# Materials and Devices Working Group Highlight: Ferroelectrics

May 16, 2023

John D. Baniecki

SLAC National Accelerator Laboratory
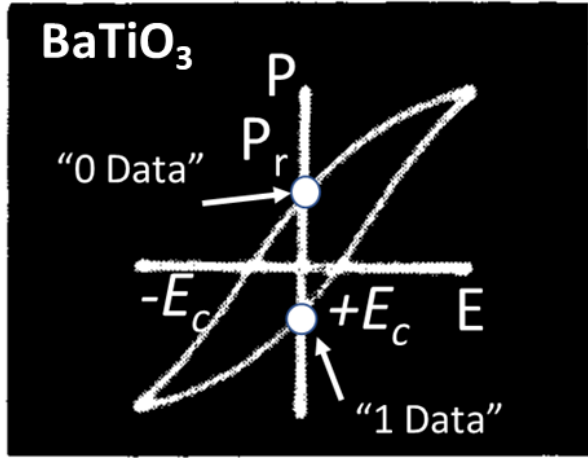
# Outline

❑ Perovskite vs fluorites (each with challenges)

❑ Ferroelectric memories: FRAM, FeFET, FTJs

❑ Comparison of memory technologies

❑ Thermal processing for BEOL integration

❑ Compute-in-memory (CIM) accelerators

SLAC

# Traditional: Perovskite-based (status quo for ~ 60 years)



$ABO_3$

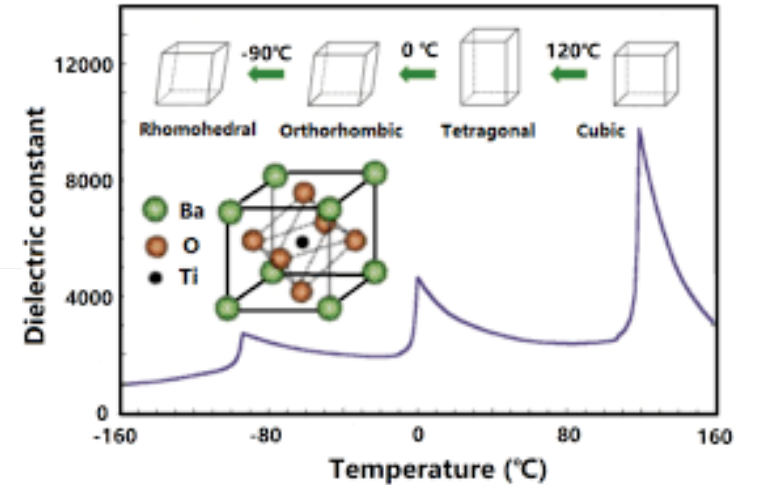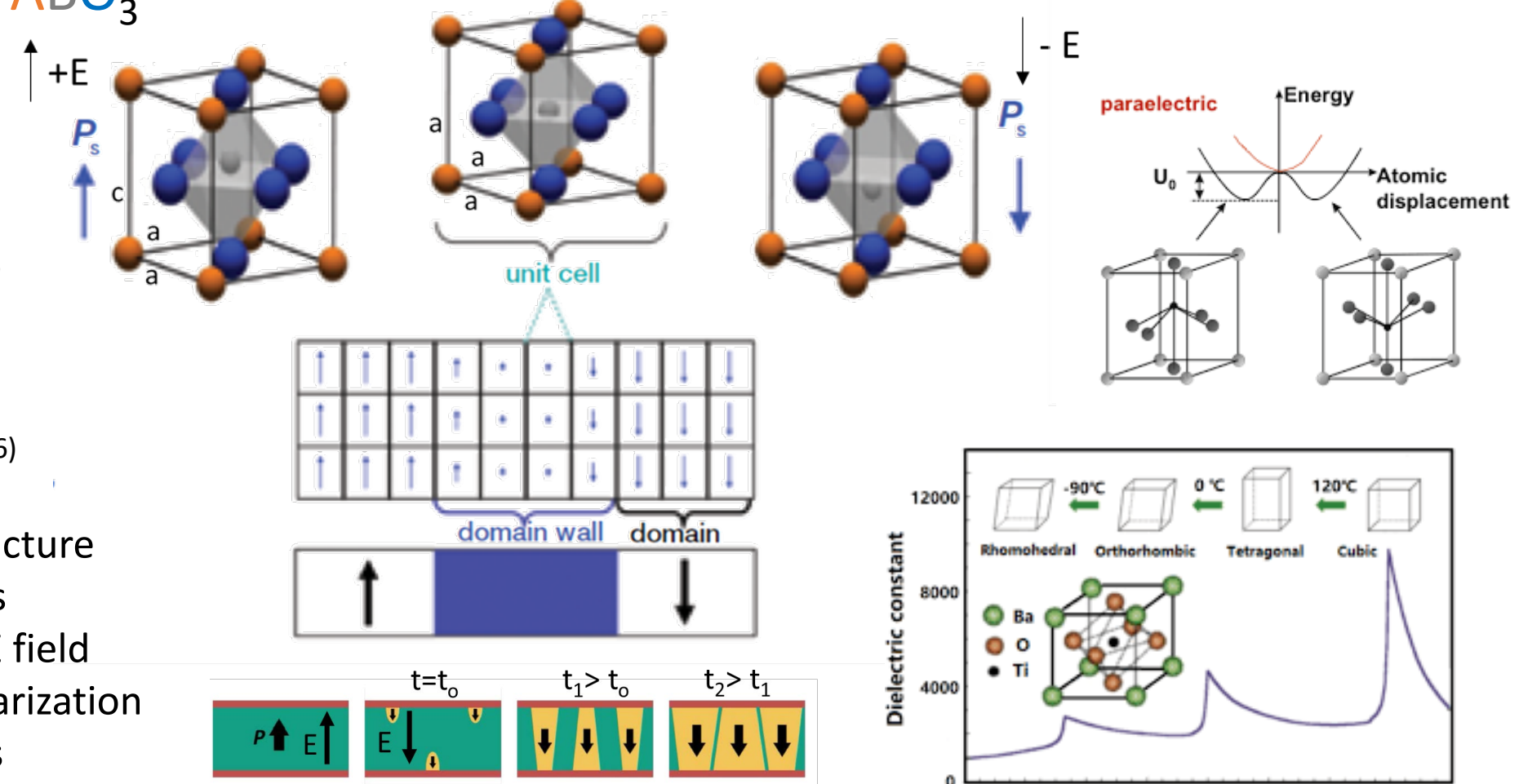S. Trolier-McKinstry et. al., Am. Ceram. Soc. Bull.(2020)

Oscillographic trace

A. Von Hippel et. al., Eng. Chem. (1946)

- Non-centrosymmetric structure
- 2 stable polarization states switchable by an applied E field
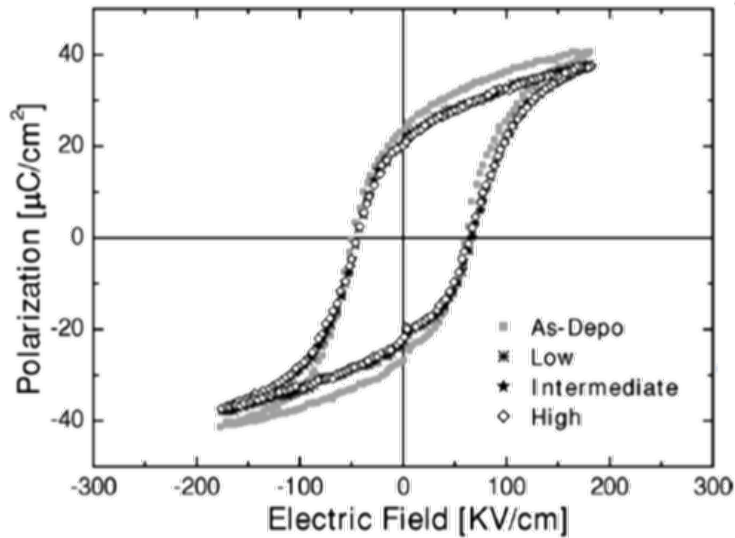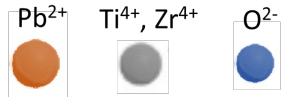- Difference in remnant polarization ($P_r$) through domain states

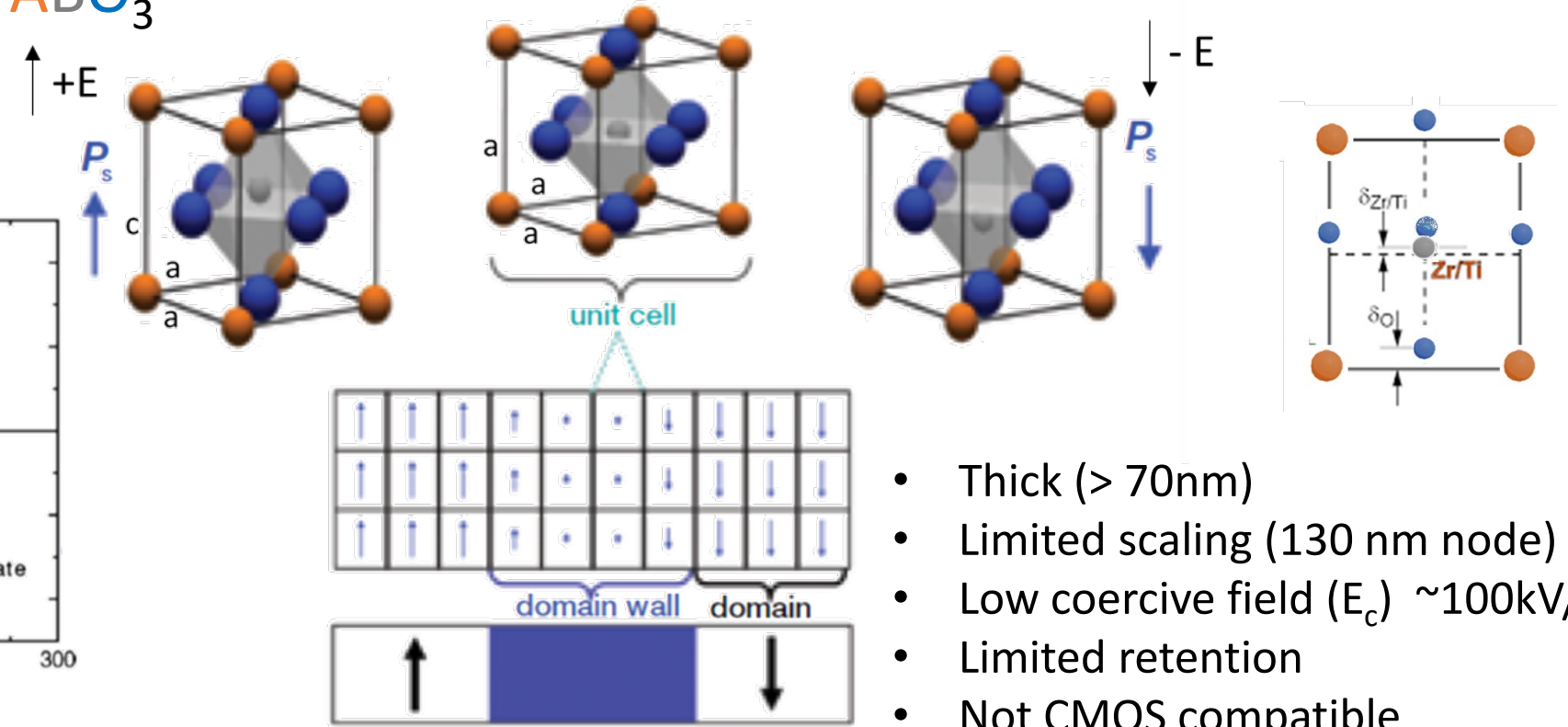J. Su et. al., J Mater Sci: Mater Electron, (2019)

SLAC

# Traditional: Perovskite-based (status quo for ~ 60 years)

PbTiO$_3$-PbZrO$_3$(PZT)

$ABO_3$

S. Trolier-McKinstry et. al., Am. Ceram. Soc. Bull.(2020)

Pb$^{2+}$   Ti$^{4+}$, Zr$^{4+}$   O$^{2-}$



+E          -E

$P_s$                    $P_s$

unit cell

$\delta_{Zr/Ti}$   Zr/Ti   $\delta_O$

domain wall   domain
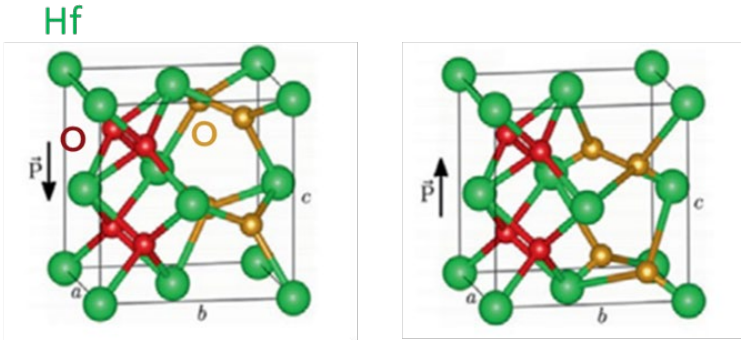
- Thick (> 70nm)
- Limited scaling (130 nm node)
- Low coercive field (E$_c$) ~100kV/cm
- Limited retention
- Not CMOS compatible
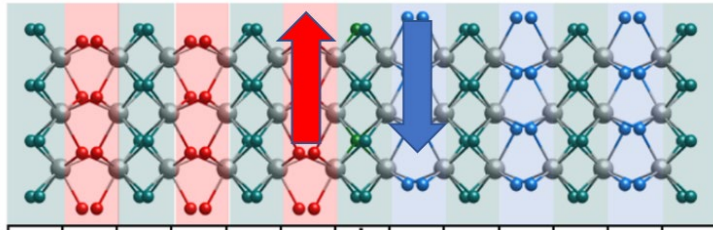
S.H. Choi et. al., Integr. Ferroelectr. (2006)

☐ PZT-based 1T1C memory in production for > 20 years FeRAM (Fujitsu), FRAM (TI), F-RAM (Infineon)

☐ Low density (4KB- 128 MB) niche applications ( IC card, robotic, automotive applications, ...)

# Fluorite-structured ferroelectrics (~2011)

## $HfO_2$-based



Hf

O

O

$\vec{P}$ ↓

$\vec{P}$ ↑

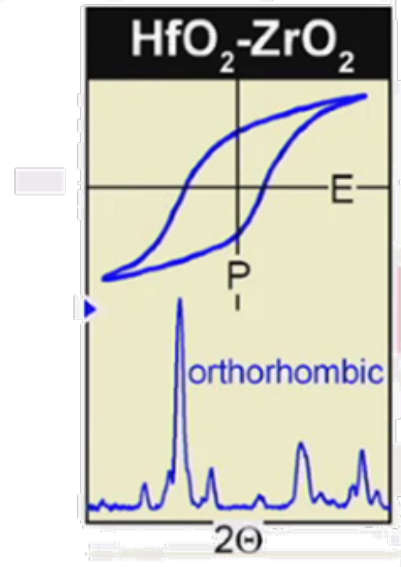Ferroelectric orthorhombic-III ($Pca2_1$)
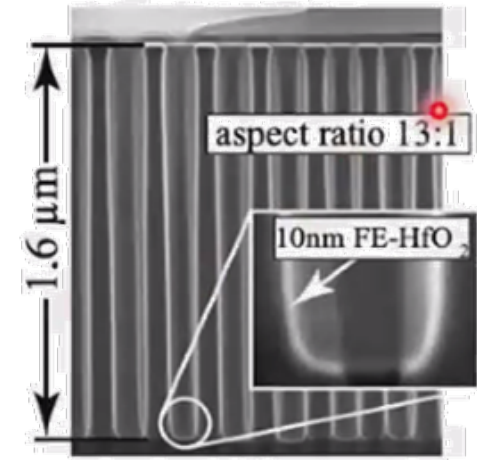(metastable phase)



H.-J. Lee et al., Science (2020)

HfO$_2$ discovered to be ferroelectric in 2006 (Tim Böscke at Qimonda, formerly Infineon), published results in 2011

Böscke, T. S. et al., Appl. Phys. Lett., (2011)

## $HfO_2$-$ZrO_2$ Binary Alloys (HZO)



HfO$_2$-ZrO$_2$

E

P

orthorhombic

2Θ

## Atomic-Layer-Deposition (ALD)



aspect ratio 13:1

10nm FE-HfO$_2$

1.6 μm
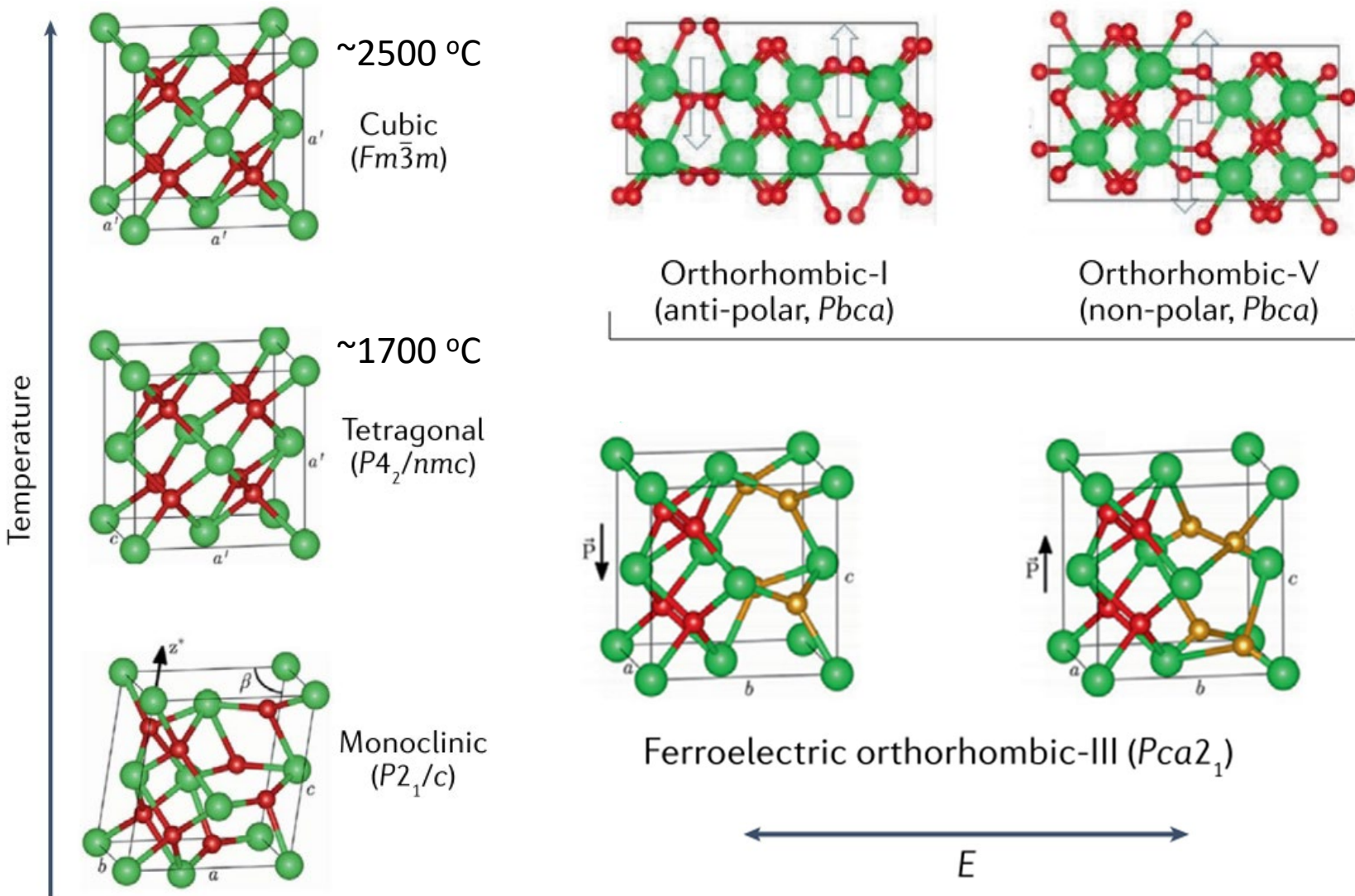
- Thin (sub 3 nm)
- Fast switching  (sub-ns)
- Scaling (22 nm & beyond)
- High  $E_c$  (~1M V/cm)

- Good retention
- CMOS compatible
- ALD growth

SLAC

# Plethora of phases present in HfO$_2$ system



~2500 °C

Cubic
(Fm$\bar{3}$m)

~1700 °C

Tetragonal
(P4$_2$/nmc)

Monoclinic
(P2$_1$/c)

Orthorhombic-I
(anti-polar, Pbca)

Orthorhombic-V
(non-polar, Pbca)

Ferroelectric orthorhombic-III (Pca2$_1$)

E

Temperature
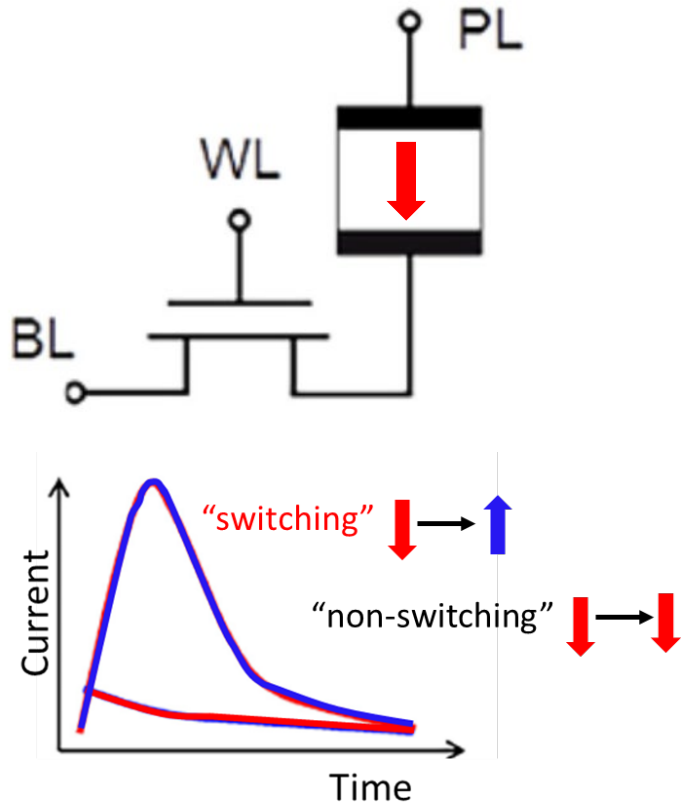
Schroeder, U. *et al.*, *Nat Rev Mater.* (2022).

- Monoclinic phase is the room-temperature bulk stable phase.

- Phases separated by energies of 10s of meV

- Rapid heating and cooling with capping layers (e.g. TiN) stabilizes ferroelectric phase

- Dopants ( Zr, Al, Gd, La, Si, Sr, and Y )

- Field induced transformations

- Reliability issues

SLAC

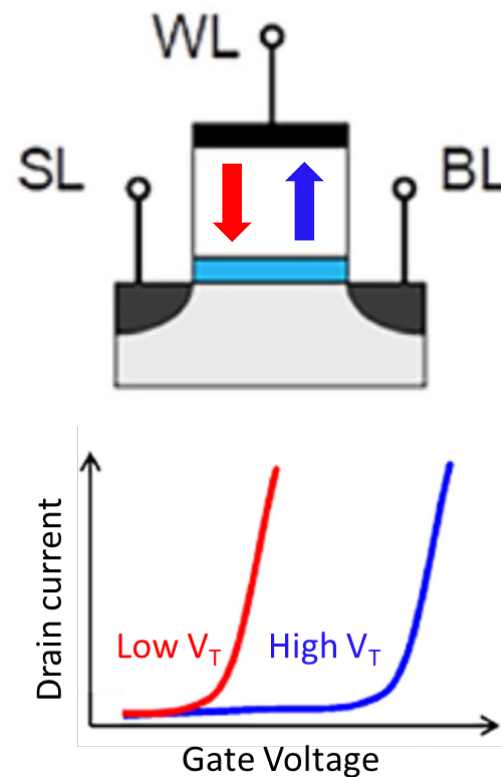# Ferroelectric memories: Operation principles

## FeRAM (1T1C)
Ferroelectric Random-Access Memory
"DRAM-like"



- Commercialized for > 20 years
- Destructive read
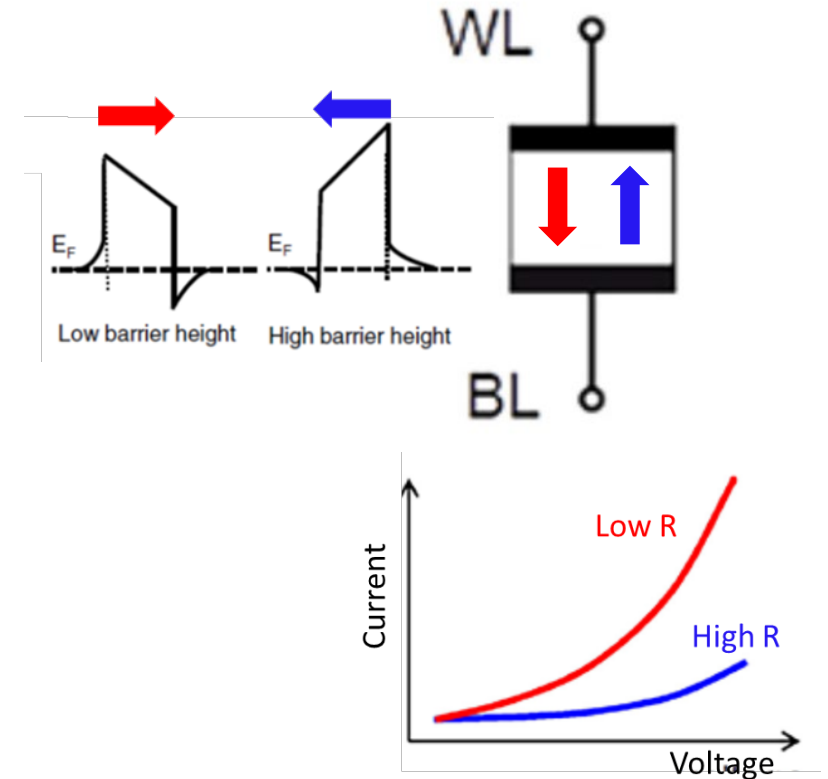- High endurance > $10^{15}$

## FeFET(1T)
Ferroelectric Field Effect Transistor
"FLASH-like"



- Intensive R&D by Semiconductor Industry
- Nondestructive read, multiple bits
- High endurance challenging

## FTJ(1R)
Ferroelectric Tunnel Junction
"Diode-like"
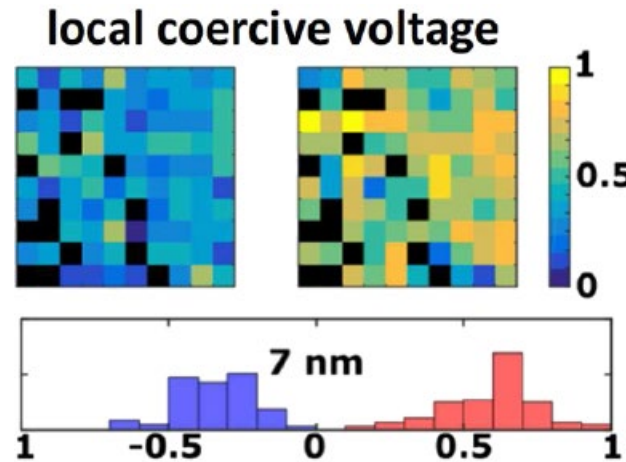


- Academia R&D
- Asymmetric free carrier screening lengths

# Multi-state $HfO_2$-based FeFET memories

**$Hf_{0.5}Zr_{0.5}O_2$ (HZO)** Ferroelectric Field Effect Transistor (FeFET)



I Stolichnov et. al., ACS Appl. Mater. Interfaces (2018)

# Scaled HfO$_2$-based FeFETs

S. Dünkel , IEDM 2017, 22 nm FDSOI CMOS



❑ Demonstrated at the 22/28 nm node

- MW of 1.5 V
- Scaled FeFET cells 0.025 μm$^2$
- Endurance cycles up to 10$^5$

❑ 3D NAND, GAA structures



K. Ni et. Al., T-ED, 2018

**SLAC**  K. Florent et. al., VLSI 2017, S-Y Lee et. al., JEDS 2021

# Challenges of ferroelectric memories

## 1. Polarization variation/Endurance

MFM (> $10^{12}$ cycles)



16 µs    0.16 s    1600 s
breakdown
wake-up
fatigue
1.6 µs

MFIS (<$10^9$ cycles)



Program
Erase
W/L = 20/2 µm

## 2. High Write Voltage

- Reduce write voltage to logic compatible level



| | eSRAM | eDRAM | FeRAM | FeFET |
|---|---|---|---|---|
| Write voltage | <1 V | <1 V | <3 V | <4 V |

## 3. Scaling and Density

Storage capacity



Stochasticity



TiN/HZO/TiN, 1x1 µm²

Multi-bit per cell



## 4. 3D Integration

- Reduced latency, energy consumption
- BEOL process temperature < 400 °C

*BEOL: Back-end-of-line*
*FEOL: Font-end-of-line*



BEOL FeFET

FEOL Access Transistor

ACS Appl. Mater. Interfaces 2018, H. Mulaosmanovic et al, EDL 2018

# Comparison of memory technologies

- ❑ Impact of tech node on energy:
  - $2P_r \sim 60\ \mu C/cm^2$, 1.5 V
- $60F^2$, 130 nm (FRAM perov.) ~ 913 fJ
- $30F^2$, 22 nm FeFET ~ 13 fJ
- Hypothetical $30F^2$, 5 nm ~ 4 fJ

- ❑ Pathways to attojoule switching not clear

- ❑ Advantages over eSRAM in energy efficiency will depend on technology node, compute to standby ratio (application specific)

- ❑ Need to reduce voltage, improve endurance

**PERSPECTIVE**     **NATURE ELECTRONICS**

**Table 1 | Key parameters and metrics**

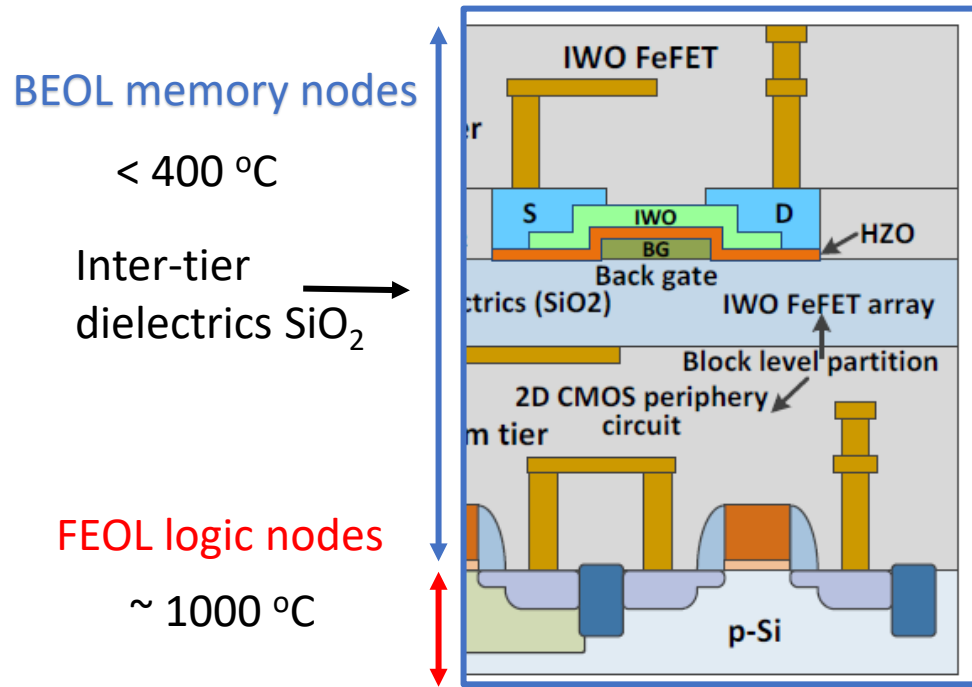| Metrics | Mainstream embedded memory | | | | | Embedded ferroelectric memory | | | |
|---|---|---|---|---|---|---|---|---|---|
| | eSRAM | eDRAM[70] | eFlash (FG) | eFlash (SG MONOS)[42] | eFlash (SONOS)[41] | FEFET (hafnia based, MFIS structure) | FEFET (hafnia based, MFMIS structure) | FRAM (hafnia based) | FRAM (perovskite based)[67] |
| Cell size | 120–150$F^2$ | 40$F^2$ | 40–60$F^2$ | 40–50$F^2$ | 50–60$F^2$ | 10–30$F^2$ | 10–30$F^2$ | 30–40$F^2$ | 50–60$F^2$ |
| Cell structure | 6T | 1T1C | 1.5T | 1.5T | 2T | 1T | 1T1FE, 1T | 1T1FE, 2T2FE | 1T1FE, 2T2FE |
| Non-volatile | No | No | Yes | Yes | Yes | Yes | Yes | Yes | Yes |
| Multi-bit operation | No | No | Yes | Yes | Yes | Yes | Yes | No | No |
| Non-destructive read | Yes | No | Yes | Yes | Yes | Yes | Yes | No | No |
| Status | Av. | Dev. | Av. | Dev. | Dev. | Res. | Res. | Res. | Av. |
| Advanced node demonstration | 7 nm FinFET | 22 nm FinFET | 40 nm | 16 nm FinFET | 28 nm HKMG | 22 nm FDSOI | N/A | N/A | 130 nm |
| Write voltage | <1 V | <1 V | ~12 V | ~12 V | ~5 V | 1.5–4 V | ~1.5 V | 1–3 V | 1.5 V |
| Write energy | ~1 fJ | ~1 pJ | ~100 pJ | ~100 pJ | 1–10 pJ | 1–10 fJ | 1–10 fJ | ~100 fJ | ~1 pJ |
| Standby power | High | Medium | Low | Low | Low | Low | Low | Low | Low |
| Write speed | <1 ns | >10 ns | ~100 ns | <100 ns | ~100 ns | 1–10 ns | 1–10 ns | 1–10 ns | 1 µs |
| Read speed | <1 ns | >10 ns | ~10 ns | <10 ns | ~10 ns | 1–10 ns | 1–10 ns | 1–25 ns | 50–100 ns |
| Endurance | >$10^{16}$ | >$10^{16}$ | ~$10^4$ | ~$10^5$ | ~$10^6$ | $10^5$–$10^9$ | >$10^{10}$ | >$10^{12}$ | >$10^{14}$ |

Key device parameters and performance metrics comparing current embedded memory candidates and ferroelectric technologies. Data for eDRAM, SG MONOS eFlash, SONOS eFlash and perovskite based FRAM are obtained from ref. [70], ref. [42], ref. [41] and ref. [67], respectively. FG, floating gate; SG MONOS, split gate metal–oxide–nitride–oxide–Si; SONOS, Si-oxide–nitride–oxide–Si; eSRAM, embedded static random-access memory; eFlash, embedded flash; eDRAM, embedded dynamic random-access memory; FRAM, ferroelectric random access memory; T, transistor; C, capacitor; FE, ferroelectric; Av., commercially available; Dev., development; Res., research.

SLAC

# Back-end-of-line (BEOL) flash thermal processing

BEOL memory nodes

< 400 °C

Inter-tier dielectrics $SiO_2$ →

FEOL logic nodes

~ 1000 °C



How not to damage BEOL components?

Intense ms pulsed light



Xenon output spectrum

- **Embedded memory in BEOL stack**
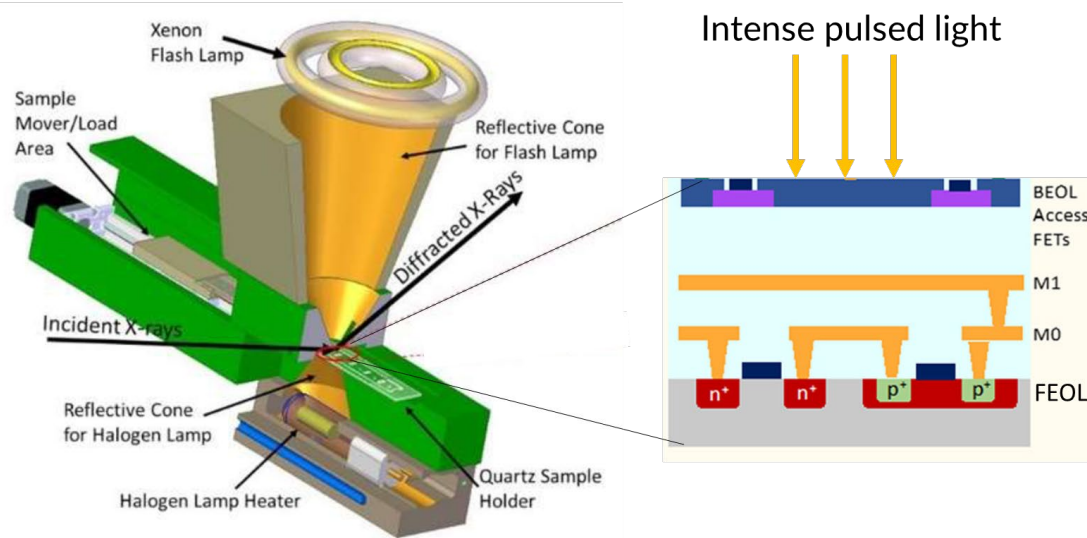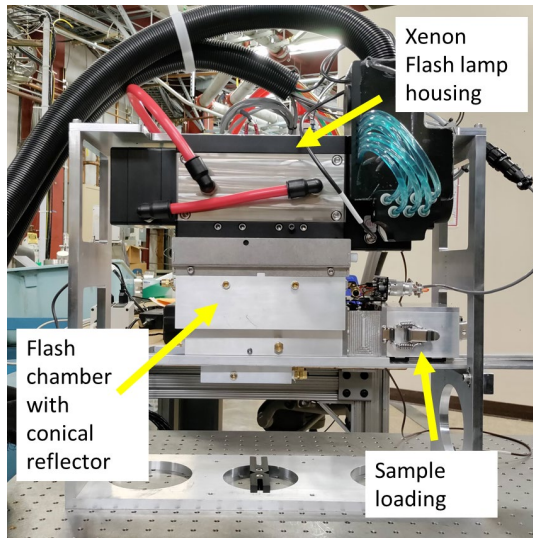
  - Can bring memory and logic closer by stacking memory on top of processing nodes

  - Decrease energy consumption and latency

Dutta, S. *et al.*, *IEEE Electron Device Lett.* **43**, 382–385 (2022).

☐ Confine thermal transients to upper layers

☐ Ultra-fast thermal treatments

# Back-end-of-line (BEOL) flash thermal processing



Xenon Flash lamp housing

Flash chamber with conical reflector

Sample loading



Xenon Flash Lamp

Sample Mover/Load Area

Reflective Cone for Flash Lamp

Diffracted X-Rays

Incident X-rays

Reflective Cone for Halogen Lamp

Halogen Lamp Heater

Quartz Sample Holder

Intense pulsed light

BEOL Access FETs

M1

M0

FEOL

$n^+$  $n^+$  $p^+$  $p^+$

High Brightness In-Vacuum Undulator Beamline

undulator

monochromator

- Static & time resolved studies
- Collect thousands of diffraction patterns during a flash time-temperature sequence
- See transformation in real-time







U.S. DEPARTMENT OF ENERGY | Energy Efficiency & Renewable Energy

ADVANCED MANUFACTURING OFFICE

# Compute-in-memory accelerators

## The Memory Demand of Modern AI Models





Shimeng Yu (Georgia Tech)

- Image classification model
- Langauge/graph model

DLRM-2022 (2TB)
DLRM-2021 (1TB)
GPT-3 (175GB)
DLRM-2020 (100GB)
Turing-NLG (17.2GB)
T5 (11GB)
Megaton-LM (8.3GB)
GPT-2 (1.5GB)
BERT-L(340MB)
AmoebaNet(155MB)
Alexnet (61MB)
VGG-19 (144MB)
SENET (146MB)
ResNet-152 (60MB)
ELMo (94MB)
EfficientNet-B7(155MB)

Number of Parameters (GB)
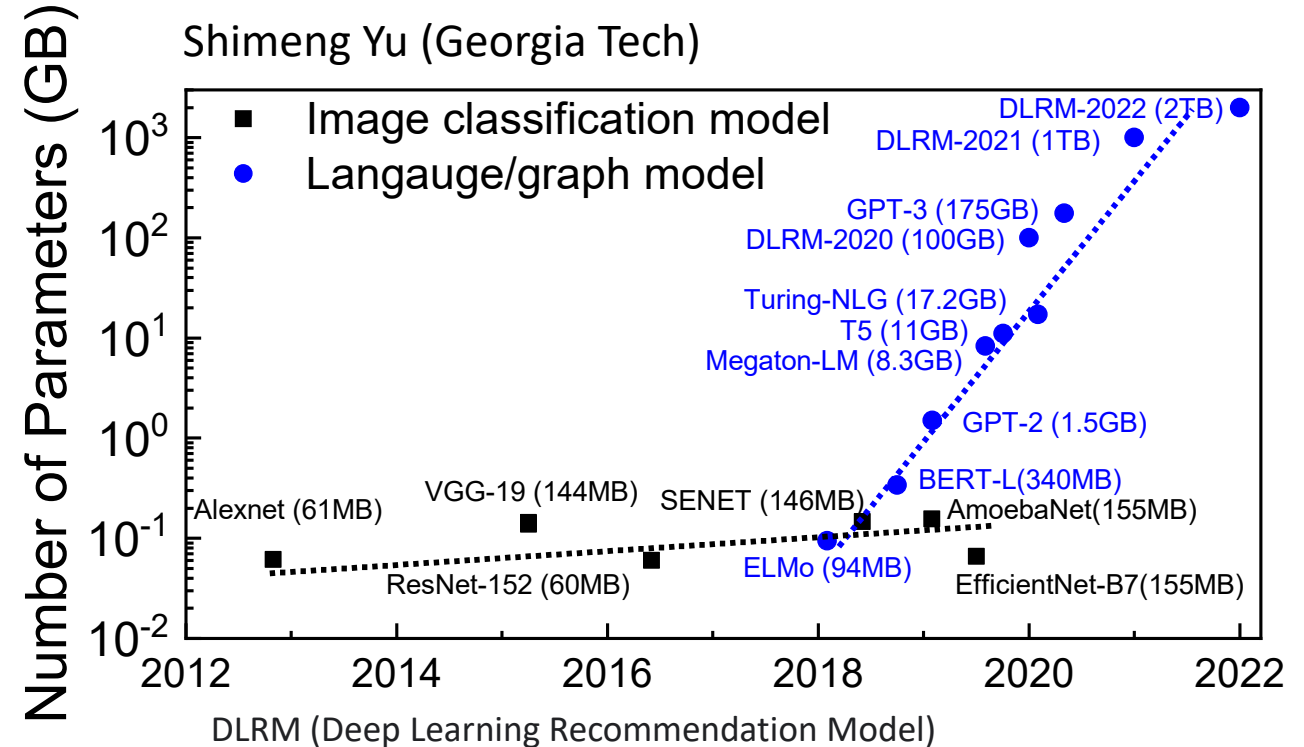
DLRM (Deep Learning Recommendation Model)

❑ Modern AI models can have more than a TB of parameters

❑ With on-chip memory limited by SRAM size, there is an extraordinary volume of data traffic between processor and off-chip memory that adds to energy consumption and latency.

❑ Compute-in-memory (CIM) is a promising approach to overcome memory bottleneck where compute is moved closer to the data residing in the memory

# FeFETS for compute-in-memory accelerators

**Synaptic Weight**

$$I_1 = \sum V_n G_{n1}$$

Input  Output

### Pseudo-crossbar array

$$I_1 = \sum V_n G_{n1}$$

### HfO$_2$-based FeFETs

Conductance Tuning

$V_{th}$ shift

## Advantages

- Nonvolatile conductance tuning
- Low switching energy
- Fast read/write
- ALD w/deeply scaled CMOS nodes

## Challenges

- Endurance (for training)
- High write voltage
- Linearity (for training), stochasticity of conductance tuning
- Density (Legacy nodes, pseudo-crossbar array)
- Area-hungry peripheral circuits (e.g., level shifters (e.g., 45%), high-precision ADC, shift-&-add and buffers)
- How to leverage multi-bit density with peripheral logic scaling (device to system co-optimization)

SLAC

# Monolithic 3D compute-in-memory



Conventional Compute-In-Memory (CIM)
Memory (Synapse) in Front-end
Level Shifter
Logic | MUX | ADC
Neurons & Peripheral CMOS Circuits in Front-end

Monolithic 3D (M3D) CIM
Stacked Memory (Synapse)
BEOL
FEOL
Peripheral CMOS under array (CUA)

BEOL Compatible FeFET Synaptic Weight Cell
M3
BEOL FeFET
M2
M1
FEOL Access Transistor

BEOL FeFET
Top Gate
Drain | $L_g=20nm$ | Source
5nm $HfO_2$
3nm IWO
10nm HZO
10nm | Bottom Gate

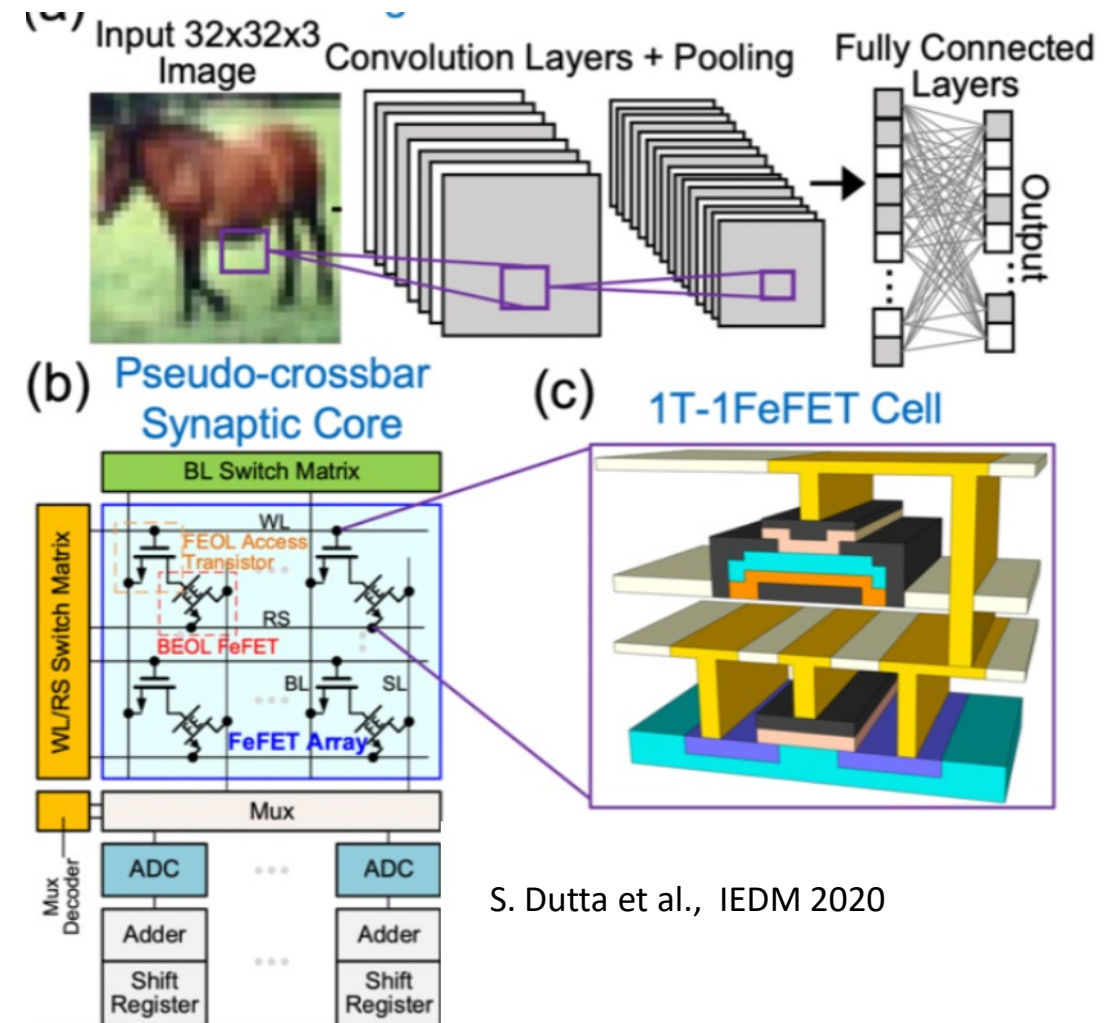$V_{GS}$
$V_{WL}$
$G_{11}$
FE

S. Dutta, (IEDM) 2020.

❑ Memory arrays in the BEOL on top of FEOL CMOS, peripheral circuits

❑ Significant advantage in terms of area, energy and latency
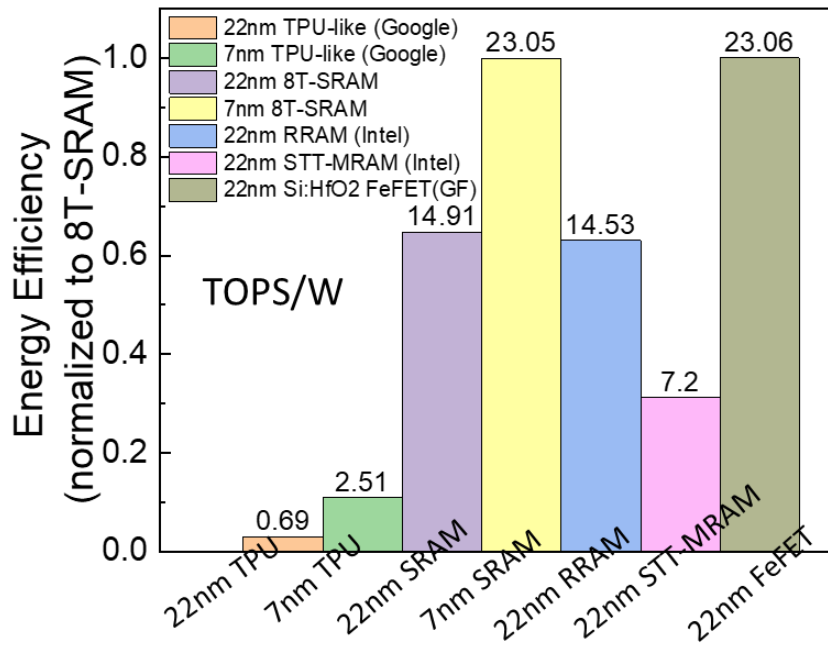
SLAC

# NeuroSim Framework: Benchmarking IMC performance

**Shimeng Yu (Georgia Tech)**   https://github.com/neurosim

❑ Hierarchical simulation framework that covers the device to algorithms to investigate design trade-offs

- Open-source simulator for "in-memory compute" interfaced with PyTorch
- Wide technology choices: SRAM, emerging NVM (RRAM, MRAM, FeFET, etc.) Periphery and interconnect accounted
- Widely used in academia worldwide (>300 registered users)
- Used by industry researchers from SRC/DARPA JUMP sponsors (Intel, TSMC, Samsung, SK Hynix, etc.)

- Validation with IMC prototype with TSMC 40nm RRAM (<2% error)

- VGG-8 model on CIFAR10 dataset (60,000 32x32 color images), with 8-bit weight and 8-bit activation precision.

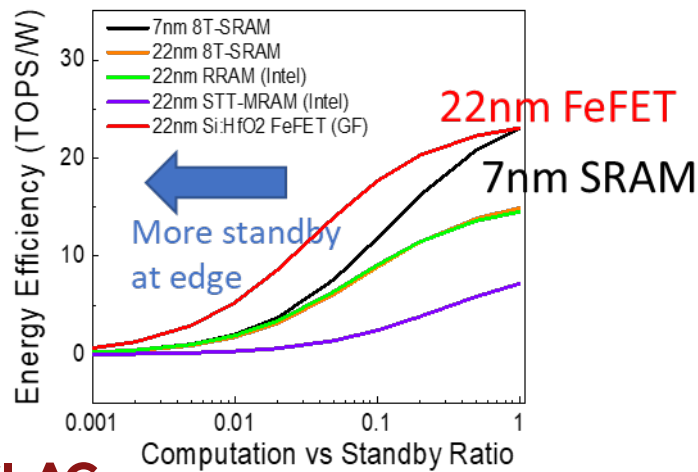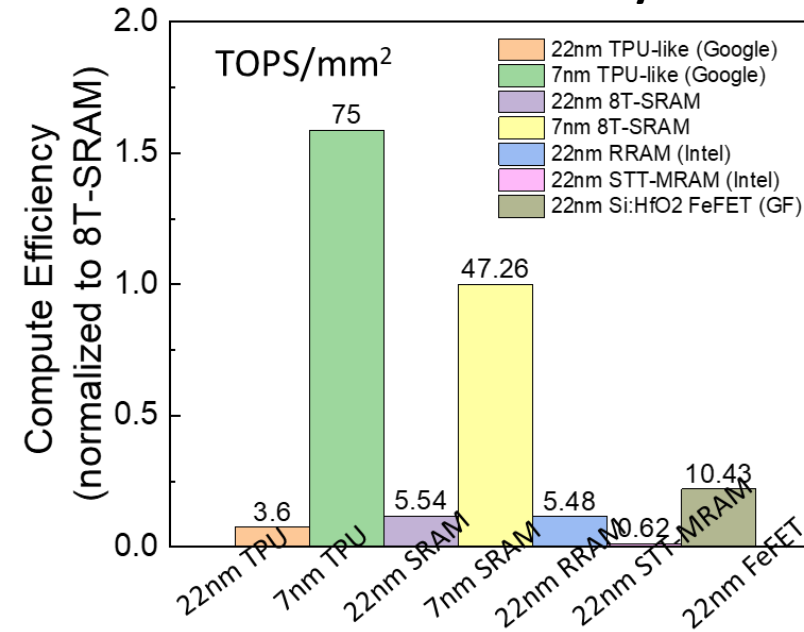  X. Peng et al., IEDM 2019 and 2020   W. Li, et al. CICC 2020

(a) Input 32x32x3 Image   Convolution Layers + Pooling   Fully Connected Layers   Output

(b) Pseudo-crossbar Synaptic Core

(c) 1T-1FeFET Cell

S. Dutta et al., IEDM 2020

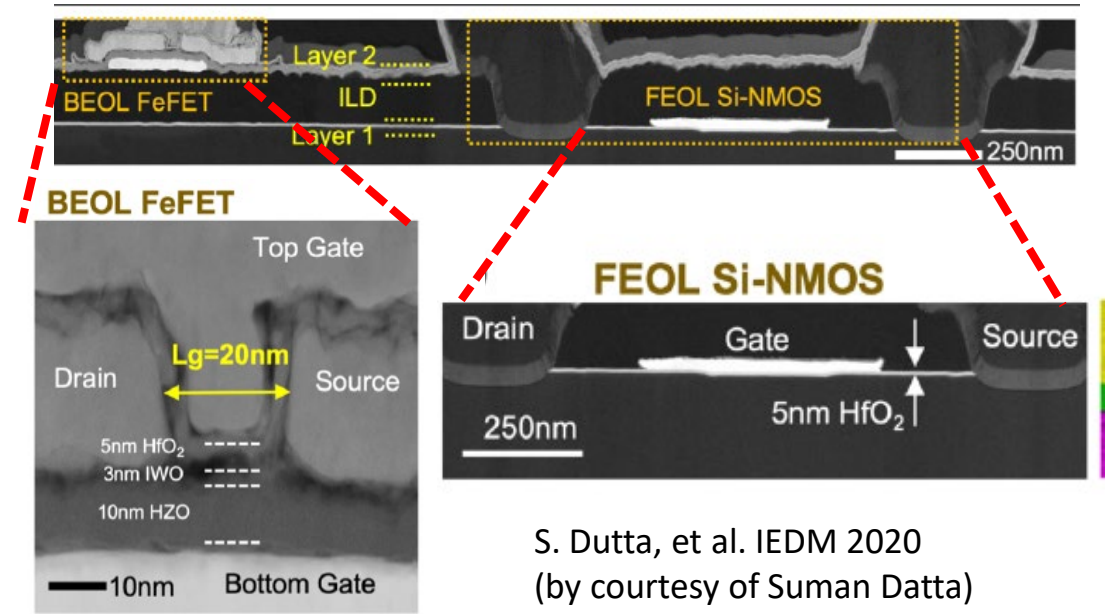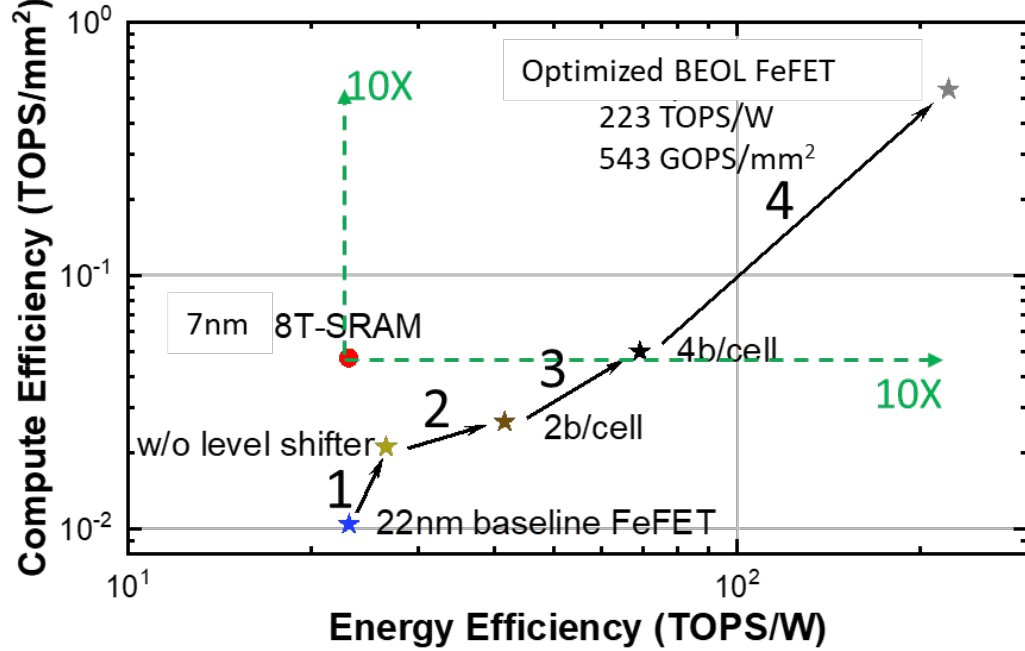# Inference H/W Benchmark Results – TOPS/W & TOPS/mm²

- Limited gain by technology scaling (e.g. TPU), need new approaches: in-memory computing.

- 7nm SRAM TOPS/W is high, but suffers from leakage when the standby is frequent at edge.
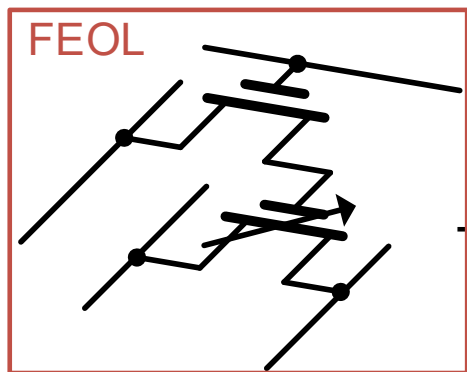- 22nm FeFET design shows superior TOPS/W, thanks to its high resistance (Ron)

# Roadmap of FeFET Improvements

**Side courtesy Shimeng Yu (Georgia Tech)**
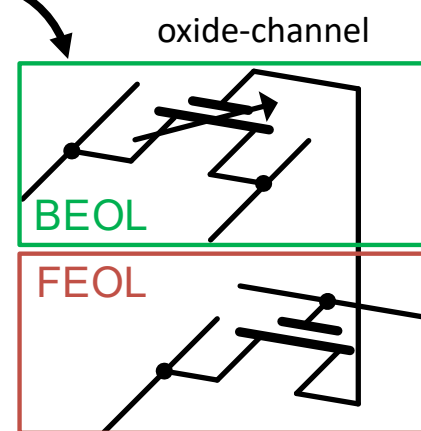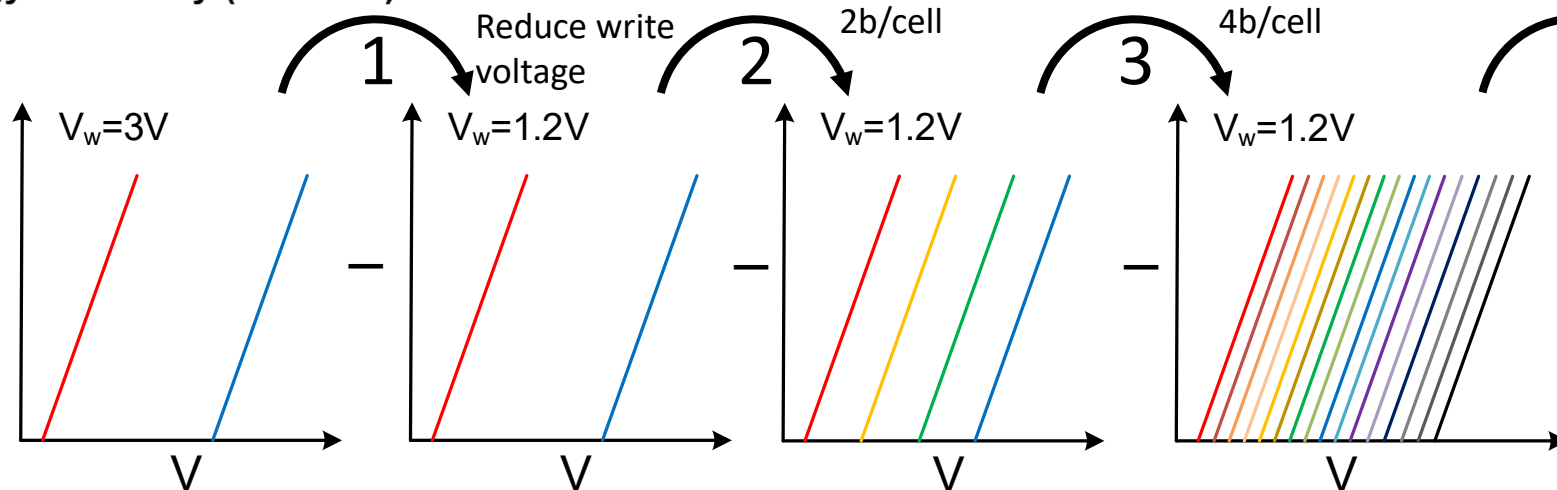
**End goal:** 10X improvement in TOPS/W and GOPS/mm² over 7nm SRAM



Optimized BEOL FeFET
223 TOPS/W
543 GOPS/mm²

7nm 8T-SRAM

w/o level shifter

22nm baseline FeFET

**Compute Efficiency (TOPS/mm²)**

**Energy Efficiency (TOPS/W)**

BEOL FeFET
ILD
Layer 2
Layer 1
FEOL Si-NMOS
250nm

**BEOL FeFET**
Top Gate
Drain
Lg=20nm
Source
5nm HfO₂
3nm IWO
10nm HZO
10nm
Bottom Gate

**FEOL Si-NMOS**
Drain
Gate
Source
250nm
5nm HfO₂

S. Dutta, et al. IEDM 2020
(by courtesy of Suman Datta)

Monolithic 3D Bit Cell at 22nm

FEOL

1 Reduce write voltage    2 2b/cell    3 4b/cell    4

$V_w=3V$    $V_w=1.2V$    $V_w=1.2V$    $V_w=1.2V$

oxide-channel

BEOL

FEOL

SLAC

Distribution Statement "A" (Approved for Public Release, Distribution Unlimited).

# Summary

❑ $HfO_2$-based ferroelectrics offer potential for deeply scaled (22 nm and beyond), low switching energy (~1fJ/bit), non-volatile, fast (sub-ns), multi-bit CMOS compatible memories

❑ FeFET needs to reduce write voltage to logic compatible level, increase cycling endurance, further increase multi-bit per cell, and manage its variability/reliability, particularly in deeply scaled structures

❑ Monolithic 3D integration of BEOL memory and transistors has to overcome thermal processing challenges while maintaining high performance of devices

## Acknowledgement

SLAC