

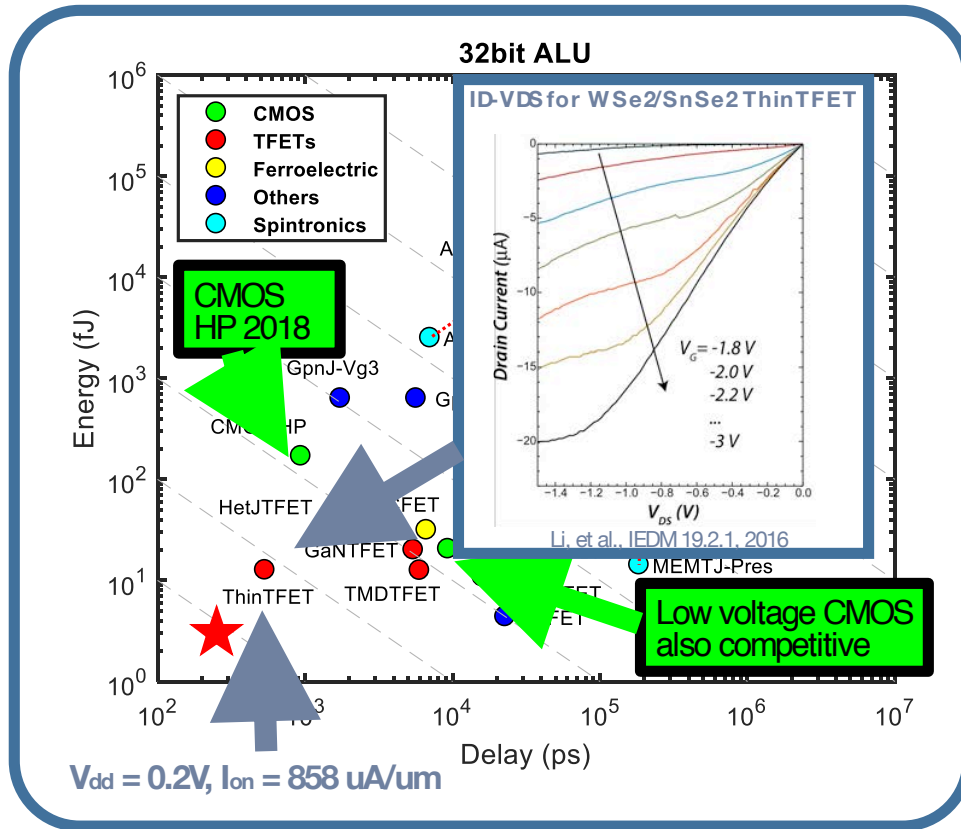
# Beyond Logic Applications for Ferroelectric Field Effect Transistors

Michael Niemier  
University of Notre Dame

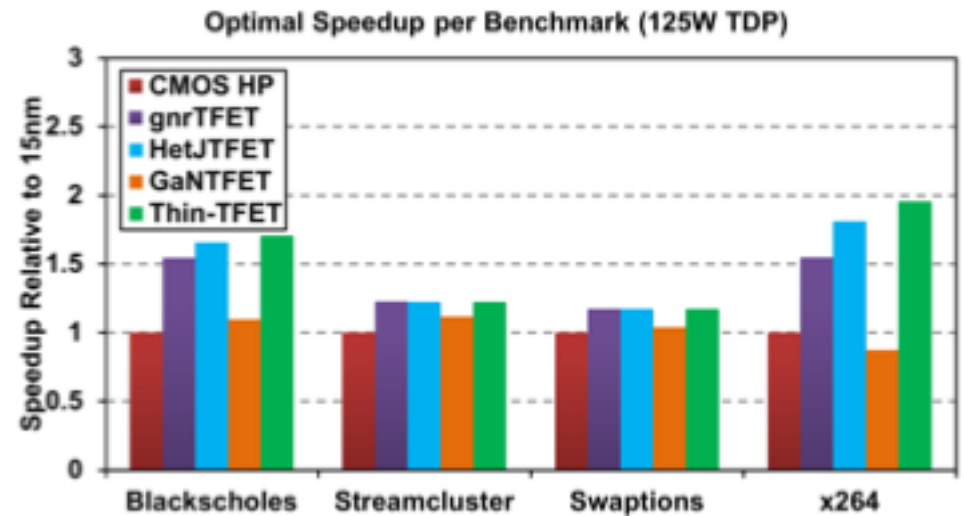


This work was supported in part by ASCENT, one of six centers in JUMP, a Semiconductor Research Corporation (SRC) program sponsored by DARPA

# How does technology scaling impact m/c scaling?



Greater speedups for highly parallelizable benchmarks...



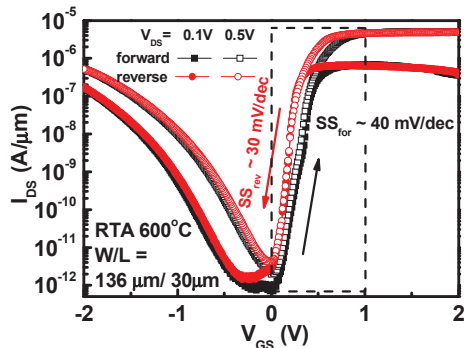
Slightly better improvements at low TDP, but still only 2X to 3X...

*"NRI research has explored a broad spectrum of beyond-CMOS devices for a 'new logic switch' to replace the current CMOS-based transistor ... a '**better switch**' has not been found. Comprehensive benchmarking of beyond-CMOS devices ... has revealed little or no advantage of these devices over CMOS for conventional Boolean logic and the von Neumann architecture."*

*"some devices demonstrate unique characteristics suitable for novel architectures or computing paradigms, e.g., non-volatility in logic devices, reconfigurability, [and/or] high computation density."*

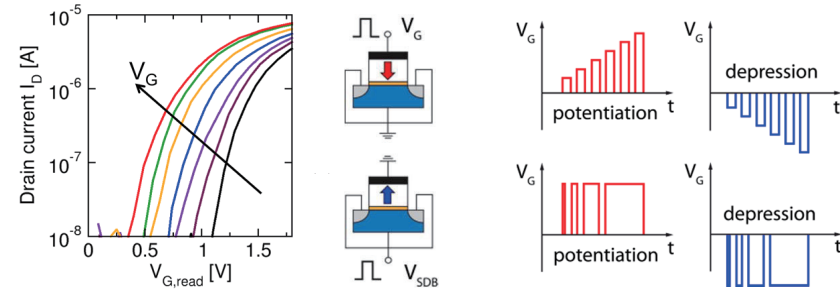
# What do ferroelectric devices offer?

## Steep subthreshold swings



Lee, et al., Prospects for Ferroelectric HfZrOx FETs with Experimentally CET=0.98nm, SS<sub>for</sub>=42mV/dec, SS<sub>rev</sub>=28mV/dec, Switch-OFF <0.2V, and Hysteresis-Free Strategies, IEDM 2015.

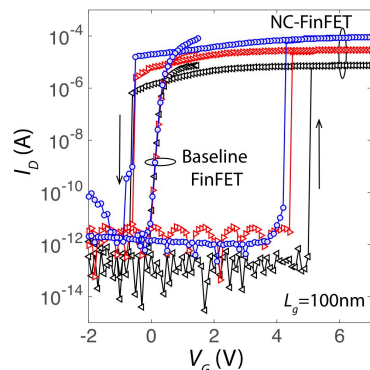
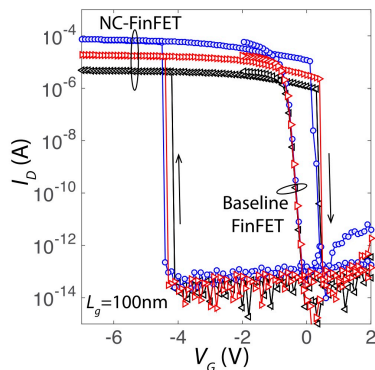
## Analog synaptic behavior



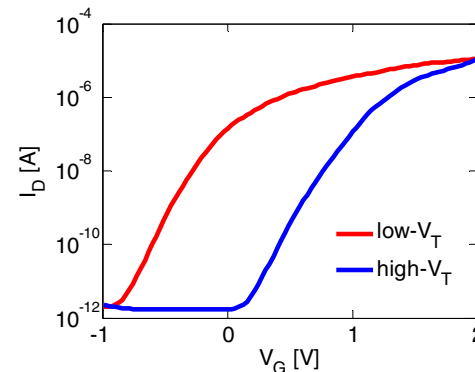
H. Mulaosmanovic, Novel ferroelectric FET based synapse for neuromorphic systems, VLSI Symposium, 2016.

## Memory functionality

○  $|V_D|=950$  mV    $\rightarrow$   $|V_D|=200$  mV    $\leftarrow$   $|V_D|=50$  mV



Kahn, et al, Negative Capacitance in Short-Channel FinFETs  
Externally Connected to an Epitaxial Ferroelectric Capacitor, IEEE ELECTRON DEVICE LETTERS, VOL. 37, NO. 1, JANUARY 2016 111



M. Trentzsch, et al, "A 28nm HKMG super low power embedded NVM technology based on ferroelectric FETs," IEDM, 11.5.1-.2, 2016

*Devices with integrated ferroelectrics are well-positioned to address aforementioned space!*



# Talk outline

- FeFET device, models
- FeFETs for **logic-in-memory (LIM)**, **compute-in-memory (CIM)**
  - Emphasis on design/benchmarking of content addressable memories (LIM)
  - Briefly discuss FeFET-based CIM
- FeFETs for **neuromorphic applications**
  - FeFET-based analog synapse
  - FeFET-based (binary) convolutional neural networks (CNNs)
- Wrap-up

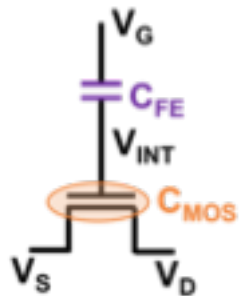
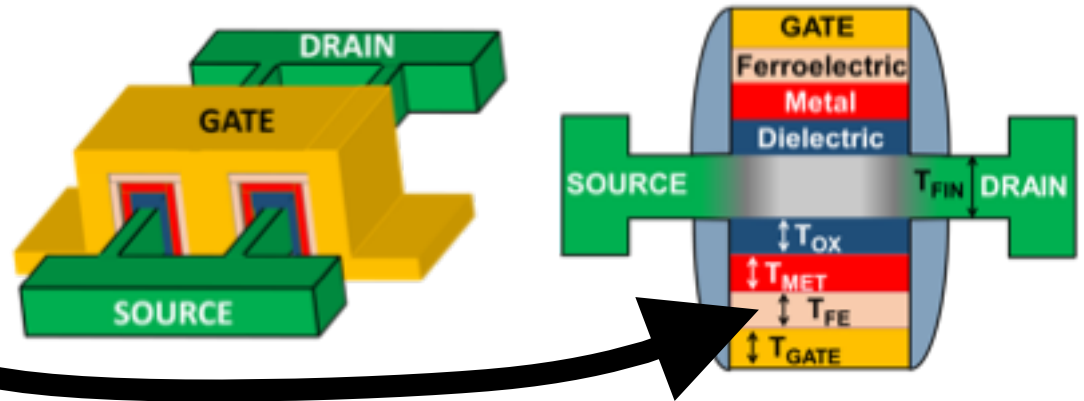
# Background

# FeFET device structure & operating modes

FeFET transistor structurally similar to bulk MOSFET or FinFET

- Ferroelectric (FE) layer integrated into gate stack

Ferroelectricity demoed in hafnium zirconium dioxide (highly compatible with CMOS)



Interplay between FE material + underlying transistor capacitance results in different modes of operation:

- **Non-volatile mode** (device can maintain state)
- **Steep switching mode** (aimed at high performance)

# Time-dependent Landau Khalatnikov (LK) model

- LK model is SPICE compatible

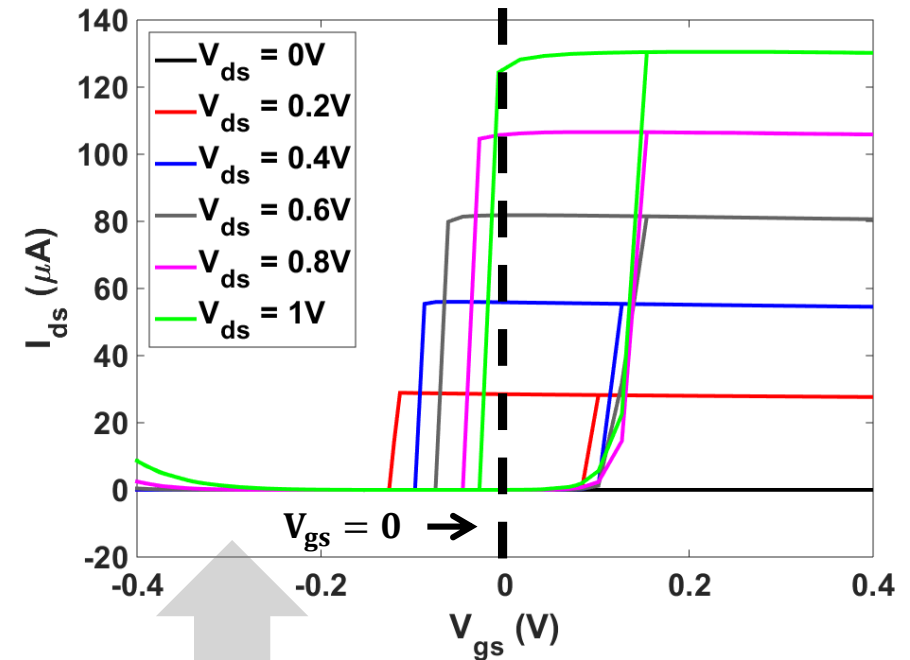
$$E = \alpha P + \beta P^3 + \gamma P^5 + \rho \frac{dP}{dt}$$

static coefficients      kinetic coefficient

Electric field      Polarization

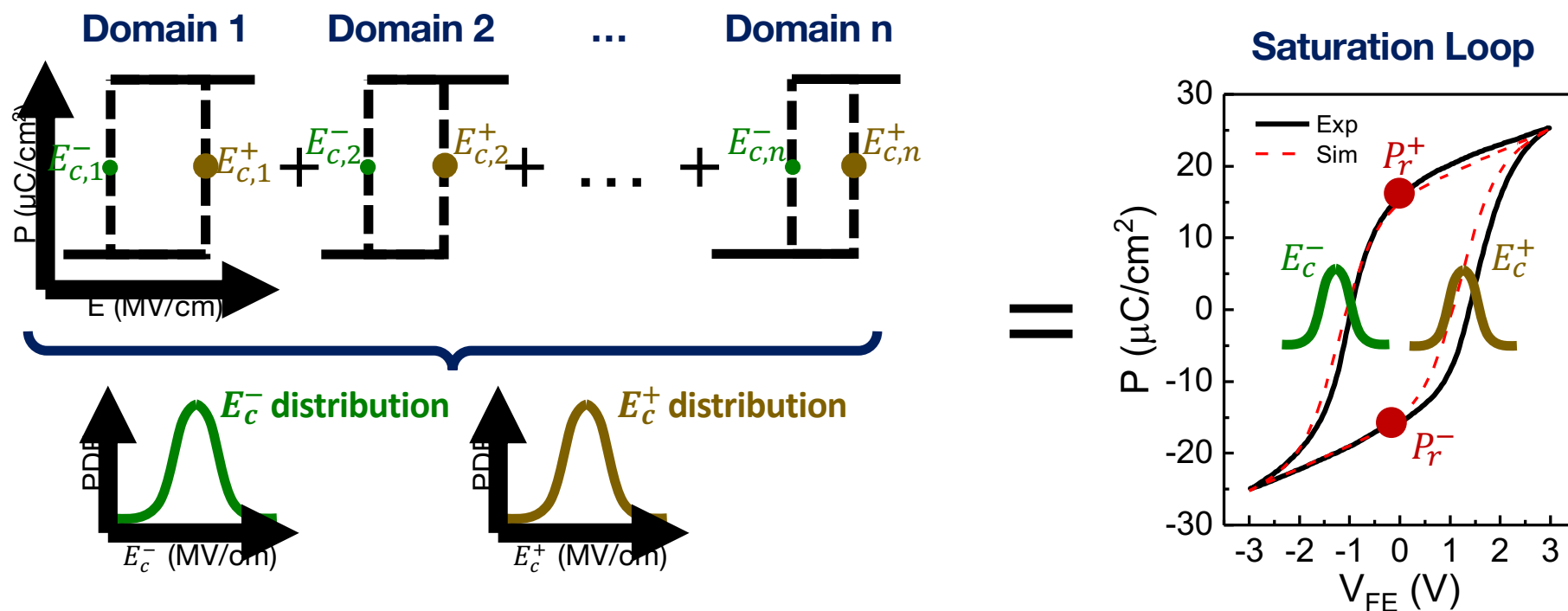
$\alpha, \beta, \gamma$  calibrated to hafnium zirconium oxide (HZO)

$\alpha$	$-7 \times 10^9 \text{ m/F}$
$\beta$	$3.3 \times 10^{10} \text{ m}^5/\text{F/coul}^2$
$\gamma$	$7 \times 10^9 \text{ m/F}$
$\rho$	0.25
$t_{\text{FE}}$	5.7 nm



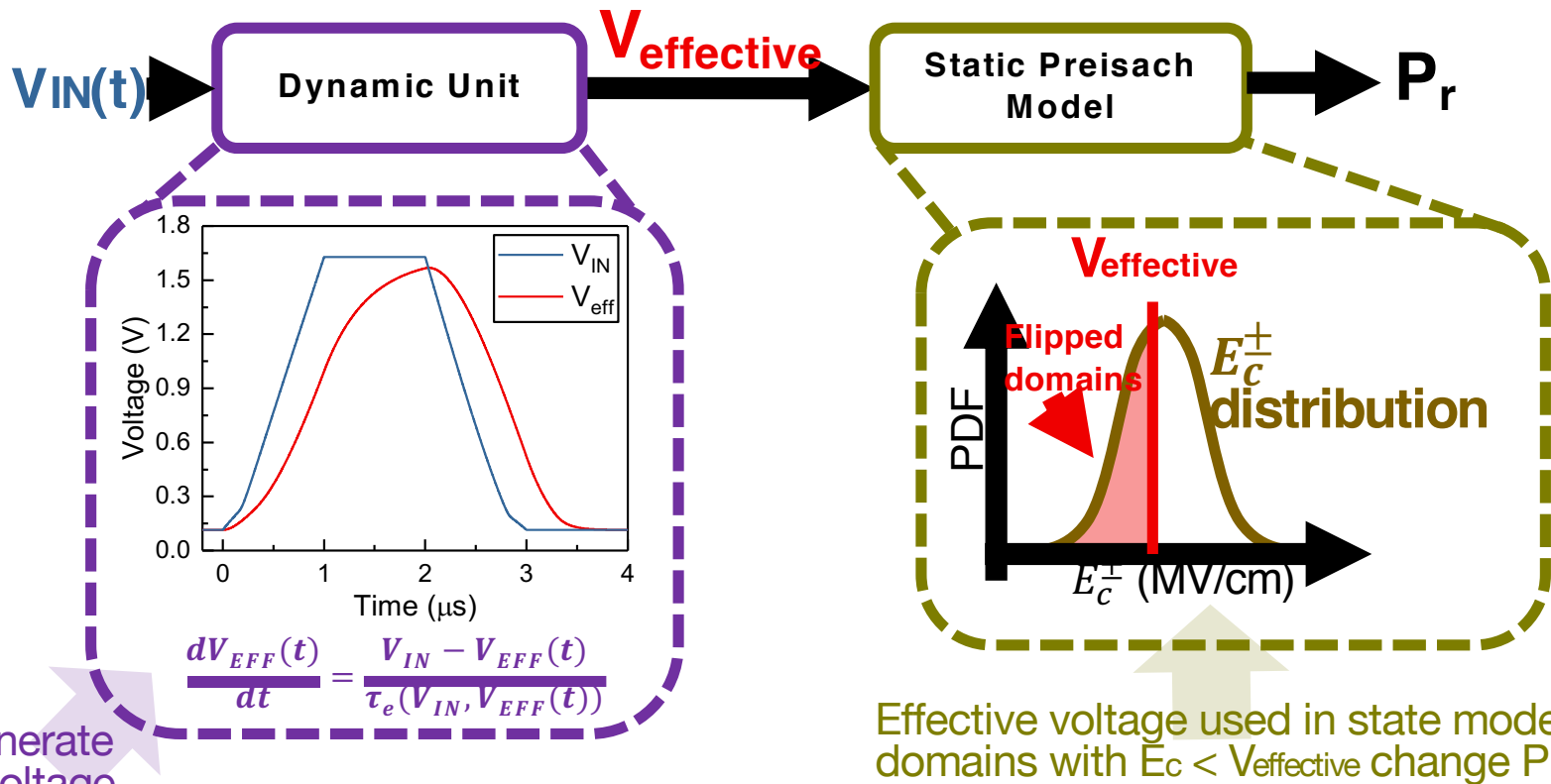
FeFET simulated by combining self-consistent LK equation with 45 nm PTM

# Multi-Domain Preisach Model



- The response of HZO film is described by the total contributions of many ideal ferroelectric domains of varying  $E_c^\pm$ .

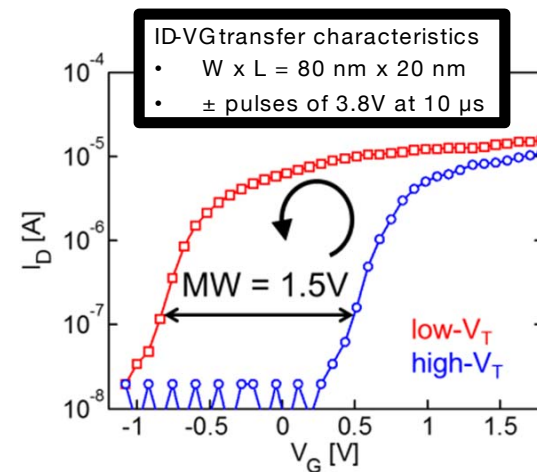
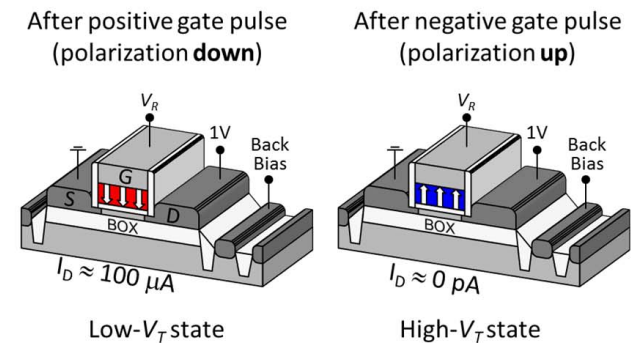
# Multi-Domain Preisach Model



Calibration to the measured data allows model to accurately capture  $P_r = f(V_{IN}, t)$

# Operation

- Can switch FeFET polarization with:
  - Positive gate voltage pulse (program)
  - Negative gate voltage pulse (erase)
- Pulse causes stable, reversible  $V_t$  shift
  - Low  $V_t$ , high  $V_t$  depends on dipole's orientation
- 2 distinguishable states = memory window
  - Sense with readout of drain current



May tradeoff pulse duration, amplitude depending on application-level figures of merit

# Logic-in-memory & Compute-in-memory

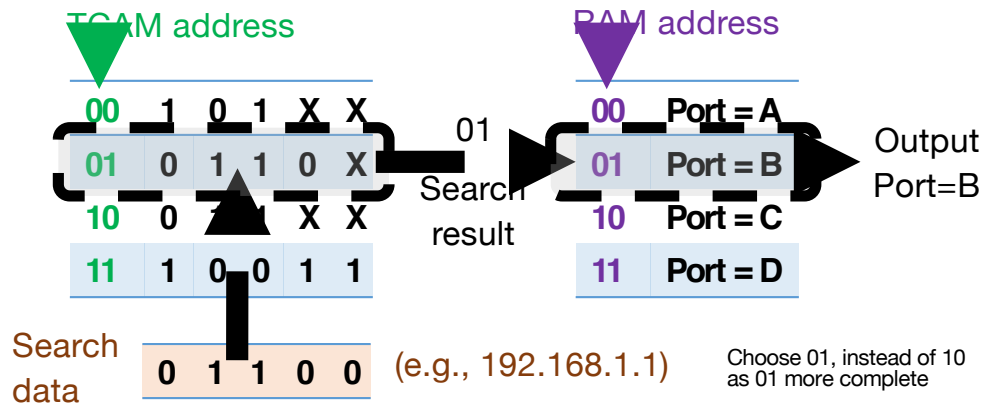


# Logic-in-memory: CAMs

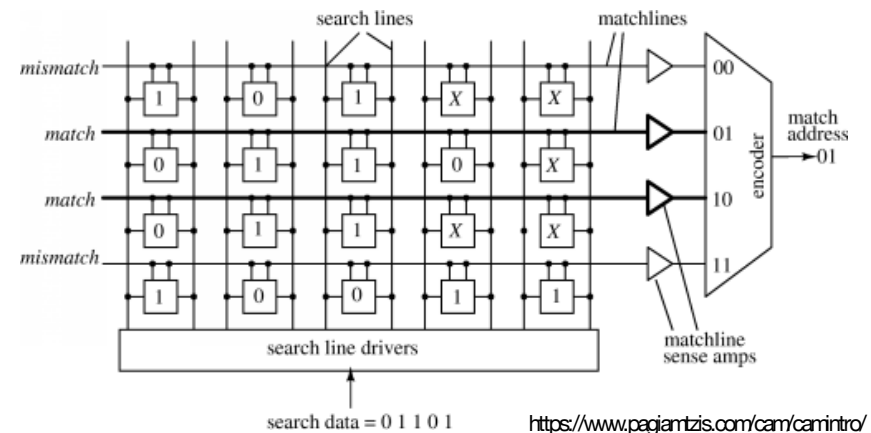
## Content Addressable Memory (CAM)

- Fast HW search  $O(1)$  for search intensive apps
- Often use **ternary CAM** (TCAM) – i.e., store 1, 0, or X (where X is “don’t care” (DC))
- TCAMs applicable to database apps, neural networks, *routers and switches*, etc.

### Address lookup with TCAM, RAM



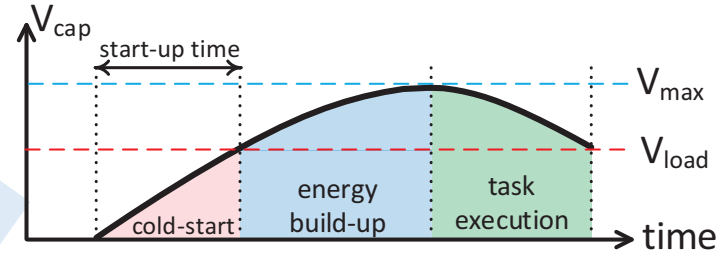
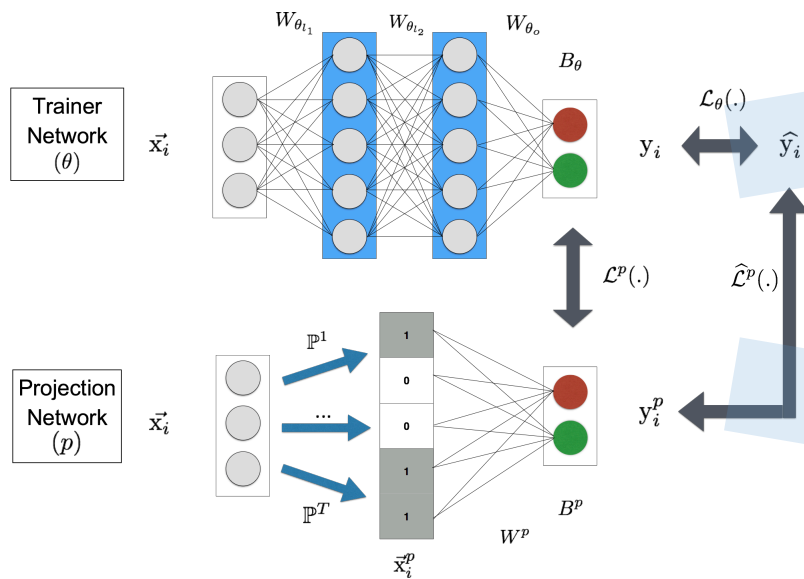
### TCAM array architecture



# Emerging neuromorphic computing models

e.g., **Projection Networks**

Train neural network, lightweight network in lockstep



**Goal:** more accurate, powerful machine learning models in resource constrained environments

The choice of the type of projection matrix  $\mathbb{P}$  as well as representation of the projected space  $\Omega_p$  in our setup has a direct effect on the computation cost and model size. We propose to leverage an efficient randomized projection method using a modified version of **locality sensitive hashing (LSH)** to define  $\mathbb{P}(\cdot)$ .

**TCAM-supported hashing again an important compute kernel**

ProjectionNet: Learning Efficient On-Device Deep Networks Using Neural Projections

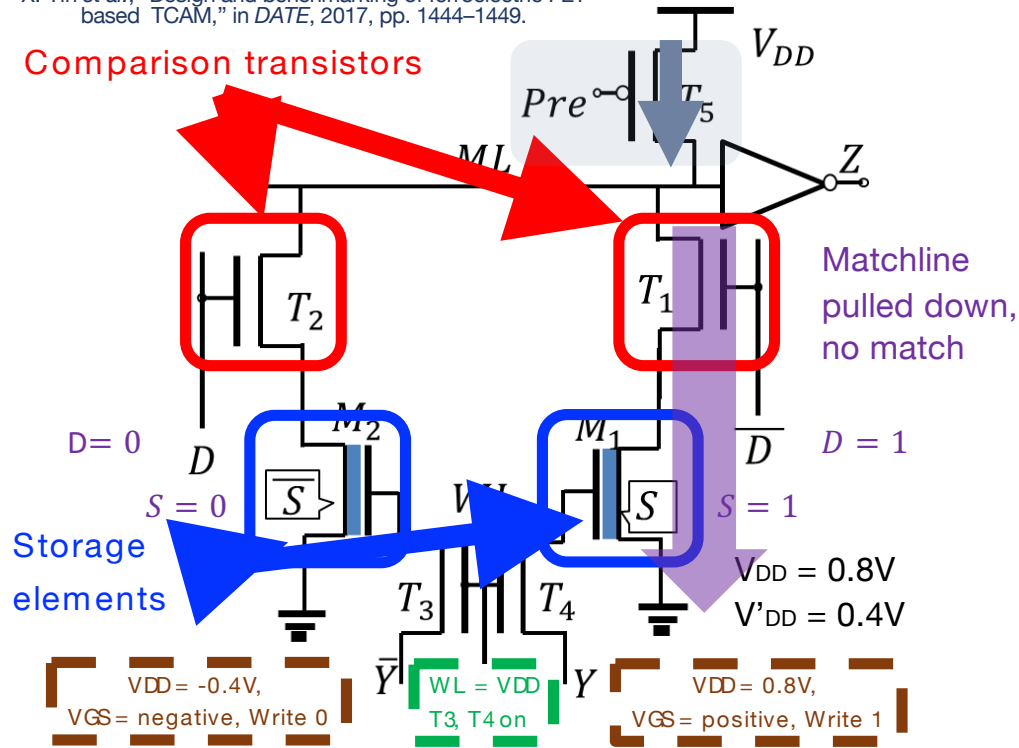
Sujith Ravi  
Google Research, Mountain View, CA, USA  
sravi@google.com

Look at TCAMs based on **ASCENT** technologies to support these models, other applications – consider FeFETs here...

# 4T, 2FeFET TCAMs (w/negative supply, LK model)

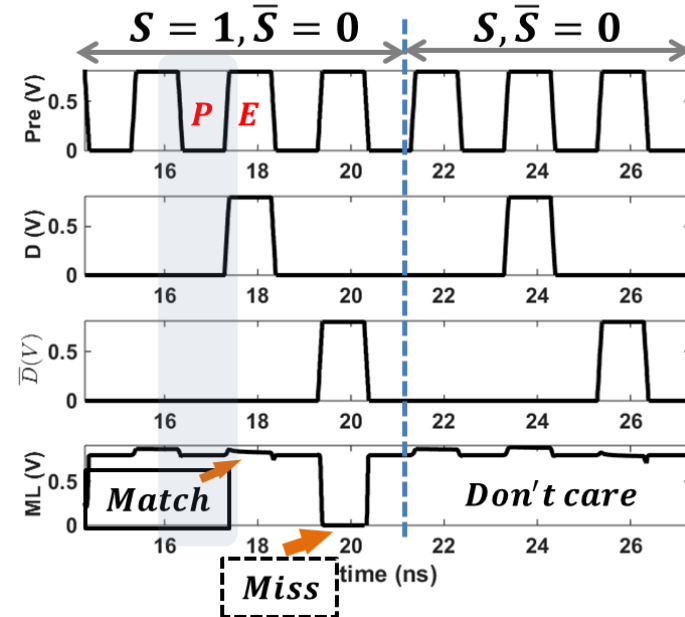
X. Yin et al., "Design and benchmarking of ferroelectric FET based TCAM," in DATE, 2017, pp. 1444-1449.

Comparison transistors



1<sup>st</sup>, pre-charge matchline

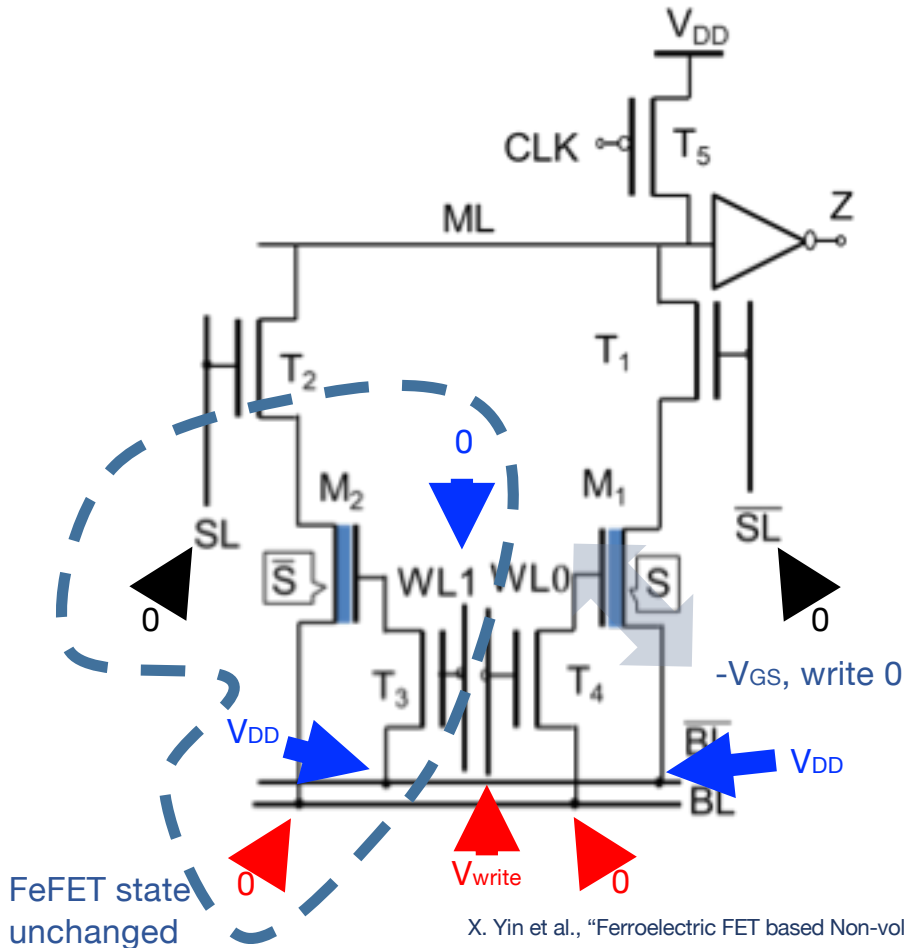
P = Pre-charge phase, E = Evaluate phase



Mode	D/D	Y/Y	WL
Write	S = 1	0	V <sub>DD</sub> / -V' <sub>DD</sub>
	S = 0	0	-V' <sub>DD</sub> / V <sub>DD</sub>

Mode	D/D	Y/Y	WL
Search	D = 1	V <sub>DD</sub> / 0	0
	D = 0	0 / V <sub>DD</sub>	0

# 4T, 2FeFET TCAMs (w/o negative supply, LK model)

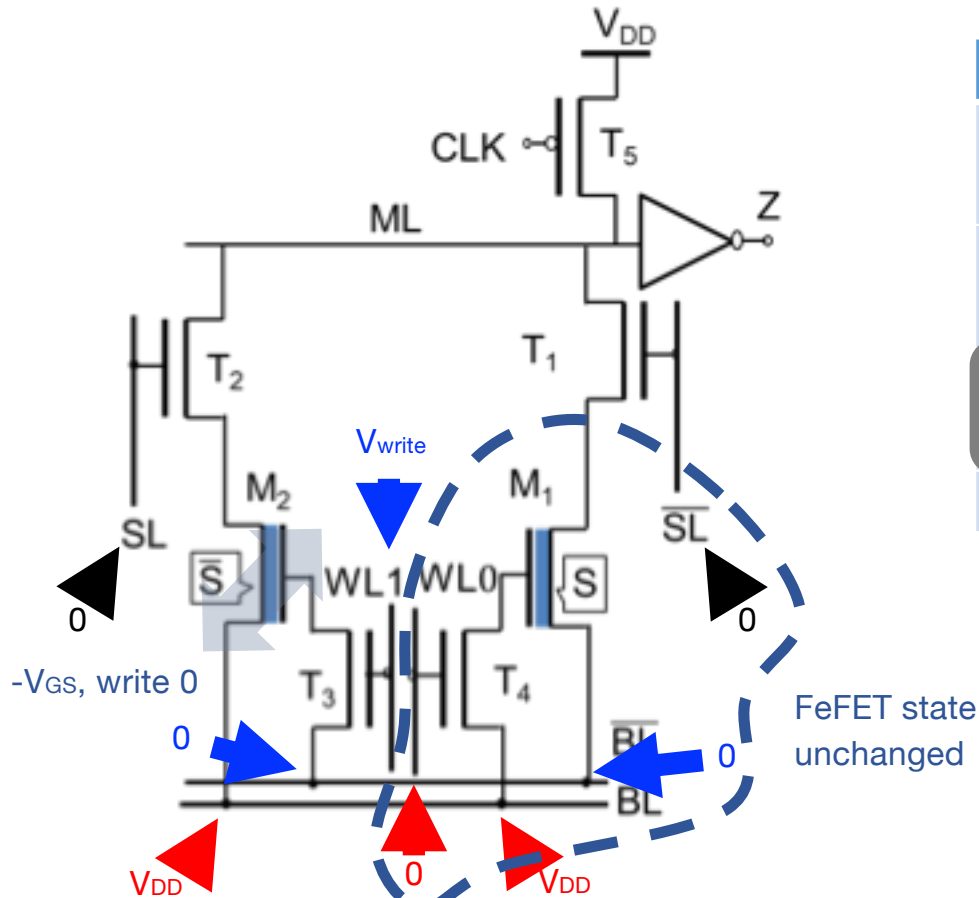


X. Yin et al., "Ferroelectric FET based Non-volatile Logic-in-Memory Circuits," IEEE TVLSI (in submission), 2018.

**Prior work** considered TCAM based on LK model with *negative* supply

	Step	WL0/WL1	BL/ $\overline{BL}$	SL/ $\overline{SL}$
Write 0	1	$V_{write}/0$	$0/V_{DD}$	0
	2	$0/V_{write}$	$V_{DD}/0$	
Write 1	1	$V_{write}/0$	$V_{DD}/0$	0
	2	$0/V_{write}$	$0/V_{DD}$	
Don't care	1	$V_{write}/0$	$0/V_{DD}$	0
	2	$0/V_{write}$	$V_{DD}/0$	
search		0/0	0/0	data

# 4T, 2FeFET TCAMs (w/o negative supply, LK model)

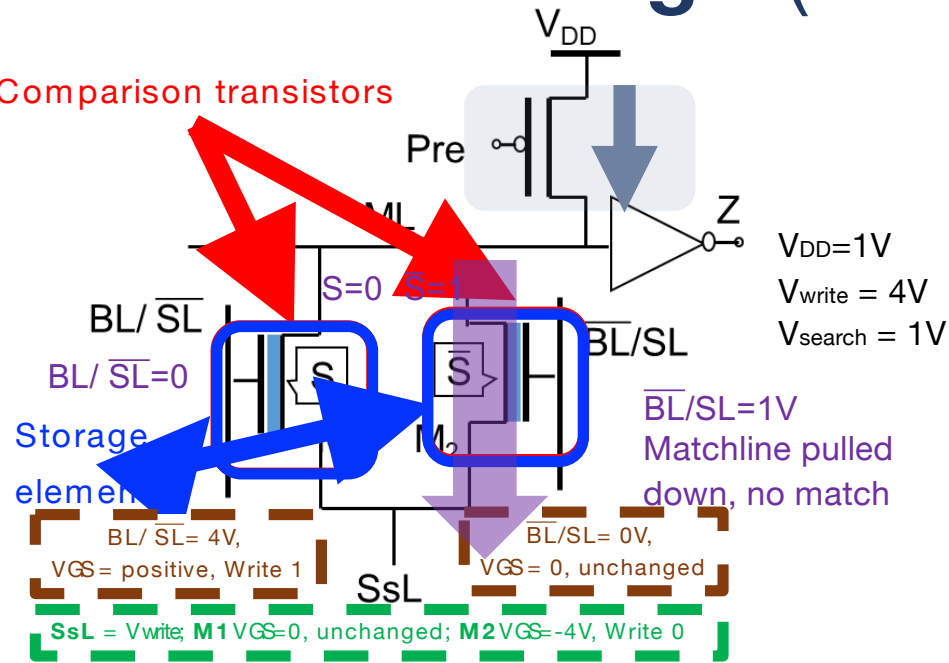


	Step	WL0/WL1	BL/ $\overline{BL}$	SL/ $\overline{SL}$
Write 0	1	$V_{write}/0$	$0/V_{DD}$	0
	2	$0/V_{write}$	$V_{DD}/0$	
Write 1	1	$V_{write}/0$	$V_{DD}/0$	0
	2	$0/V_{write}$	$0/V_{DD}$	
Don't care	1	$V_{write}/0$	$0/V_{DD}$	0
	2	$0/V_{write}$	$V_{DD}/0$	
search		$0/0$	$0/0$	data

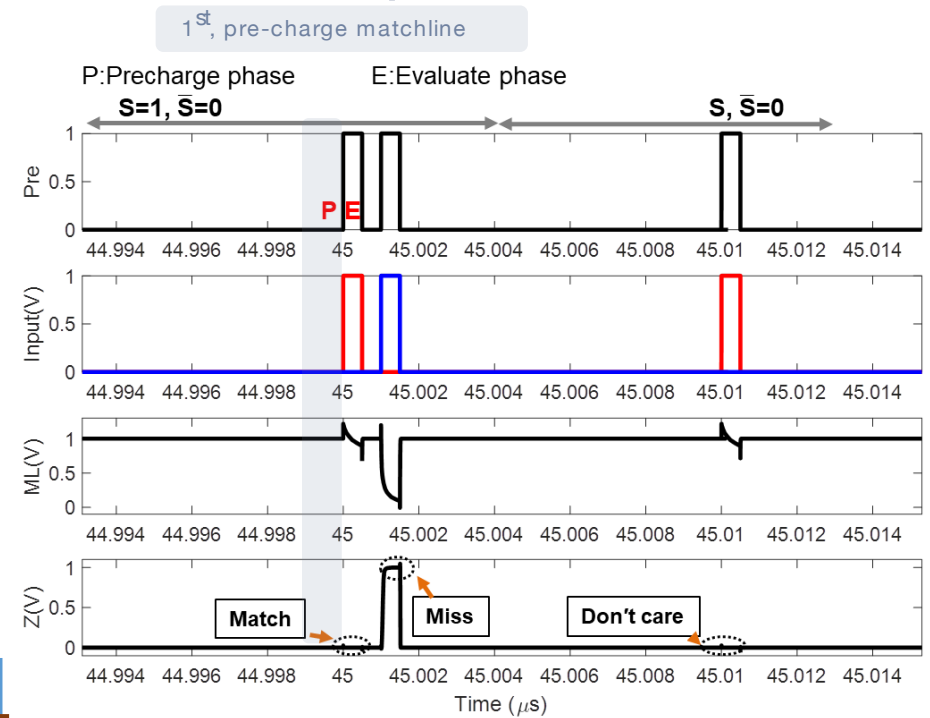
X. Yin et al., "Ferroelectric FET based Non-volatile Logic-in-Memory Circuits," IEEE TVLSI (in submission), 2018.

# 2T FeFET design (Preisach model)

Comparison transistors

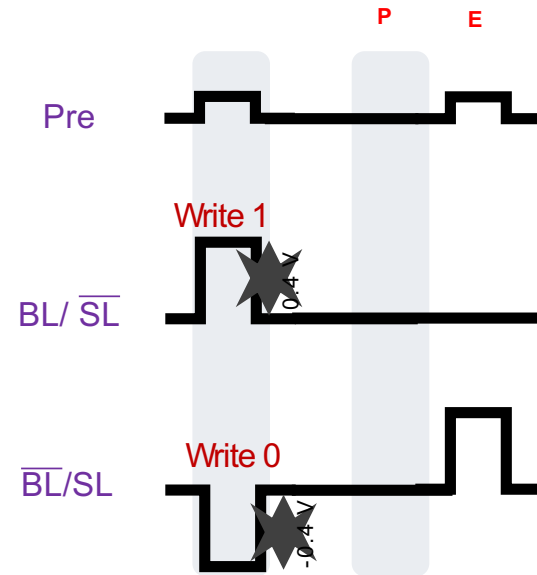
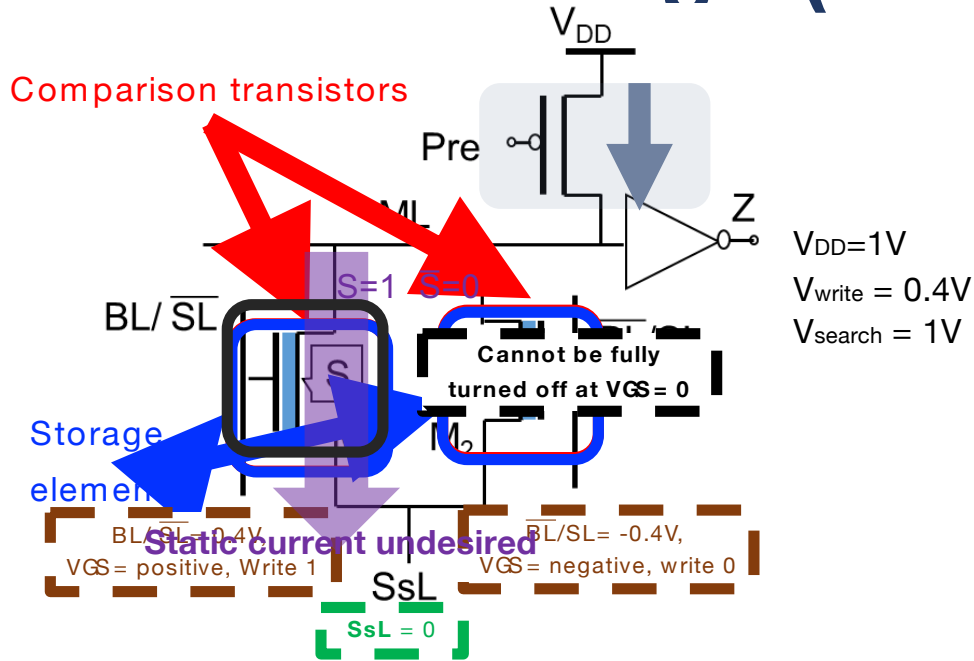


Mode			$BL/\bar{S}L$	$\bar{BL}/SL$	$SsL$
Write	$S = 1$	Step 1	$V_{write}$	0	0
		Step 2	$V_{write}$	0	$V_{write}$
	$S = 0$	Step 1	0	$V_{write}$	$V_{write}$
		Step 2	0	$V_{write}$	0

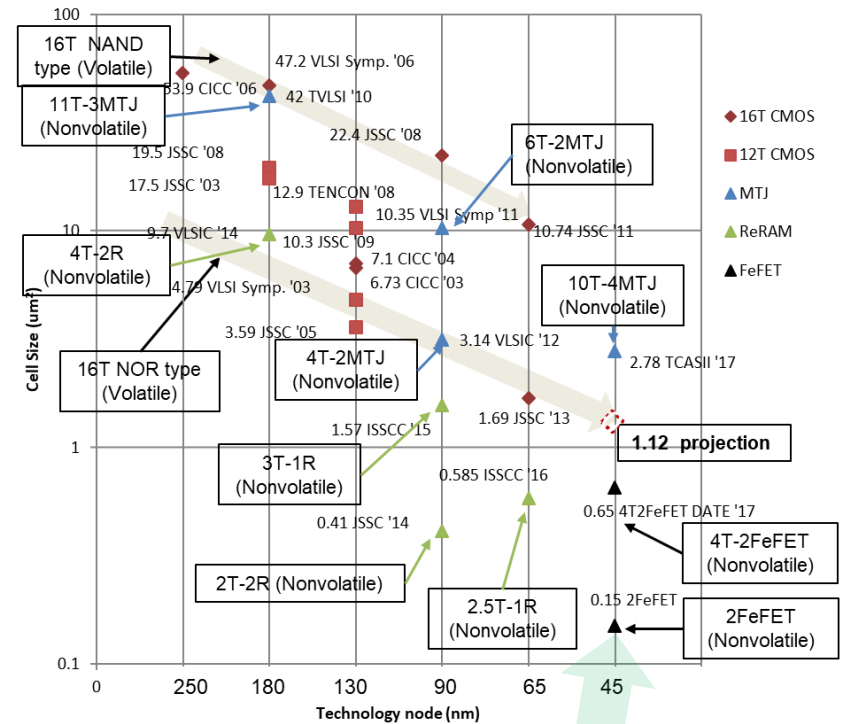
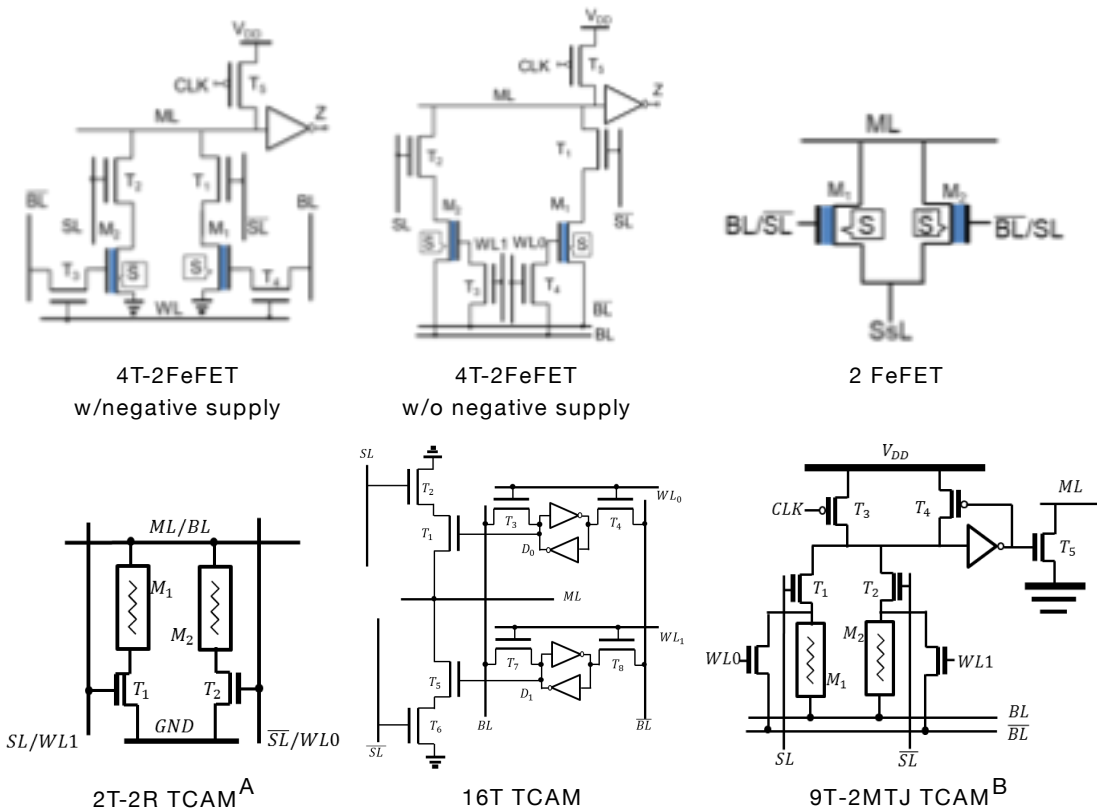


Mode		$BL/\bar{S}L$	$\bar{BL}/SL$	$SsL$
Search	1	0	$V_{search}$	0
	0	$V_{search}$	0	0

# 2T FeFET design (LK model)



# Benchmarking (area comparisons)



Current projections suggest competitive density

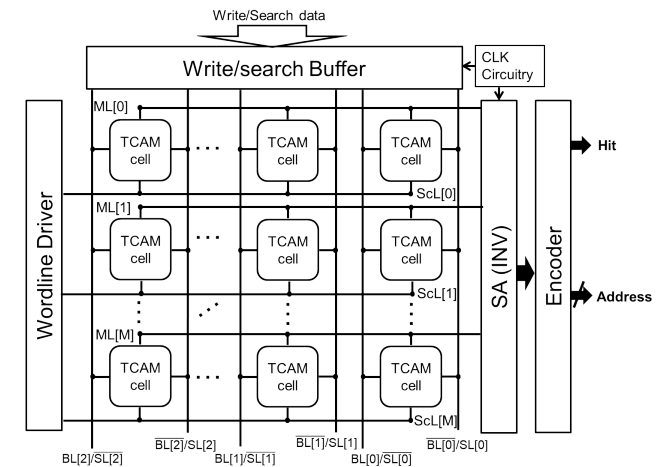
<sup>A</sup>J. Li, R. K. Montoye, M. Ishii, and L. Chang, "1 mb 0.41  $\mu\text{m}^2$  2t-2r cell nonvolatile TCAM with two-bit encoding and clocked self-referenced sensing," *IEEE Journal of Solid-State Circuits*, vol. 49, no. 4, pp. 896–907, 2014.

<sup>B</sup>S. Matsunaga, A. Katsumata, M. Natsui, T. Endoh, H. Ohno, and T. Hanyu, "Design of a nine-transistor/two-magnetic-tunnel-junction-cell-based low-energy nonvolatile ternary content-addressable memory," *Japanese J. of Applied Physics*, vol. 51, no. 2S, p. 02BM06, 2012

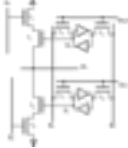
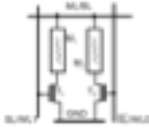
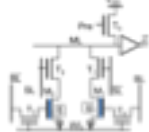



# Benchmarking methodology

- All designs evaluated in context of 64x64 array
- Assume
  - 45 nm PTM
  - Inverter-based SA
  - Minimum sized transistors for TCAM cell, SA
- Extract wiring parasitics from DESTINY
  - M. Poremba, et al., “Destiny: A tool for modeling emerging 3D NVM and EDRAM caches,” in *DATE*, 2015, pp. 1543–6.
- Delay assumes worst case
  - (i.e., 1-bit mismatch...)

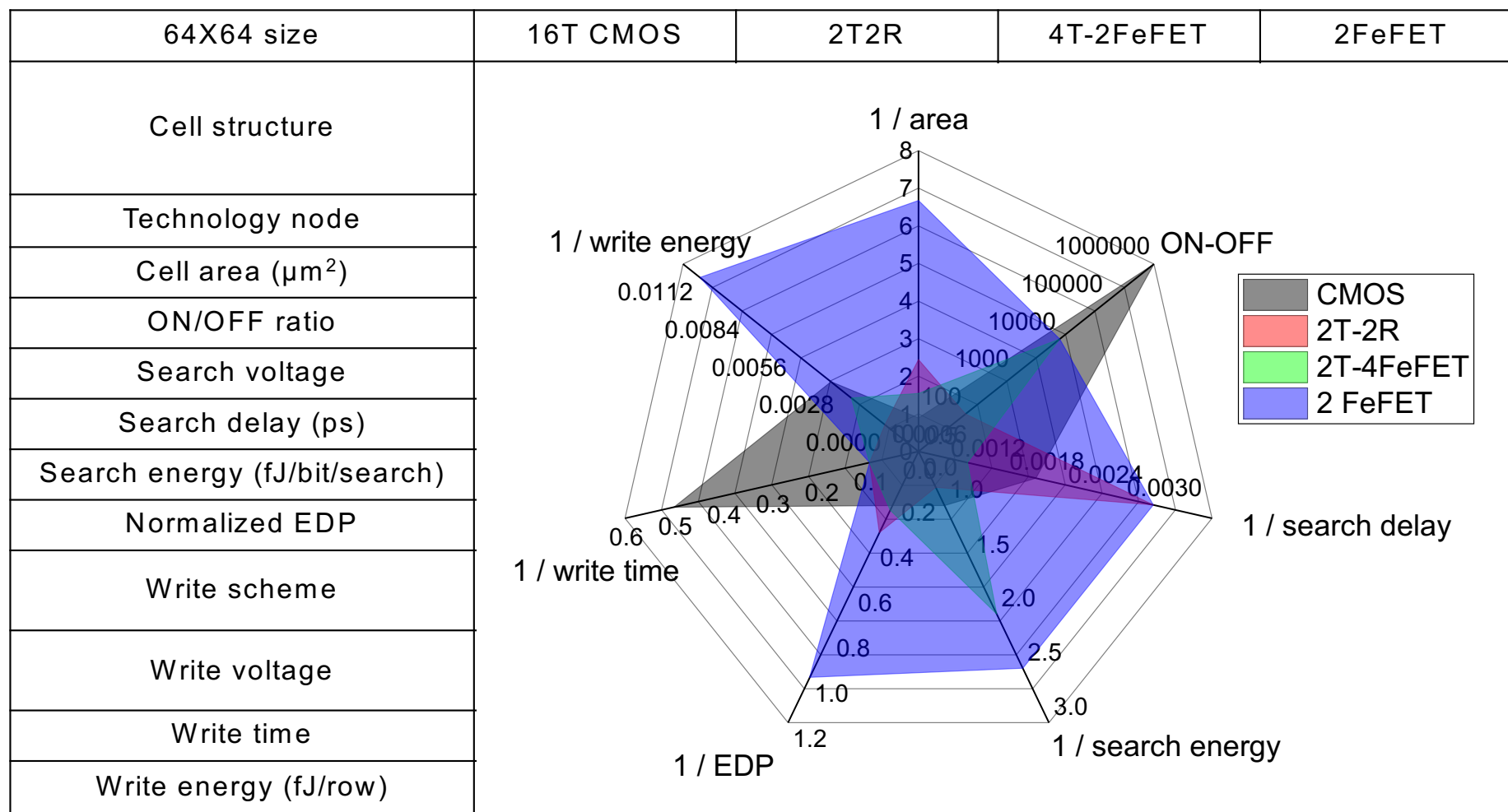


# Benchmarking (other figures of merit)

64X64 size	16T CMOS	2T2R	4T-2FeFET	2FeFET
Cell structure				
Technology node	45nm	45nm	45nm	45nm
Cell area ( $\mu\text{m}^2$ )	1.12 (7.5x)	0.41 <sup>[1]</sup> (2.7x)	0.65 (4.3x)	0.15 (1x)
ON/OFF ratio	$\sim 10^6$	$\sim 100$ <sup>[2]</sup>	$\sim 10^4$	$\sim 10^4$
Search voltage	1V	1V	1V	1V
Search delay (ps)	582 (1.7x)	350 (1.03x)	1013 (3.0x)	341 (1x)
Search energy (fJ/bit/search)	1.0 (2.4x)	1.2 (2.7x)	0.5 (1.3x)	0.4 (1x)
Normalized EDP	4.1x	2.8x	3.8x	1x
Write scheme	Voltage driven dynamic switching	Current driven	Voltage driven	Voltage driven
Write voltage	1V	Set 1.8V <sup>[3]</sup> Reset 1.2V <sup>[3]</sup>	$\pm 4\text{V}$	$\pm 4\text{V}$
Write time	< 2ns	$\sim 10$ ns	10 ns	10 ns
Write energy (fJ/row)	309 (3.5x)	288000 (3225x) <sup>[3]</sup>	512 (5.7x)	89 (1x)

1. Li, Jing, et al. "1Mb 0.41  $\mu\text{m}^2$  2T-2R cell nonvolatile TCAM with two-bit encoding and cloaked self-referenced sensing." IEEE VLSI/C, 2013.  
 2. Lastras-Montano, Miguel Angel, et al. "Architecting energy efficient crossbar-based memristive random-access memories." IEEE NANARCH 2015  
 3. Li, Shuangchen, et al. "Nvsm-cam: a circuit-level simulator for emerging nonvolatile memory based content-addressable memory." ICCAD, 2016.

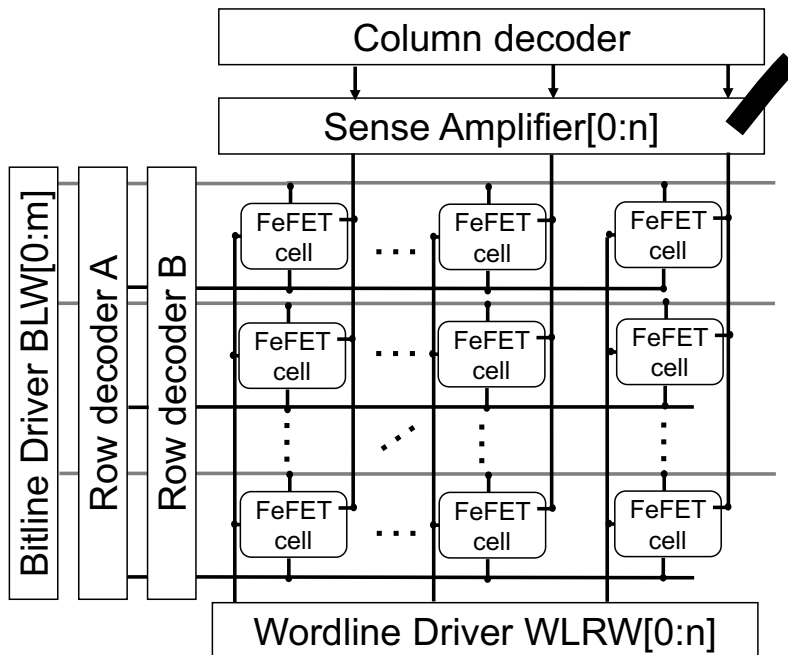
# Benchmarking (other figures of merit)



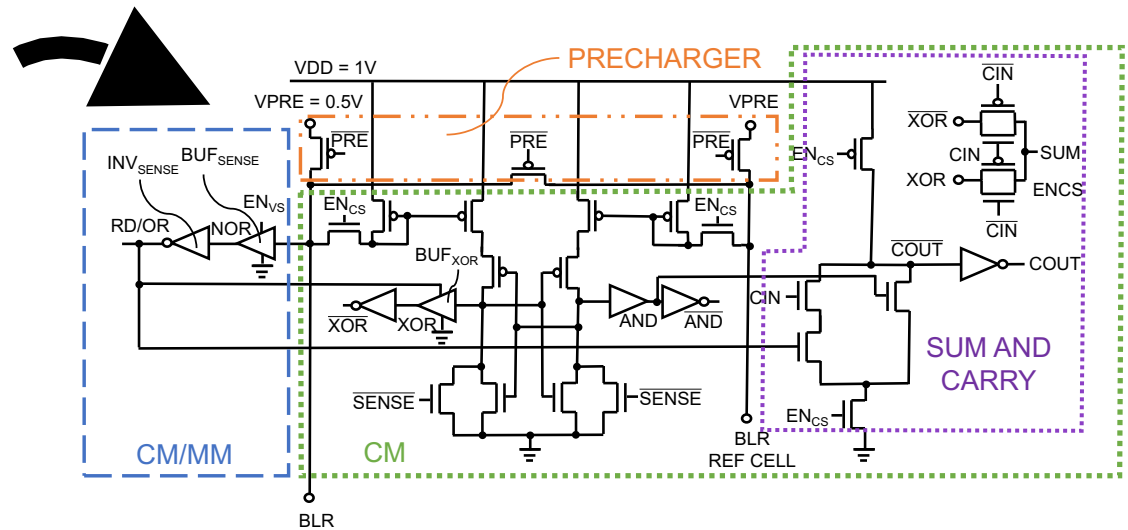
1. Li, Jing, et al. "1Mb 0.41  $\mu\text{m}^2$  2T-2R cell nonvolatile TCAM with two-bit encoding and cloaked self-referenced sensing." IEEE VLSI, 2013.  
 2. Lastras-Montano, Miguel Angel, et al. "Architecting energy efficient crossbar-based memristive random-access memories." IEEE NANARCH 2015  
 3. Li, Shuangchen, et al. "Nysim-cam: a circuit-level simulator for emerging nonvolatile memory based content-addressable memory." ICCAD, 2016.

# FeFET-based CIM: architecture

## FeFET-based CIM architecture

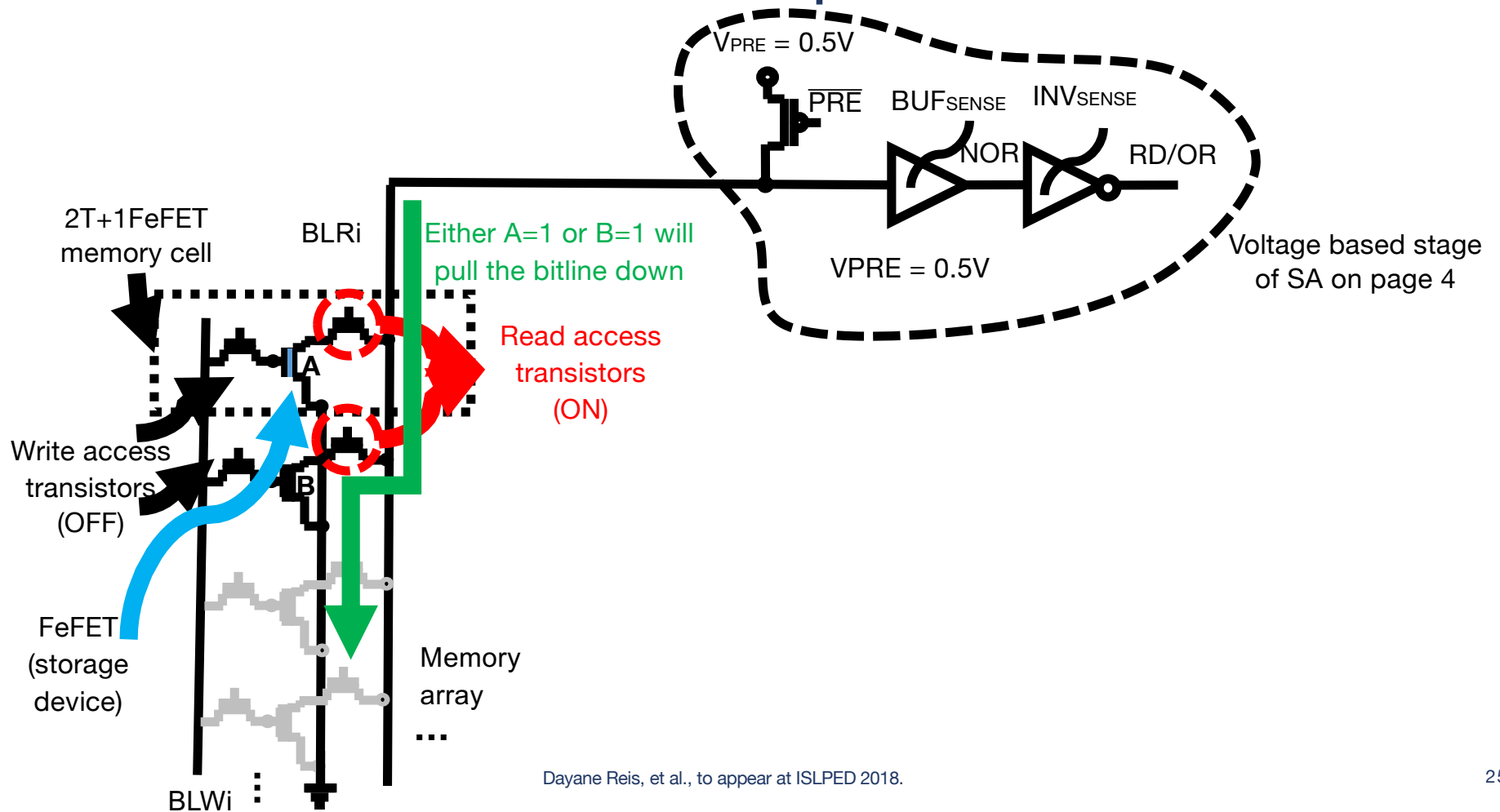


## FeFET-CIM customized sense amp (SA)



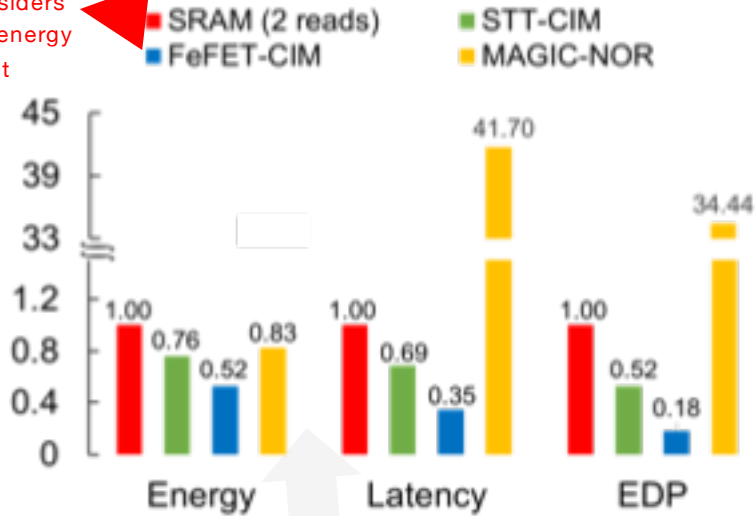
- **CM/MM** is voltage-based sense scheme responsible for (N)OR logic and reads
- **CM** is current-based sense scheme used for Boolean (N)AND, X(N)OR, and ADD; also leverages voltage scheme
- **SUM** and **CARRY** is additional circuitry for carry and sum

# FeFET-based CIM: OR operations



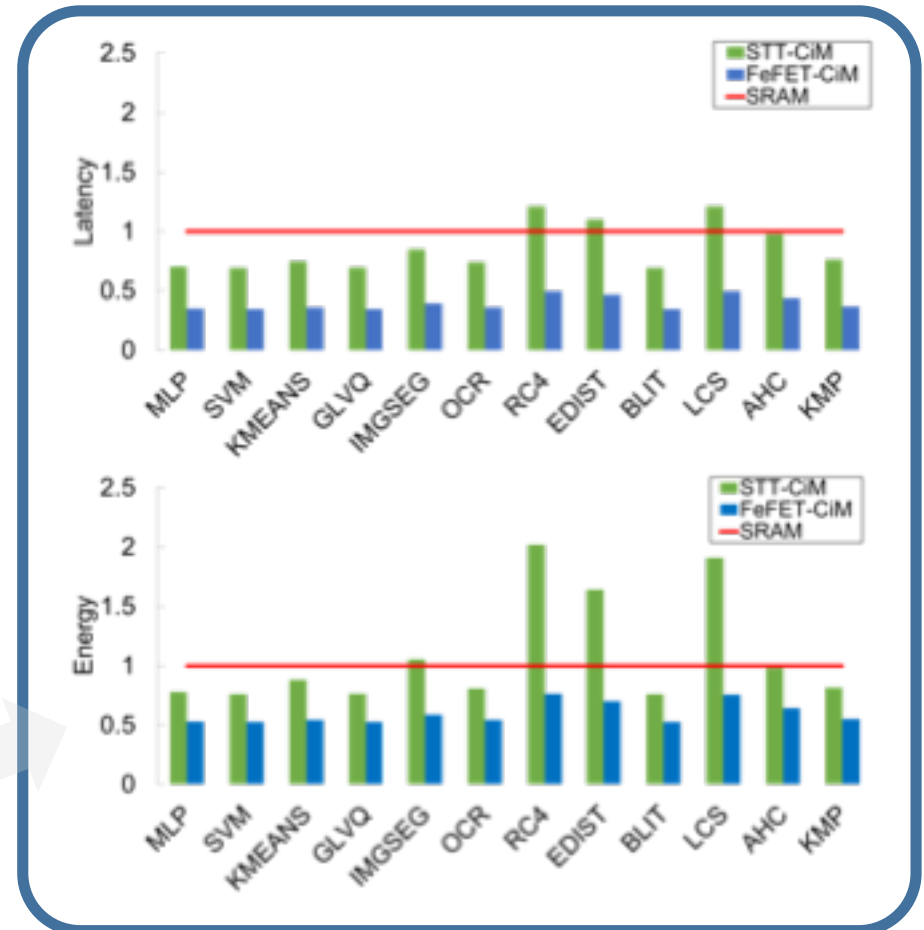
# FeFET-based CIM: benchmarking

Only considers memory energy at present



FeFET-CIM has speed-ups (energy reductions) of  $\sim 119X$  ( $\sim 1.6X$ ) and  $\sim 1.97X$  ( $\sim 1.5X$ ) over ReRAM and STT-RAM CIM for in-memory addition of 32-bit words

FeFET-CIM approach offers an average speedup of  $\sim 2.5X$  and energy reduction of  $\sim 1.7X$  when compared to a conventional (not in-memory) approach.



Computing in Memory with Spin-Transfer Torque Magnetic RAM

Shubham Jain, Ashish Ranjan, Kaushik Roy, Anand Raghunathan  
School of Electrical and Computer Engineering, Purdue University  
{jain130,aranjan,kaushik,raghunathan}@purdue.edu

STT-CIM

Logic Design Within Memristive Memories Using Memristor-Aided loGIC (MAGIC)

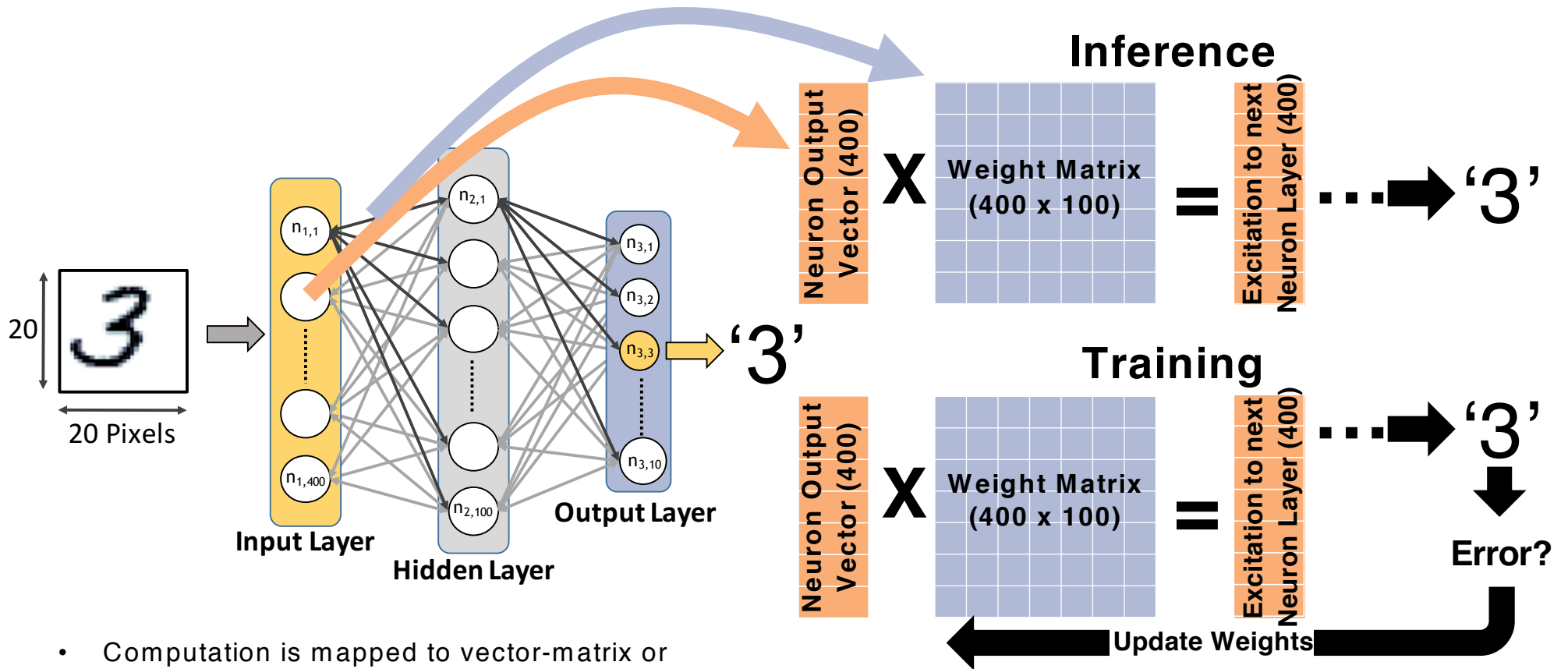
MAGIC-NOR

Nishil Talati, Saransh Gupta, Pravin Mane, and Shahar Kvatinsky, Member, IEEE

Dayane Reis, et al., to appear at ISLPED 2018.

# Neuromorphic applications

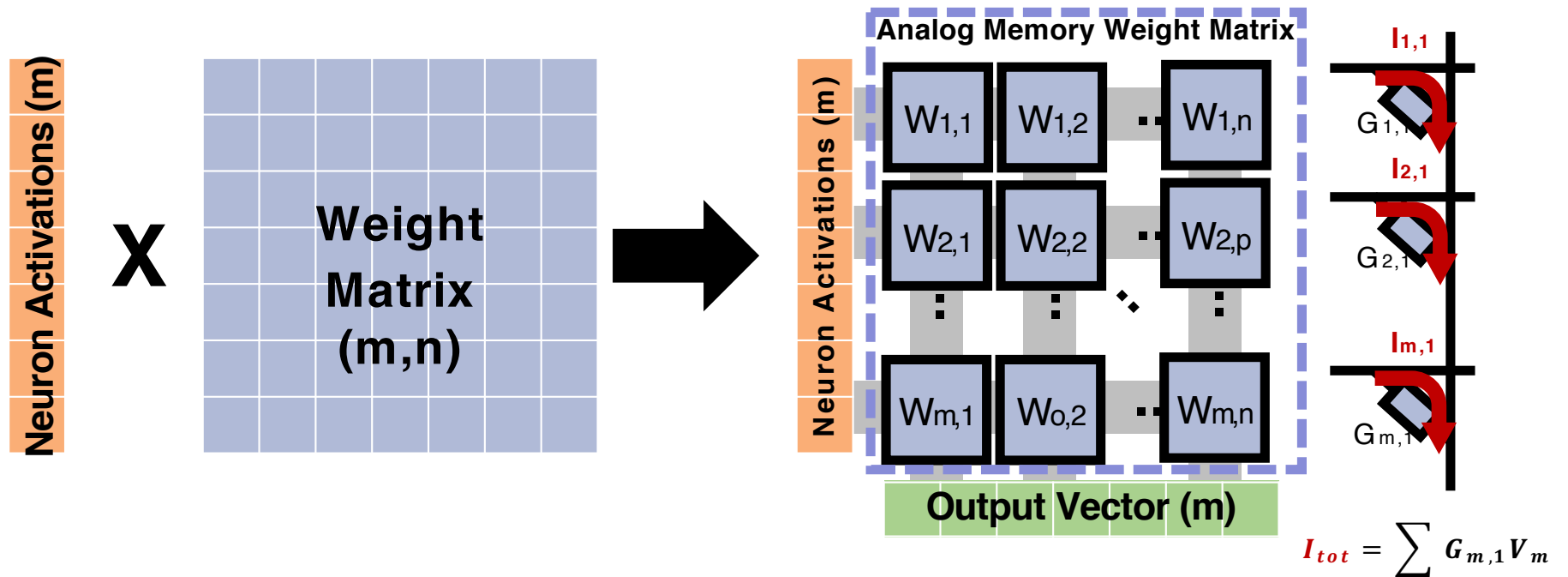
# Inference & training



- Computation is mapped to vector-matrix or matrix-matrix multiplication

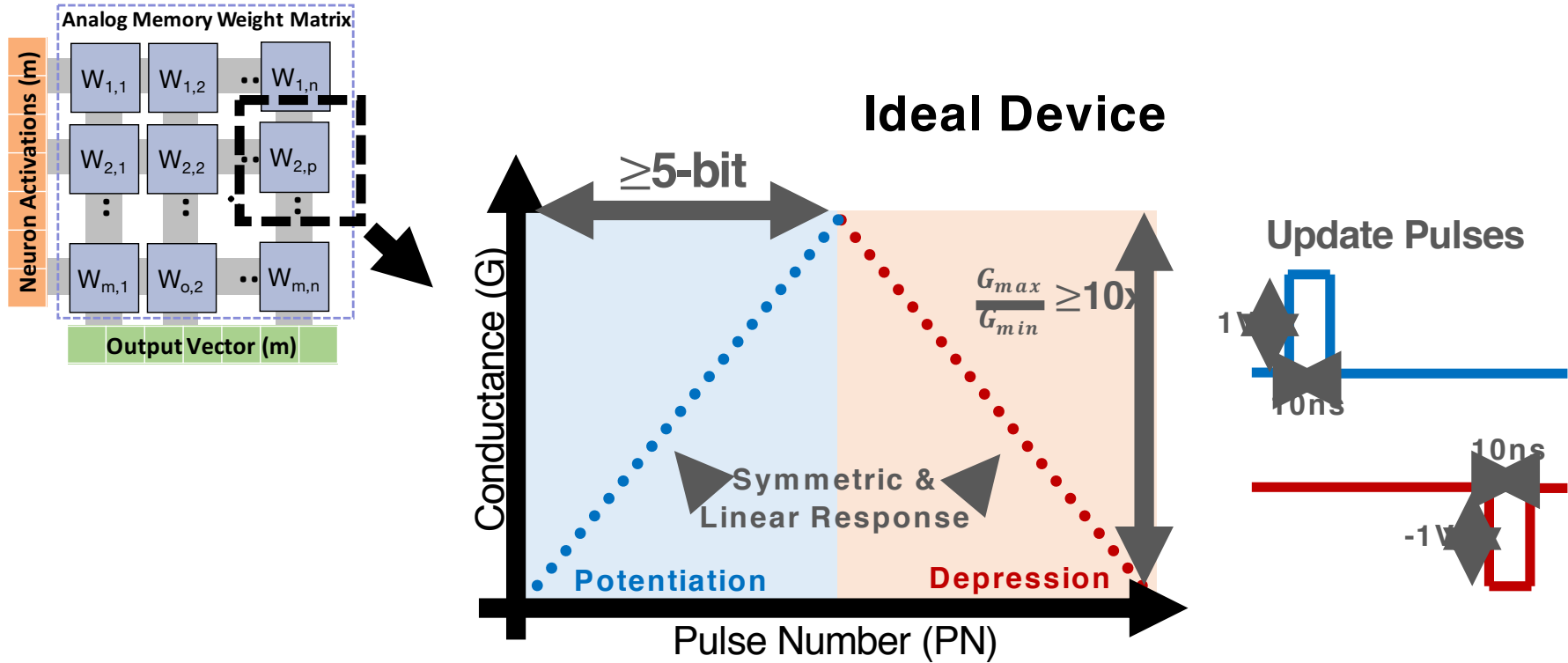


# Vector-matrix multiplication with crossbars



- Dense analog synaptic memory arrays perform MACs and update at the location of the data

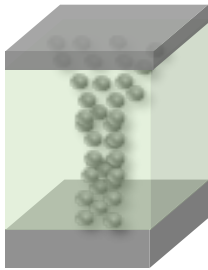
# Analog synaptic device characteristics



- Synaptic memory needs to be high density, low latency, energy efficient, and preserve high network accuracies.

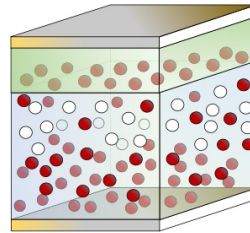
# Analog synapses

Filamentary O-RRAM



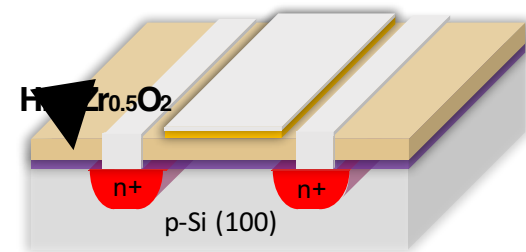
- ✓ **High density**
- Electro-thermal switching
- ✓ **<100ns pulse widths**
- Low  $G_{\max}/G_{\min}$  ratios demonstrated thus far

Non-Filamentary RRAM



- ✓ **High density**
- Electro-thermal switching
- <100ns pulse widths to be demonstrated
- Low  $G_{\max}/G_{\min}$  ratios demonstrated thus far

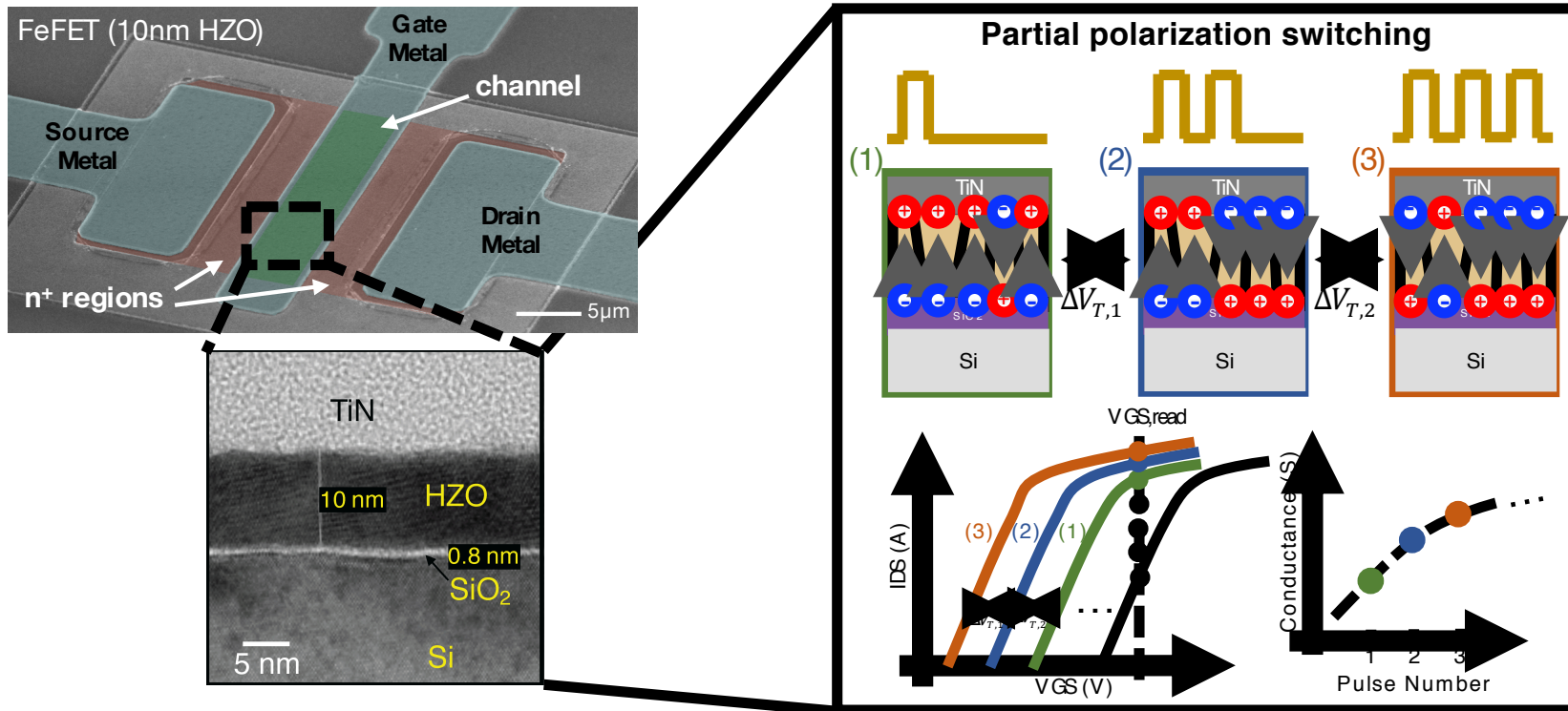
Ferroelectric FET



- 2T design proposed
- ✓ **Electric-field switching**
- ✓ **75ns pulse widths**
- ✓ **Large and tunable  $G_{\max}/G_{\min}$**

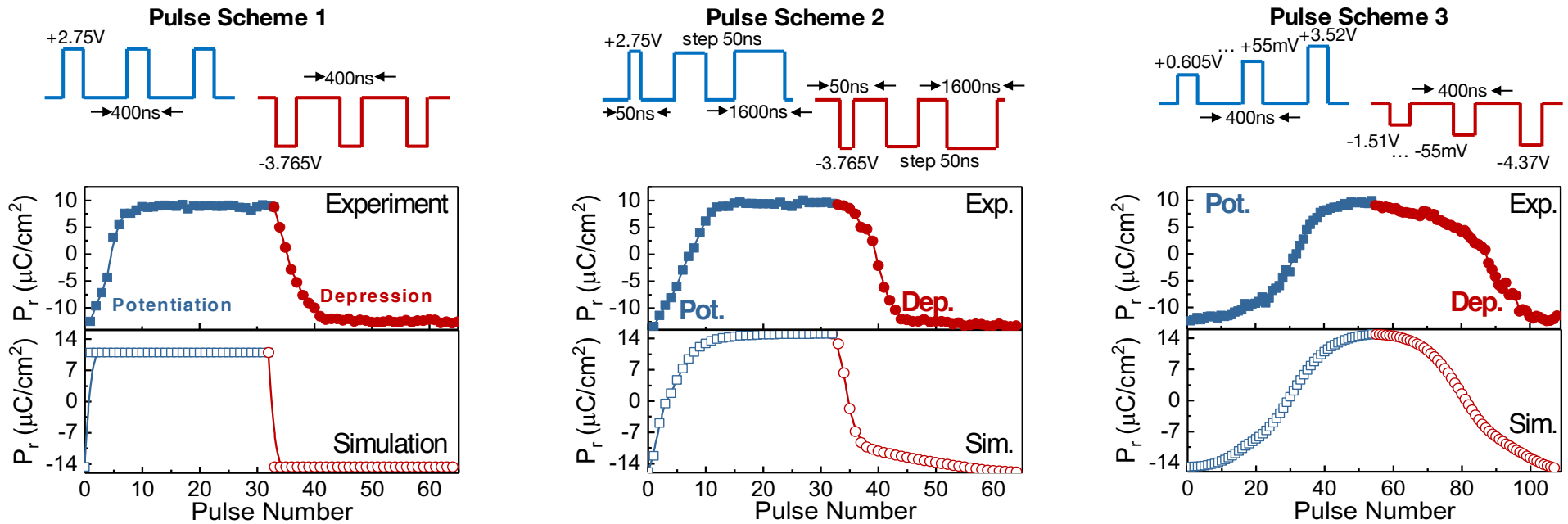
- Ferroelectric FET is a promising candidate for an analog synaptic memory device.

# Ferroelectric FET analog synapse



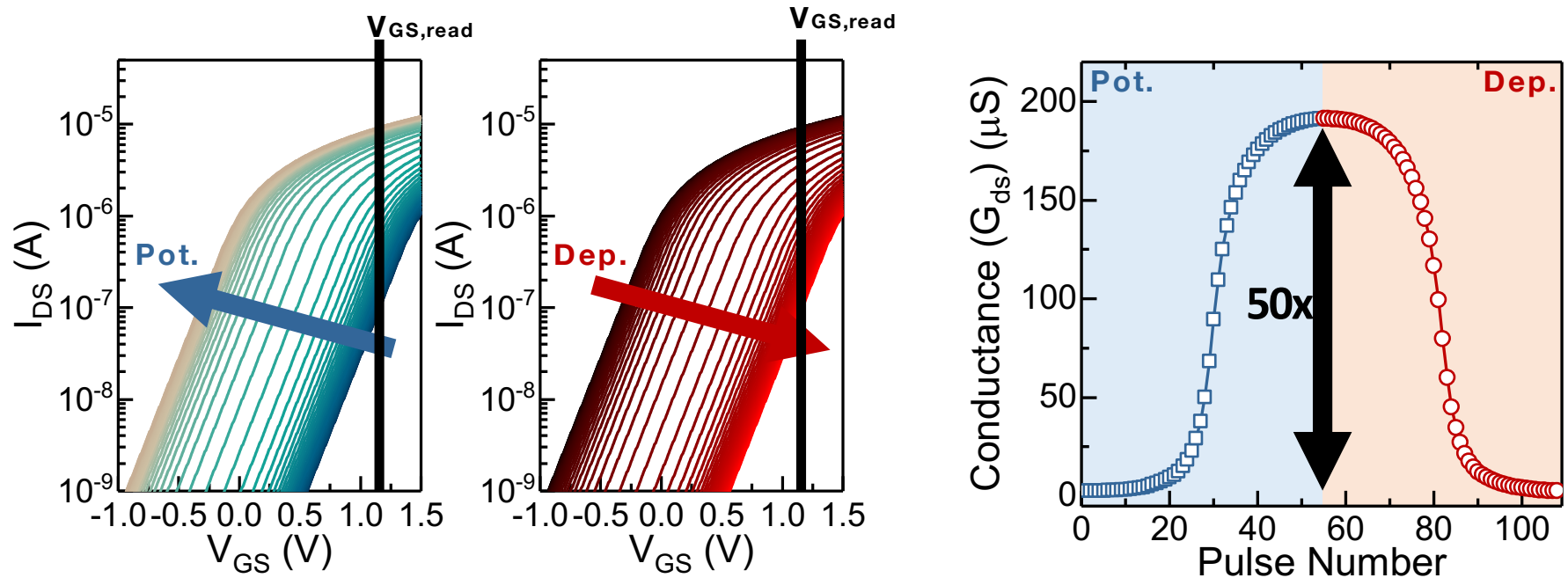
- Electric-field controlled partial polarization switching in ferroelectrics FETs can be harnessed for synaptic memory with nanosecond updates.

# Effect of Pulse Scheme on $P_r$



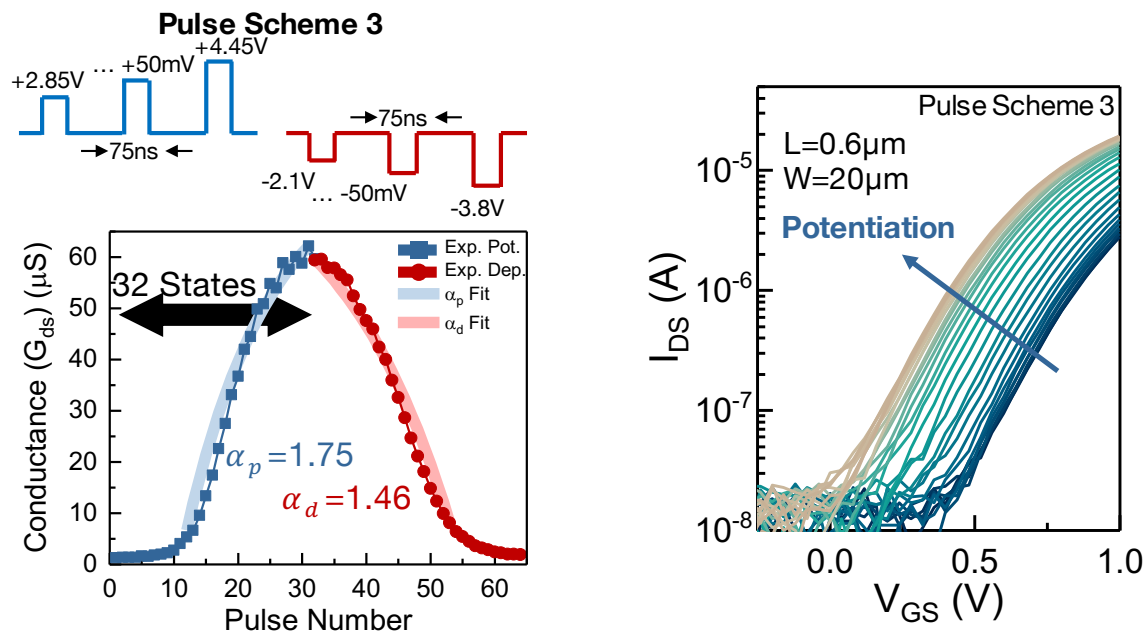
- Multi-domain Preisach model accurately captures the the response of the remnant polarization

# Simulated G vs Pulse Number: Scheme 3



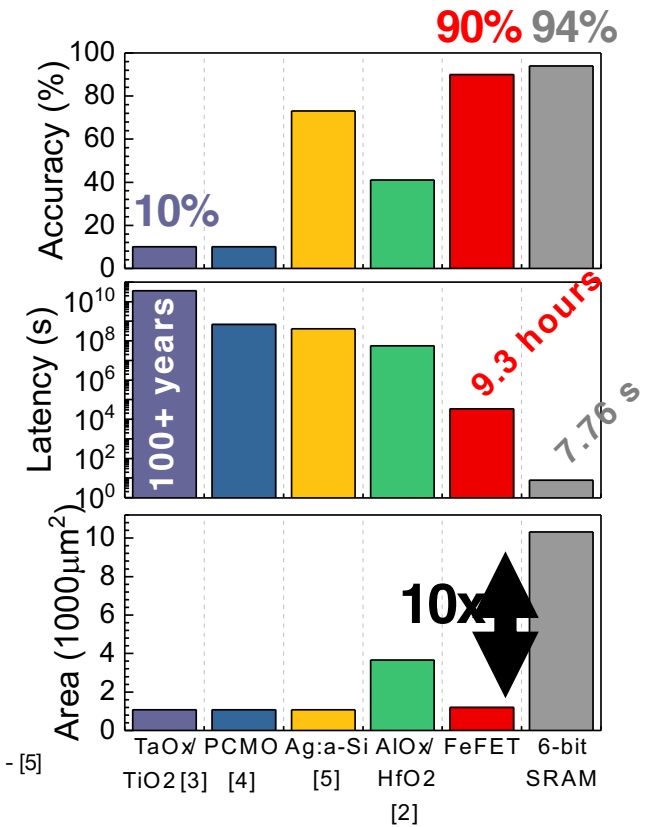
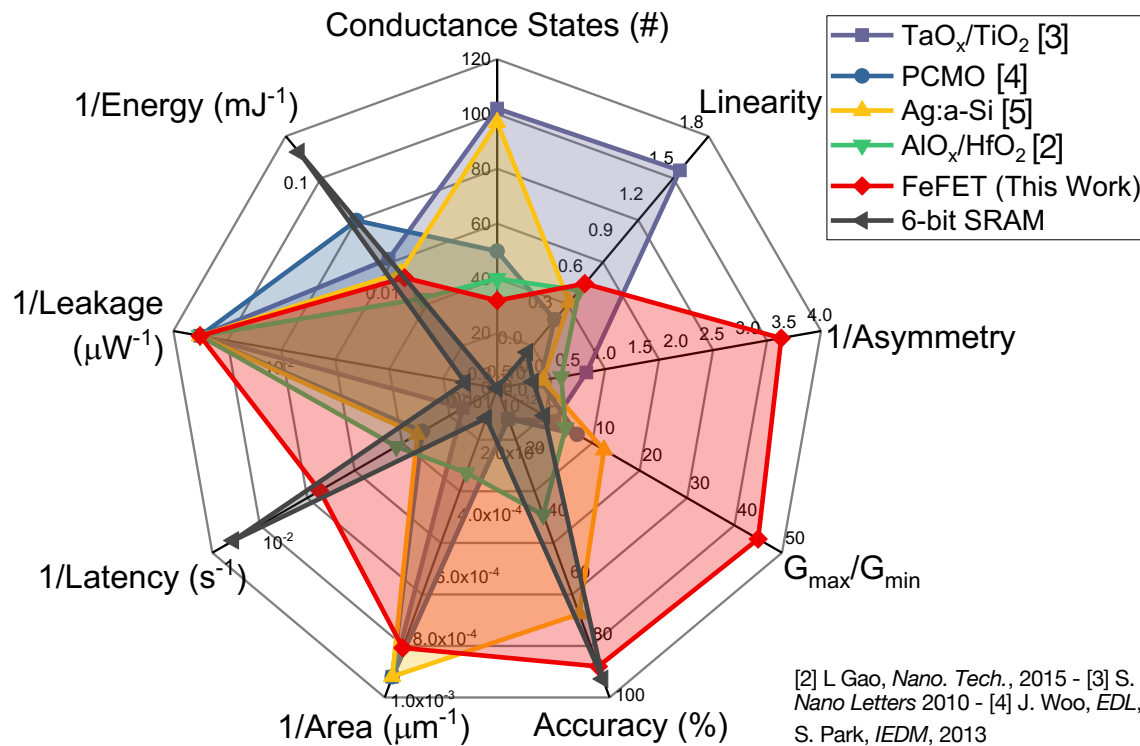
- FeFET synapse response from simulated  $P_r$  in programming scheme 3.

# FeFET Analog Synapse: Scheme 3



- Partial polarization switching within the ferroelectric gate oxide results in a gradual decrease/increase (potentiation/depression) in  $V_T$ .

# Analog Synapse Benchmarking

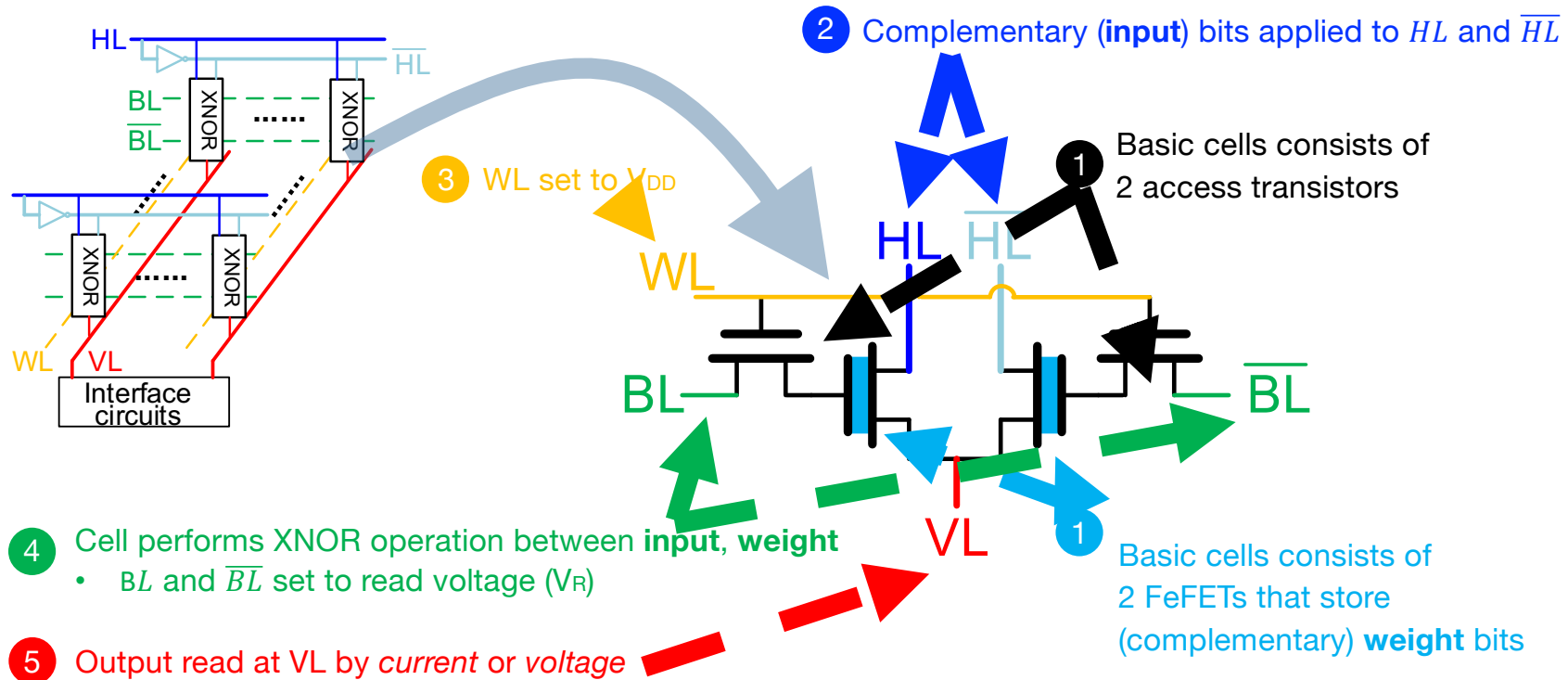


- FeFET under pulse scheme 3 exhibits the reduced footprint, high accuracy, and low latency

Jerry, Matthew, et al. "Ferroelectric FET analog synapse for acceleration of deep neural network training." in *IEDM*, p. 6-2, 2017.



# FeFET-based binary crossbars: circuits



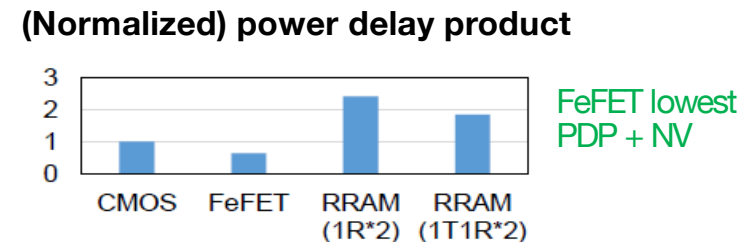
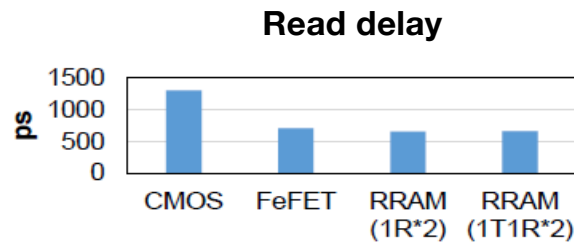
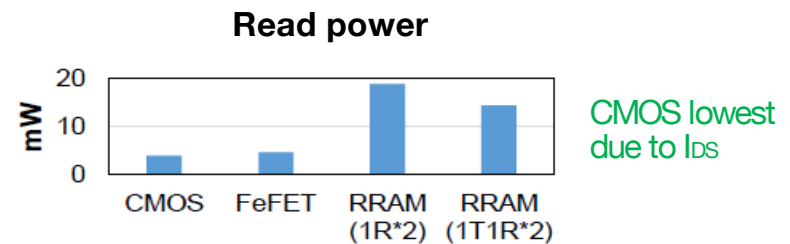
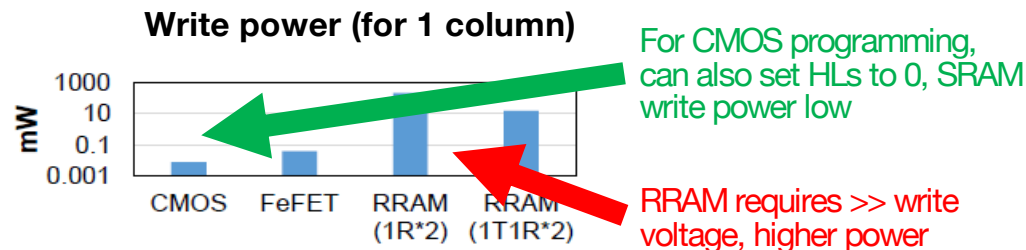
# FeFET-based crossbars: benchmarking

## Benchmarking assumptions for 64x64 crossbar array

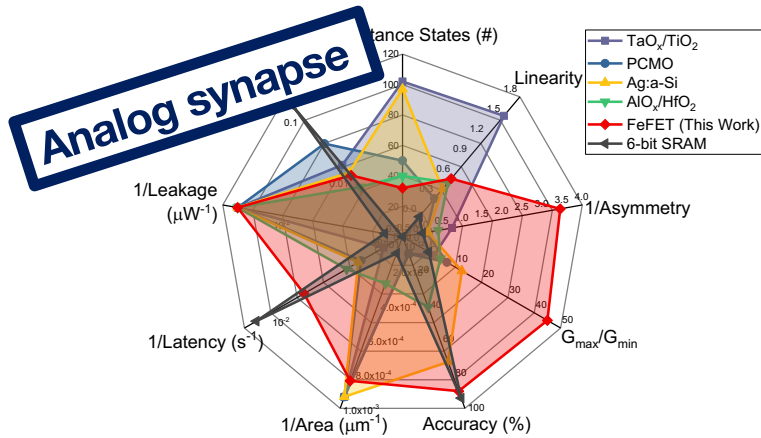
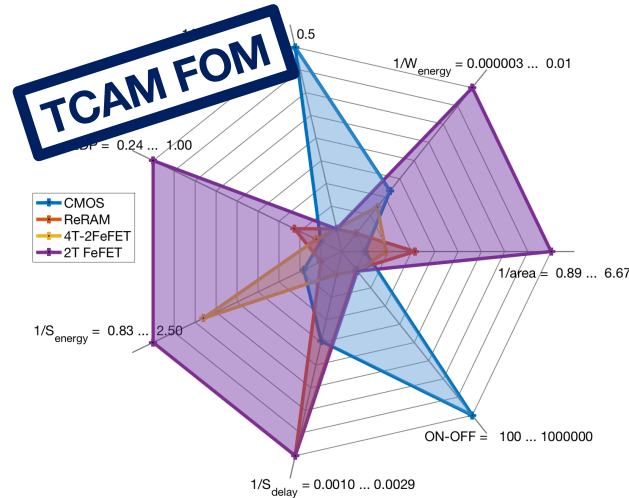
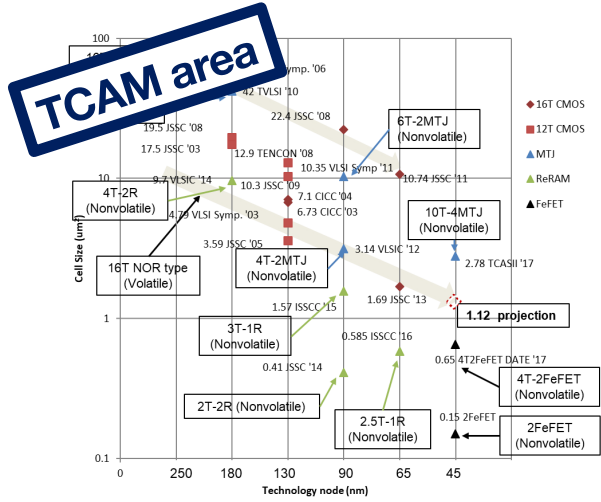
- **FeFET:** 10nm FinFET,  $T_{fe}=10.5\text{nm}$ ,  $V_{WL}=0.6$ ,  $V_W=0.6$ ,  $V_R=-0.55$
- **RRAM:**  $R_{on}=10\text{K}\Omega$ ,  $R_{off}=1\text{M}\Omega$ ,  $V_W=2$
- **For both:**  $V_{HL}=0.3$
- **Average case:** half input bits and half weights are 1

## RC parasitic parameters used in simulations

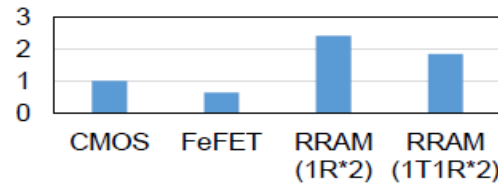
	Cell area ( $F^2$ )	$R_{\text{wire}}$ ( $\Omega$ )	$C_{\text{wire}}$ (fF)
CMOS	150	0.245	0.059
FeFET	60	0.155	0.037
RRAM (1R*2)	4 ( $\times 2$ ) [30]	0.04	0.0096
RRAM (1T1R*2)	20 ( $\times 2$ ) [30]	0.09	0.022



# Takeaways ... promising metrics!



## Binary CNN (normalized) power delay product



FeFET lowest PDP + NV



***JUMP***

Joint University Microelectronics Program

[www.src.org/program/jump](http://www.src.org/program/jump)



Semiconductor Research Corporation



@srcJUMP