Digital Integrated Circuits (83-313)

Lecture 11: Other Memories

Emerging Nanoscaled Integrated Circuits and Systems Labs Prof. Adam Teman 27 May 2022

The Alexander Kofkin Faculty of Engineering



Bar-Ilan University

Disclaimer: This course was prepared, in its entirety, by Adam Teman. Many materials were copied from sources freely available on the internet. When possible, these sources have been cited; however, some references may have been cited incorrectly or overlooked. If you feel that a picture, graph, or code example has been copied from you and either needs to be cited or removed, please feel free to email <u>adam.teman@biu.ac.il</u> and I will address this as soon as possible.

Lecture Content





The Memory Hierarchy





Semiconductor Memory Classification



Memory Hierarchy of a Personal Computer



[©] Adam Teman, 2022

Another View







READ Only Memory (ROM)







SRAM Drawbacks

- The 6T SRAM is by far the dominant on-chip (embedded) memory solution
 - Static operation
 - High noise margins
 - Differential (→Fast) readout
 - Symmetric structure → dense layout
- But six transistors is still a lot...
 - On-chip SRAM is (usually) limited to a <30 MB
 - Smallest Linux distro requires >64 MB to run
 - Windows 10 requires >2 GB to run
- Can we achieve higher density?





- One option is to define it in the RTL...
- But what if we want it to be nicely structured and easily modified for bug fixes and future versions of our product?
- The solution: Mask ROM

Reminder: NOR vs. NAND Pull-Down Networks

 In a NOR gate, we <u>only need one</u> conducting path to pull down the output

• In a NAND gate, we <u>need all paths</u> to be conducting to pull down the output





Read-Only Memory Cells



Diode ROM

MOS ROM 1

MOS ROM 2 dam Teman, 2022

NOR ROM

- 4-word x 6-bit NOR-ROM
 - Selected word-line high
- Represented with dot diagram

Word 3: 101010









NAND ROM

4-word 4-bit NAND ROM



NAND ROM Layout



One-Time-Programmable Memory (OTP)

Mask ROM is set during fabrication

- An actual lithography mask encodes the ROM content.
- A re-spin of the chip is needed in order to change the ROM.
- Programmable Read-Only Memories (PROM) are an alternative
 - Content of the PROM can be programmed (blown) after manufacturing
 - Generally, One-Time-Programmable (OTP) devices used today.
 - Often used to store boot code, encryption keys and configuration parameters.

Source: Chung

- Achieved by integrating an "antifuse" within the CMOS process
 - Apply a high-voltage pulse (~6V) across the gate
 → oxide breakdown.





Dynamic RAM (DRAM)





ROMs are very limited in usage

• ROMs can be only used for data that is:

- Written during manufacturing
- Can only be read.
- OTPs:
 - Can be configured after manufacturing
 - But still, only once.

 How can we make a high-density (→ high-capacity) memory that can be used for both reading and writing?



Let's go dynamic

- The Intel 1103 used a 3T DRAM cell
 - Store data on capacitor (C_s)
 - Write through M_1
 - Read through M_2 and M_3
- Advantages:
 - No constraints on device ratios.
 - Reads are non-destructive.

Disadvantages

- Data leaks away over time
- "Weak 1" stored on node X.
- Single-ended read operation



3T DRAM Layout

- Route word-lines (WWL, RWL) in Poly
- Route bit-lines (BL1, BL2) and GND in M2
- Disadvantage:
 - Contact-dominated layout





BLWLCan we go even smaller? bit0 bit1 bit511 Dennard (IBM, 1967) patented a 1T cell Ę • Data is dynamically stored on C_{s} . • Write and read through M_1 $C_{BL} \neq$ Read achieved through charge sharing between C_s and C_{BL} : $\Delta V = V_{\rm BL} - V_{\rm pre} = V_{\rm BIT} - V_{\rm pre} \frac{C_{\rm S}}{C_{\rm S} + C_{\rm BI}}$ Write 1 Read 1 WL • Drawbacks: X GND Weak '1' stored on C_s

- Charge sharing is slow, small delta
- Read is destructive. Write back is necessary

 V_{DD}

BL

 $V_{DD}/2$



- Popular in 70s and 80s
- Stacked Capacitor
 - Make a "radiator" on top of the cell
- Trenched Capacitor
 - Dig a ditch under the cell
 - Fill it with conductor, insulator, conductor

Latch-Based Sense Amplifier (DRAM)

Sense amp activated

Word line activated

- Initialized in its meta-stable point with EQ
- Once adequate voltage gap created, sense amp enabled with SE
- Positive feedback quickly forces output to a stable operating point.

V_{BL} ,

V_{PRE}



DRAM Architecture

- To get a 16GB DRAM DIMM (Dual-Inline Memory Module)
 - DIMM has two ranks
 - Each rank (64-bit I/O) has 8 chips (1 GB/chip)
 - Each chip (8-bit I/O) has 8 banks (128 MB/bank)
 - Each bank (1-bit I/O) is divided into 8 arrays (16 MB/array)
 - Each array is divided into ~128 1024x1024 bit MATs (sub-arrays)
- So, to access DRAM:
 - First provide Row Access Strobe (RAS) 10 bits to MAT
 - Readout latches rows (~128 x 1024 bits) in Sense Amps
 - ~7 bits of RAS select 1024 bits into bank Row Buffer
 - Then Column Access Strobe (CAS)
 - Select 1-bit from Row Buffer to bank output.





Address

Embedded DRAM

- DRAM requires a special (dedicated) process (Fab)
 - Manufactured by only a handful of providers (predominantly Micron, Samsung and SK Hynix)
 - Is supplied on an external standalone chip.
- However, it is possible to make an on-chip "embedded" version
 - Known as "eDRAM"
 - Requires special process cost-adders to fabricate trench capacitors (IBM) or stacked MIM (Intel)
- Recent chips with L4 eDRAM
 - Intel Broadwell 128 MB
 - IBM z15 SC 960 MB



Copper m 2

Cu

Bordered Bitline Contact

Poly PWL

STI

Dielectric

Poly AWL

P-Well

Buried N Plate

Buried Strap

Contact Stud

Cu

Bitline

BP SG

rench Poly-Si

Oxide Collar

N + Polv-Si

J.Barth, IBM

SRDC 2009

Gain-Cell Embedded DRAM

WBL

eDRAM is not "Logic Compatible"

- Process adders are expensive
- Not available for many processes
- Yet to be demonstrated under 14nm

Why not go back to the old Intel DRAM?

- Use parasitic capacitance for storage node
- Manipulate current gain of read transistor
- Known as "Gain-Cell embedded DRAM"
- 2T, 3T, 4T, and others demonstrated
- Most-scaled demonstrated GC-eDRAM
 - "DynOR" 28nm FD-SOI (Giterman, TCAS-I 2017)
 - "DAFNA" 28nm (Giterman, JSSC 2018)
 - "Negev" 16nm (Giterman, SSC-L 2020)





Non-Volatile Memory (NVM)





Memory Volatility

• **SRAM** and **DRAM** are examples of *Volatile* memories

- They only store their data as long as a power source is connected
- Mask ROM and OTP are *non-Volatile*
 - Their content is saved regardless of the state of the battery

What about non-Volatile Read/Write Memories?

The Floating-gate transistor (EPROM)



Floating-Gate Transistor Programming



Avalanche injection

Removing programmingProgramming results involtage leaves charge trappedhigher V_T .



FLOTOX EEPROM



FLOTOX transistor

Fowler-Nordheim *I-V* characteristic

EEPROM Cell

Absolute threshold control is hard

- Unprogrammed transistor might be depletion
- Therefore, use a 2 transistor cell:
 - Program it such that VT>VDD
 - Erase it, such that VT is really low (even VT<0)



Flash EEPROM



Many other options ...

Cross-sections of NVM cells



Flash

Courtesy Intel

© Adam Teman, 2022

EPROM

Two major architectures: NOR and NAND

Source: Micron



NOR Flash Operation

Program

- Run current through and apply vertical field
- Takes a long time (µsecs)
- Readout and see if programming successful.

• Erase

- First program all cells in block
- Then Erase entire block and verify

Read

- Apply median voltage on the gate (e.g., 5V)
- '0' State BL will discharge
- '1' State BL will stay charged



NAND Flash Operation

- Erase:
 - Source Line = 20V
 - BL=20V WL=0
 - All devices normally on
- Program (Write '1'):
 - Source Line disconnected
 - BL=0V
 - WL_{selected}=20V WL_{unselected}=5V
- Read:
 - Source Line connected
 - BL connected
 - WL_{selected}=0V WL_{unselected}=5V





NAND vs. NOR

• NOR architectures are large, but fast.

- Actually, NOR is fast for READs, but slow for WRITEs and ERASEs, as they need precise control over threshold voltages.
- NORs are used for applications such as program code storage that are read a lot and at high speeds.
- NAND architectures are slow, but small.
 - Actually, just READs are slow. Programming and Erasure are relatively fast.
 - Used for video and audio file storage, where we need high density, as well as fast write/erase times.
 - NAND is about 40% smaller than NOR and generally uses Fowler Nordheim tunneling for program and erase.

Increasing Flash Density

- Add levels to every bit
 - MLC = Multi-level Cell
- Or perhaps go in the third dimension...



Source: Kingston





1	2	3	4	5	
Hierarchy	ROM	DRAM	NVM	Emerging NVM	

Emerging Memories





The "Universal Memory"

- A single memory technology with the advantages of all existing technologies and without their limitations.
- The requirements of a Universal Memory are:
 - High Density → Flash
 - Scalability → SRAM
 - Unlimited Retention → Flash/HDD
 - High Performance → SRAM
 - Unlimited endurance → SRAM/DRAM
 - Process Integration → SRAM/eDRAM
 - Low Power → NVM + Low Voltage

F² (F-squared)

Minimum bitcell size measured as F², where F is the minimum feature size.

The smallest possible structure is 4F² achieved through a crossbar



bitcell

2F



• The "memristor" – the missing fourth element?

L.O. Chua, "Memristor – The Missing Circuit Element," IEEE Trans., 1971

The Missing Memristor Found!

• HP Labs, 2008



D.B. Strukov et al, "The missing memristor found," Nature, 2008



Charge vs. Resistive memories

Charge Memory (e.g., DRAM)

- 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
- Often constructed as a resistive crossbar.



Source: Teman, et al., Wiley 2022 © Adam Teman, 2022

Resistive Memory Array Architectures



Emerging Nonvolatile 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
- FeRAM
 - Utilize the ferroelectric effect
 - Resistance determined by dipole

D. Ielmini and G. Pedretti, Adv. Intell. Syst. 2000040 (2020)

Phase Change Memory (PCM)

- Chalcogenide glass (e.g., CD-ROM) exists in two states:
 - Amorphous: High Resistance
 - Crystalline: Low Resistance
- Advantages
 - Scalable (sub 10nm)
 - Dense (4F², MLC, 3D Stacking)
 - Non-Volatile
 - Read/Write Performance (between DRAM and Flash)
- Current availability:
 - Intel Optane (3D Xpoint)
 - STMicro, Samsung

Magnetoresistance RAM (STT MRAM)

- Resistance controlled by spin across Magnetic Tunnel Junction (MTJ)
 - Parallel spin state: Low Resistance
 - Anti-parallel spin state: High Resistance
- Advantages
 - Scalable (sub 10nm)
 - Non-Volatile
 - Read/Write Performance (comparable to DRAM)
 - Endurance/Retention
- Current availability:
 - Everspin (1Gb, 28nm, standalone)
 - eMRAM: TSMC, GF, Samsung, Intel, UMC
 - Avalanche, Renesas, IBM

Spin-Transfer Orbit MRAM (SOT-MRAM)

Decouple read and write operations of MRAM

Higher write speeds

© Adam Teman, 2022

Source: Teman, et al., Wiley 2022

Resistive RAM (ReRAM/RRAM)

Change resistance by creating a conductive filament:

- Conductive Bridge (CBRAM)
- Oxygen Vacancy (OxRAM)
- Advantages
 - Scalable (sub 10nm)
 - Dense (4F², MLC, 3D Stacking)
 - Non-Volatile
 - Read/Write Performance (between DRAM and Flash)
- Current availability:
 - Adesto (now Dialog) CBRAM licensed to GF
 - Cerfe Labs (ARM) CeRAM licensed from Symetrix
 - eReRAM offered by TSMC, GlobalFoundries
 - Mitsubishi, Fujitsu, Panasonic Winbond
 - WeeBit-Nano promising Israeli start-up

Source: Teman, et al., Wiley 2022

Carbon-Nanotube RAM (CNT-RAM/NRAM)

BL

- Create conductive bridge by connecting or disconnecting stochastic arrays of carbon nanotubes (CNTs)
 - Held together by strong Van der Waals binding forces.
 - 5nm bitcell expected to have approximately 1000 switchable CNT junctions
 - Also demonstrated 3D Stacking
- Commercialized by Nantero

© Adam Teman, 2022

Source: Wikichip, Nantero

Ferroelectric Memory (FeRAM, FeFET)

• Ferroelectricity:

• A material with an electric polarization that can be reversed by application of an external electric field.

• FRAM:

- A very old technology (invented 1952)
- Classic bitcell: PZT-based 1T-1C (130nm)

• FeFET:

Based on ferroelectric HfO₂ (discovered in 2007)

High-V_T

P-type

Low Vth state

Oxide

- Add HfO2 layer → 1T bitcell
- Being developed by FMC

<u>* * 1</u>

High Vth state

Source: S. Mueller, EMF, 2015.

Source: FMC

Comparison between technologies

	SRAM	DRAM	NAND Flash	NOR Flash	РСМ	RRAM	STT- MRAM	SOT- MRAM	FeRAM	FeFET
Circuit	6T	1T1C	1T	1T	1T1R or 1D1R	1T1R or 1D1R	1T1MTJ	2T1MTJ	1T1C	1T
Cell Size	>150 F ²	6 F ²	<4 F ² (3D)	10 F ²	4-12 F ²	4-10 F ²	6-50 F ²	12-100 F ²	6-50 F ²	6-50 F ²
Voltage	<1V	<1V	>10V	>10V	<3V	<3V	<1V	<1V	<2V	<3V
Stackable	No	Yes	Yes	No	Yes	Yes	No	No	No	No
MLC (bits/cell)	No	No	Yes (4)	Yes (2)	Yes (2)	Yes	No	No	No	Yes
Scaling	< 3 nm	~10nm	~14nm	~45nm	<10nm	<10nm	<10nm	<10nm	<10nm	<10nm
Read Latency	~1ns	~10ns	~10us	~50ns	<50ns	<50ns	~10ns	<10ns	<100ns	<50ns
Write Latency	~1ns	~10ns	>100us	>100us	<100ns	<100ns	<20ns	<3ns	<100ns	<100ns
Write Energy	~fJ	~10fJ	~10fJ	~100pJ	~10pJ	~pJ	~pJ	~pJ	~100fJ	~fJ
Endurance	> 1016	>1016	10 ⁴ -10 ⁵	10 ⁵	10 ⁶ -10 ⁸	10 ⁶ -10 ¹¹	10 ⁶ -10 ¹⁵	~1012	10 ⁹ -10 ¹²	10 ⁶ -10 ⁹
Retention	Volatile	Volatile	>10yr	>10yr	>10yr	>10yr	>10yr	>10yr	>10yr	>10yr

Source: Teman, et al., Wiley 2022 © Adam Teman, 2022

RRAM demos

[©] Adam Teman, 2022

Further Reading

- Rabaey, et al. "Digital Integrated Circuits" (2nd Edition)
- Elad Alon, Berkeley ee141 (online)
- Weste, Harris, "CMOS VLSI Design (4th Edition)" S. Kvatinsky "Memristors: Not Just Memory", ChipEx 2013
- S. Kvatinsky "Emerging Non-volatile Memories: Opportunities and Challenges", DevelopEx 2015
- Samira Khan, U. Virginia CS6354
- Daniele lelmini "Resistive switching memory for in-memory computing applications", ESSCIRC Tutorial 2020