

# Fundamental Limits on the Computational Accuracy of Resistive Crossbar-based In-memory Architectures

Saion K. Roy<sup>1</sup>, Ameya Patil<sup>2</sup>, and Naresh R. Shanbhag<sup>1</sup>

1: University of Illinois at Urbana-Champaign

2: Amazon Lab126

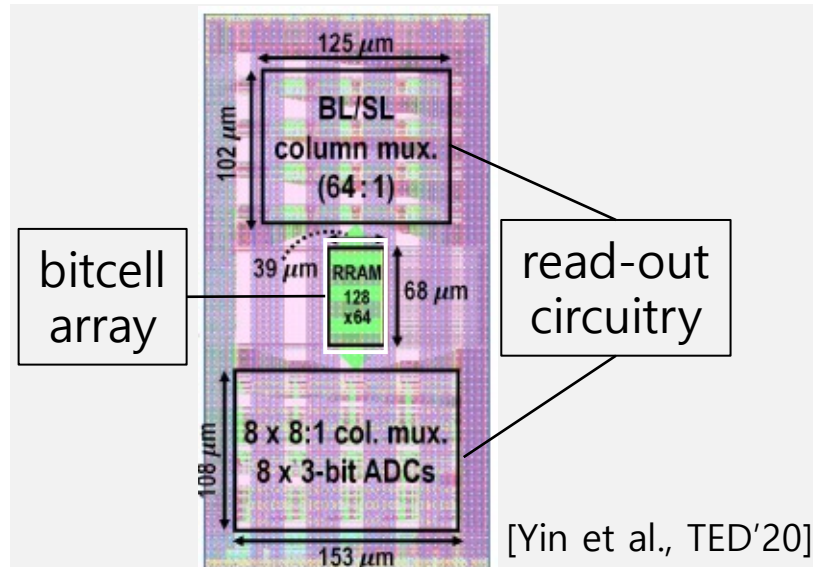
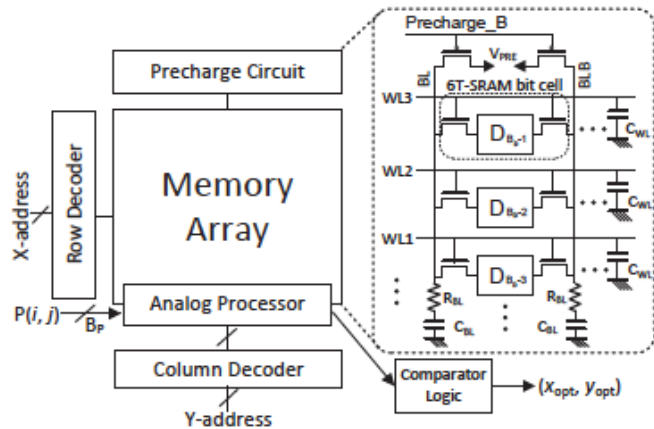
**2022 IEEE International Symposium on Circuits and Systems**  
**May 28- June 1, 2022 Hybrid Conference**

# Outline

- Introduction
- Resistive Crossbar Architecture
- Behavioral Modeling
- Simulation Results
  - Model validation
  - Compute SNR analysis for MRAM, ReRAM, and FeFET crossbars
  - System level accuracy of ResNet-20 on CIFAR-10
- Conclusion

# In-memory Computing (IMC)

## compute memory



## first IMC concept paper (ICASSP 2014)

### AN ENERGY-EFFICIENT VLSI ARCHITECTURE FOR PATTERN RECOGNITION VIA DEEP EMBEDDING OF COMPUTATION IN SRAM

Mingu Kang\*, Min-Sun Keel\*, Naresh R. Shanbhag\*, Sean Eilert†, and Ken Curewitz†

\*Dept. Electrical and Computer Engineering, University of Illinois at Urbana-Champaign  
†Micron Technology, Inc

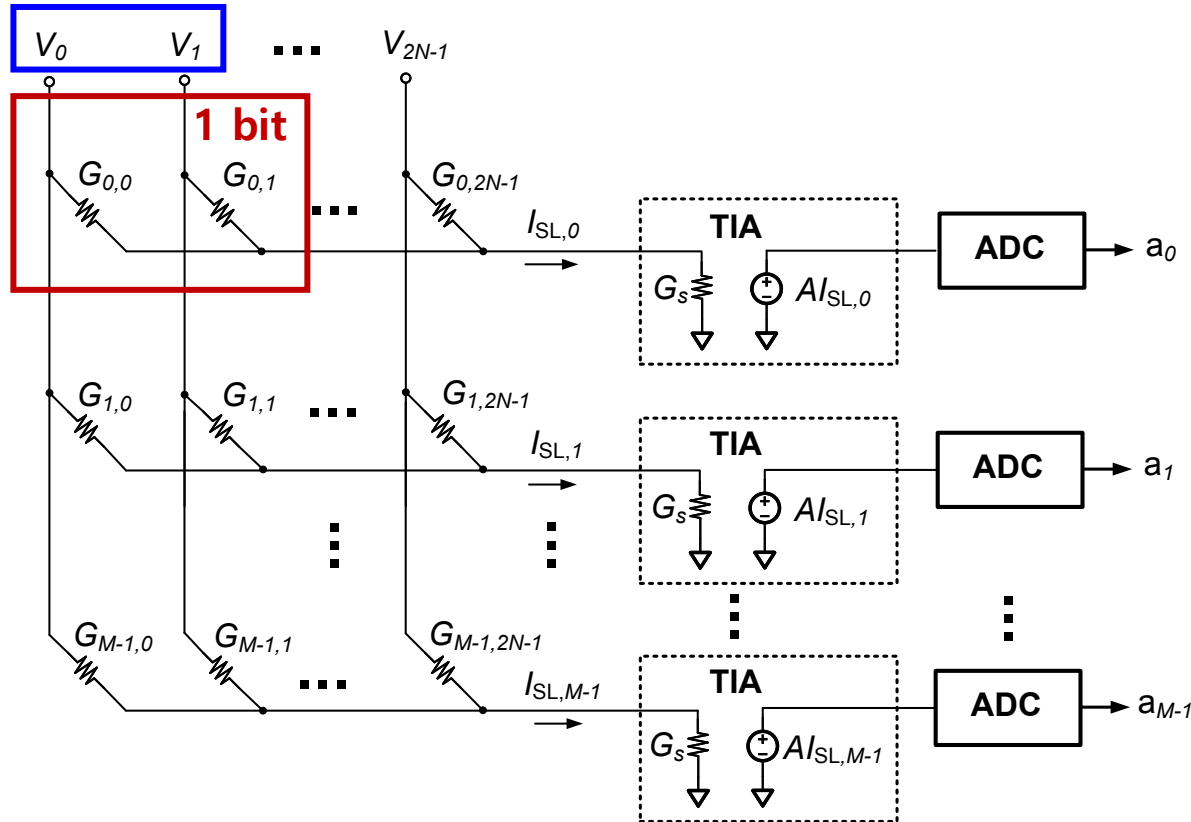
- computes a  $M \times N$  matrix-vector multiply (MVM)
- SRAM-based IMC banks are mature  $\rightarrow$   $20\times$  lower energy +  $9\times$  higher compute density than digital<sup>1</sup>
- eNVM-based (MRAM/ReRAM) IMCs have potential for high compute density but lags digital due to low compute SNR  $\rightarrow$  **this work explains why**

<sup>1</sup>N. R. Shanbhag and S. K. Roy, "Comprehending In-memory Computing Trends via Proper Benchmarking," CICC 2022 (invited)

# Resistive Crossbar Architecture

## voltage-drive current-sensing crossbar

differential inputs



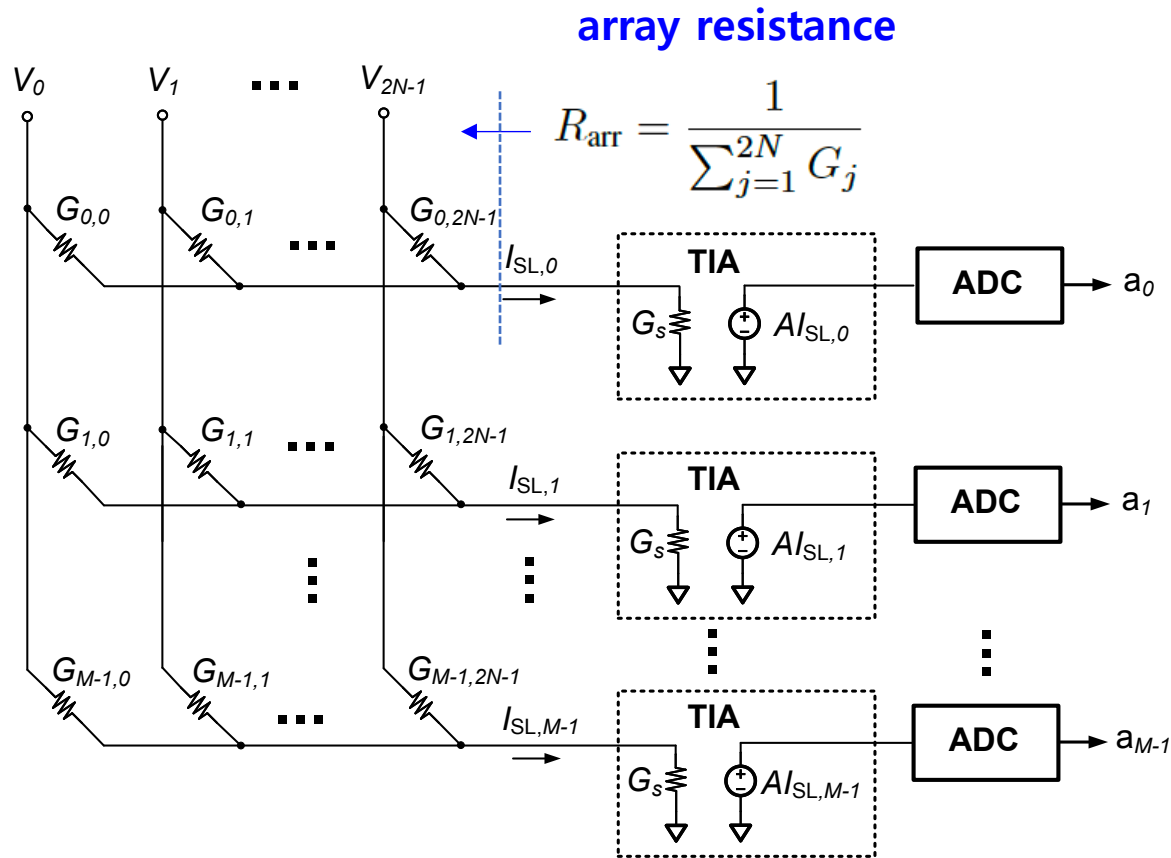
[A. Patil et al., ISCAS 2019]

- computes a  $M \times N$  matrix-vector multiply (MVM)
- V-DACs provide differential inputs on BLs ( $V_{2k} = -V_{2k-1}$ )
- two BCs ( $G_{2k-1}, G_{2k}$ ) store 1 bit
- current summing & sensing on SLs
- device **resistive contrast**

$$\rho = \frac{R_{\text{off}}}{R_{\text{on}}}$$

2 (MRAM); 12 (ReRAM);  $10^3$  (FeFET);

# Behavioral Modeling



- signal current in SL

$$I_{sig} = \underbrace{\left[ \frac{R_{arr}}{R_{arr} + R_s} \right]}_{S_I \text{ (current scaling factor)}} \underbrace{\left( \sum_{k=1}^N V_{2k} \Delta G_{2k} \right)}_{I_{ideal}} = S_I I_{ideal}$$

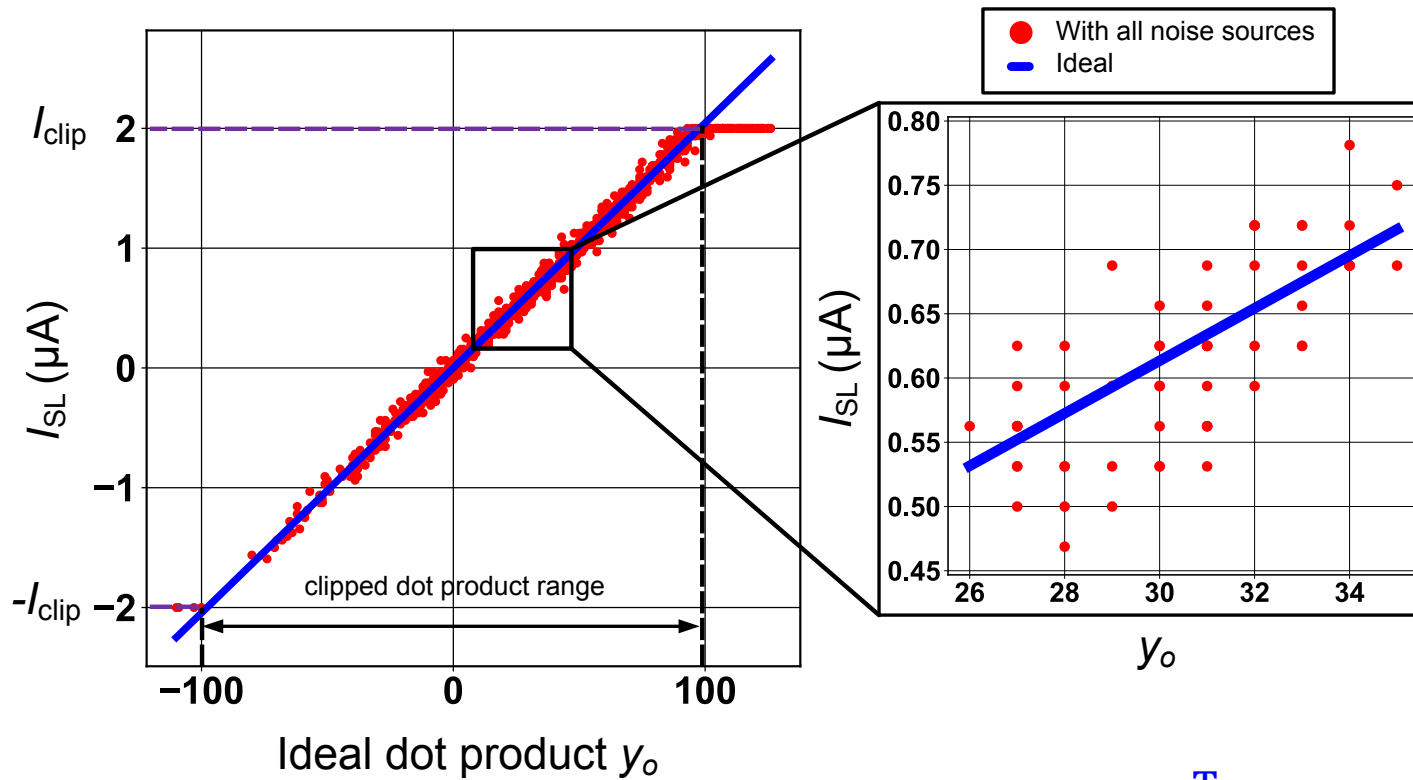
- total current in SL

$$I_{SL} = I_{sig} + \underbrace{I_{nb}}_{\text{conductance variation}} + \underbrace{I_{nd}}_{\text{DAC mismatch}} + \underbrace{I_{nc}}_{\text{clipping noise}} + \underbrace{I_{nq}}_{\text{quantization noise}}$$

# Results – Model Validation

## SPICE simulations in a 22nm process

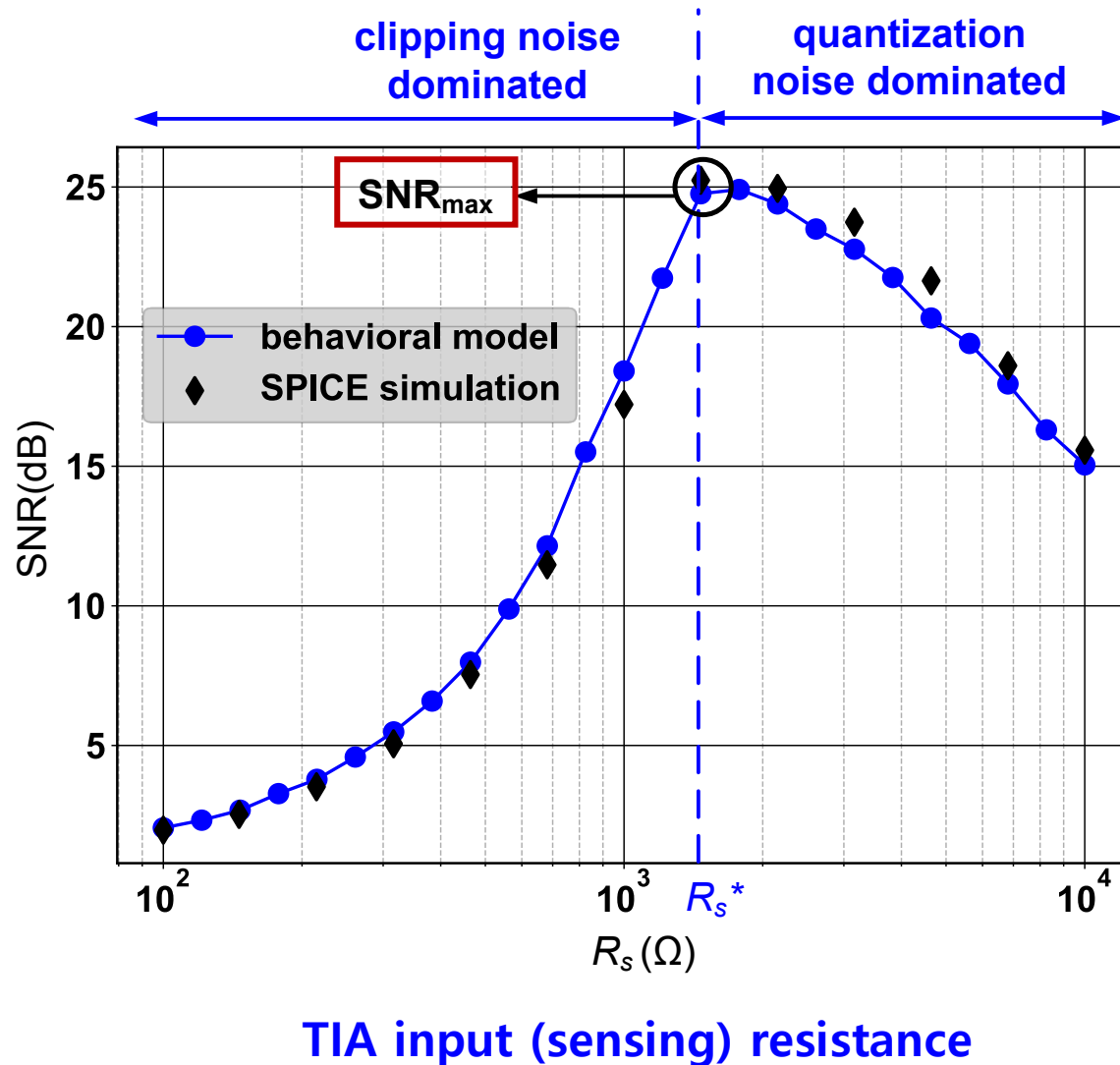
$N = 512, R_s = 316\Omega, 6b$  ADC



$$y_o = \mathbf{w}^T \mathbf{x}$$

- DAC input: signed 5b with  $V_{Isb} = 3mV$
- DAC mismatch: 4%
- Conductance variation: 4%
- ADC clipping range:  $[-2\mu A, +2\mu A]$
- SL current varies due to analog non-idealities

# Results – Compute SNR Analysis



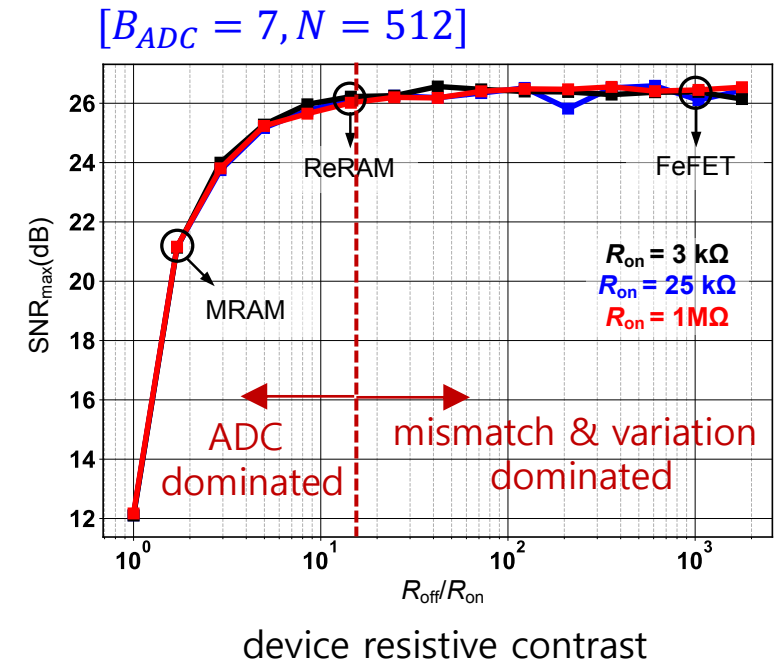
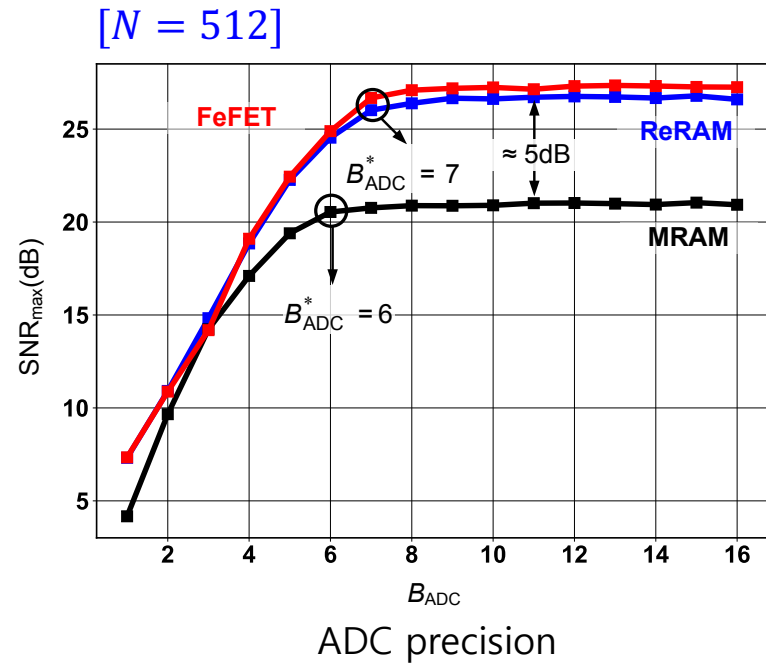
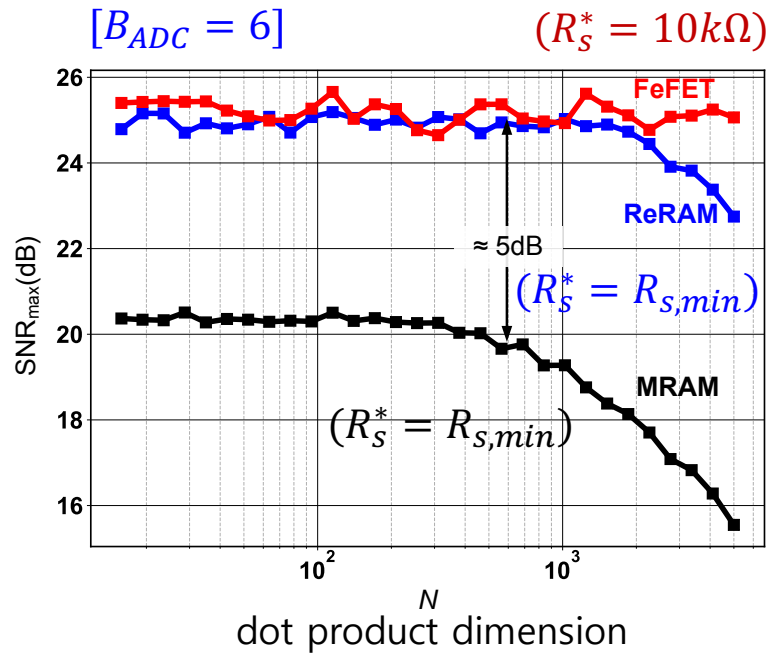
$$\text{SNR} = \frac{\mathbb{E}[I_{\text{sig}}^2]}{\mathbb{E}[I_{\text{nb}}^2] + \mathbb{E}[I_{\text{nd}}^2] + \mathbb{E}[I_{\text{nc}}^2] + \mathbb{E}[I_{\text{nq}}^2]}$$

- model and simulations match → further validates model
- ADC clipping vs. quantization noise trade-off:

$$S_I = \left[ \frac{R_{\text{arr}}}{R_{\text{arr}} + R_s} \right]$$

- compute SNR maximized if  $R_s = R_s^*$   
→ clipping noise & quantization noise are equal

# Results – SNR Dependence



- $\text{SNR}_{\max}$  roll-off with higher DP dimension  $\rightarrow$  small  $R_{\text{arr}}$   $\rightarrow$  small  $R_s = R_{s,\min} (= 1\text{k}\Omega) \neq R_s^*$
- higher absolute  $R_{\text{on}}, R_{\text{off}}$  critical for high DP dimension

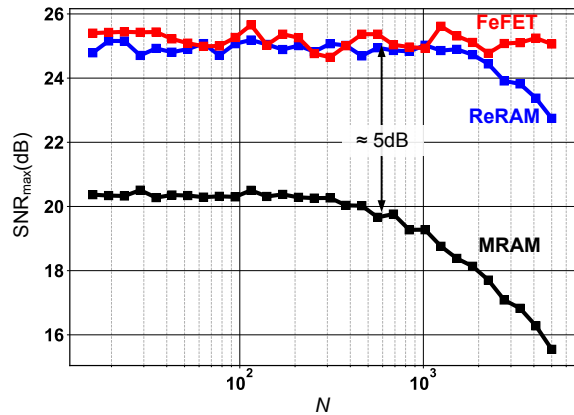
- minimum ADC precision  $B_{\text{ADC}}^*$  increases with  $\text{SNR}_{\max}$
- $\text{SNR}_{\max}$  saturates for  $B_{\text{ADC}} > B_{\text{ADC}}^* \rightarrow$  DAC mismatch and  $G$  variations dominate

- $\text{SNR}_{\max}$  improves with higher resistive contrast but..
- increasing device resistive contrast beyond ( $\sim 12$ -to- $15$ ) is futile



# System Level Accuracy Prediction Set-up

## Proposed SNR analysis

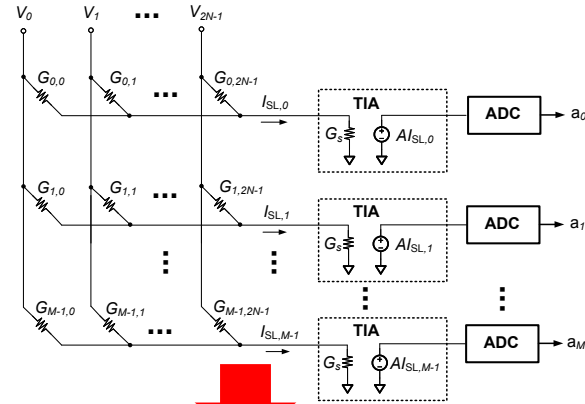


Crossbar Design Parameters

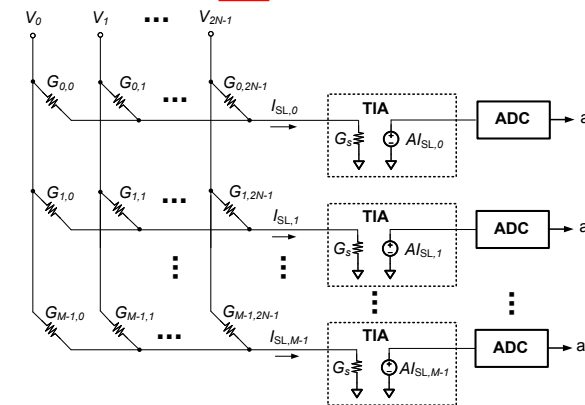
$$R_S, B_{ADC}, N, M, V_{Isb}$$

Exhaustive search

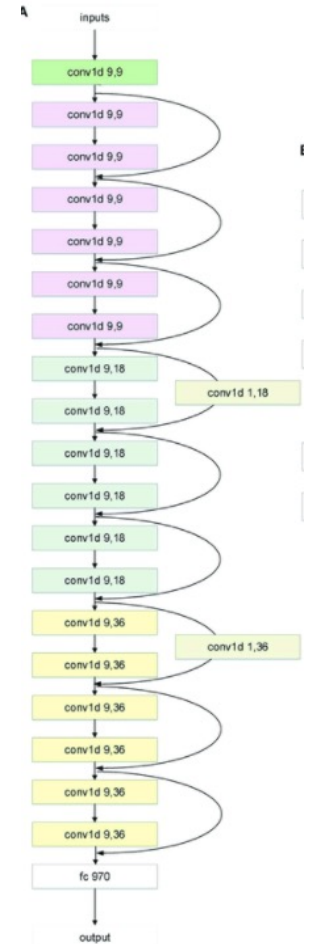
## SNR maximizing Crossbar Designs



network accuracy comparison

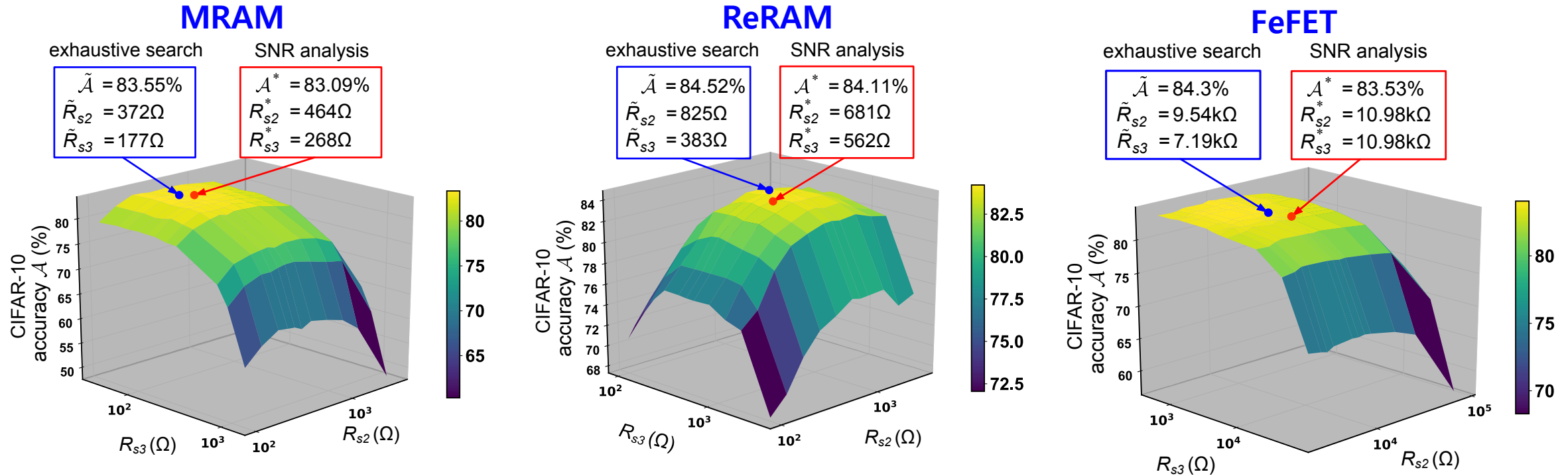


## ResNet-20 on CIFAR-10



DNN mapping

# Results – System Level Accuracy Prediction



- Baseline ResNET-20 on CIFAR-10: 5b input, ternary weights, accuracy = 84.94% (3-layer network)
- 3 Crossbars:  $N = 144, 288, 576$ ;  $R_{s1} = R_{s1}^*$ ; sweep  $(R_{s2}, R_{s3})$
- SNR analysis predicted crossbar design achieves system-level accuracy to within 1% (exhaustive search) to within 2% of digital baseline value (84.94%)
- bank-level SNR is a good proxy for network level accuracy → SNR analysis bypasses trial & error

# Conclusion

- proposed an analytical framework to obtain SNR-optimal resistive crossbar parameters → avoids expensive trial and error
- insights provided by the framework:
  - SNR-optimal sensing resistance  $R_S^*$  exists which equalizes the clipping and quantization noise in the column ADCs
  - system level inference accuracy is maximized when bank-level compute SNR is maximized
  - increasing device resistive contrast improves SNR up to a point (~12-15). Diminishing returns due to mismatch (input DACs) and variations (device conductance)
- proposed framework can be extended to other resistive IMC and devices

# Thank You

[saionkr2@illinois.edu](mailto:saionkr2@illinois.edu)

## **Acknowledgement**

Work supported by the Defense Advanced Research Agency (DARPA) and the Semiconductor Research Corporation (SRC) via the FRANC Program, and the Center for Brain-inspired Computing (C-BRIC).