

Computer Architecture: Memory and Storage Technologies

Hossein Asadi (asadi@sharif.edu)

Department of Computer Engineering

Sharif University of Technology

Spring 2024



Copyright Notice

- Some Parts (text & figures) of this Lecture adopted from following:
 - "Computer Architecture", Prof. Onur Mutlu, ETH Zurich, Fall 2022.

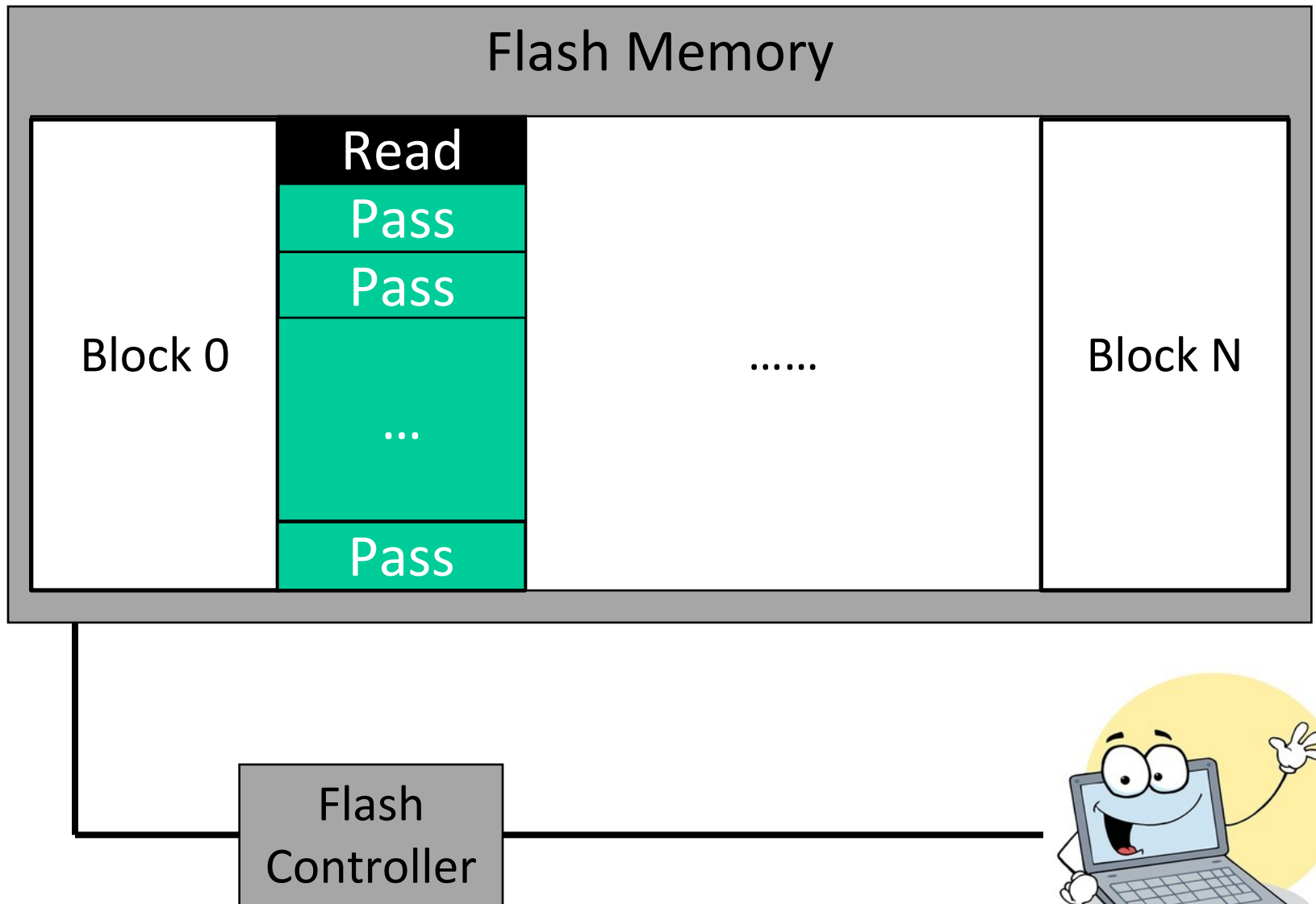


Topics Covered in This Lecture

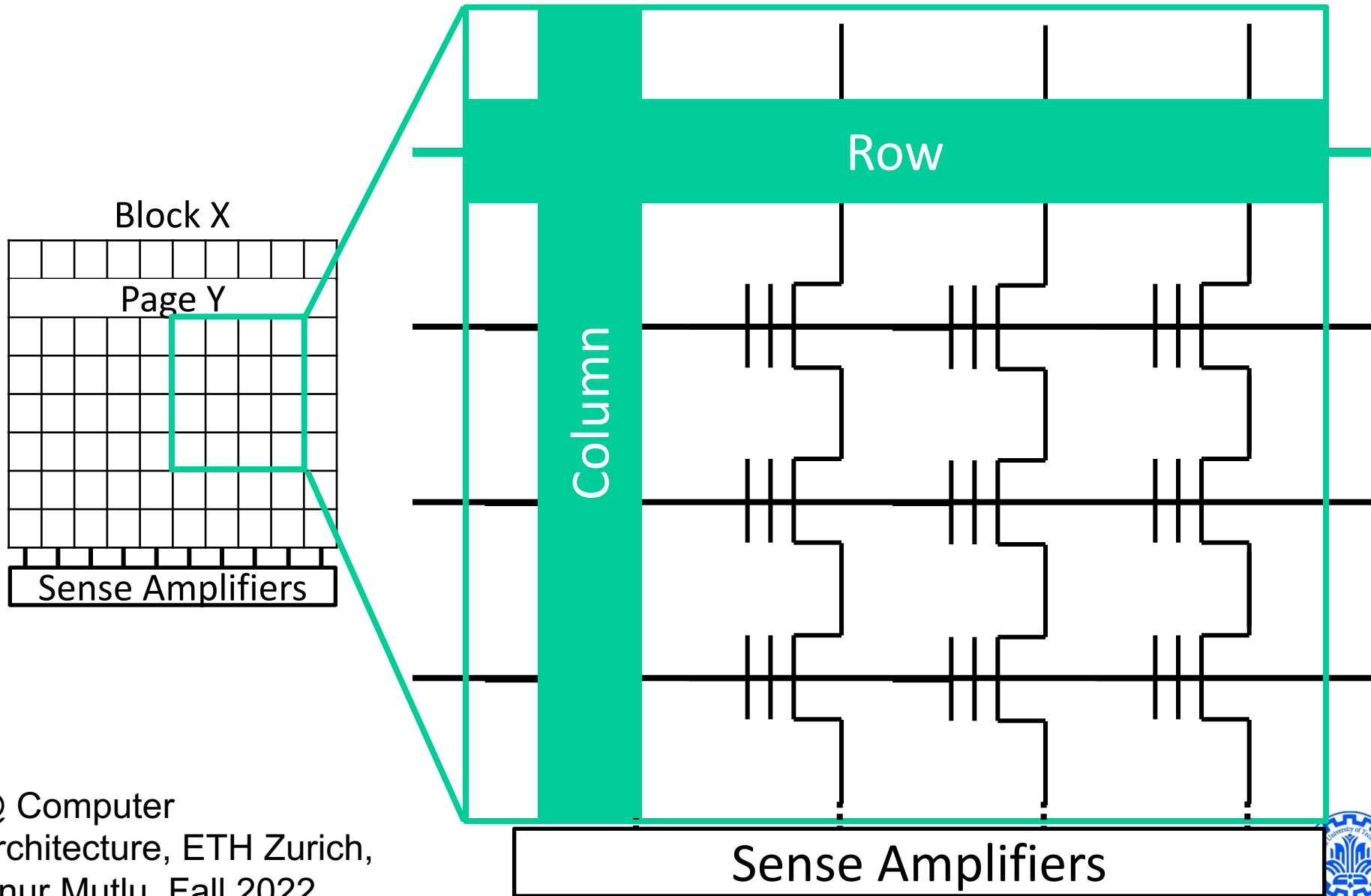
- **Flash Technology**
- **Solid-State Drives**
- **NAND Flash**
- **NOR Flash**
- **Phase Change Memory**



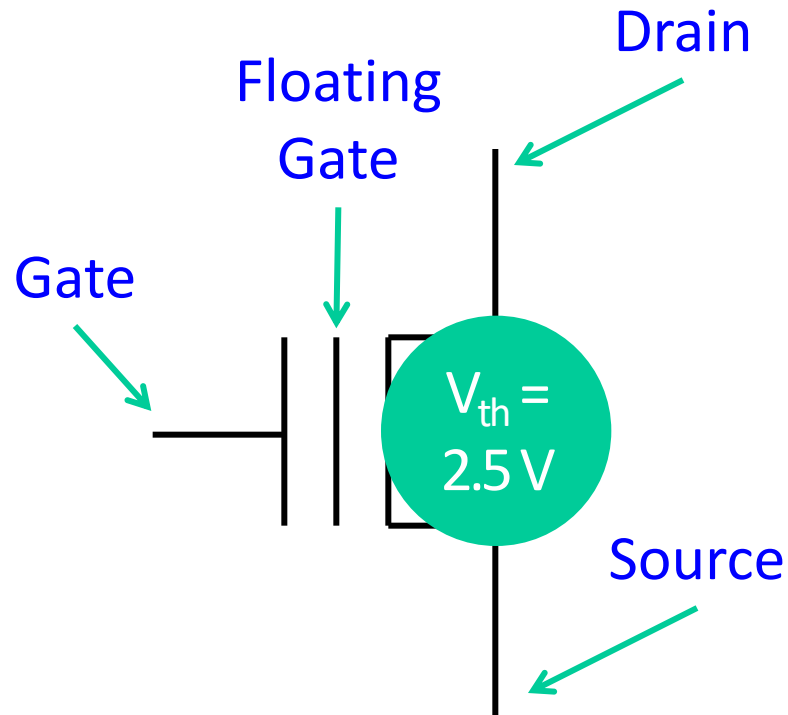
NAND Flash Memory Background



Flash Cell Array

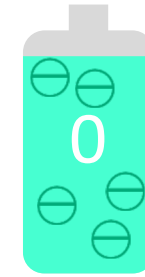
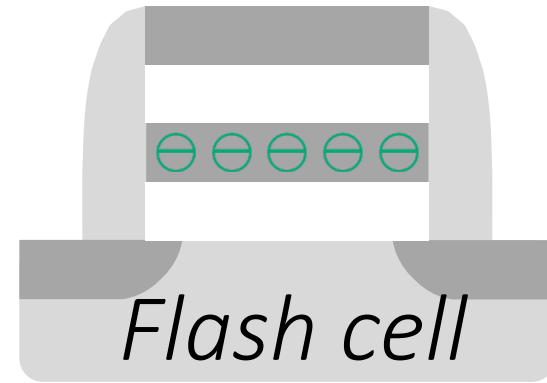
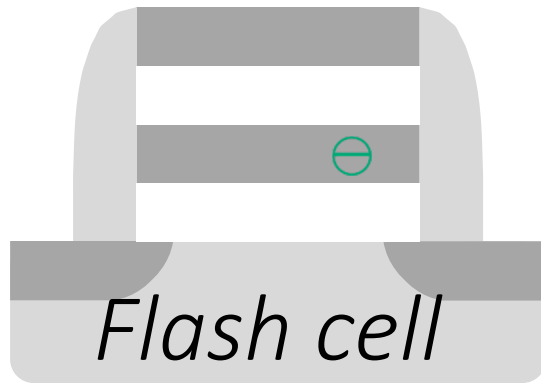


Flash Cell



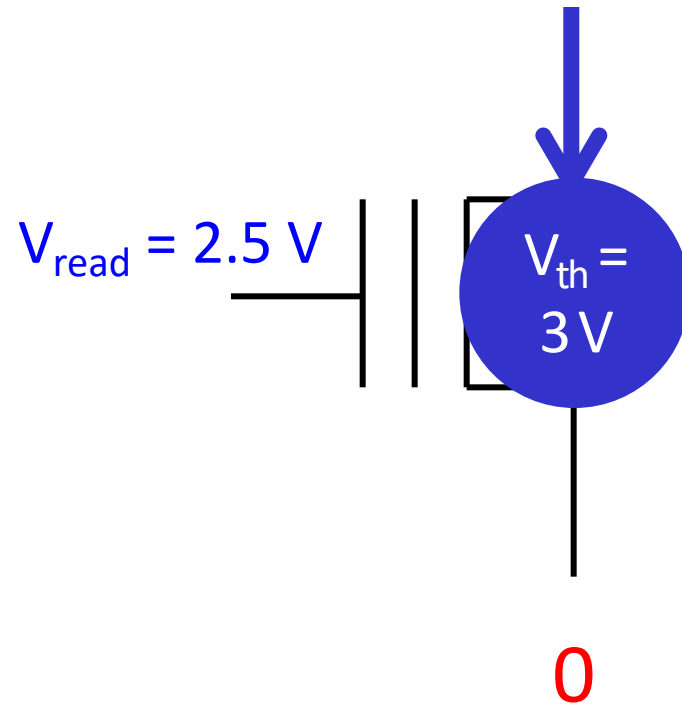
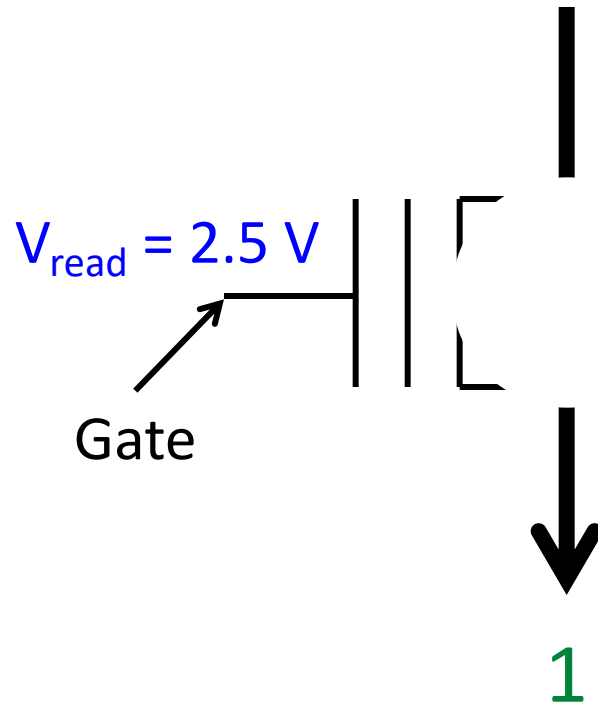
Floating Gate Transistor
(Flash Cell)

Threshold Voltage (V_{th})

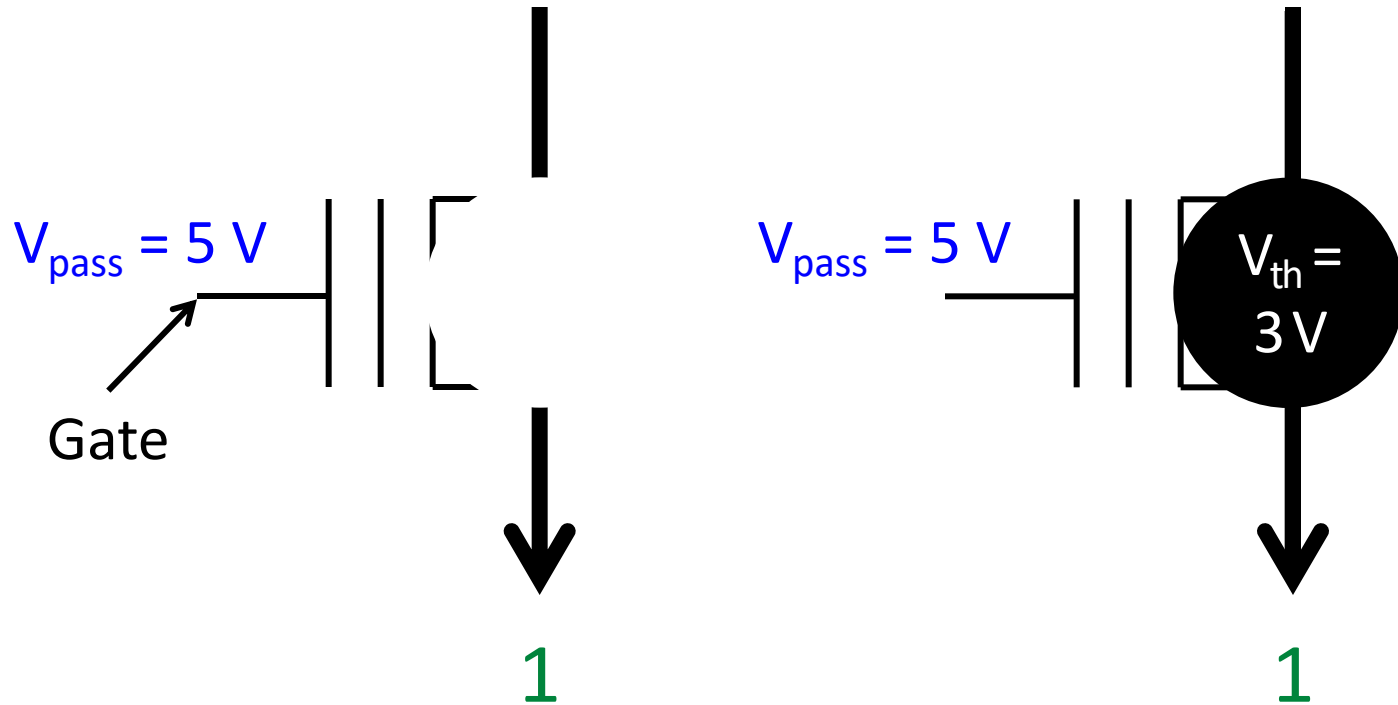


Normalized V_{th}

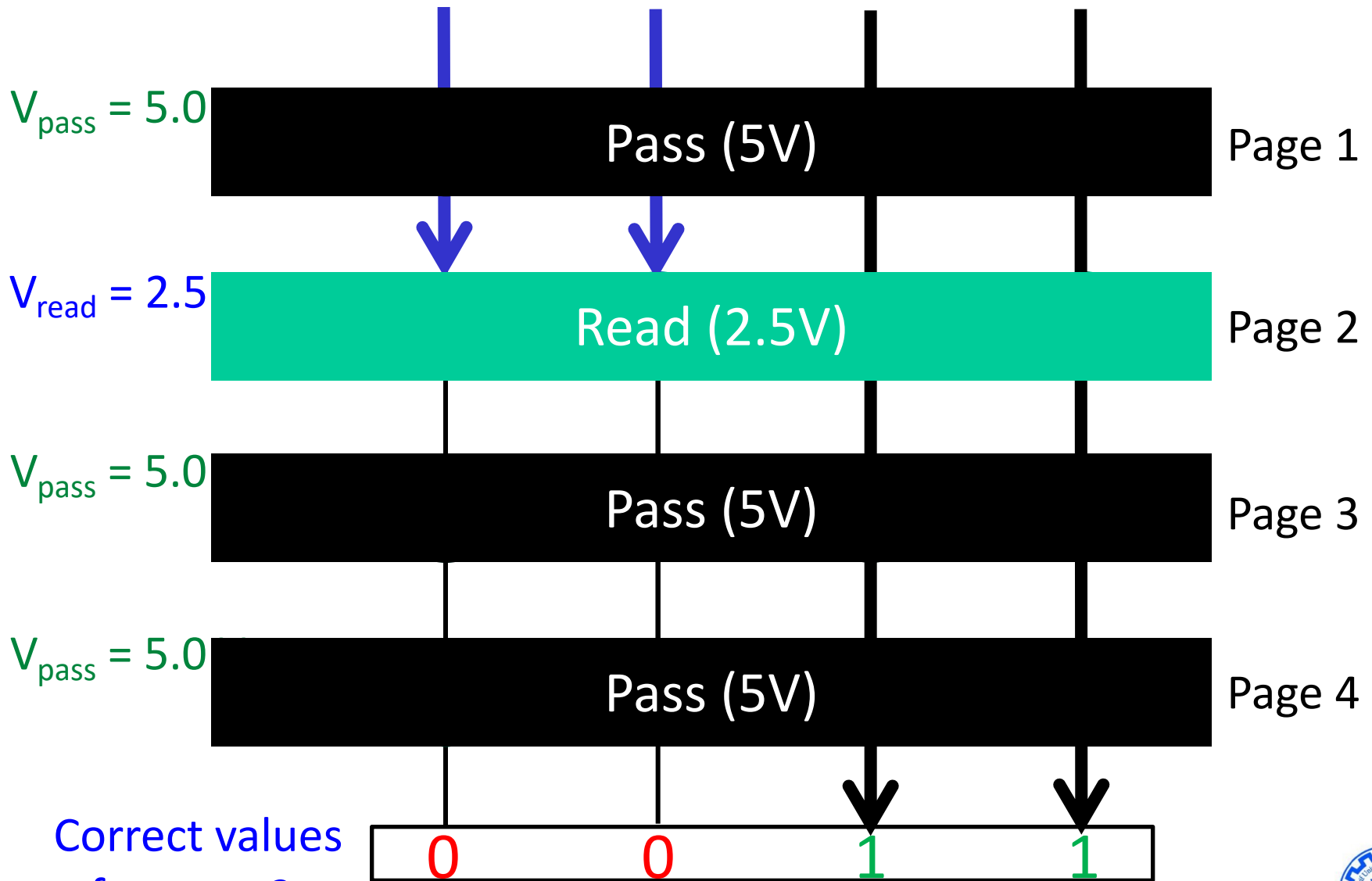
Flash Read



Flash Pass-Through

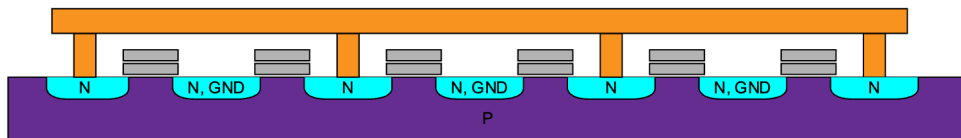
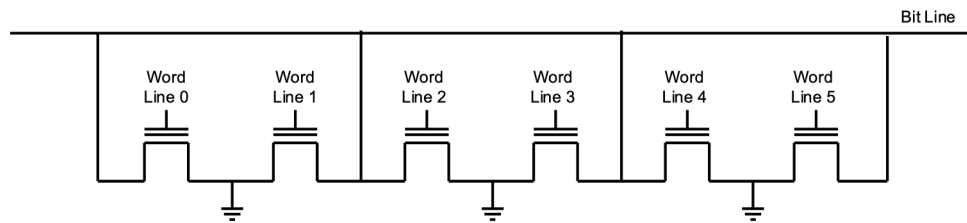
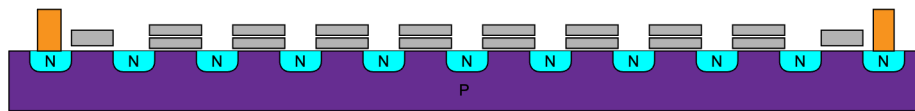
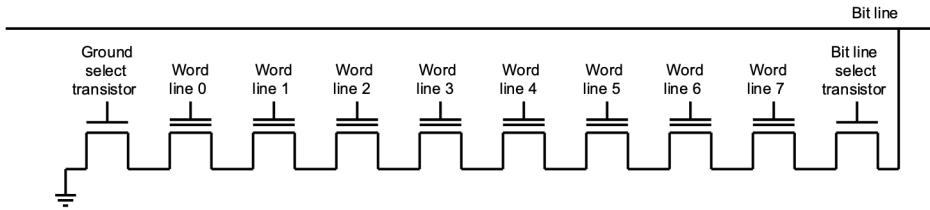


Read from Flash Cell Array

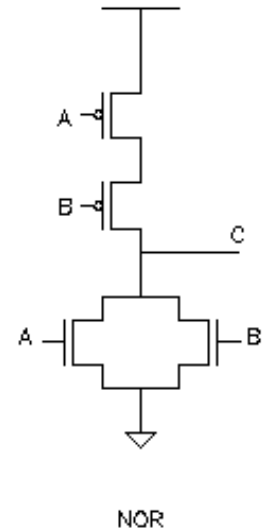
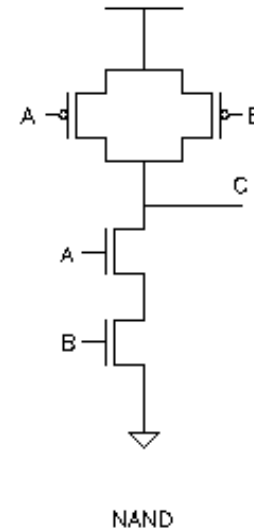


Aside: NAND vs. NOR Flash Memory

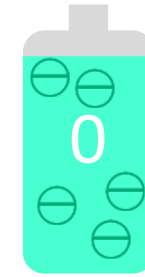
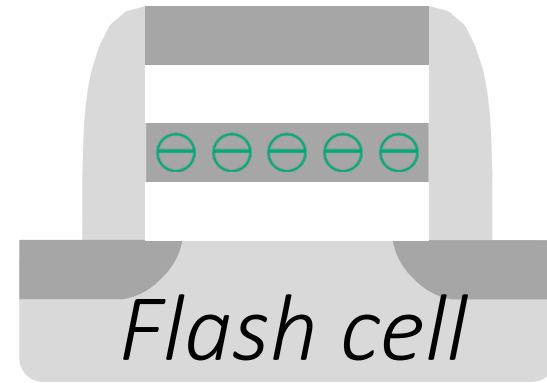
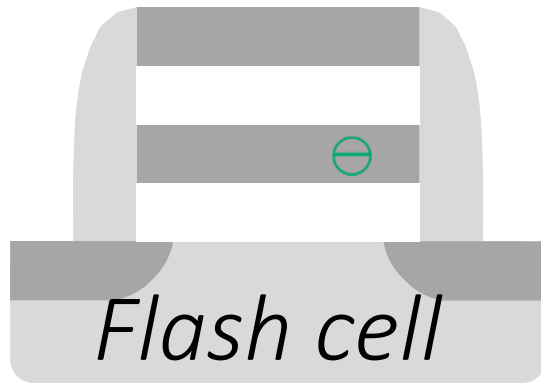
NAND



NOR



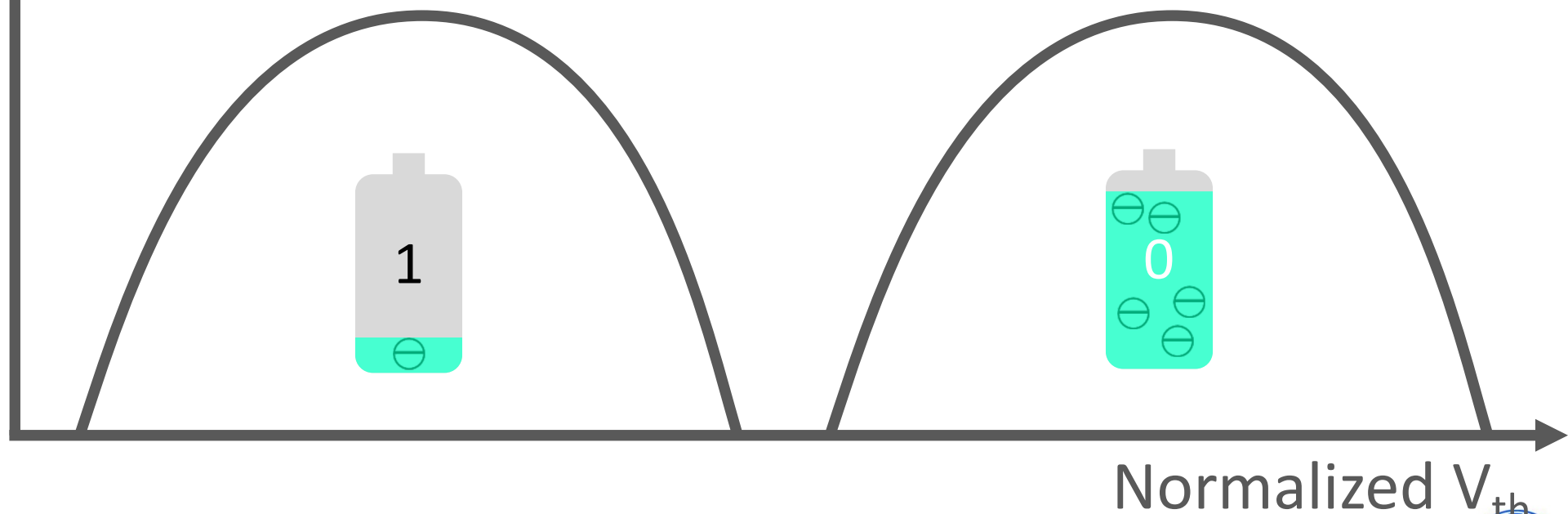
Threshold Voltage (V_{th})



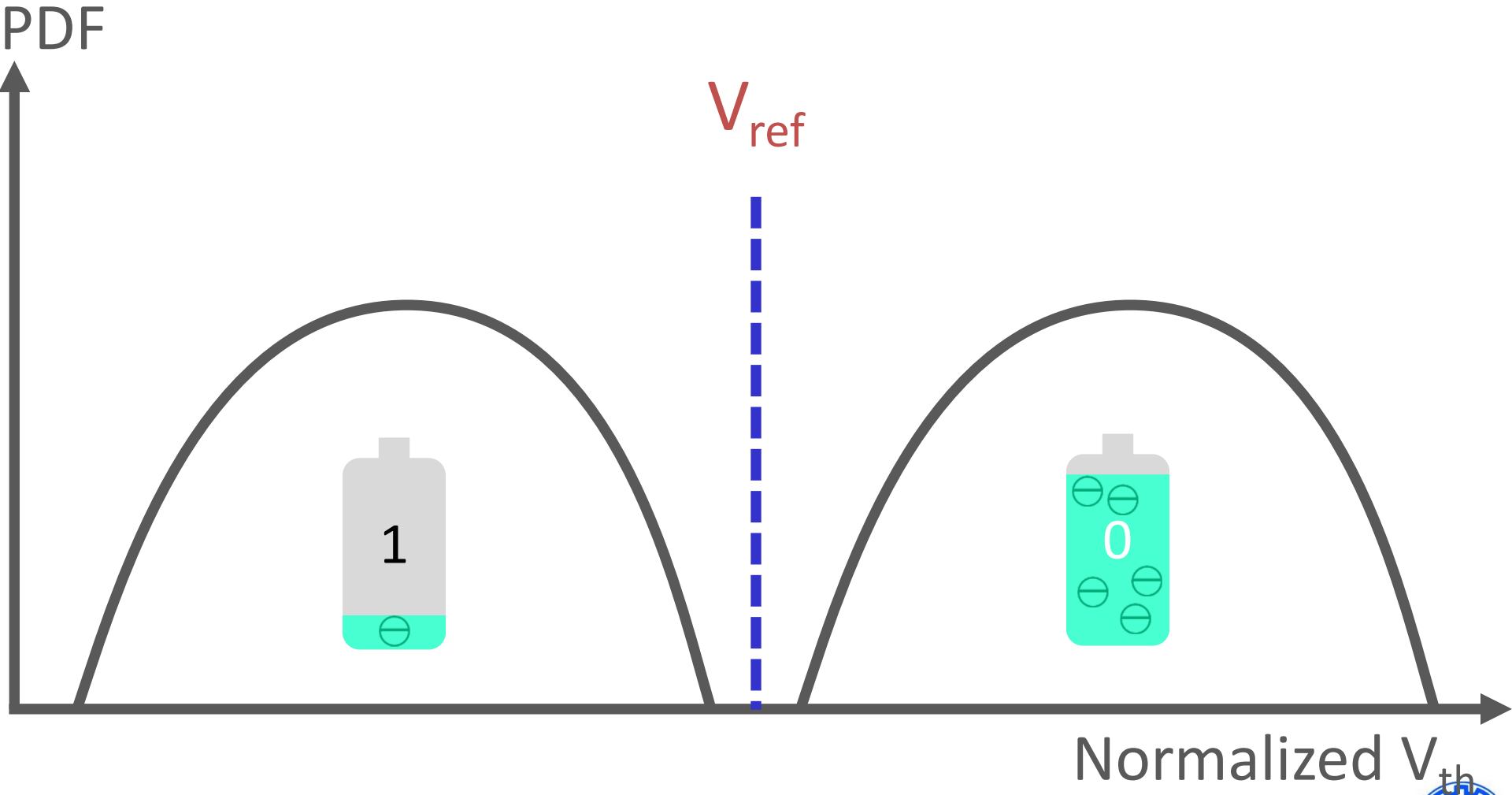
Normalized V_{th}

Threshold Voltage (V_{th}) Distribution

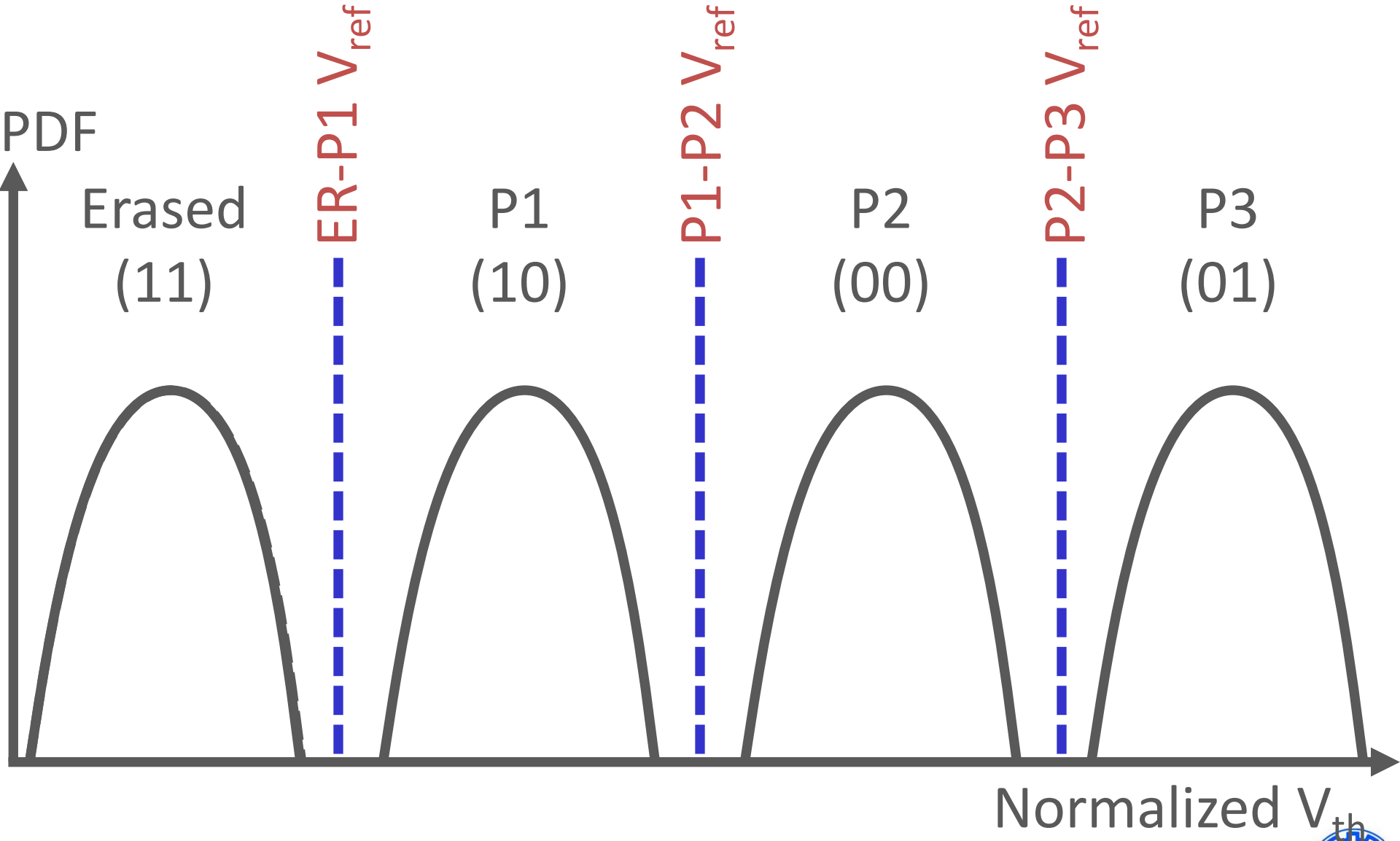
Probability Density
Function (PDF)



Read Reference Voltage (V_{ref})

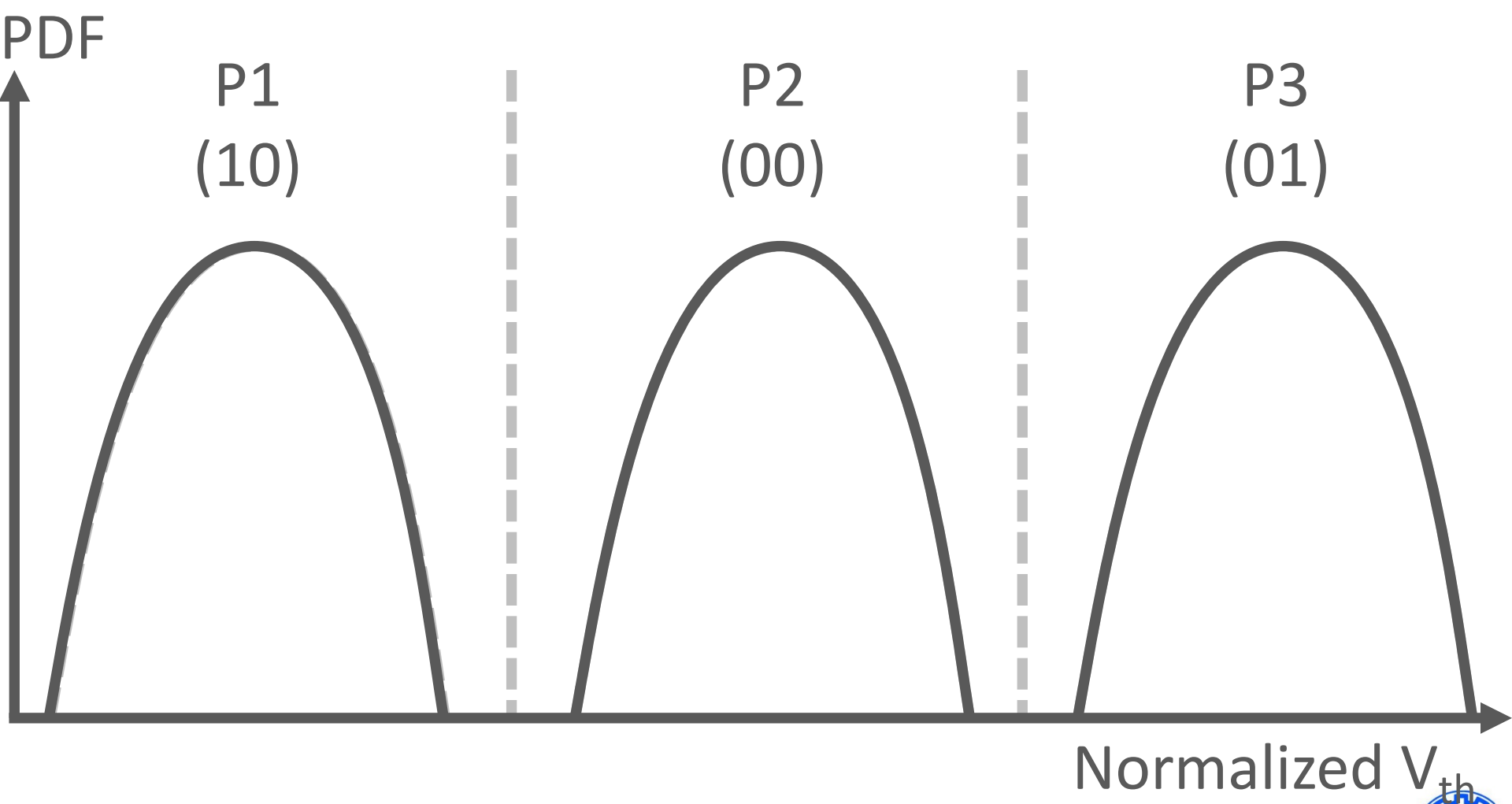


Multi-Level Cell (MLC)



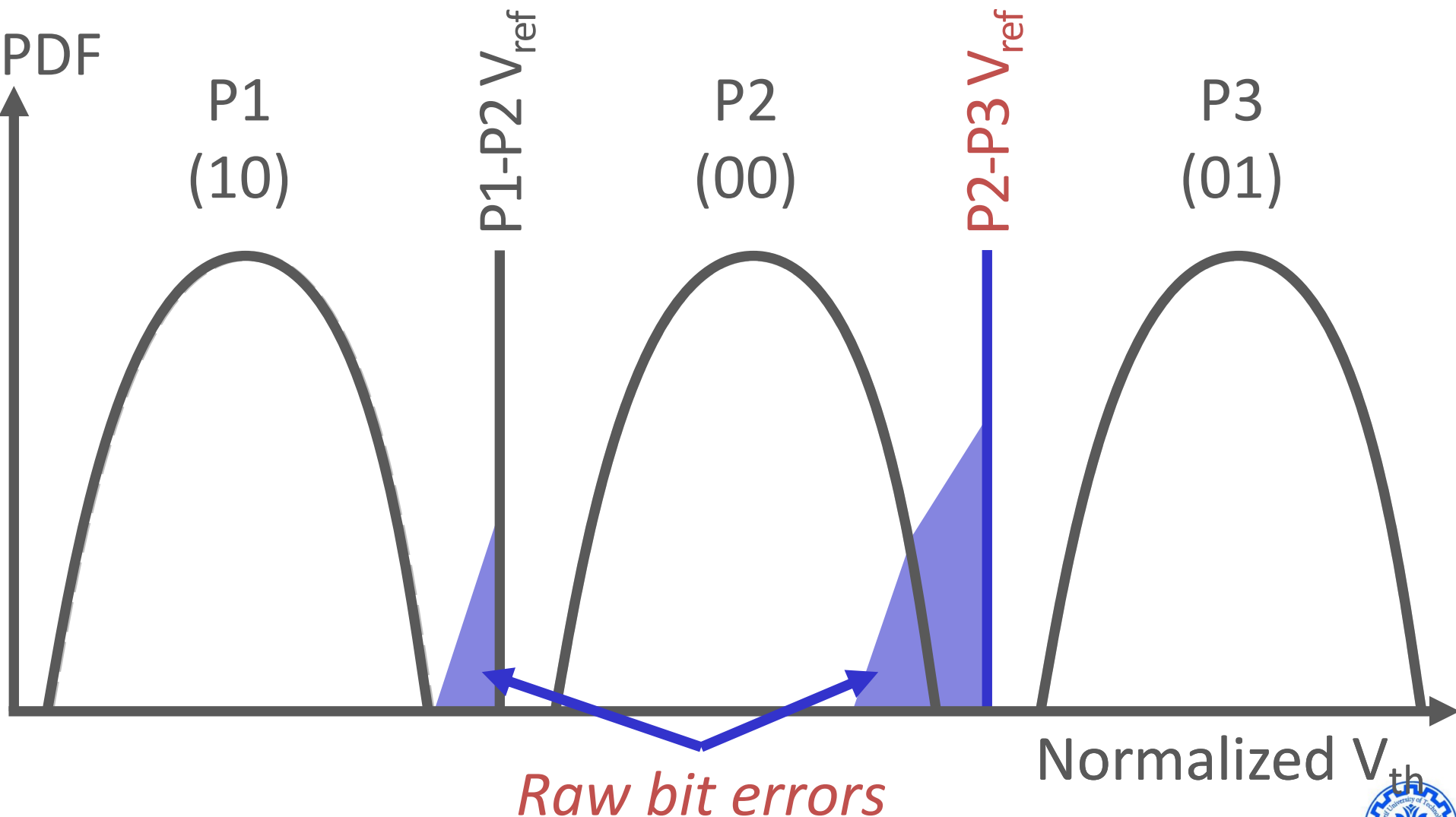
Threshold Voltage Reduces Over Time

After some retention loss:



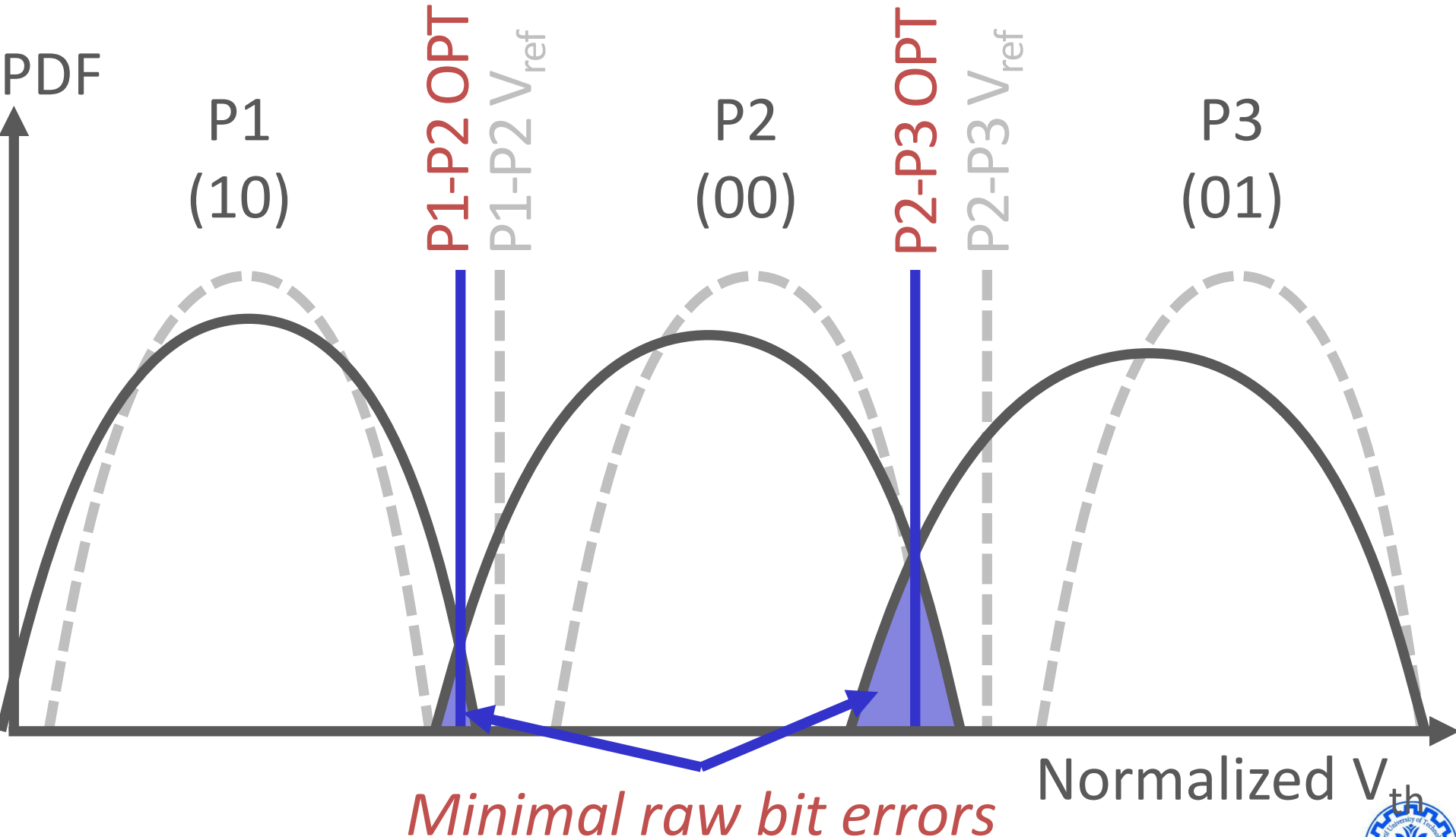
Fixed Read Reference Voltage Becomes Suboptimal

After some retention loss:



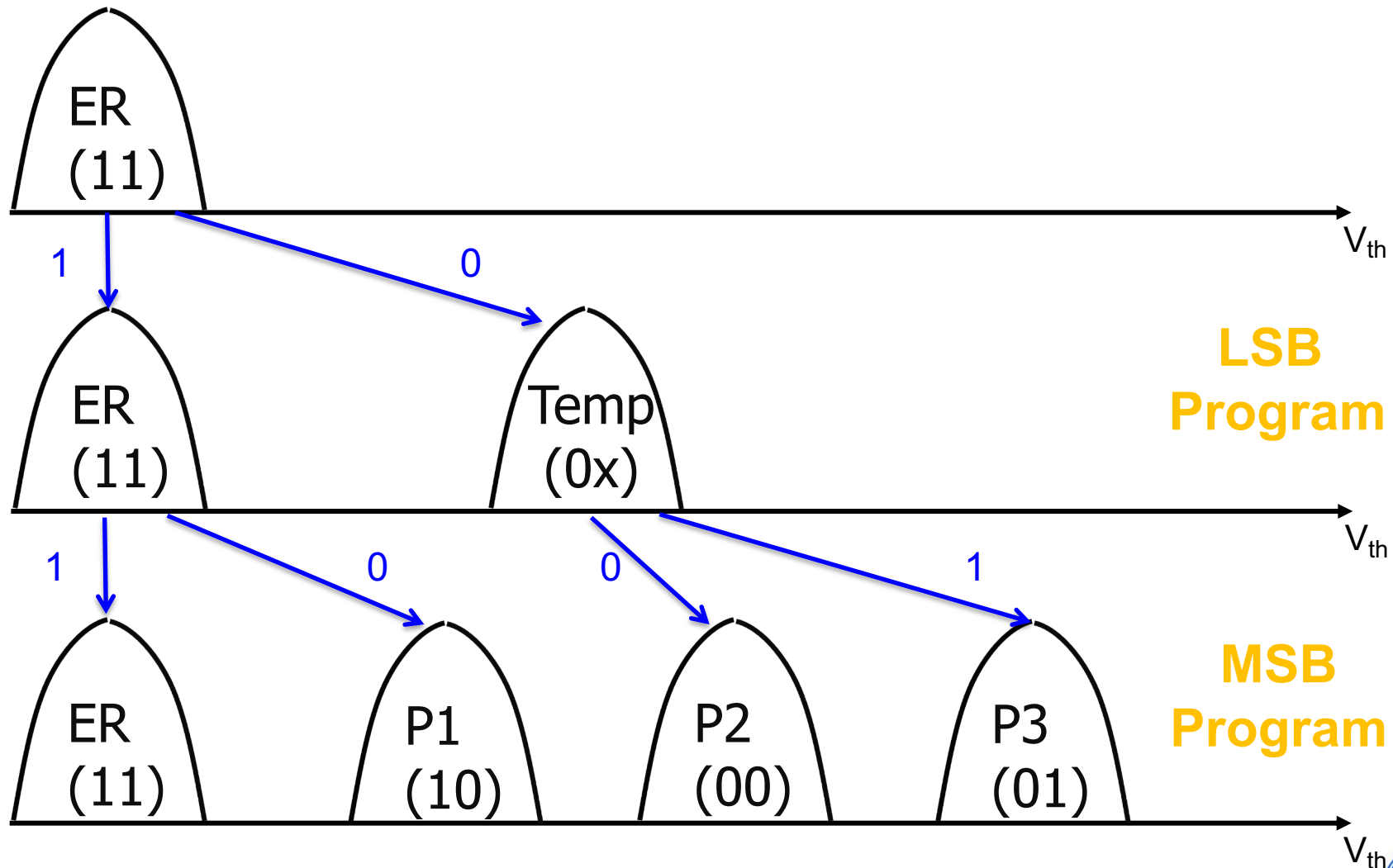
Optimal Read Reference Voltage (OPT)

After some retention loss:

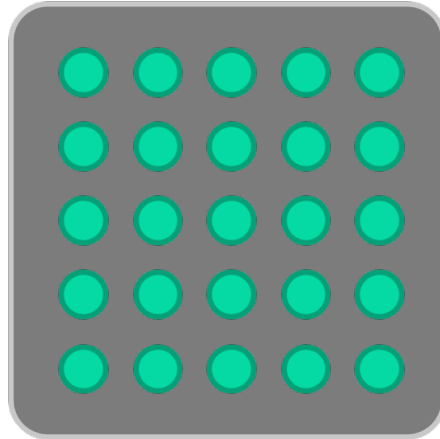


How Current Flash Cells are Programmed

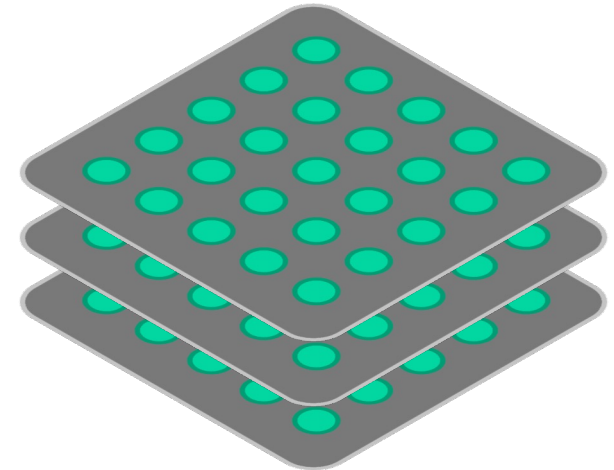
- Programming 2-bit MLC NAND flash memory in two steps



Planar vs. 3D NAND Flash Memory



**Planar NAND
Flash Memory**



**3D NAND
Flash Memory**

Scaling

Reduce flash cell size,
Reduce distance b/w cells

Increase # of layers

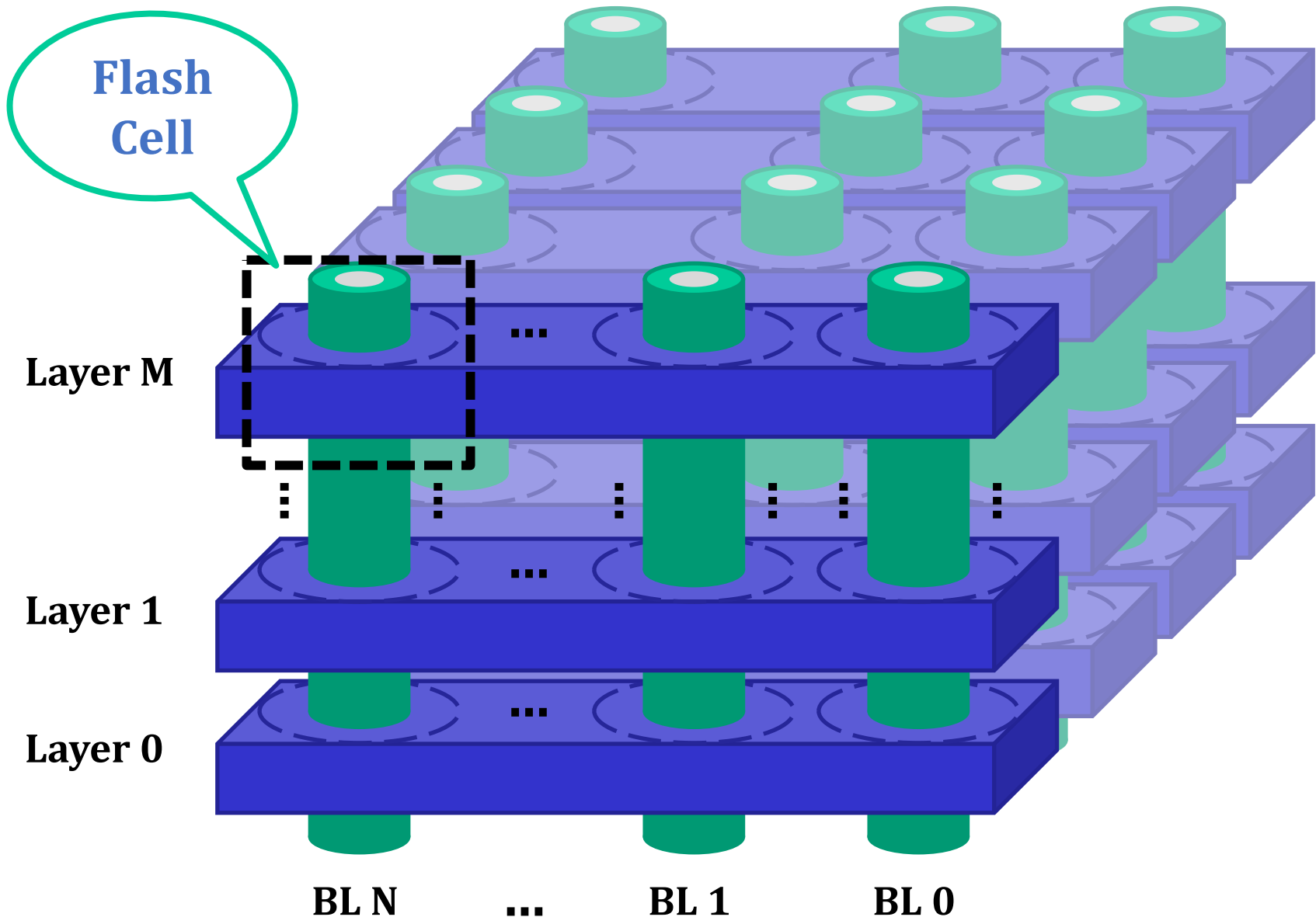
Reliability

Scaling hurts reliability

Not well studied!

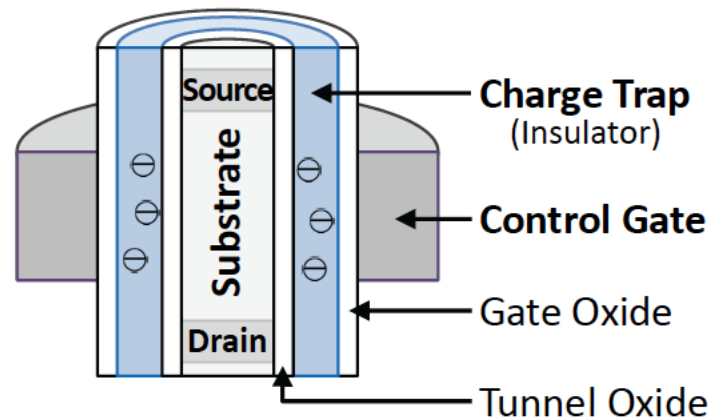


3D NAND Flash Memory Structure



Charge Trap Based 3D Flash Cell

- Cross-section of a charge trap transistor



3D NAND Flash Memory Organization

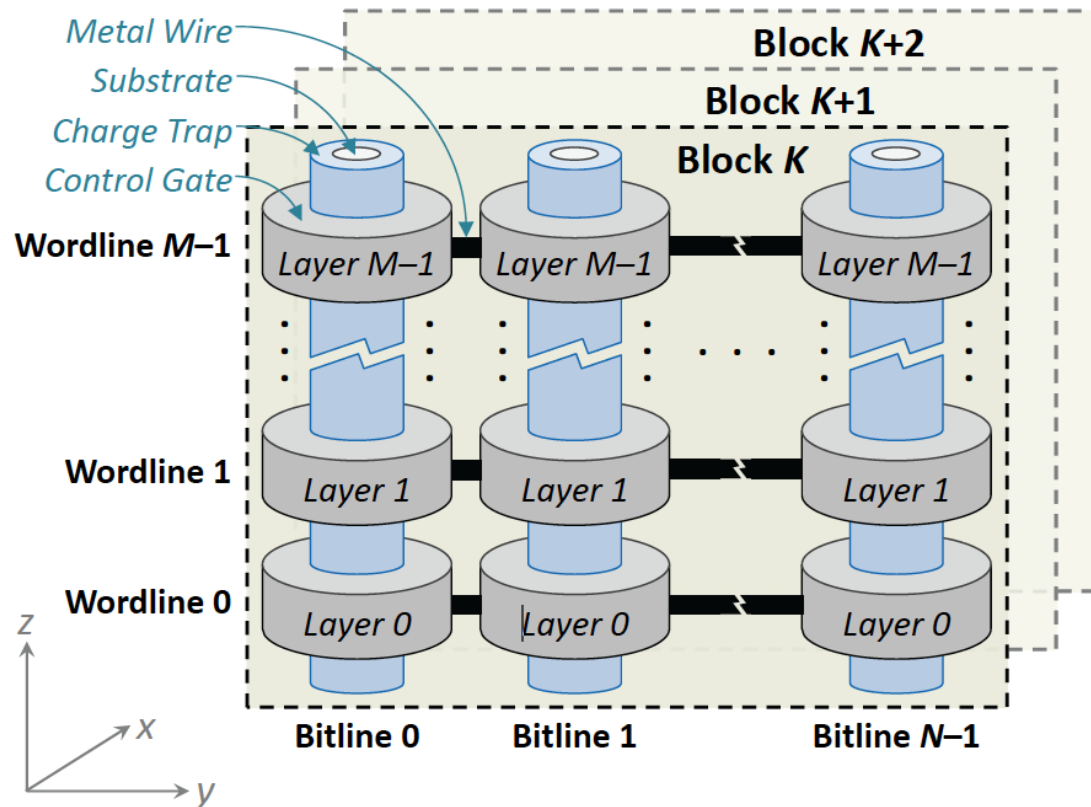
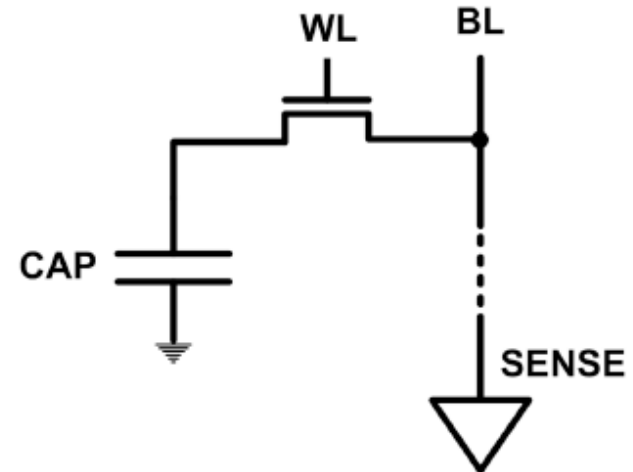
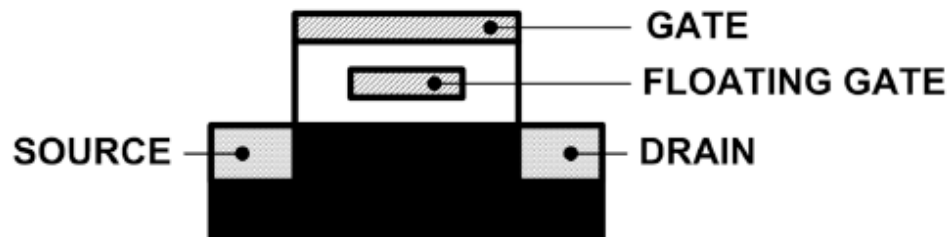


Fig. 43. Organization of flash cells in an M -layer 3D charge trap NAND flash memory chip, where each block consists of M wordlines and N bitlines.

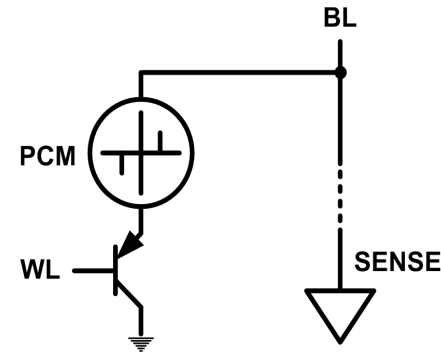
Limits of Charge Memory

- Difficult charge placement and control
 - Flash: floating gate charge
 - DRAM: capacitor charge, transistor leakage
- Reliable sensing becomes difficult as charge storage unit size reduces



Solution: Emerging Memory Technologies

- Some emerging **resistive** memory technologies seem more scalable than DRAM (and they are non-volatile)
- Example: Phase Change Memory
 - Data stored by changing phase of material
 - Data read by detecting material's resistance
 - Expected to scale to 9nm (2022 [ITRS 2009])
 - Prototyped at 20nm (Raoux+, IBM JRD 2008)
 - Expected to be denser than DRAM: can store multiple bits/cell
- But, emerging technologies have (many) shortcomings
 - Can they be enabled to replace/augment/surpass DRAM?



Intel Optane Persistent Memory (2019)

- Non-volatile main memory
- Based on 3D-XPoint Technology



PCM as Main Memory: Idea in 2009

- Benjamin C. Lee, Engin Ipek, Onur Mutlu, and Doug Burger,
"Architecting Phase Change Memory as a Scalable DRAM Alternative"
Proceedings of the [36th International Symposium on Computer Architecture \(ISCA\)](#), pages 2-13, Austin, TX, June 2009. [Slides \(pdf\)](#)
One of the 13 computer architecture papers of 2009 selected as Top Picks by IEEE Micro.
Selected as a CACM Research Highlight.

Architecting Phase Change Memory as a Scalable DRAM Alternative

Benjamin C. Lee[†] Engin Ipek[†] Onur Mutlu[‡] Doug Burger[†]

[†]Computer Architecture Group
Microsoft Research
Redmond, WA
{blee, ipek, dburger}@microsoft.com

[‡]Computer Architecture Laboratory
Carnegie Mellon University
Pittsburgh, PA
onur@cmu.edu



Charge vs. Resistive Memories

- Charge Memory (e.g., DRAM, Flash)
 - Write data by capturing charge Q
 - Read data by detecting voltage V
- Resistive Memory (e.g., PCM, STT-MRAM, memristors)
 - Write data by pulsing current dQ/dt
 - Read data by detecting resistance R



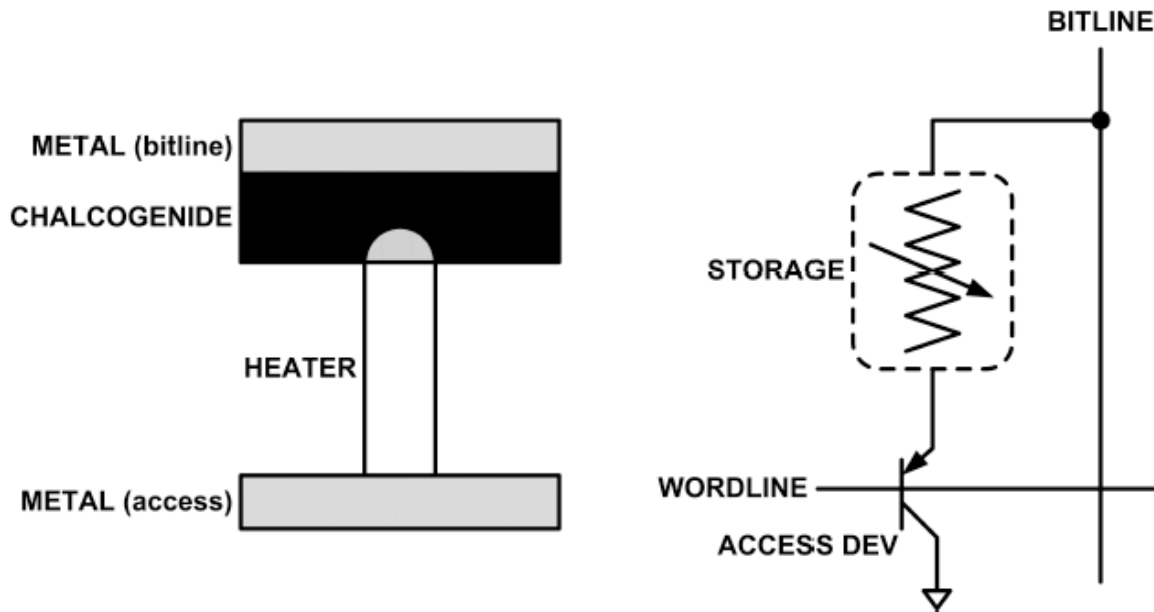
Promising Resistive Memory Technologies

- PCM
 - Inject current to change **material phase**
 - Resistance determined by phase
- STT-MRAM
 - Inject current to change **magnet polarity**
 - Resistance determined by polarity
- Memristors/RRAM/ReRAM
 - Inject current to change **atomic structure**
 - Resistance determined by atom distance



What is Phase Change Memory?

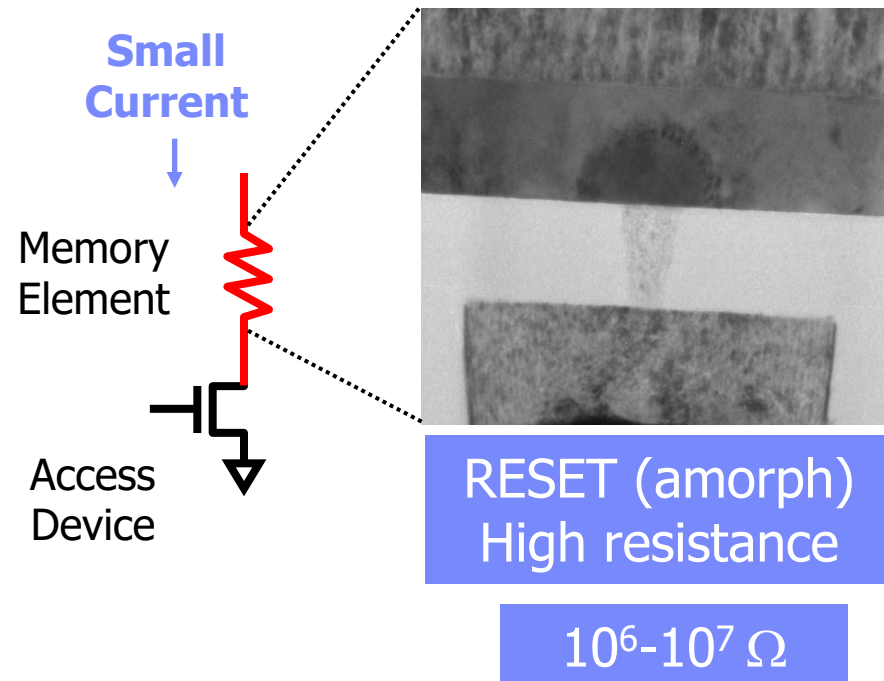
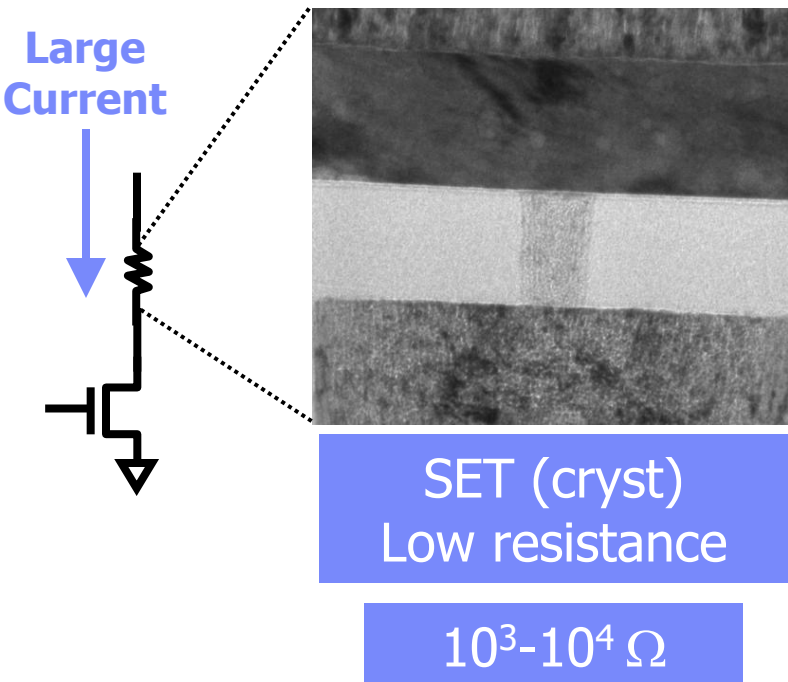
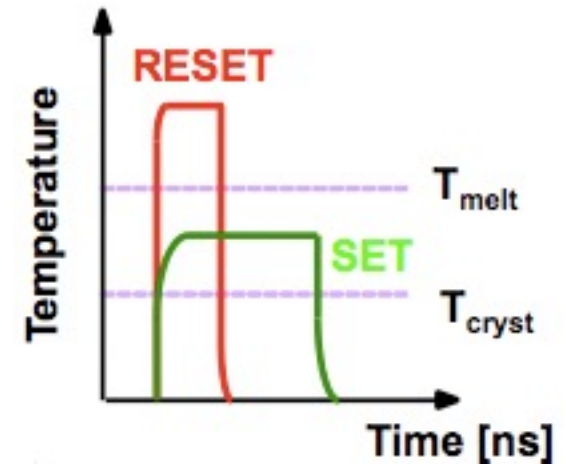
- Phase change material (chalcogenide glass) exists in two states:
 - Amorphous: Low optical reflexivity and high electrical resistivity
 - Crystalline: High optical reflexivity and low electrical resistivity



PCM is resistive memory: High resistance (0), Low resistance (1)
PCM cell can be switched between states reliably and quickly

How Does PCM Work?

- Write: change phase via current injection
 - SET: sustained current to heat cell above T_{cryst}
 - RESET: cell heated above T_{melt} and quenched
- Read: detect phase via material resistance
 - amorphous/crystalline

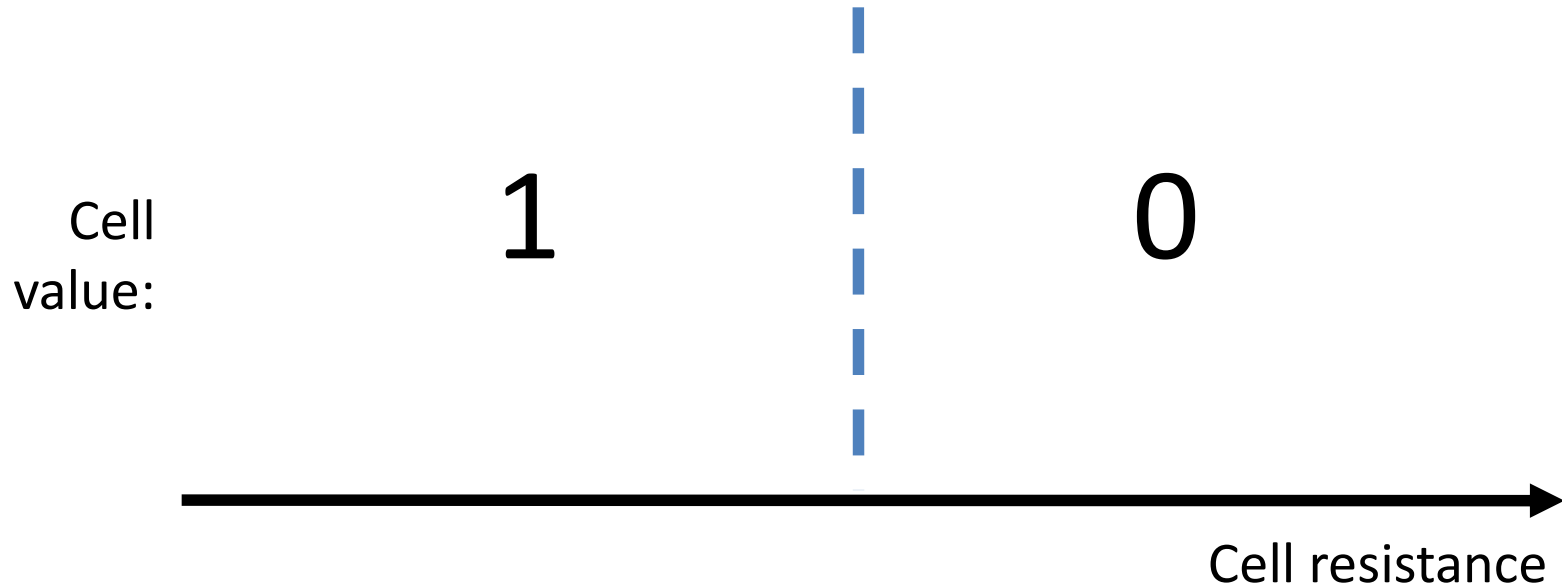


Opportunity: PCM Advantages

- Scales better than DRAM, Flash
 - Requires current pulses, which scale linearly with feature size
 - Expected to scale to 9nm (2022 [ITRS])
 - Prototyped at 20nm (Raoux+, IBM JRD 2008)
- Can be denser than DRAM
 - Can store multiple bits per cell due to large resistance range
 - Prototypes with 2 bits/cell in ISSCC' 08, 4 bits/cell by 2012
- Non-volatile
 - Retain data for >10 years at 85C
- No refresh needed, low idle power



PCM Resistance \rightarrow Value

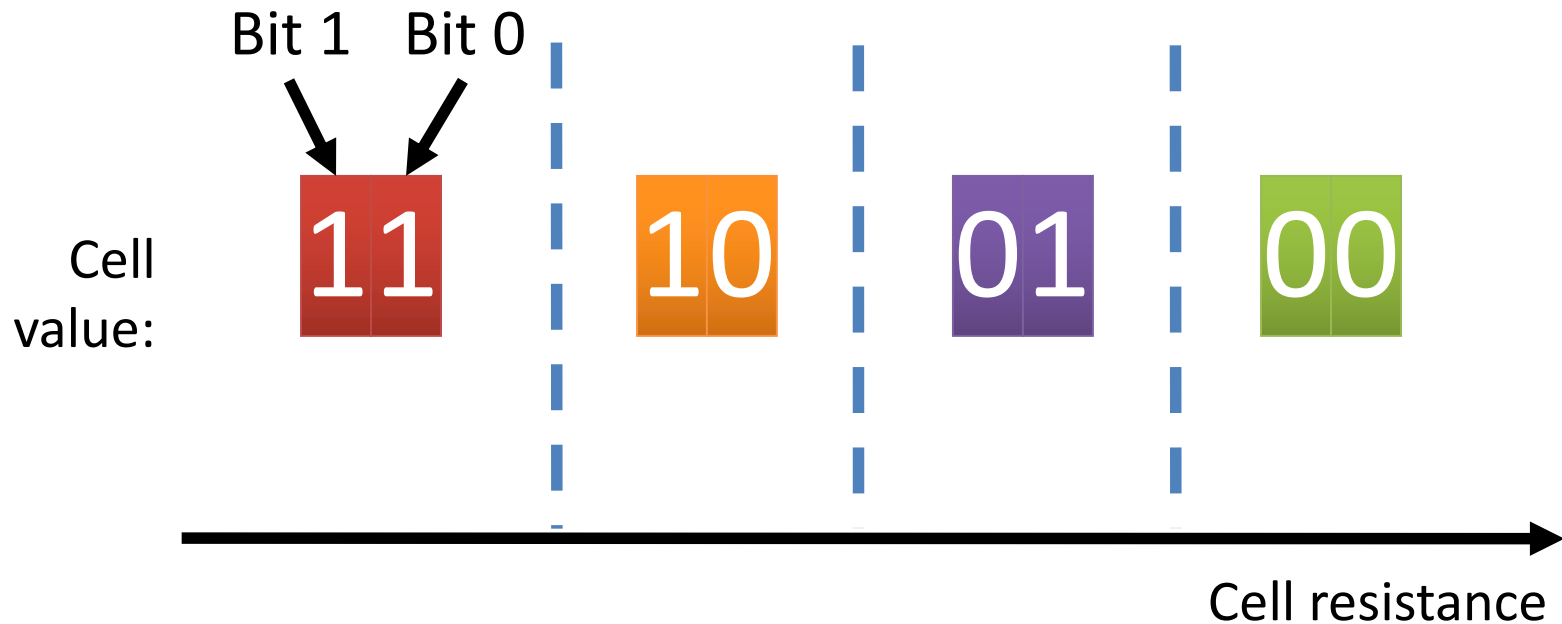


Multi-Level Cell PCM

- Multi-level cell: more than 1 bit per cell
 - Further increases density by 2 to 4x [Lee+,ISCA'09]
- But MLC-PCM also has drawbacks
 - Higher latency and energy than single-level cell PCM



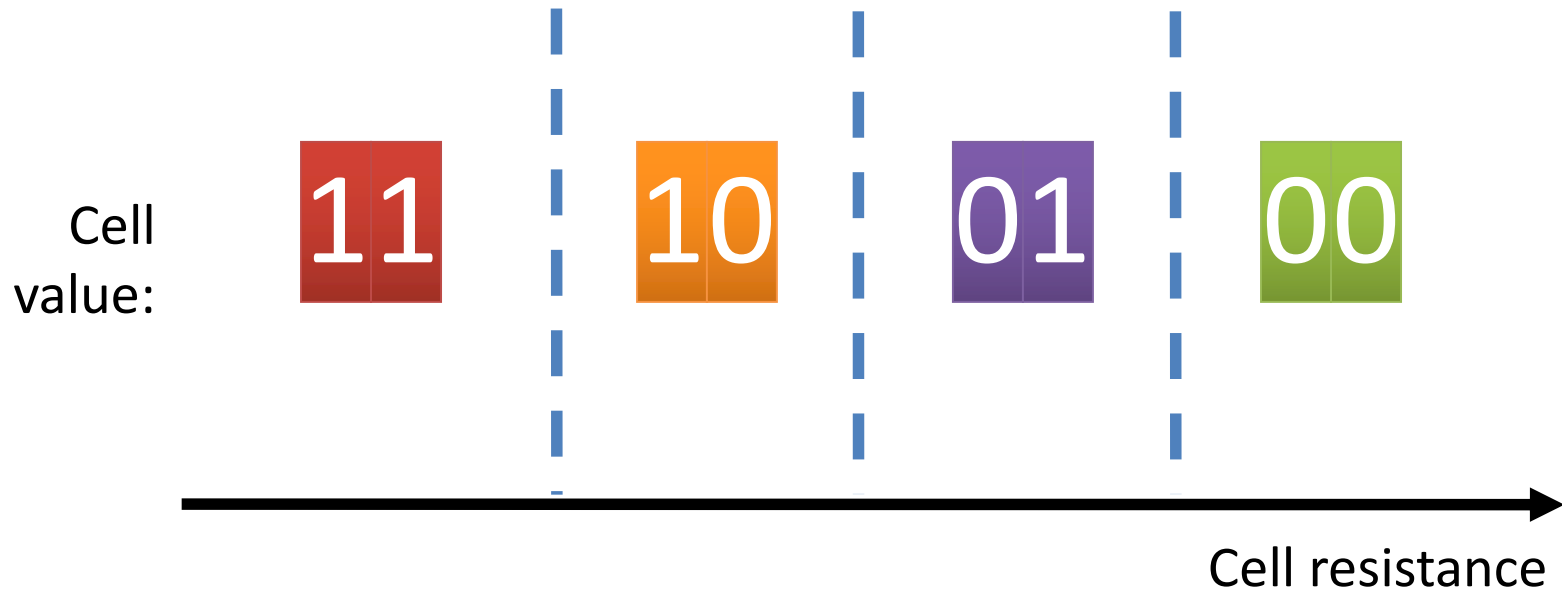
MLC-PCM Resistance \rightarrow Value



MLC-PCM Resistance → Value

Less margin between values

- need more precise sensing/modification of cell contents
- higher latency/energy (~2x for reads and 4x for writes)



Phase Change Memory Properties

- Surveyed prototypes from 2003-2008 (ITRS, IEDM, VLSI, ISSCC)
- Derived PCM parameters for $F=90\text{nm}$
- Lee, Ipek, Mutlu, Burger, “Architecting Phase Change Memory as a Scalable DRAM Alternative,” ISCA 2009.
- Lee et al., “Phase Change Technology and the Future of Main Memory,” IEEE Micro Top Picks 2010.



Table 1. Technology survey.

Published prototype

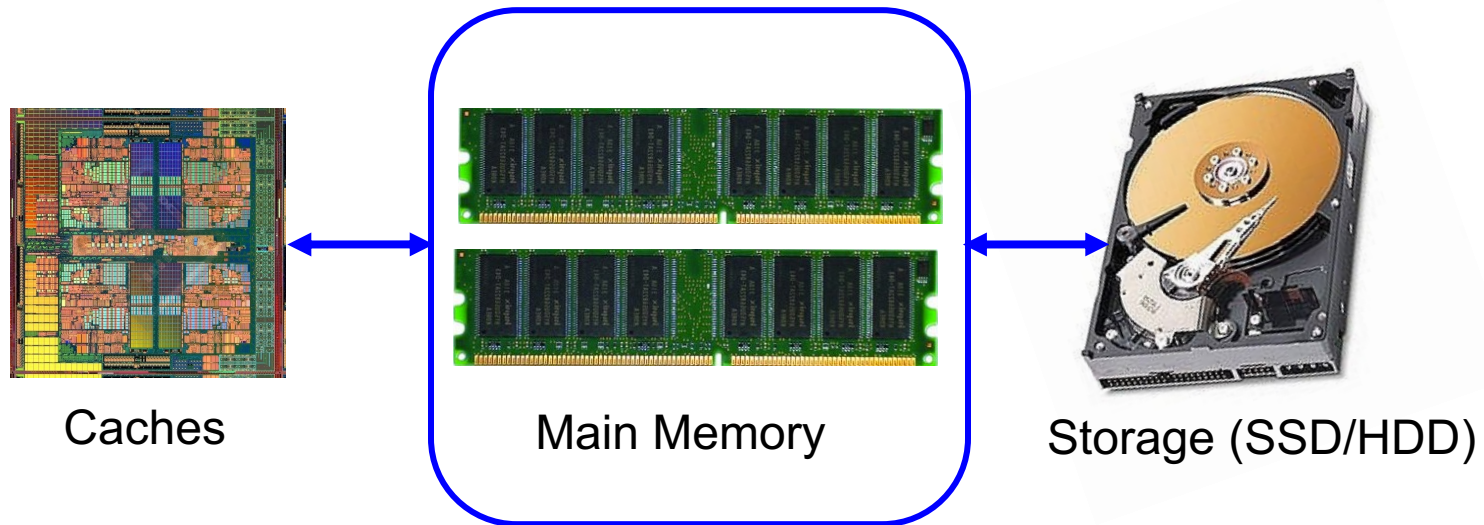
Parameter*	Horri ⁶	Ahn ¹²	Bedeschi ¹³	Oh ¹⁴	Pellizer ¹⁵	Chen ⁵	Kang ¹⁶	Bedeschi ⁹	Lee ¹⁰	Lee ²
Year	2003	2004	2004	2005	2006	2006	2006	2008	2008	**
Process, F (nm)	**	120	180	120	90	**	100	90	90	90
Array size (Mbytes)	**	64	8	64	**	**	256	256	512	**
Material	GST, N-d	GST, N-d	GST	GST	GST	GS, N-d	GST	GST	GST	GST, N-d
Cell size (μm^2)	**	0.290	0.290	**	0.097	60 nm ²	0.166	0.097	0.047	0.065 to 0.097
Cell size, F^2	**	20.1	9.0	**	12.0	**	16.6	12.0	5.8	9.0 to 12.0
Access device	**	**	BJT	FET	BJT	**	FET	BJT	Diode	BJT
Read time (ns)	**	70	48	68	**	**	62	**	55	48
Read current (μA)	**	**	40	**	**	**	**	**	**	40
Read voltage (V)	**	3.0	1.0	1.8	1.6	**	1.8	**	1.8	1.0
Read power (μW)	**	**	40	**	**	**	**	**	**	40
Read energy (pJ)	**	**	2.0	**	**	**	**	**	**	2.0
Set time (ns)	100	150	150	180	**	80	300	**	400	150
Set current (μA)	200	**	300	200	**	55	**	**	**	150
Set voltage (V)	**	**	2.0	**	**	1.25	**	**	**	1.2
Set power (μW)	**	**	300	**	**	34.4	**	**	**	90
Set energy (pJ)	**	**	45	**	**	2.8	**	**	**	13.5
Reset time (ns)	50	10	40	10	**	60	50	**	50	40
Reset current (μA)	600	600	600	600	400	90	600	300	600	300
Reset voltage (V)	**	**	2.7	**	1.8	1.6	**	1.6	**	1.6
Reset power (μW)	**	**	1620	**	**	80.4	**	**	**	480
Reset energy (pJ)	**	**	64.8	**	**	4.8	**	**	**	19.2
Write endurance (MLC)	10^7	10^9	10^6	**	10^8	10^4	**	10^5	10^5	10^8

* BJT: bipolar junction transistor; FET: field-effect transistor; GST: $\text{Ge}_2\text{Sb}_2\text{Te}_5$; MLC: multilevel cells; N-d: nitrogen doped.

** This information is not available in the publication cited.

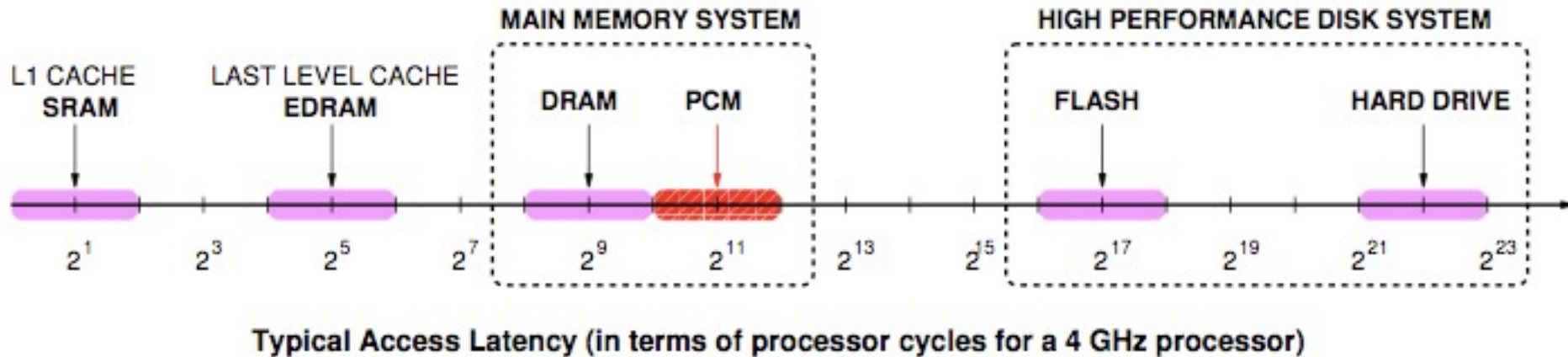


Where Can PCM Fit in the System?



Phase Change Memory Properties: Latency

- Latency comparable to, but slower than DRAM



- Read Latency
 - 50ns: 4x DRAM, 10^{-3} x NAND Flash
- Write Latency
 - 150ns: 12x DRAM
- Write Bandwidth
 - 5-10 MB/s: 0.1x DRAM, 1x NAND Flash



Phase Change Memory Properties

- Dynamic Energy
 - 40 μA Rd, 150 μA Wr
 - 2-43x DRAM, 1x NAND Flash
- Endurance
 - Writes induce phase change at 650C
 - Contacts degrade from thermal expansion/contraction
 - 10^8 writes per cell
 - 10^{-8}x DRAM, 10^3x NAND Flash
- Cell Size
 - 9-12F² using BJT, single-level cells
 - 1.5x DRAM, 2-3x NAND (will scale with feature size, MLC)



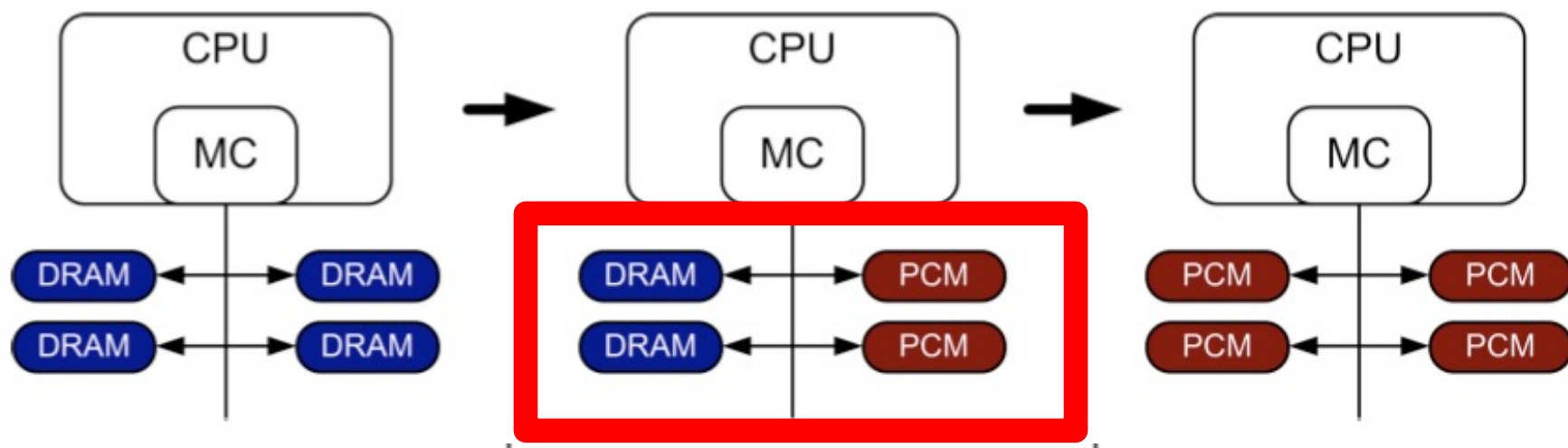
Phase Change Memory: Pros and Cons

- Pros over DRAM
 - Better technology scaling (capacity and cost)
 - Non volatile → Persistent
 - Low idle power (no refresh)
- Cons
 - Higher latencies: $\sim 4\text{-}15\times$ DRAM (especially write)
 - Higher active energy: $\sim 2\text{-}50\times$ DRAM (especially write)
 - Lower endurance (a cell dies after $\sim 10^8$ writes)
 - Reliability issues (resistance drift)
- Challenges in enabling PCM as DRAM replacement/helper:
 - Mitigate PCM shortcomings
 - Find the right way to place PCM in the system



PCM-based Main Memory (I)

- How should PCM-based (main) memory be



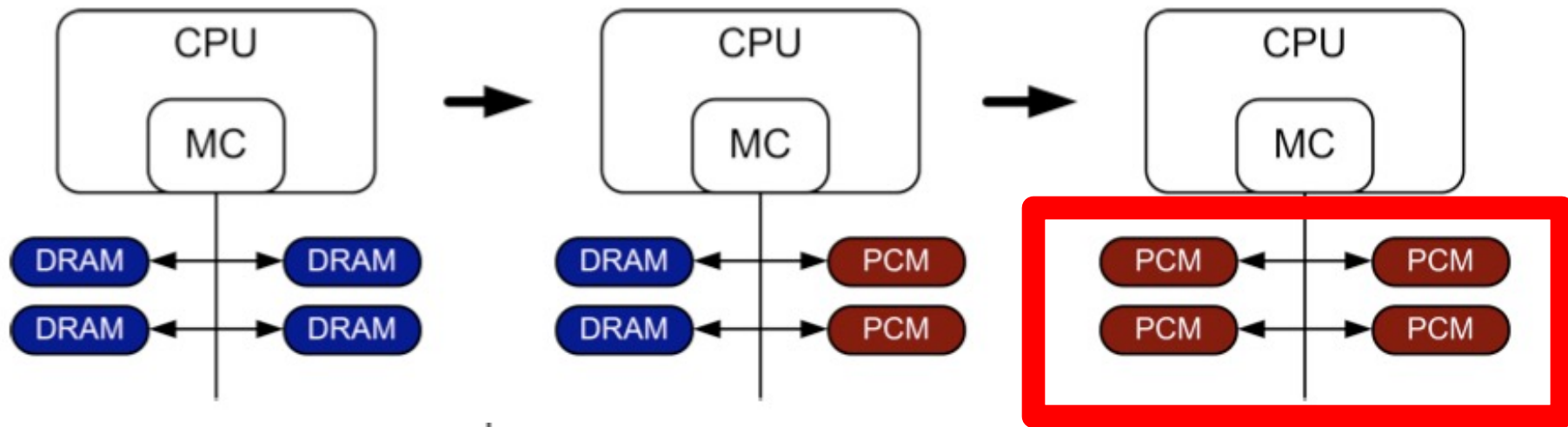
- Hybrid PCM+DRAM** [Qureshi+ ISCA'09, Dhiman+ DAC'09]:
 - How to partition/migrate data between PCM and DRAM

M. Tarihi, H. Asadi, A. Haghdoost, M. Arjomand, and H. Sarbazi-Azad,
“A Hybrid Non-Volatile Cache Design for Solid-State Drives Using
Comprehensive I/O Characterization,”
IEEE Transactions on Computers (TC), Vol. 65, Issue 6, 2016



PCM-based Main Memory (II)

- How should PCM-based (main) memory be



- Pure PCM main memory** [Lee et al., ISCA'09, Top Picks'10]:
 - How to redesign entire hierarchy (and cores) to overcome PCM shortcomings

An Initial Study: Replace DRAM with PCM

- Lee, Ipek, Mutlu, Burger, “Architecting Phase Change Memory as a Scalable DRAM Alternative,” ISCA 2009.
 - Surveyed prototypes from 2003-2008 (e.g. IEDM, VLSI, ISSCC)
 - Derived “average” PCM parameters for F=90nm

Density

- ▷ $9 - 12F^2$ using BJT
- ▷ $1.5\times$ DRAM

Latency

- ▷ 50ns Rd, 150ns Wr
- ▷ $4\times, 12\times$ DRAM

Endurance

- ▷ $1E+08$ writes
- ▷ $1E-08\times$ DRAM

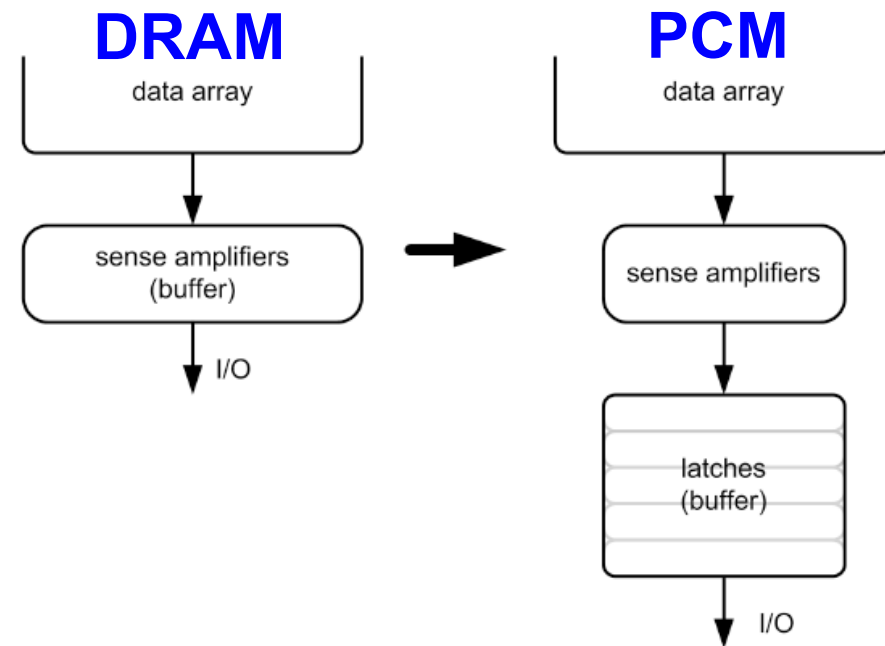
Energy

- ▷ $40\mu A$ Rd, $150\mu A$ Wr
- ▷ $2\times, 43\times$ DRAM



Architecting PCM to Mitigate Shortcomings

- Idea 1: Use multiple narrow row buffers in each PCM chip
→ Reduces array reads/writes → better endurance, latency, energy
- Idea 2: Write into array at cache block or word granularity
→ Reduces unnecessary wear

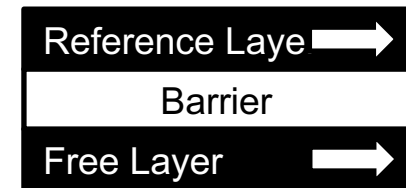


STT-RAM as Main Memory

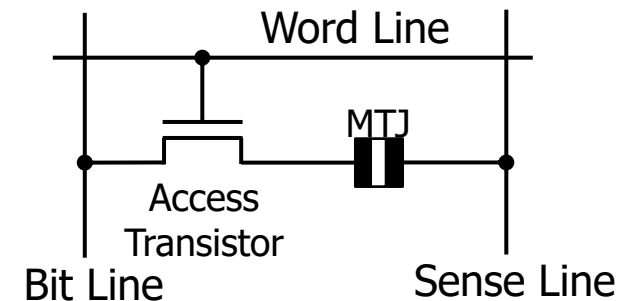
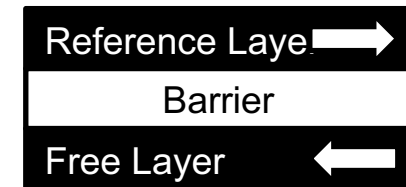
STT-MRAM as Main Memory

- Magnetic Tunnel Junction (MTJ) device
 - Reference layer: Fixed magnetic orientation
 - Free layer: Parallel or anti-parallel
- Magnetic orientation of the free layer determines logical state of device
 - High vs. low resistance
- Write: Push large current through MTJ to change orientation of free layer
- Read: Sense current flow
- Kultursay et al., “Evaluating STT-RAM as an Energy-Efficient Main Memory Alternative,” ISPASS 2013.

Logical 0



Logical 1

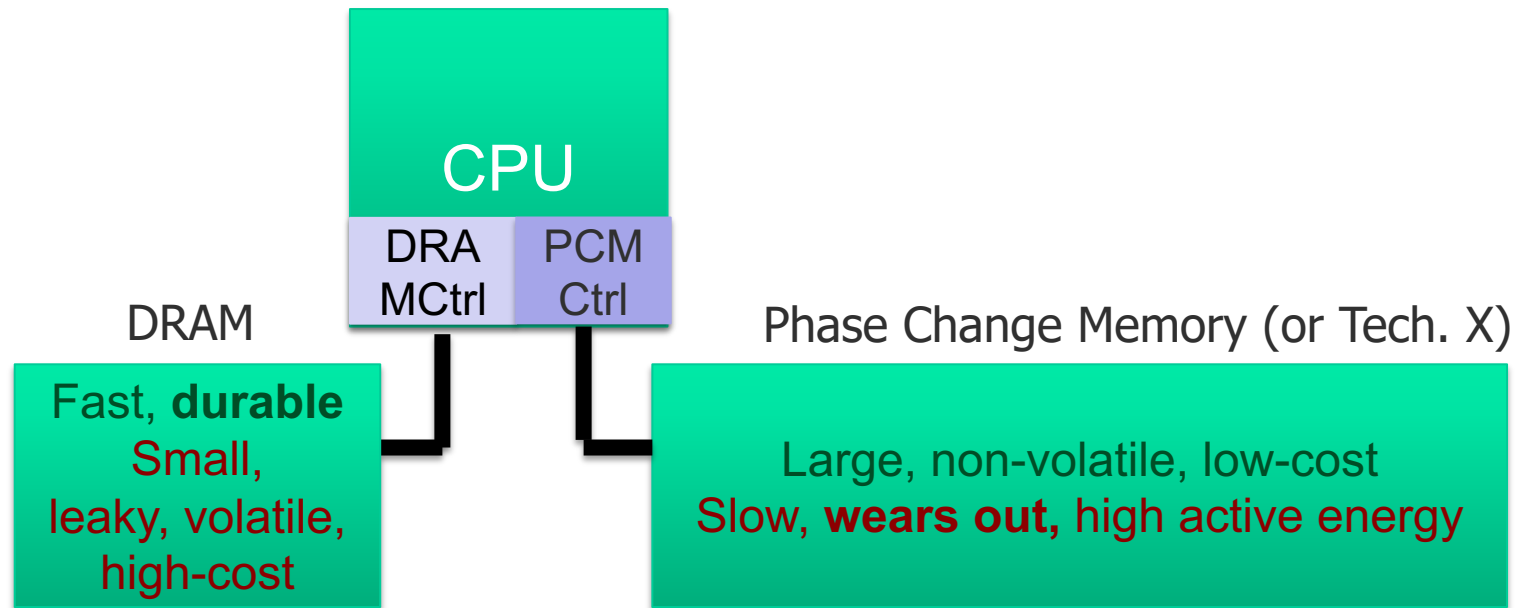


STT-MRAM: Pros and Cons

- Pros over DRAM
 - Better technology scaling (capacity and cost)
 - Non volatile → Persistent
 - Low idle power (no refresh)
- Cons
 - Higher write latency
 - Higher write energy
 - Poor density (currently)
 - Reliability?
- Another level of freedom
 - Can trade off non-volatility for lower write latency/energy (by reducing the size of the MTJ)



A More Viable Approach: Hybrid Memory Systems



Hardware/software manage data allocation and movement
to achieve the best of multiple technologies

Meza+, "[Enabling Efficient and Scalable Hybrid Memories](#)," IEEE Comp. Arch. Letters, 2012.
Yoon+, "[Row Buffer Locality Aware Caching Policies for Hybrid Memories](#)," ICCD 2012 Best Paper Award.



Providing the Best of
Multiple Metrics
with

Multiple Memory Technologies



Challenge and Opportunity

Heterogeneous,
Configurable,
Programmable
Memory Systems



Hybrid Memory Systems: Issues

- Cache vs. Main Memory
- Granularity of Data Move/Management: Fine or Coarse
- Hardware vs. Software vs. HW/SW Cooperative
- When to migrate data?
- How to design a scalable and efficient large cache?



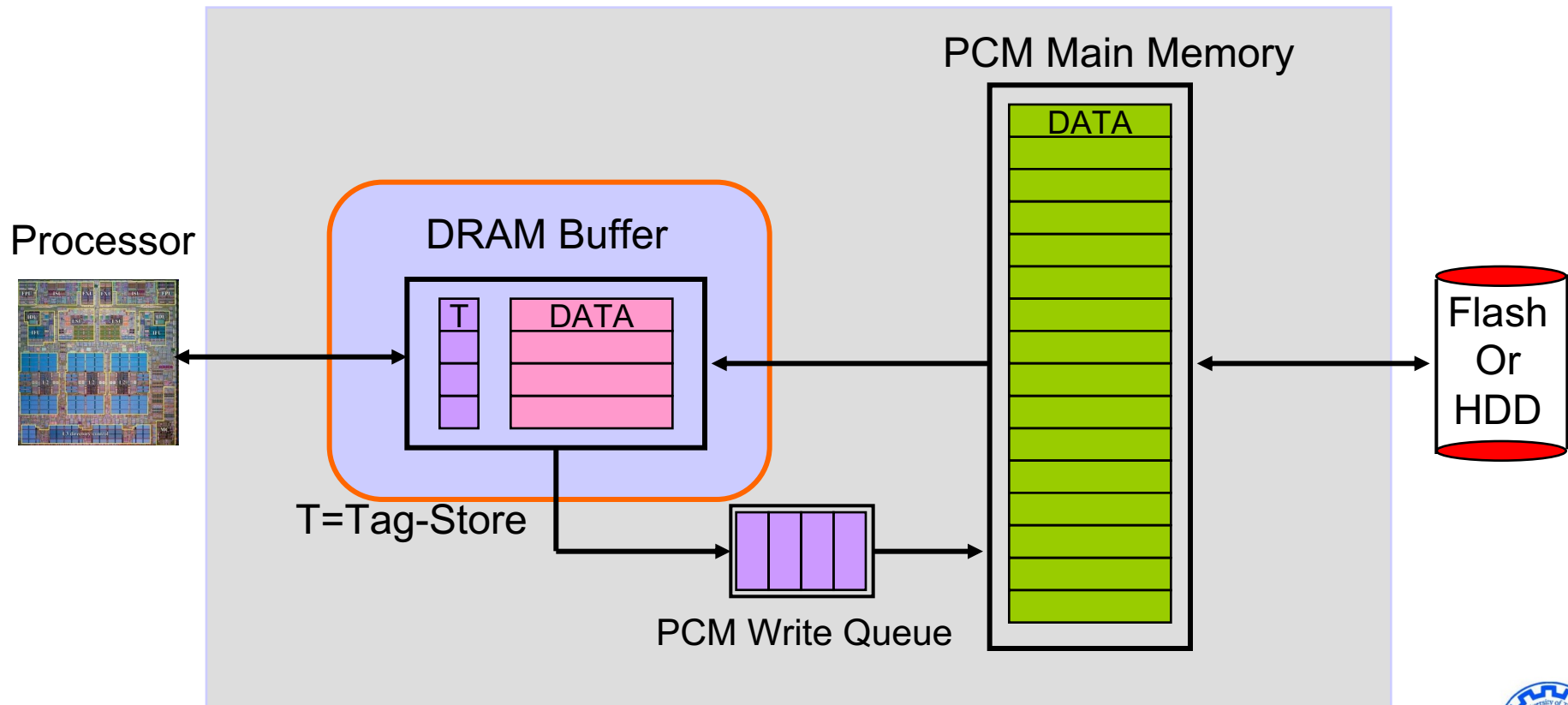
One Option: DRAM as a Cache for PCM

- PCM is main memory; DRAM caches memory rows/blocks
 - Benefits: Reduced latency on DRAM cache hit; write filtering
- Memory controller hardware manages the DRAM cache
 - Benefit: Eliminates system software overhead
- Three issues:
 - What data should be placed in DRAM versus kept in PCM?
 - What is the granularity of data movement?
 - How to design a low-cost hardware-managed DRAM cache?



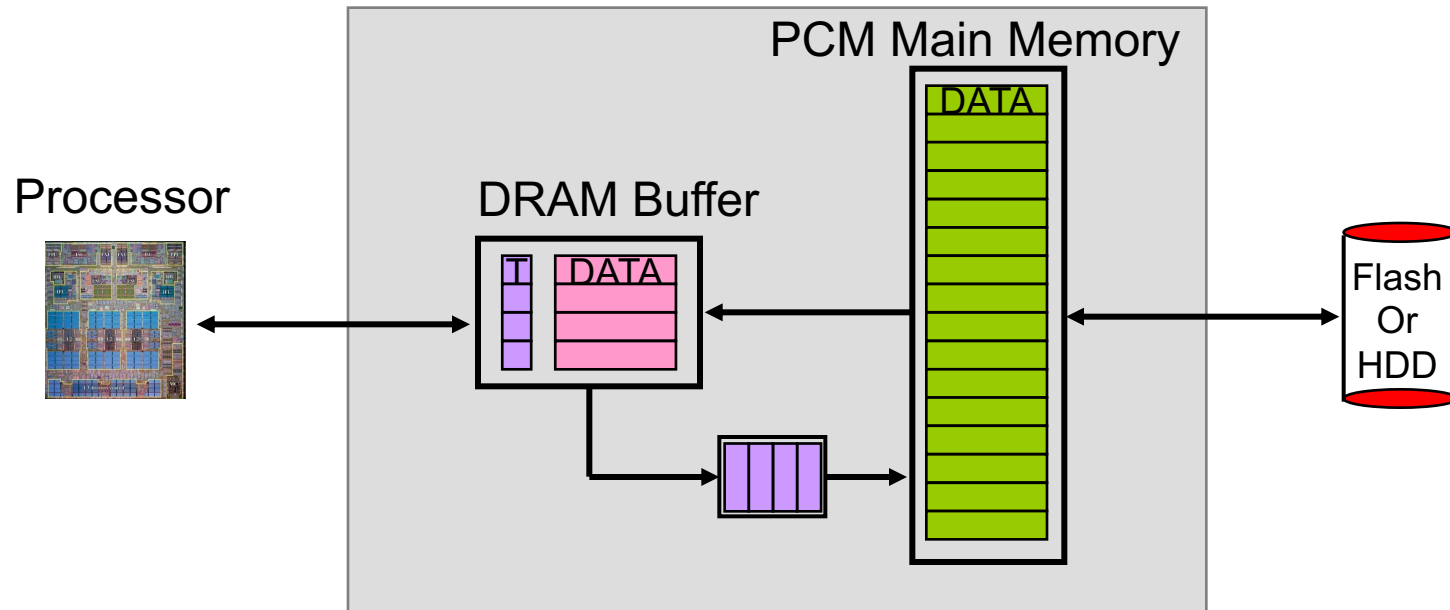
DRAM as a Cache for PCM

- Goal: Achieve the best of both DRAM and PCM/NVM
 - Minimize amount of DRAM w/o sacrificing performance, endurance
 - DRAM as cache to tolerate PCM latency and write bandwidth
 - PCM as main memory to provide large capacity at good cost and power



Write Filtering Techniques

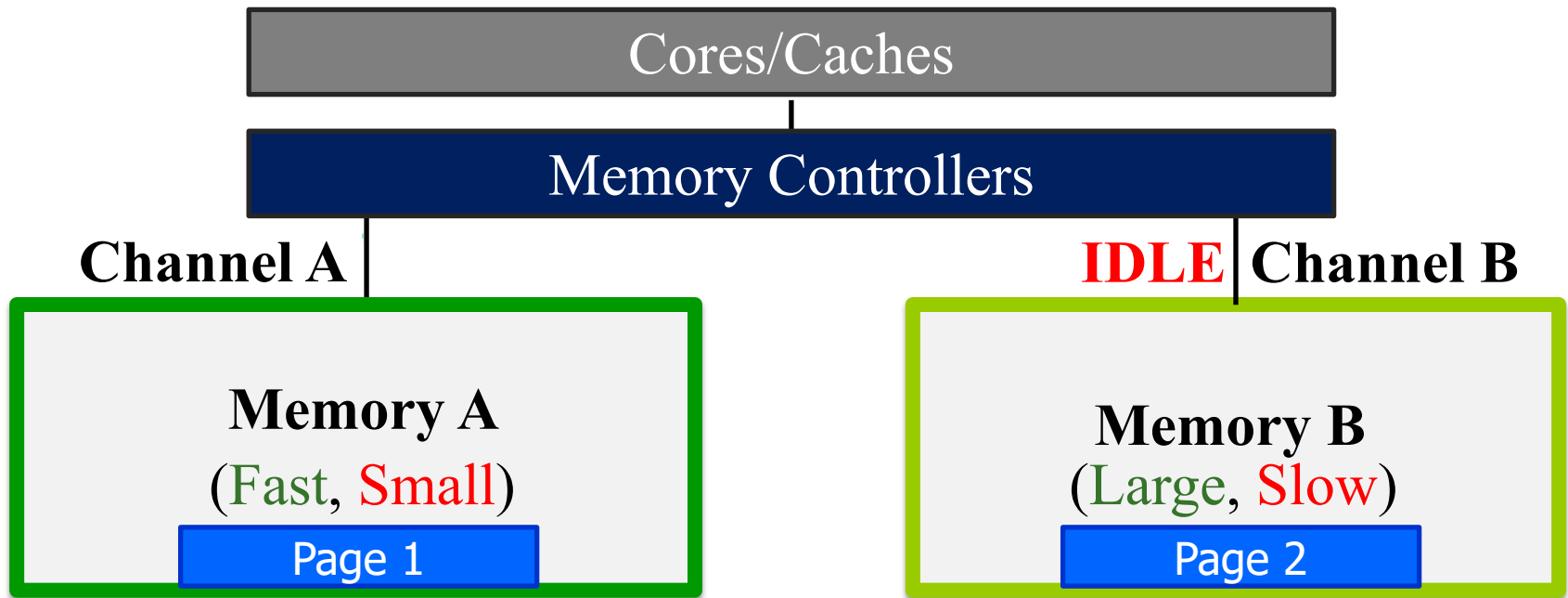
- Lazy Write: Pages from disk installed only in DRAM, not PCM
- Partial Writes: Only dirty lines from DRAM page written back
- Page Bypass: Discard pages with poor reuse on DRAM eviction



- Qureshi et al., “Scalable high performance main memory system using phase-change memory technology” ISCA 2009



Data Placement in Hybrid Memory



Which memory do we place each page in,
to **maximize system performance**?

- Memory A is fast, but small
- Load should be balanced on both channels?
- Page migrations have performance and energy overhead



Key Observation & Idea

- Row buffers exist in both DRAM and PCM
 - Row **hit** latency **similar** in DRAM & PCM [Lee+ ISCA'09]
 - Row **miss** latency **small** in DRAM, **large** in PCM
- Place data in DRAM which
 - is likely to miss in the row buffer (**low row buffer locality**) → miss penalty is smaller in DRAM
 - AND
 - is **reused many times** → cache only the data worth the movement cost and DRAM space



