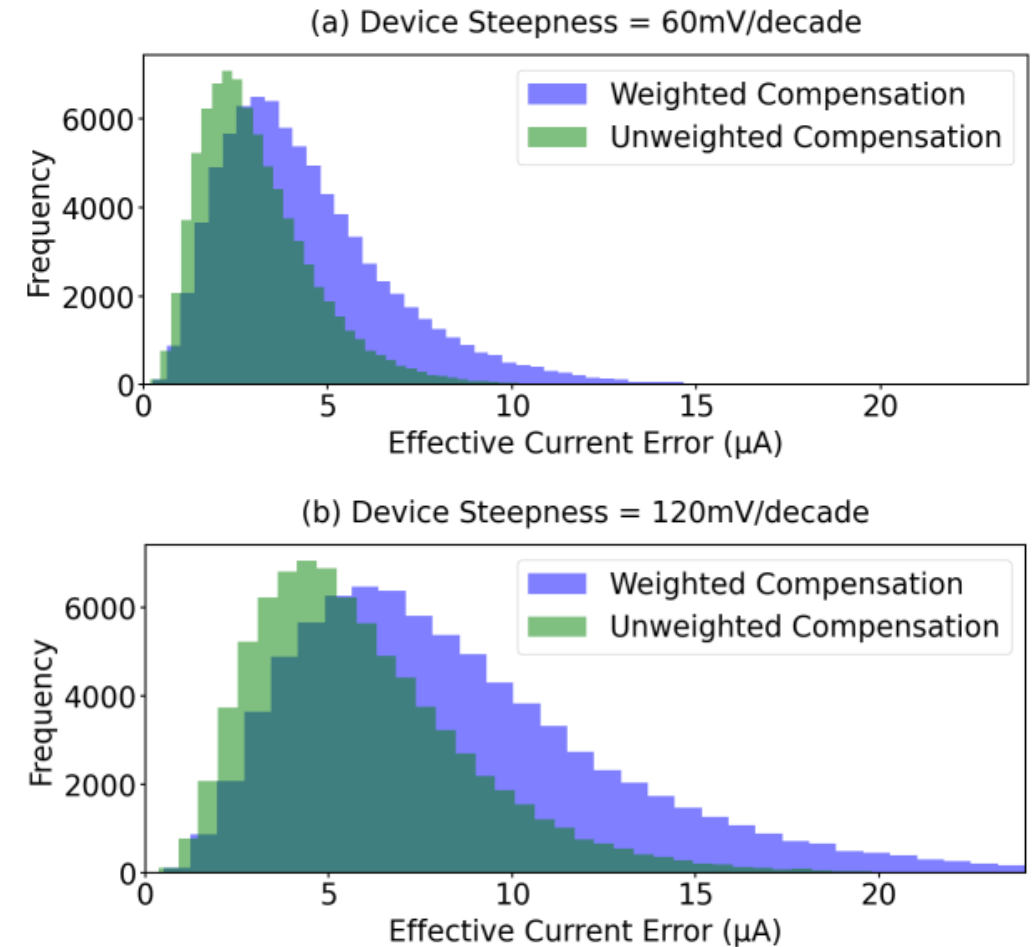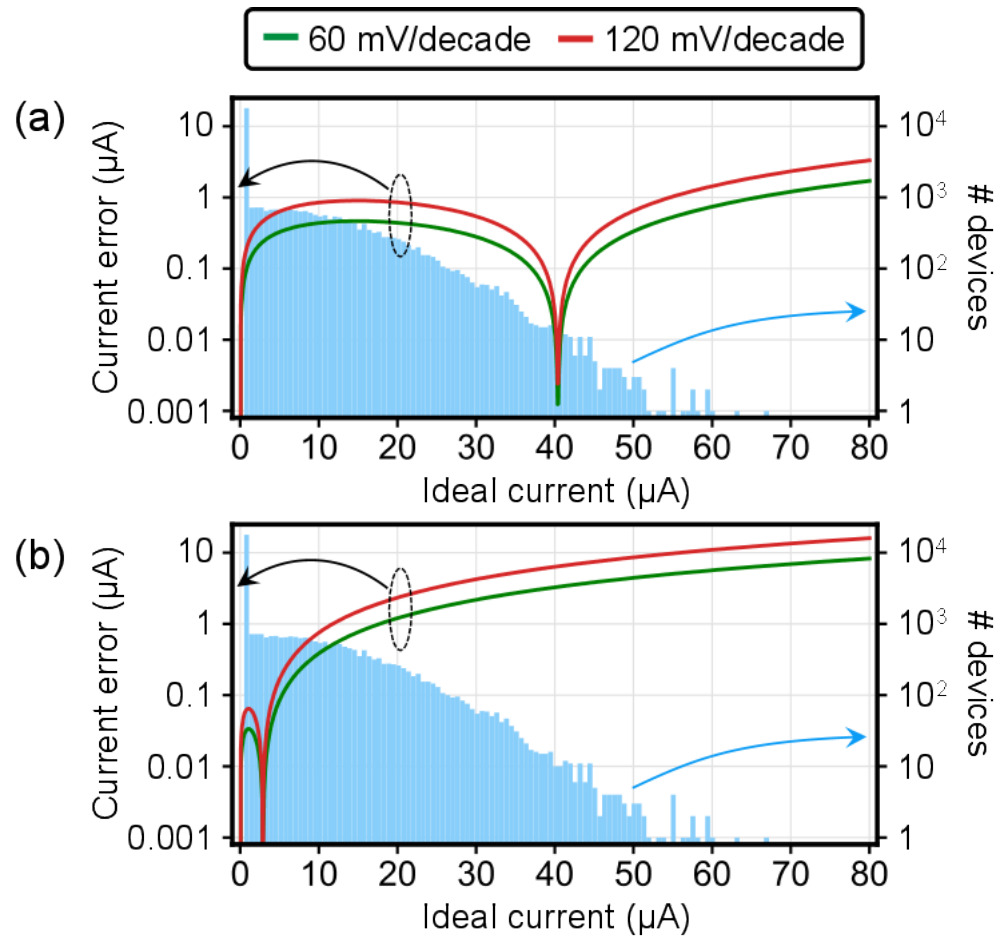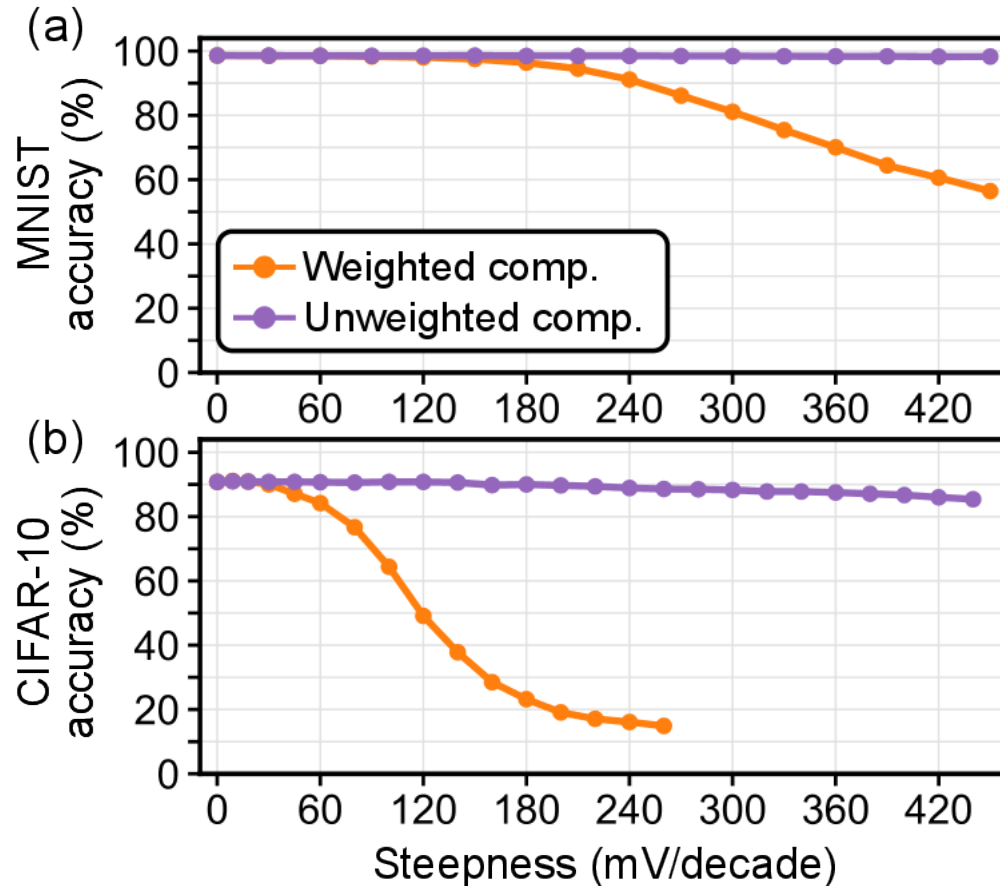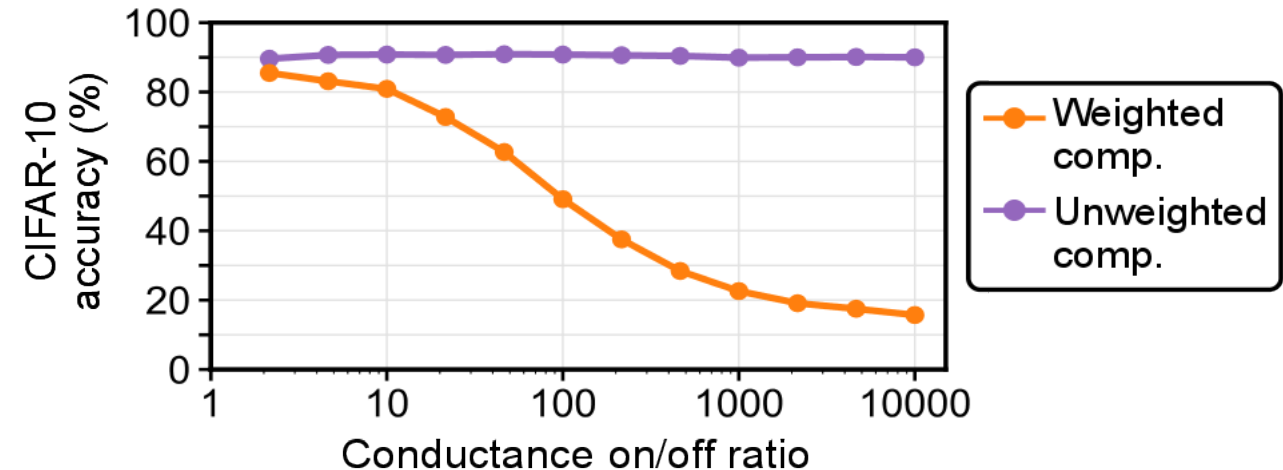
# Weighted vs Unweighted Compensation



Fig. 5.  (a) unweighted compensation  (b) weighted compensation

(a)

(b)

- Unweighted compensation accuracy on both MNIST and CIFAR-10 is minimally affected by device steepness and conductance on/off ratio
- Assuming a realistic 60mV/dec, unweighted compensation achieves an accuracy of 90.29%
- This is only 0.44% below ideal floating-point results

# Conclusions

- Non-volatile memory arrays with a 1S1R topology can be made compatible with accurate neural network inference if the errors induced by the select device are appropriately compensated
- Showed that a single compensation voltage, applied uniformly across the entire system, can effectively reduce these errors to enable accurate inference
- With this compensation, a CIFAR-10 accuracy that is within 0.44% of the floating-point digital result can be achieved using a realistic selector with 60 mV/decade steepness
- The accuracy is insensitive to the memory device On/Off ratio
- These results are promising for the use of dense 1S1R arrays for analog neural network inference
- Future work should investigate how selector-induced errors interact with other sources of analog error, and how these different errors can be mitigated together
      - These include parasitic resistance, process variations, other non-idealities, etc.

# Acknowledgements

- This work was funded by:

- Sandia National Laboratories Laboratory Directed Research and Development (LDRD)