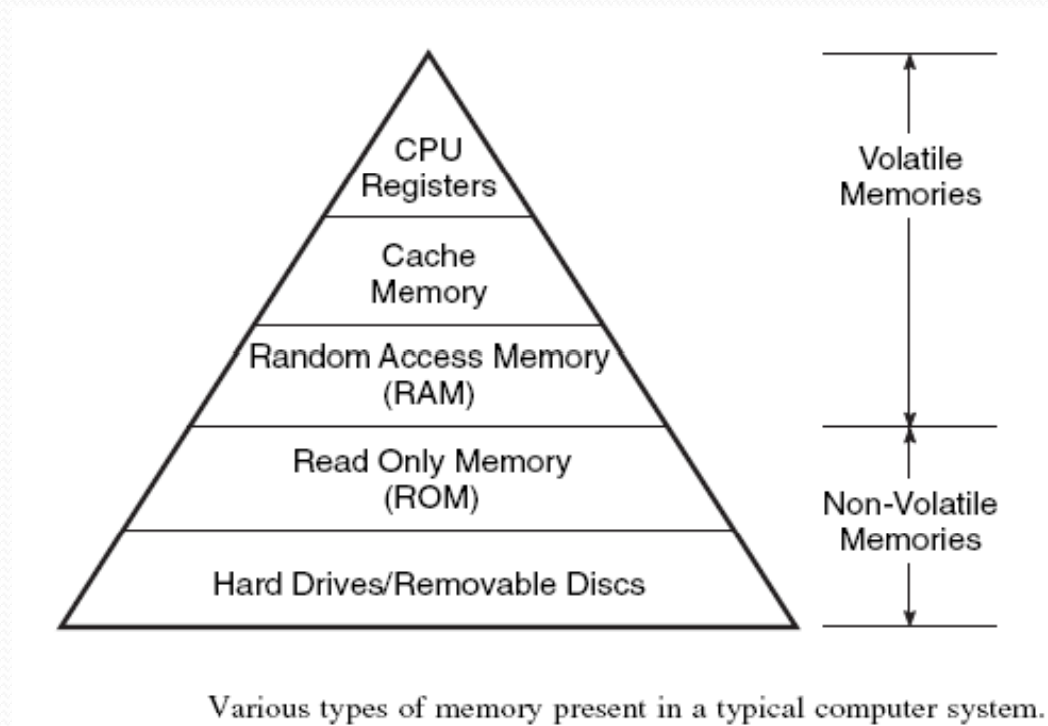
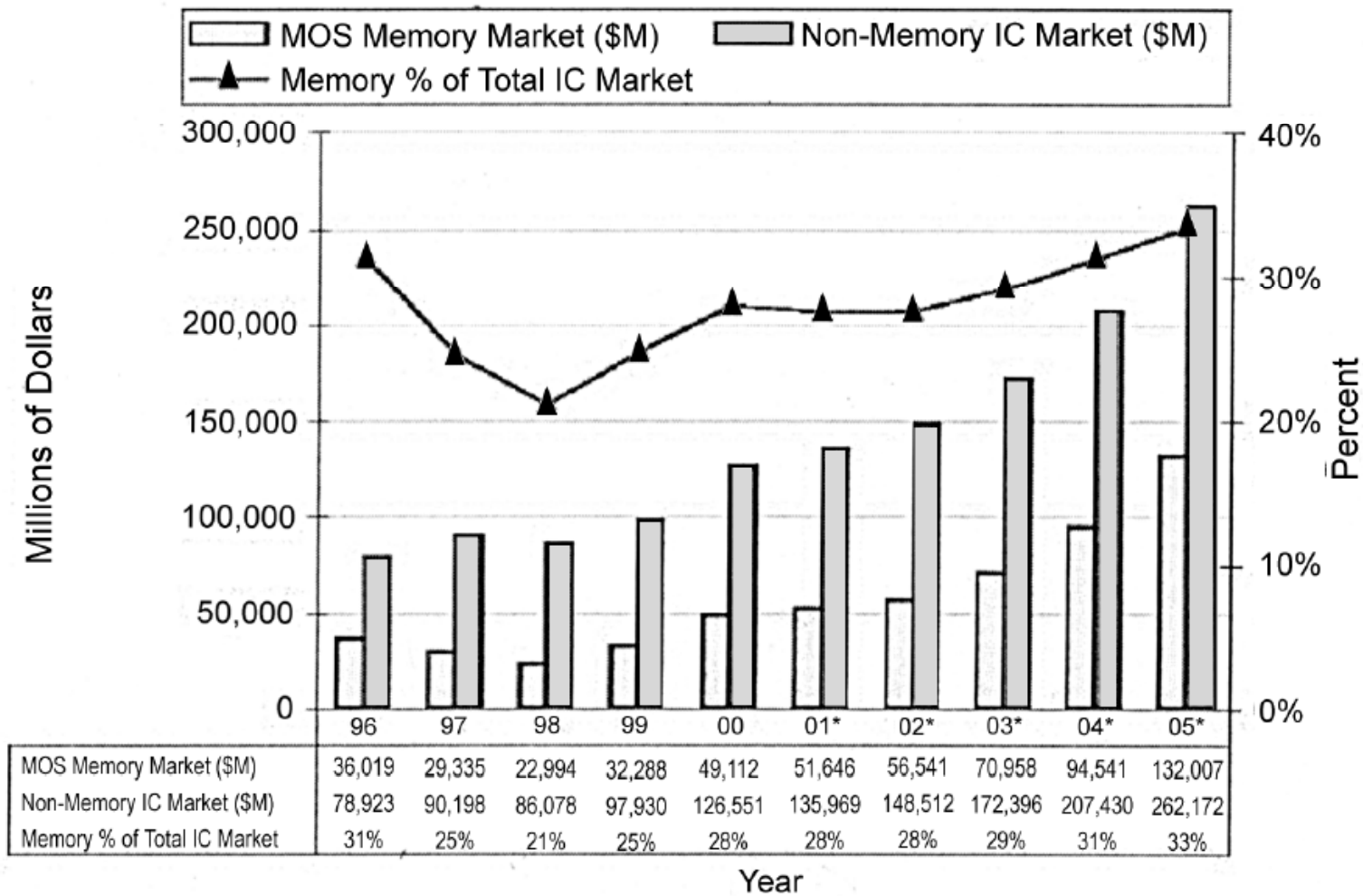
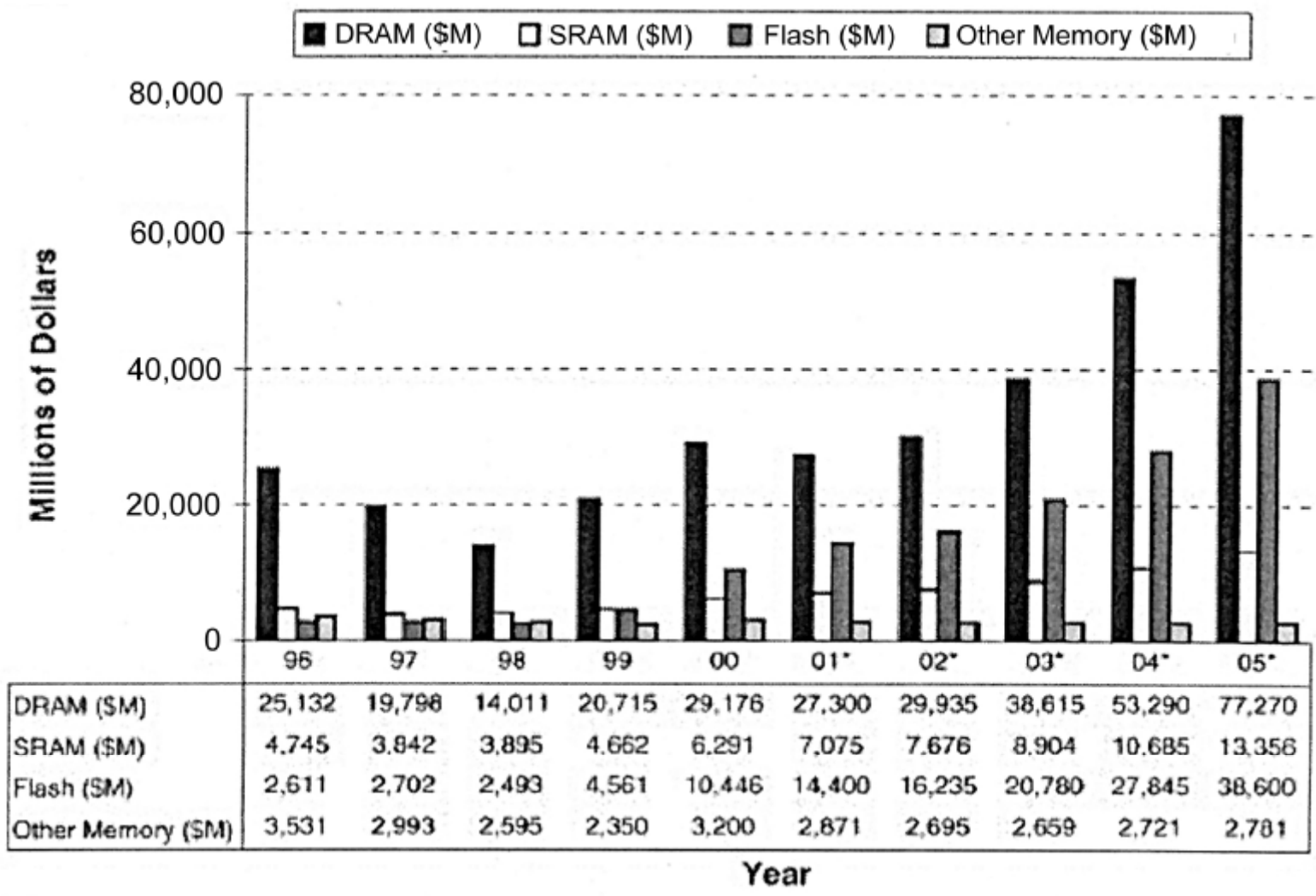


Memories

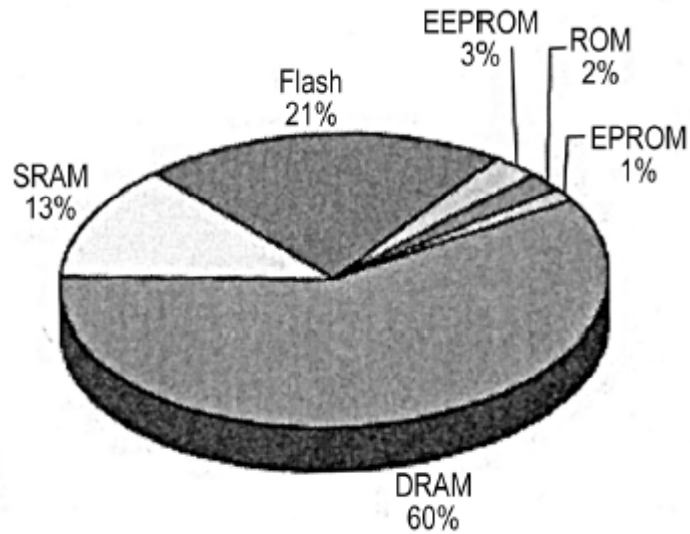




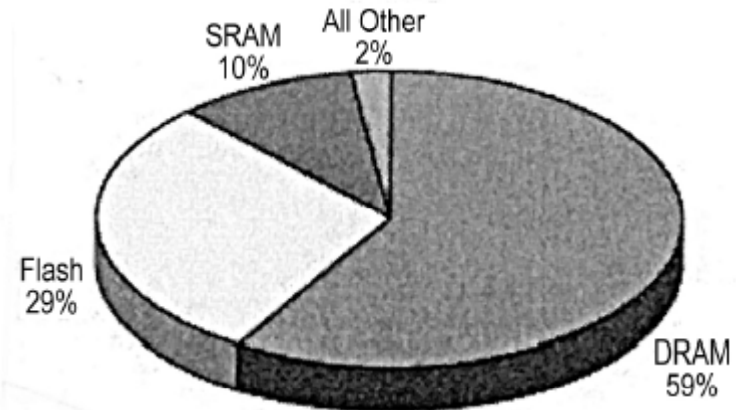
Semiconductor memory market as a percentage of the total IC market



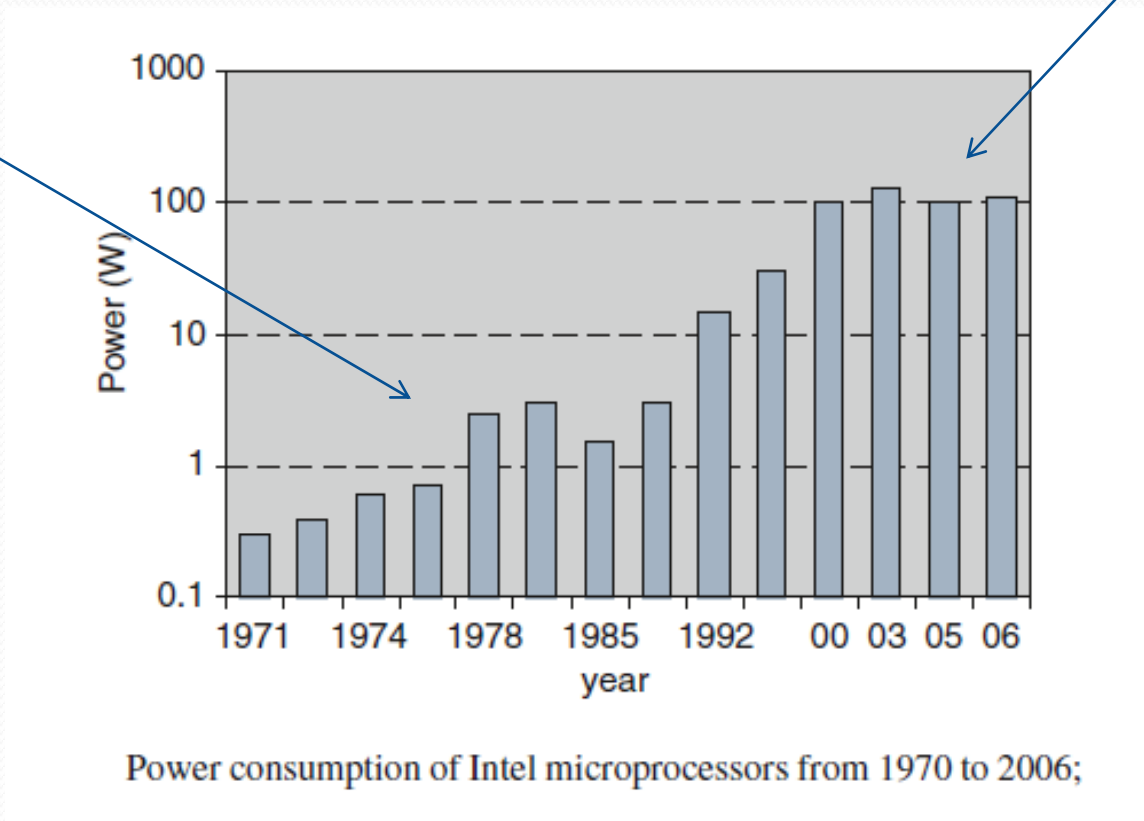
2000 MOS Memory Market
(\$49.1B)



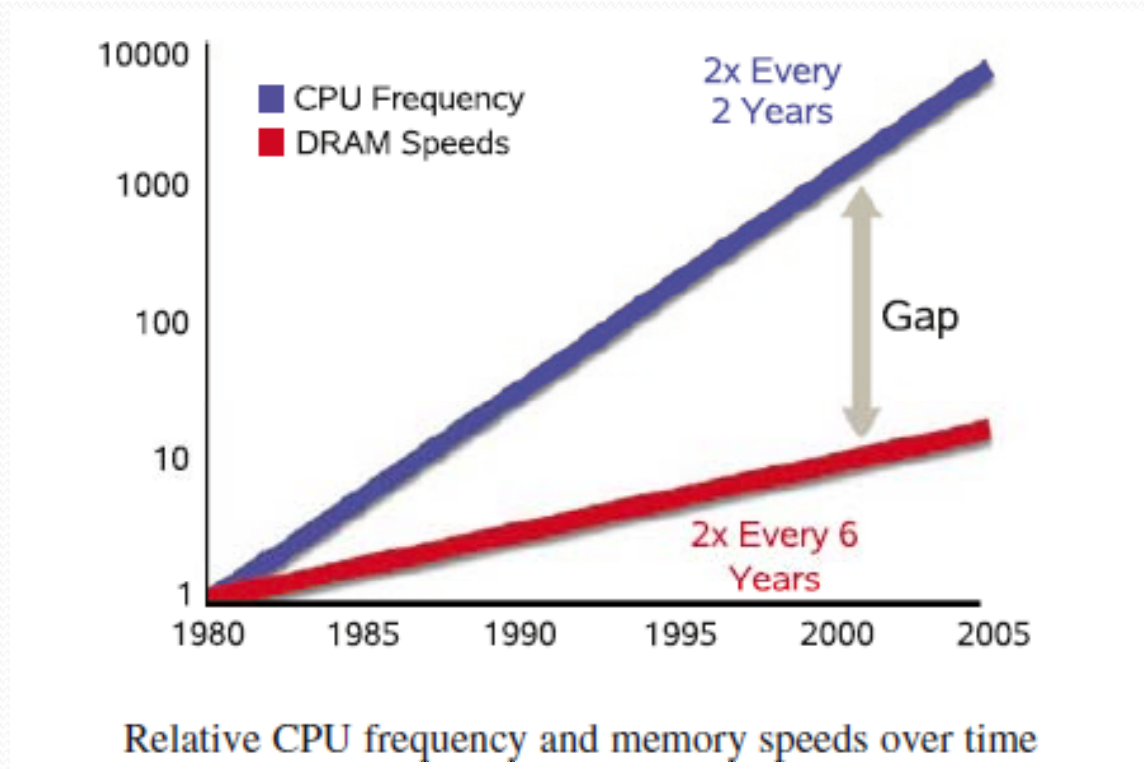
2005 MOS Memory Market
(\$132.0B)

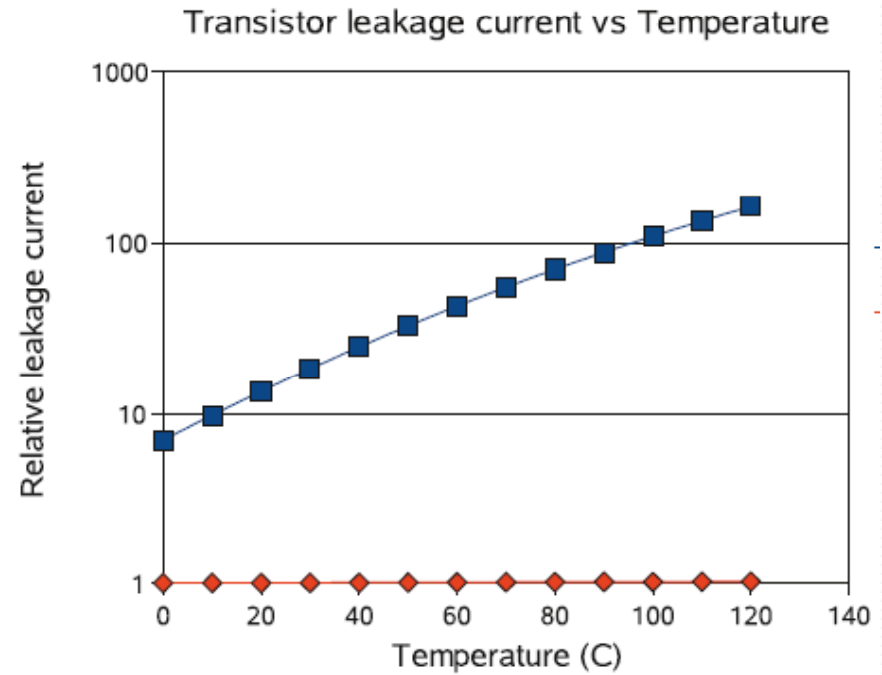
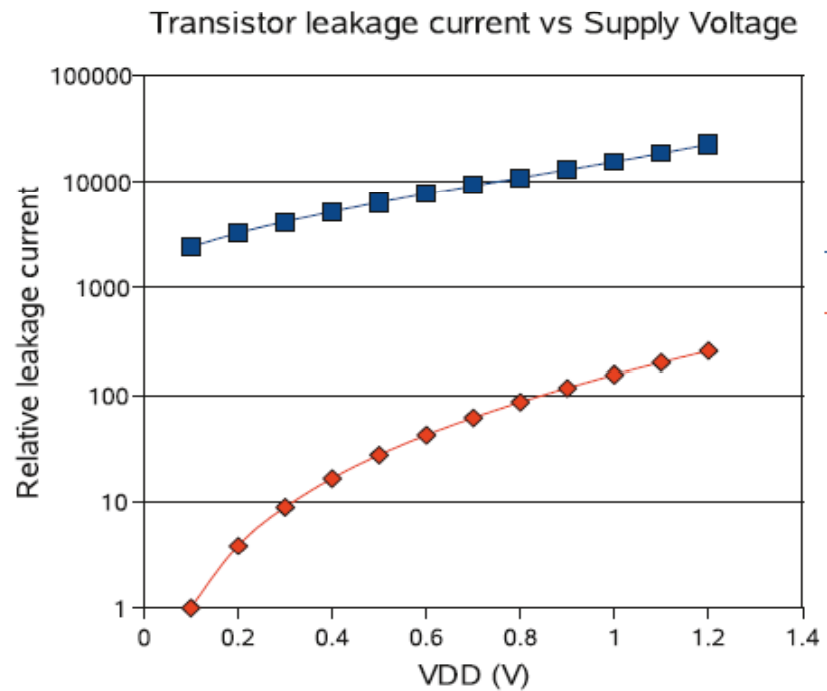


Moore's law

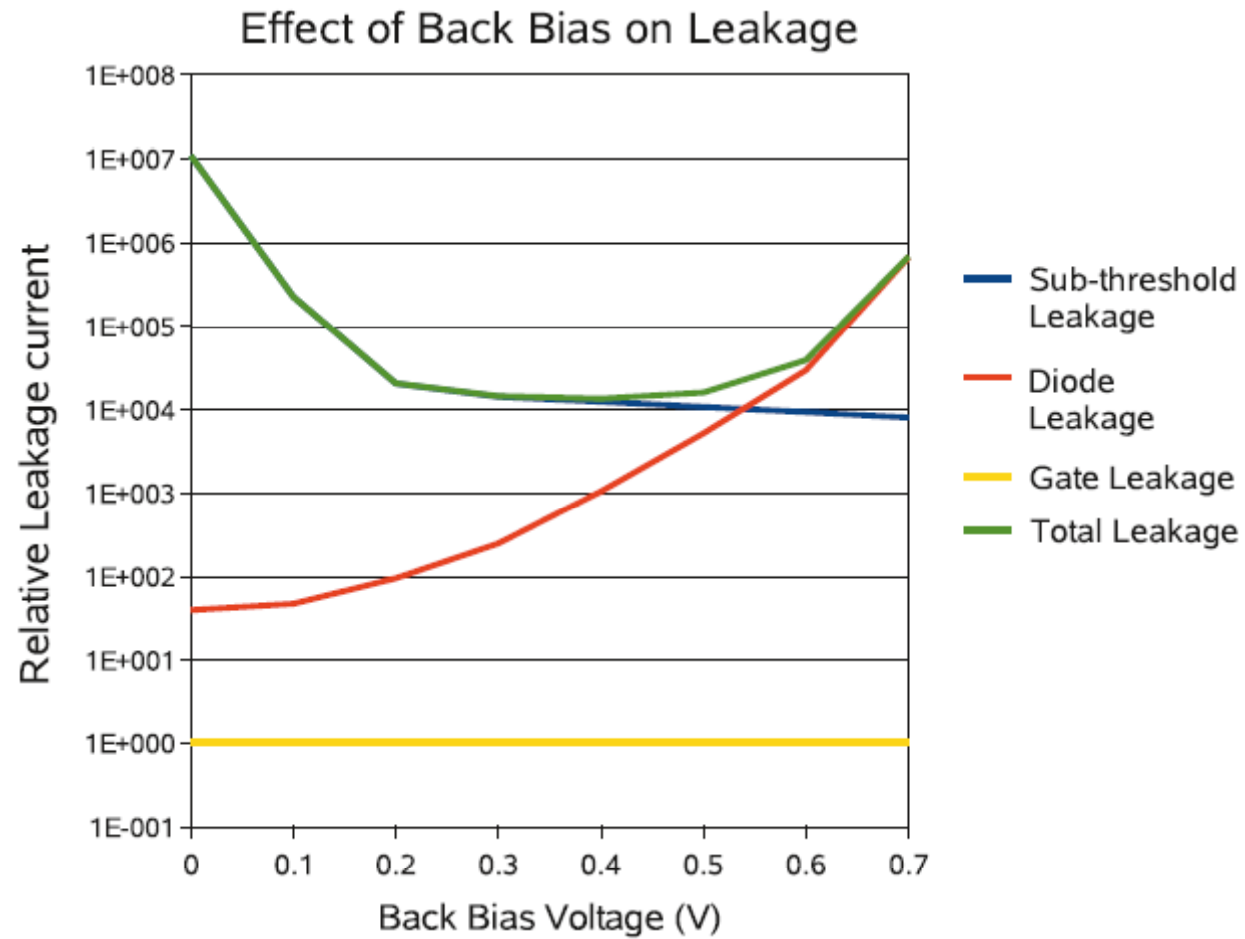


LV LP techniques





- Relative Drain Leakage
- ◆ Relative Gate Leakage



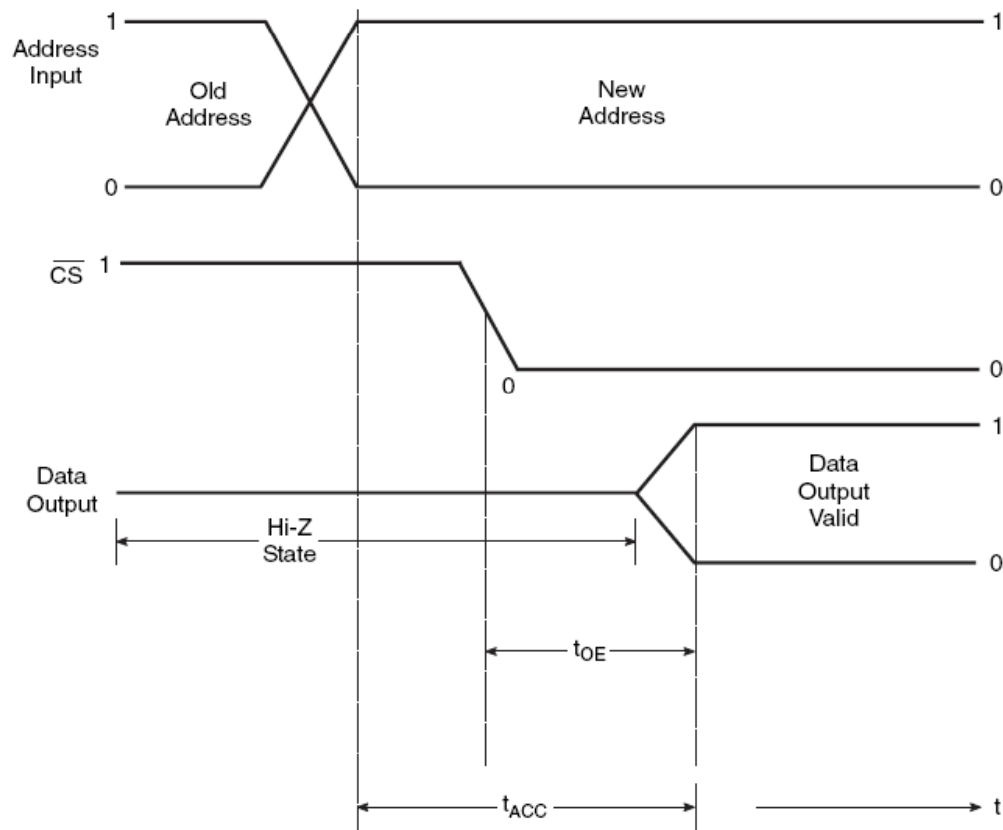
$$V_{th} = V_{th0} \pm \gamma * (|V_{SB}|)^{1/2}$$

Non-Volatile Memories (NVM's)

	Cell Phone	Consumer	Automotive	Computer & Communication
EPROM	Analog, Residential	Games, Set Top Box	Engine Mgt	HDD, Copiers, Fax, Switching
FLASH	Digital (GSM)	Set Top Box PDA	All Power Train Car Navigation, ABS, GPS	HDD, PC Bios, CDROM
EEPROM	Digital (GSM)	Audio, Video	All Car Body	PC SPD, Graphic boards, Printers



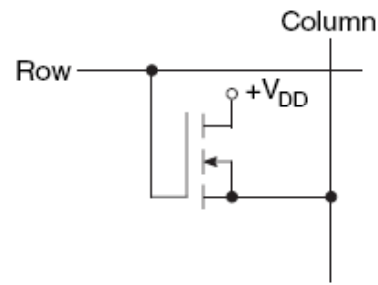
Read-Only Memory (ROM)



Typical timing diagram of a ROM READ operation.

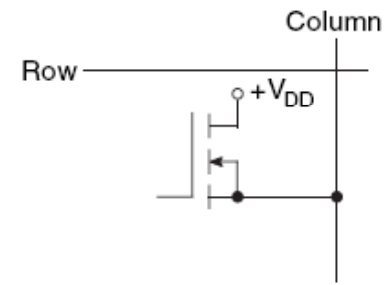
Mask-programmed ROM

- ROM is programmed at the manufacturer's site according to the specifications of the customer
- economical only when manufactured in large quantities
- once programmed, it cannot be reprogrammed.
- The basic storage element is an NPN bipolar transistor, connected in common-collector configuration, or a MOSFET in common-drain configuration.



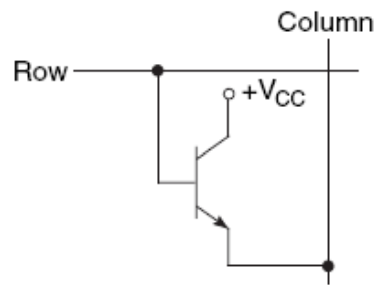
(a)

stored 1

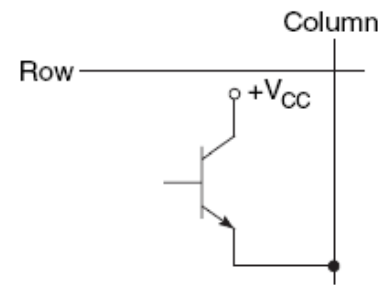


(b)

stored 0

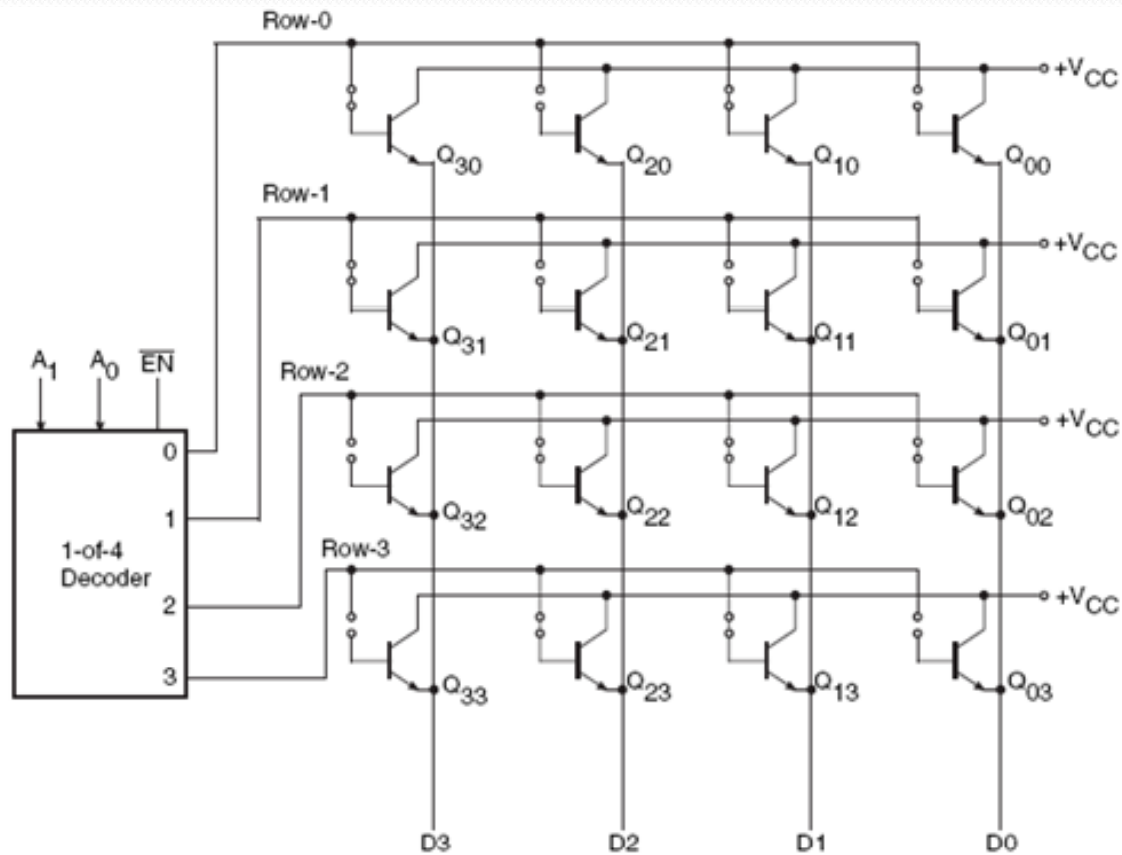


(c)



(d)

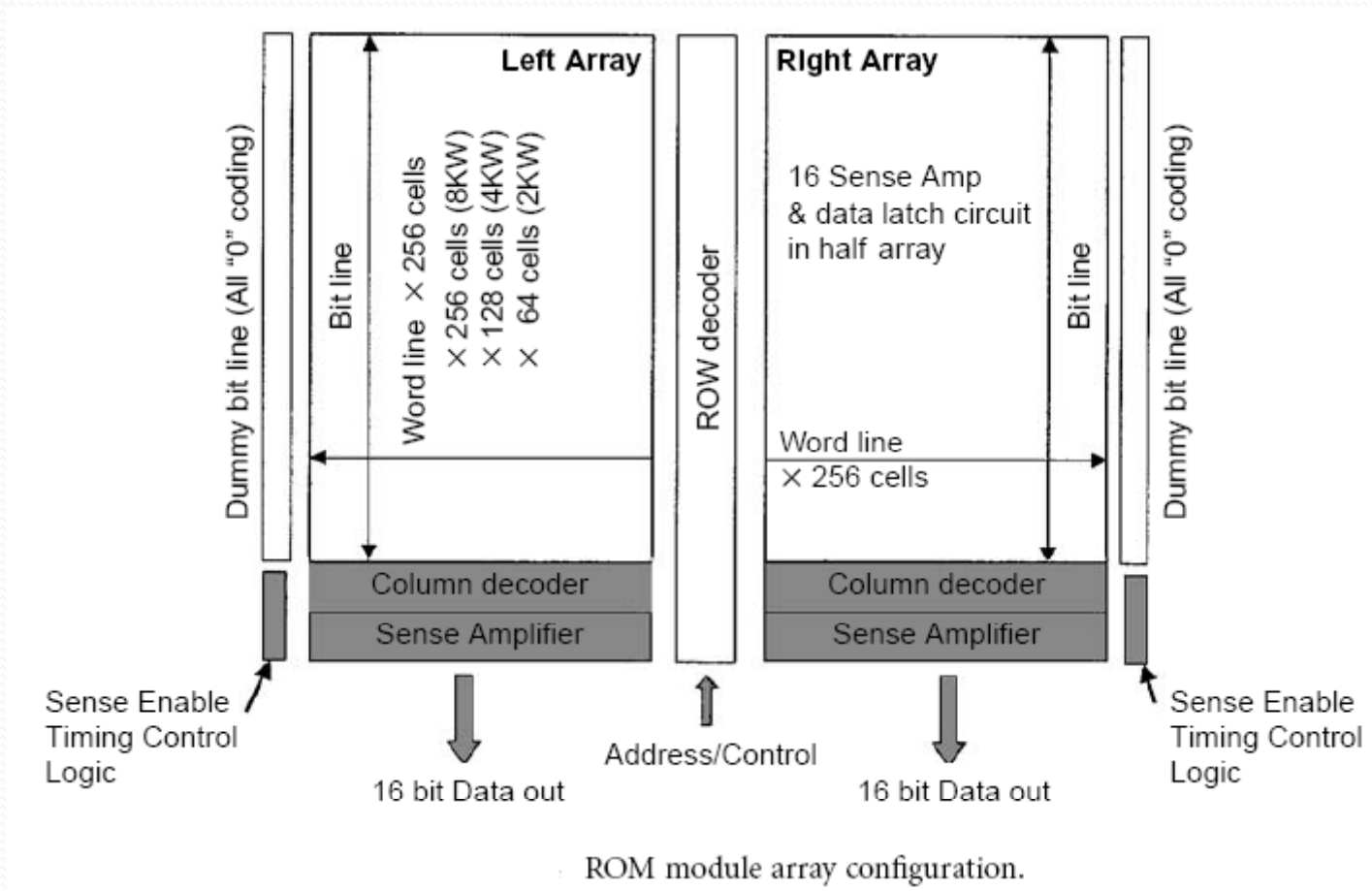
Basic cell connection of a mask-programmed ROM.

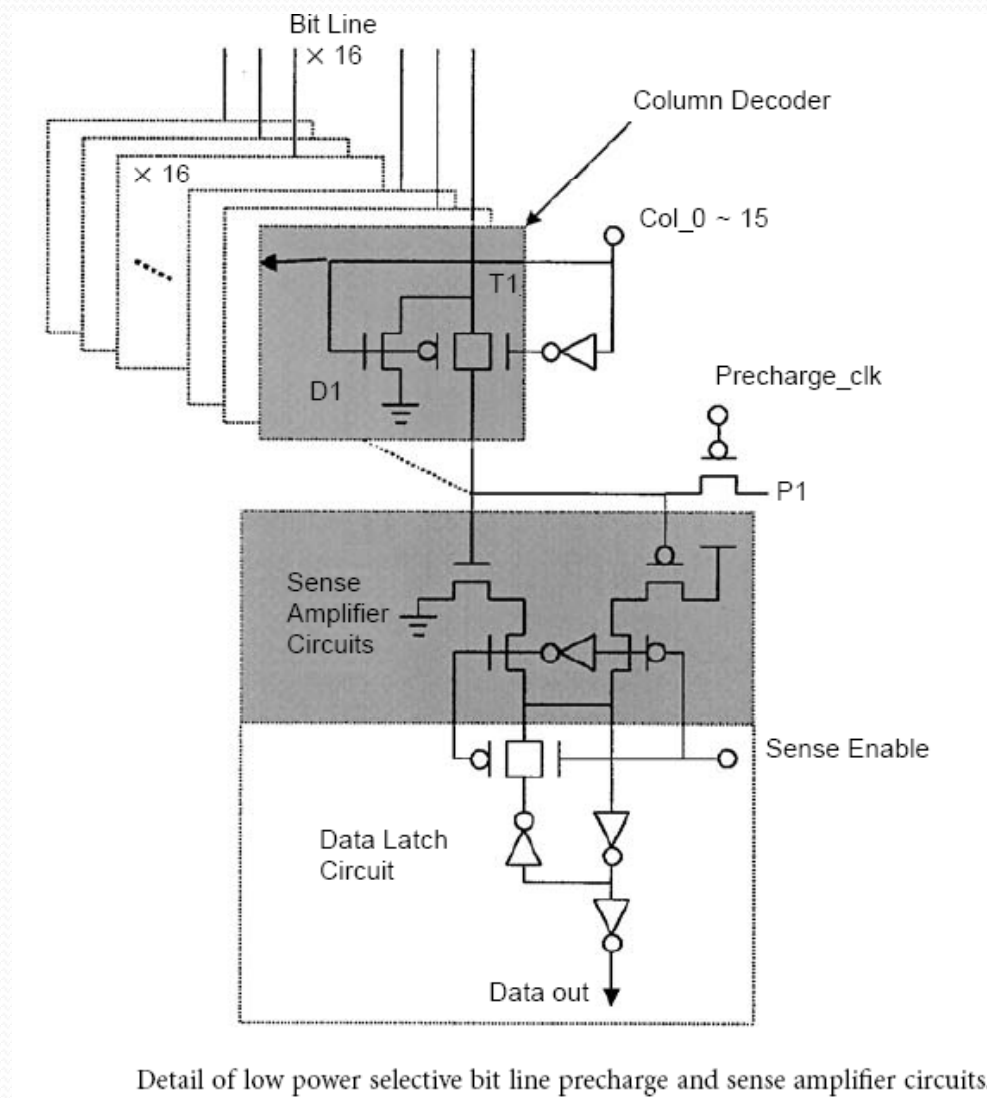


Internal structure of a 4×4 bipolar mask-programmed ROM.

Truth Table

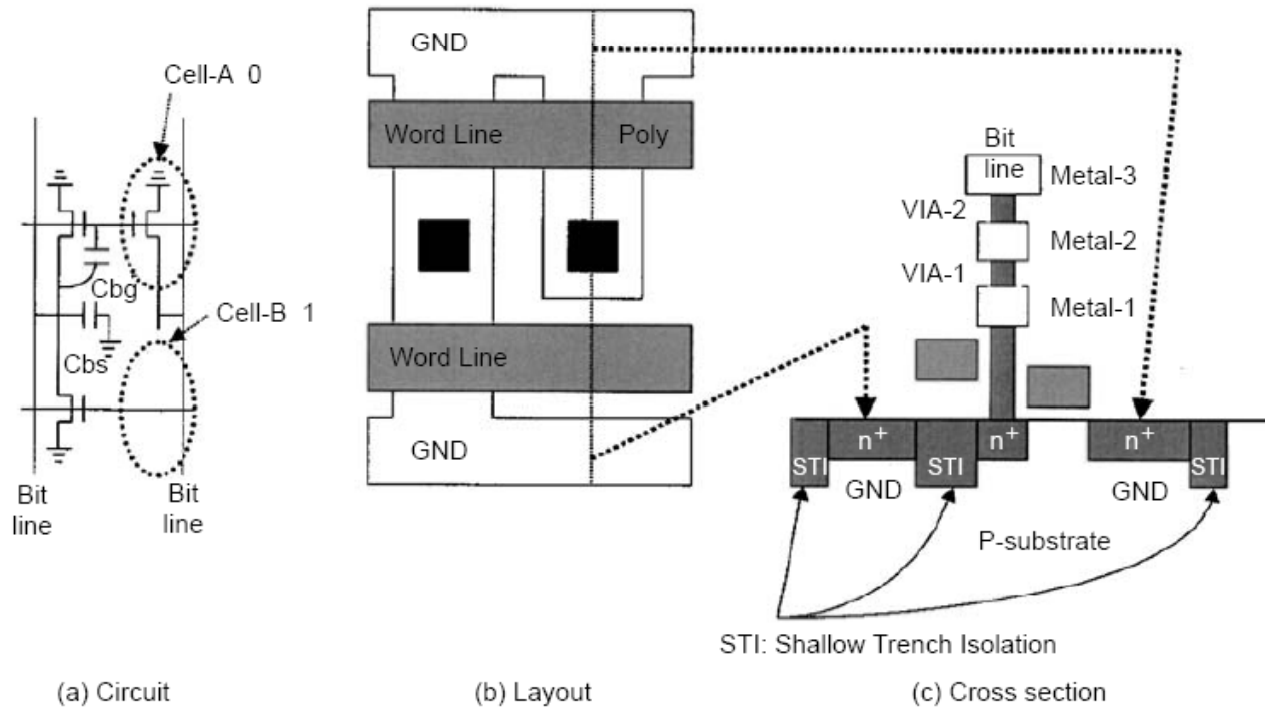
Address		Data			
A_1	A_0	D_3	D_2	D_1	D_0
0	0	1	0	1	0
0	1	1	0	0	0
1	0	0	1	1	1
1	1	0	1	1	0





Detail of low power selective bit line precharge and sense amplifier circuits.

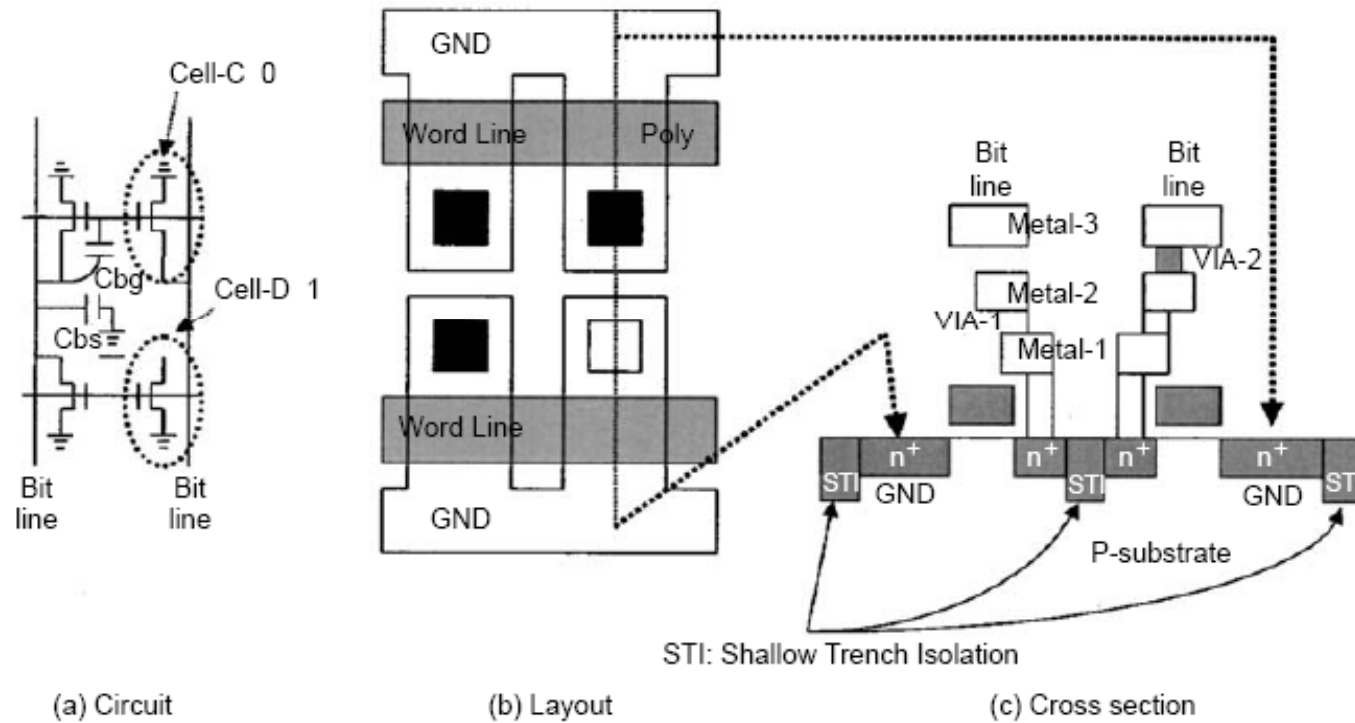
Diffusion-programming ROM




Diffusion programming ROM.

- highest density: bit line contact to discharge transistor can be shared by two-bit cells
- very long fabrication cycle time: diffusion programming at early process stage

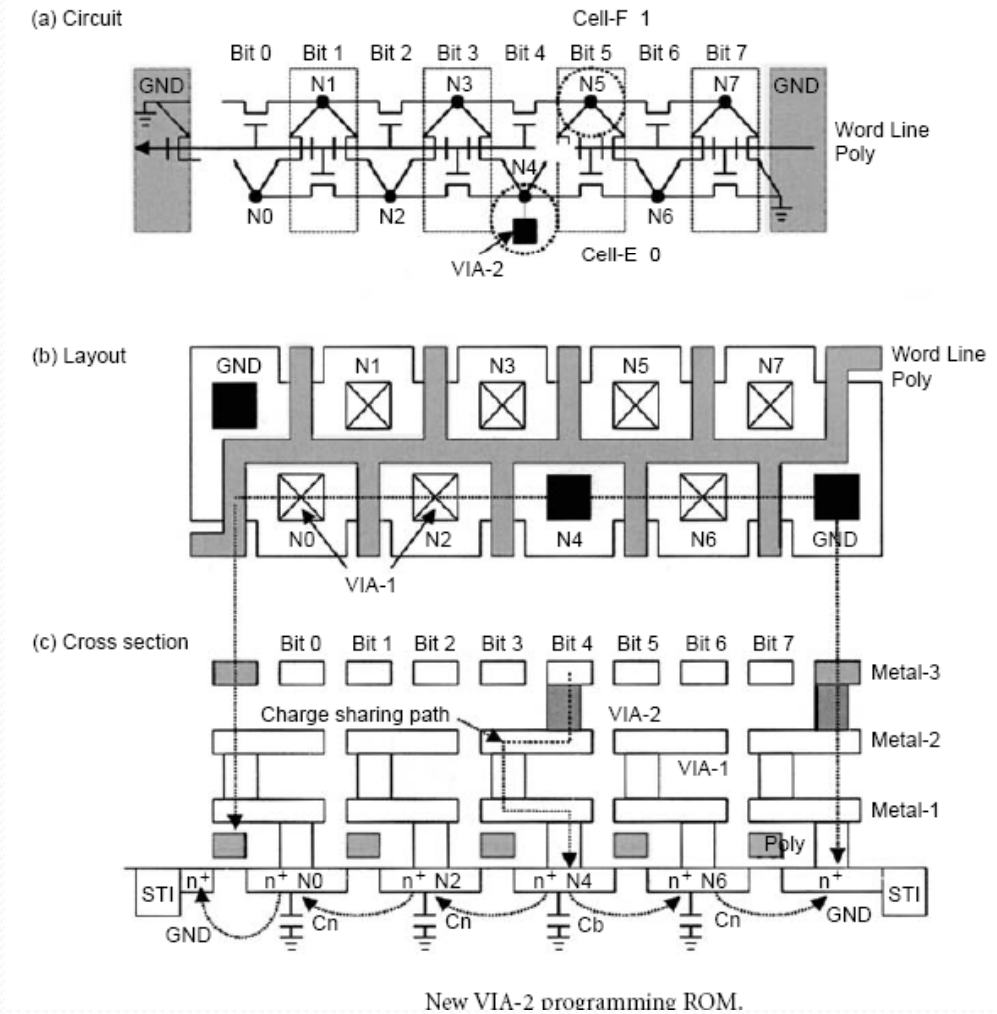
VIA-2 Contact Programming ROM (conventional)



Conventional VIA-2 programming ROM.

- 
- The VIA-2 is final stage of process and base process can be completed just before VIA-2 etching and remaining process steps are quite few
 - VIA-2 ROM fabrication cycle time is about 1/5 of the diffusion ROM
 - Drawback: poorer density \Rightarrow diffusion area and contact must be separated in each ROM bit cell

VIA-2 Contact Programming ROM (new)



- 8-bit block with GND on each side
- higher density

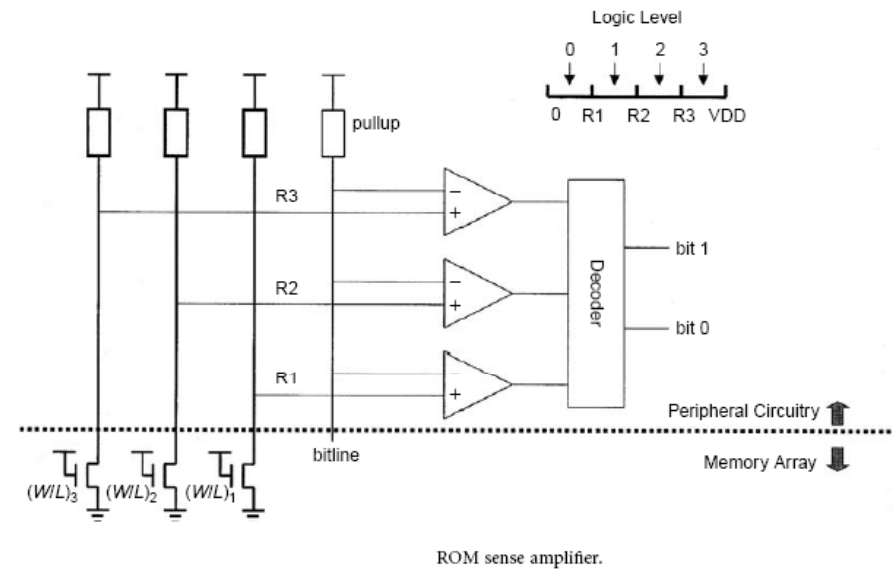
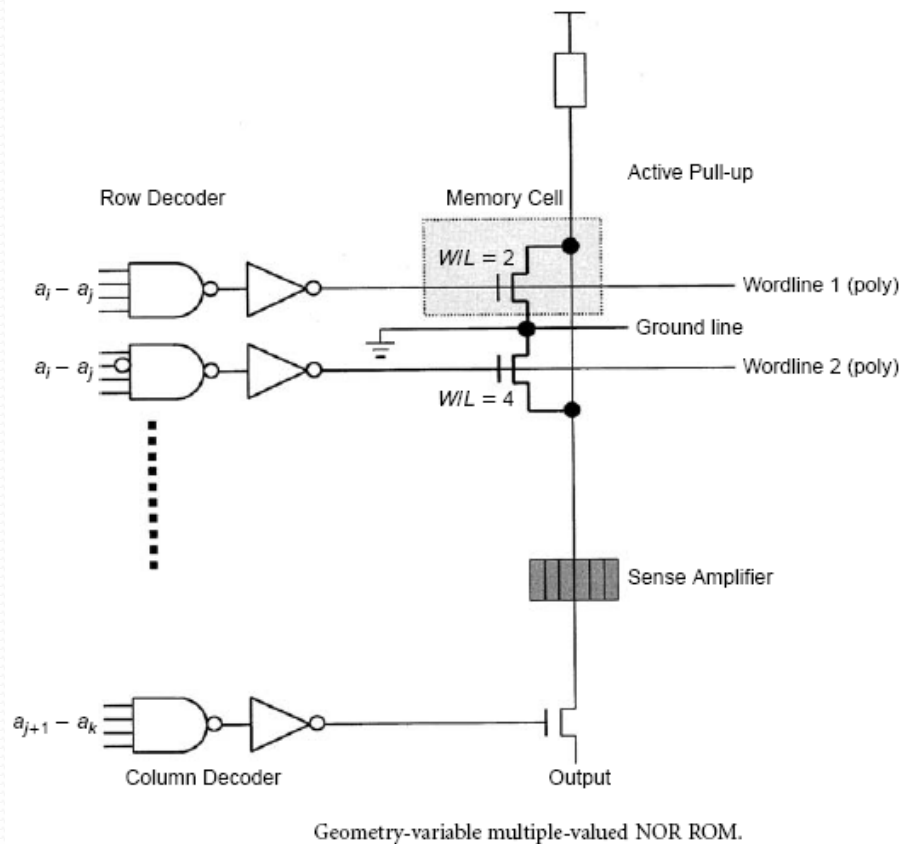
TABLE 51.1 Comparison of ROM Performance

Comparison Item	Diffusion ROM	Conventional VIA-2 ROM	New VIA-2 ROM
8KW (Area ratio)	1.0	1.2	1.04
TAT (Day ratio)	1.0	0.2	0.2
Speed @ 2.13 V, 125dc. Weak.	83 MHz	86 MHz	123 MHz
Speed @ 2.13 V, 125dc. Typical.	166 MHz	98 MHz	149 MHz
Speed @ 2.81 V, -40dc. Strong.	277 MHz	179 MHz	294 MHz
Speed @ 1.66 V, 125dc, Typical.	103 MHz	75 MHz	106 MHz
Power@2.81 V, -40dc. Strong. 100 MHz. (16-bit single access)	15.6 mW	19.3 mW	2 UrnW
Power@2.81 V@40dc. Strong. 100 MHz. (32-bit dual access)	29.6 mW	37.1 mW	401 mW

Performance was measured with worst coding (all coding "1").

Multi-level ROM

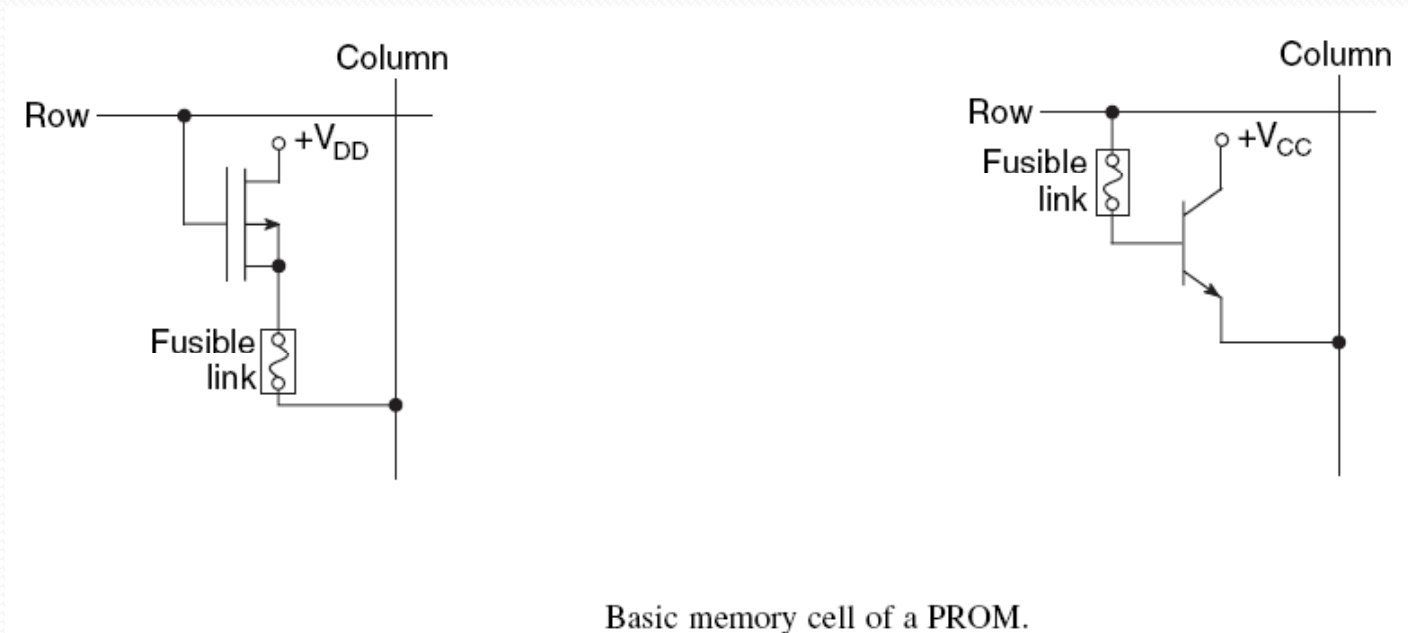
- transistor-cell (W/L) adjusted according to its logic state



- 2 bits per cell
- Intel (1980's)

Programmable ROM (PROM)

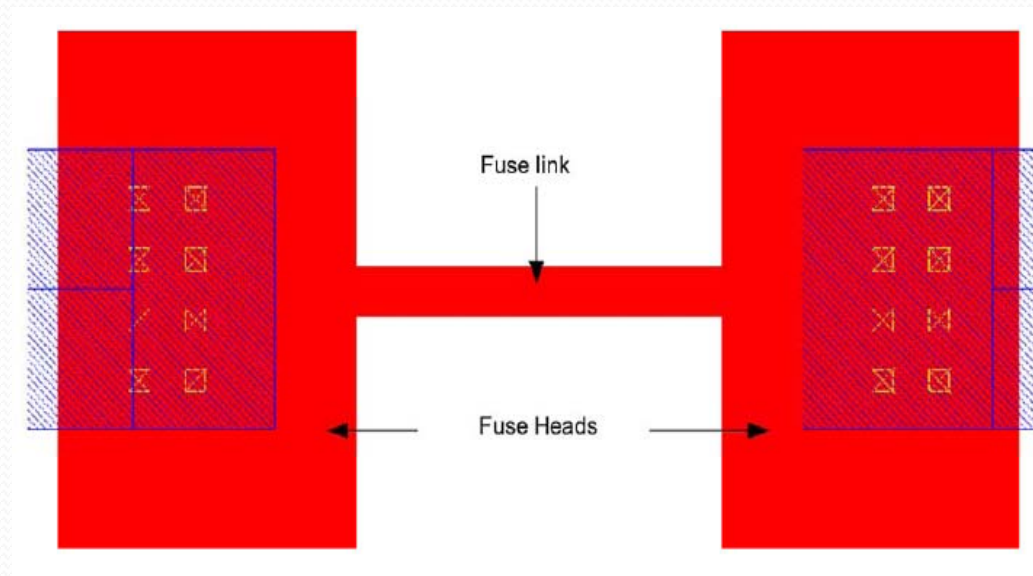
- programmed by the customer
- fuse links: metal or poly-Si



Basic memory cell of a PROM.

poly-Si/metal fuses

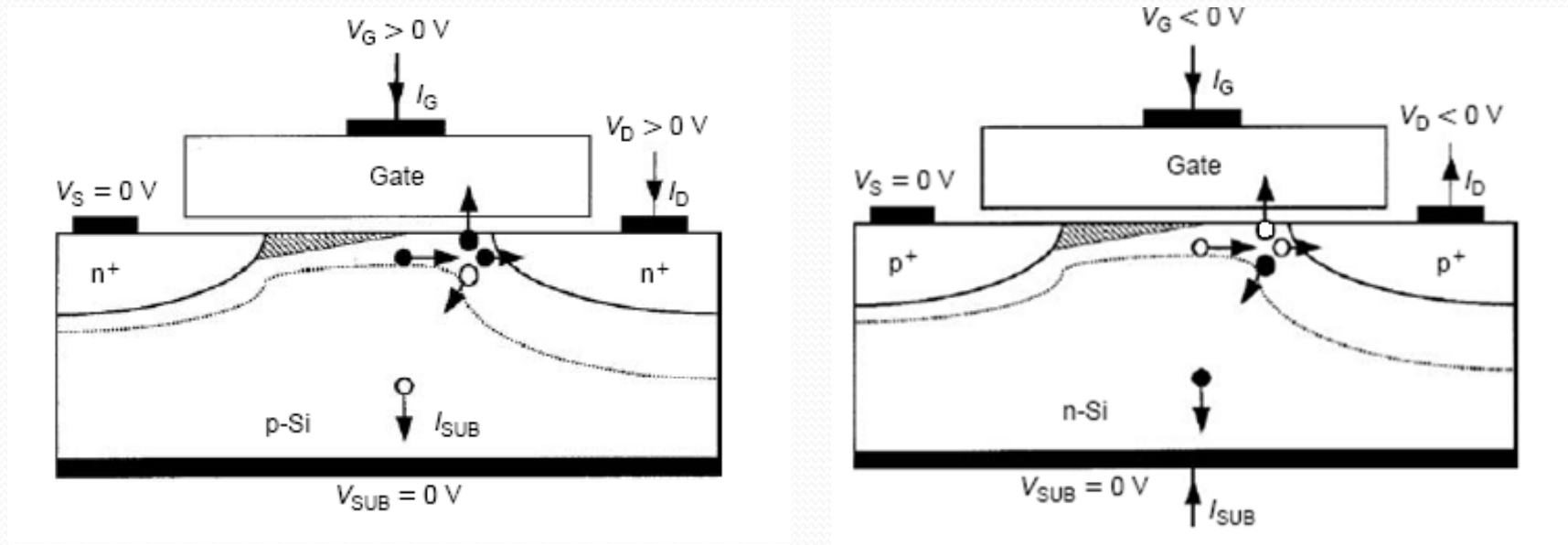
- line heats up to fusion by self-heating mechanism due to current intensity during programming mode
- fusion increases its resistance and eventually opens the link



polyfuse/metalfuse layout

Stacked-Gate Non-Volatile Memory (EPROM/EEPROM)

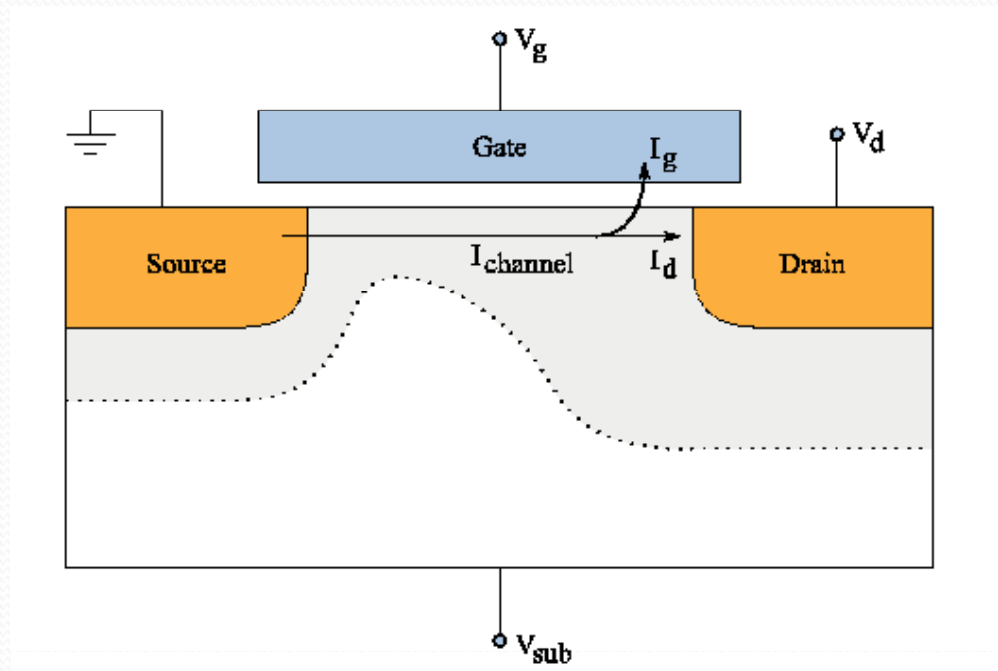
- **DAHC: Drain-Avalanche Hot Carrier**



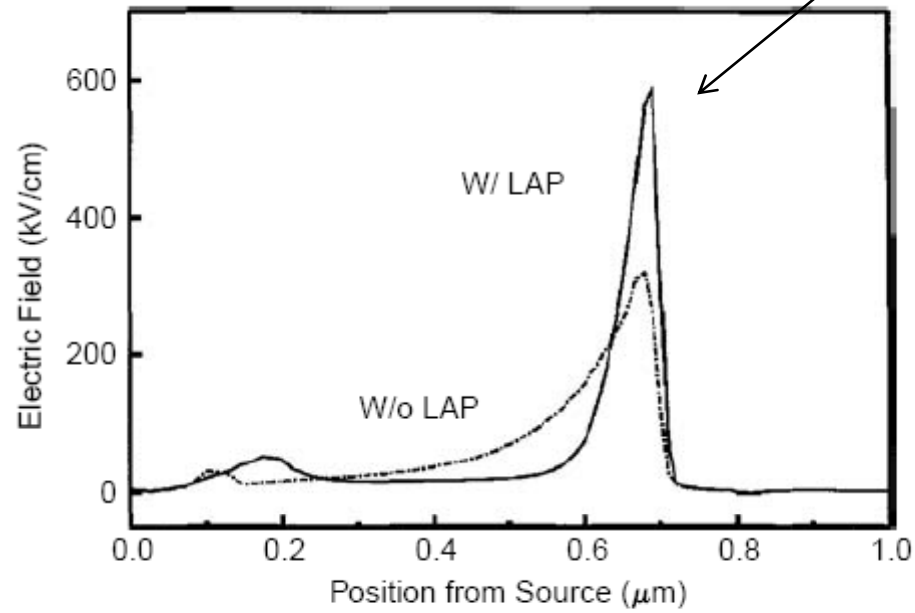
Schematic illustration of the channel hot carrier effect in (a) n-channel MOSFET, and (b) p-channel MOSFET.

• HCI: Hot Carrier Injection

- either an electron or hole gains sufficient kinetic energy to overcome a potential barrier necessary to break an interface state
- can be injected into the gate dielectric, where they can get trapped



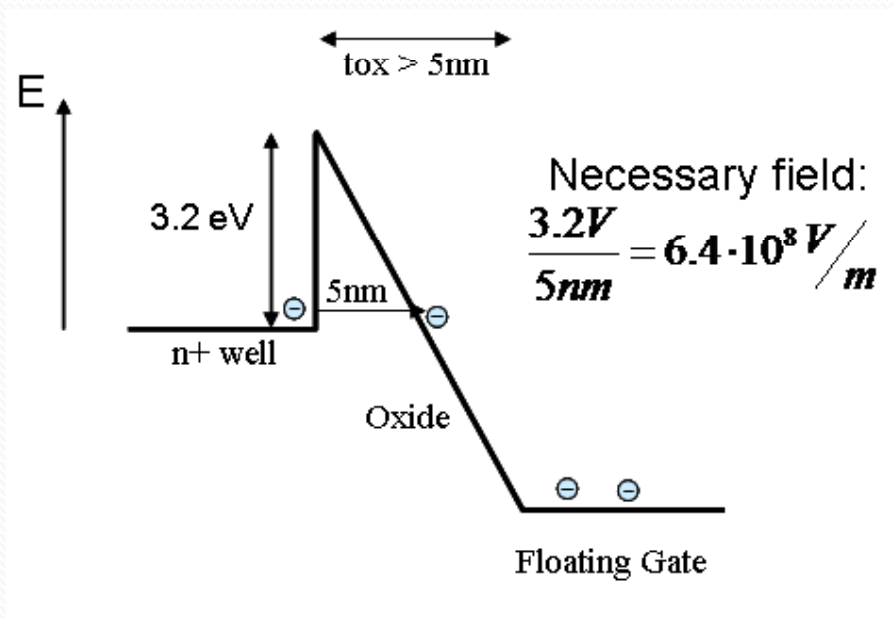
pinch-off region



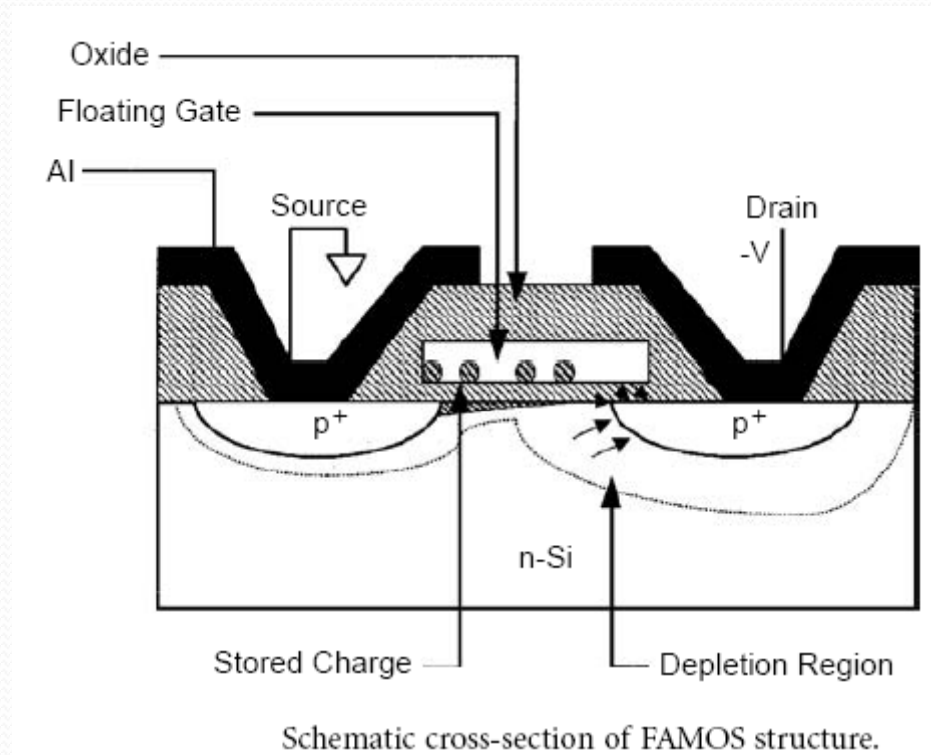
Simulated electric field along the channel in the n-channel MOSFET.

• Fowler-Nordheim tunneling

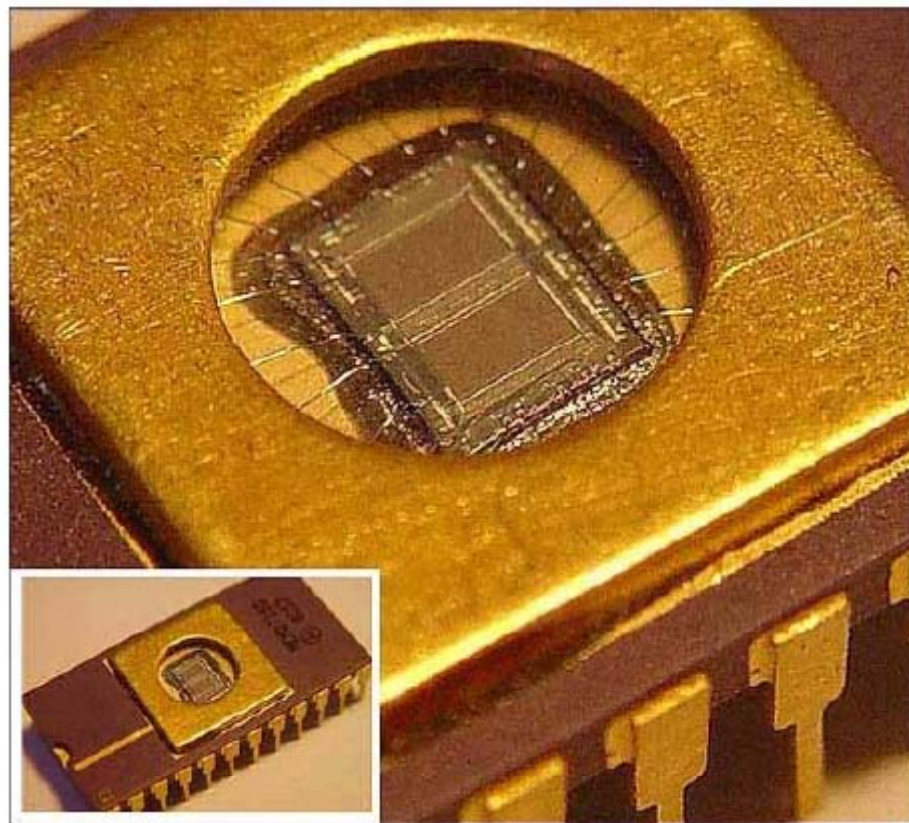
- current flowing across MOS structure at the high electric field in the oxide; electrons tunnel from semiconductor conduction band into oxide conduction band through part of the potential barrier at the semiconductor-oxide interface (oxide 5-10 nm thick)



- **floating-gate MOS (FAMOS) transistor**



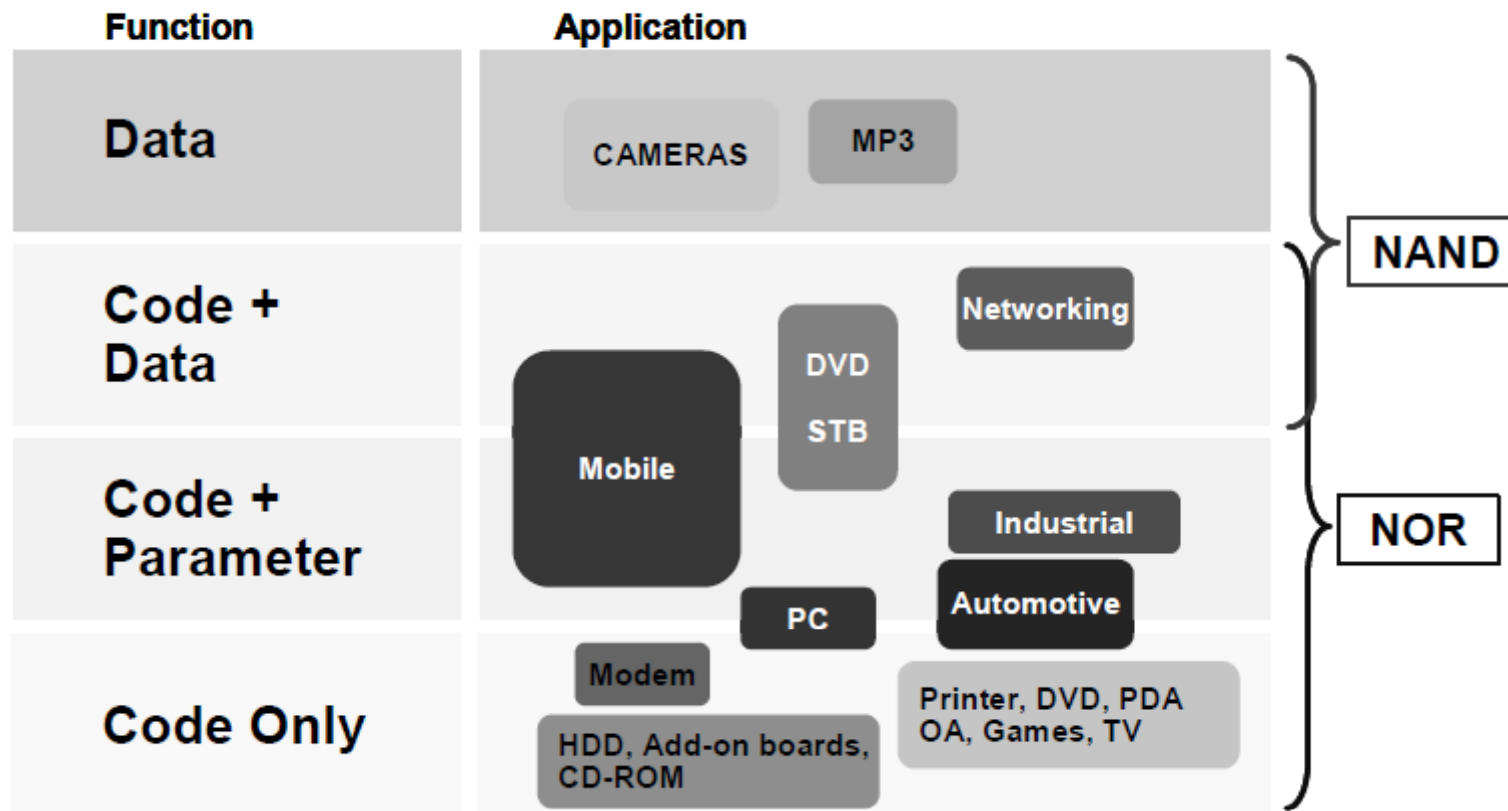
- not electrically erasable (only by time-consuming UV-irradiation)
- compatible with double-poly processes



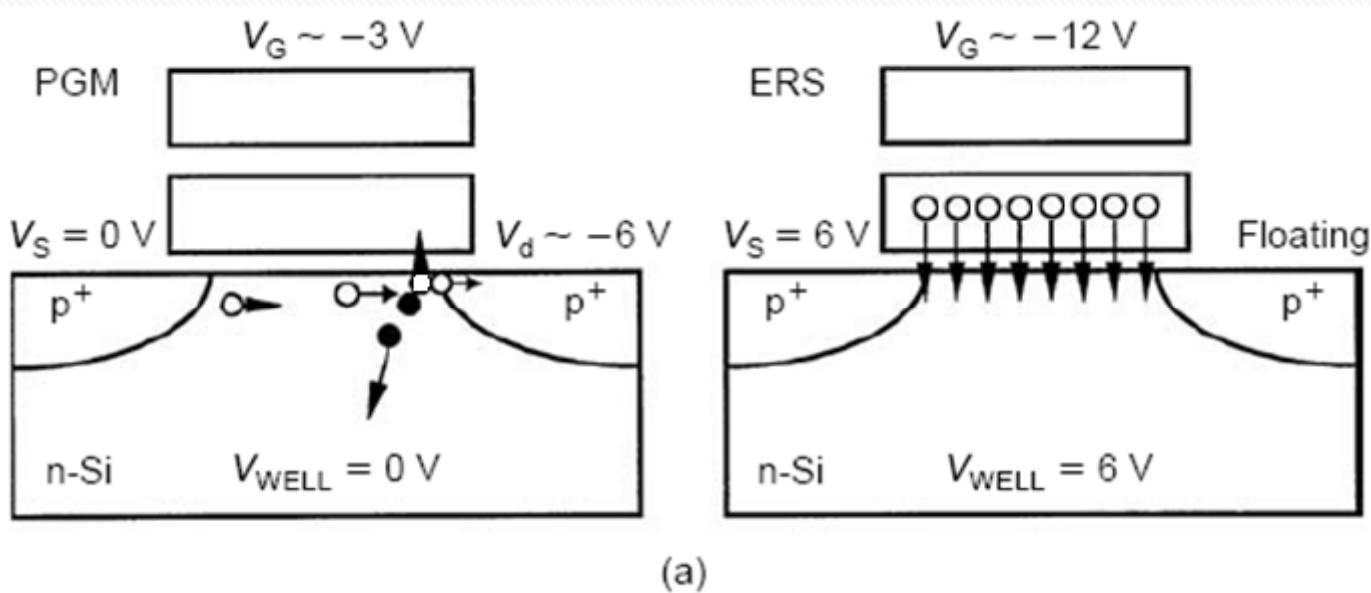
EPROM with UV-irradiation erasing

Flash Memories

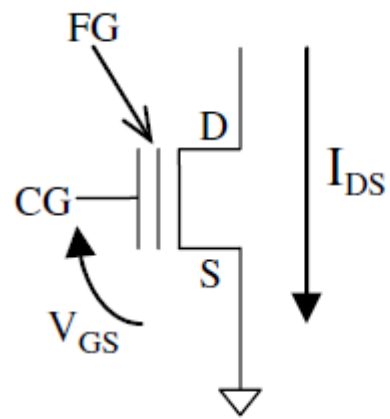




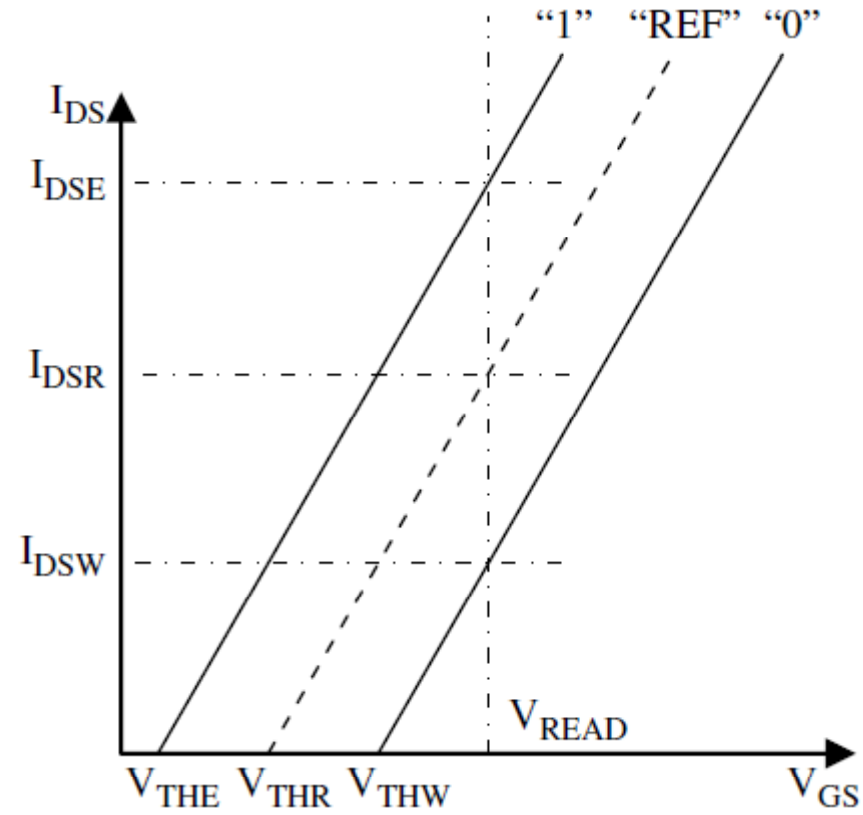
Main applications of flash memories



Different p-channel Flash write/erase operations (a) programming with DAHC and erase with FN tunneling action



CG: Control Gate
 FG: Floating Gate
 D: Drain
 S: Source



Threshold voltage distribution
for the erased cells:

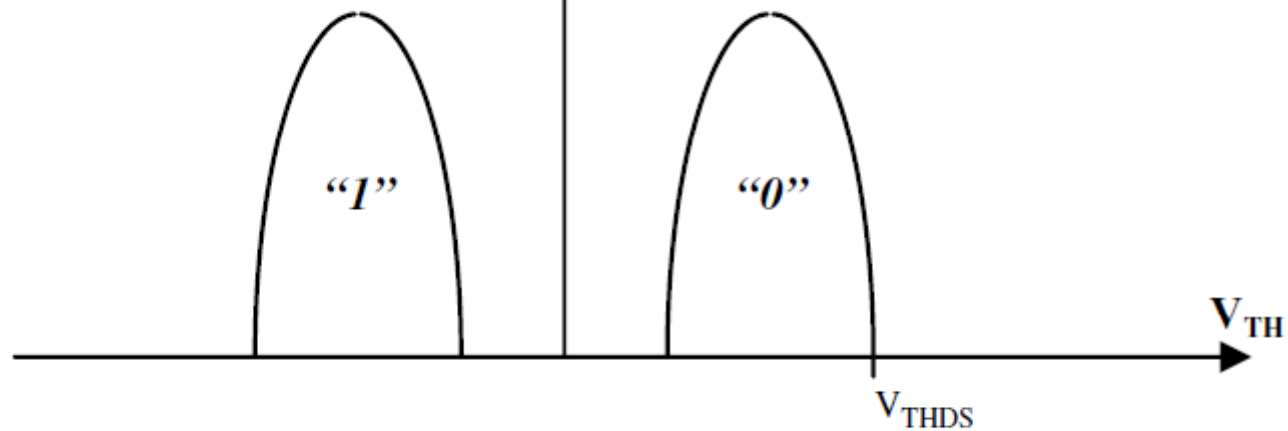
$$-V_{TH} < 0$$

$$-I > 0 @ WL = 0V$$

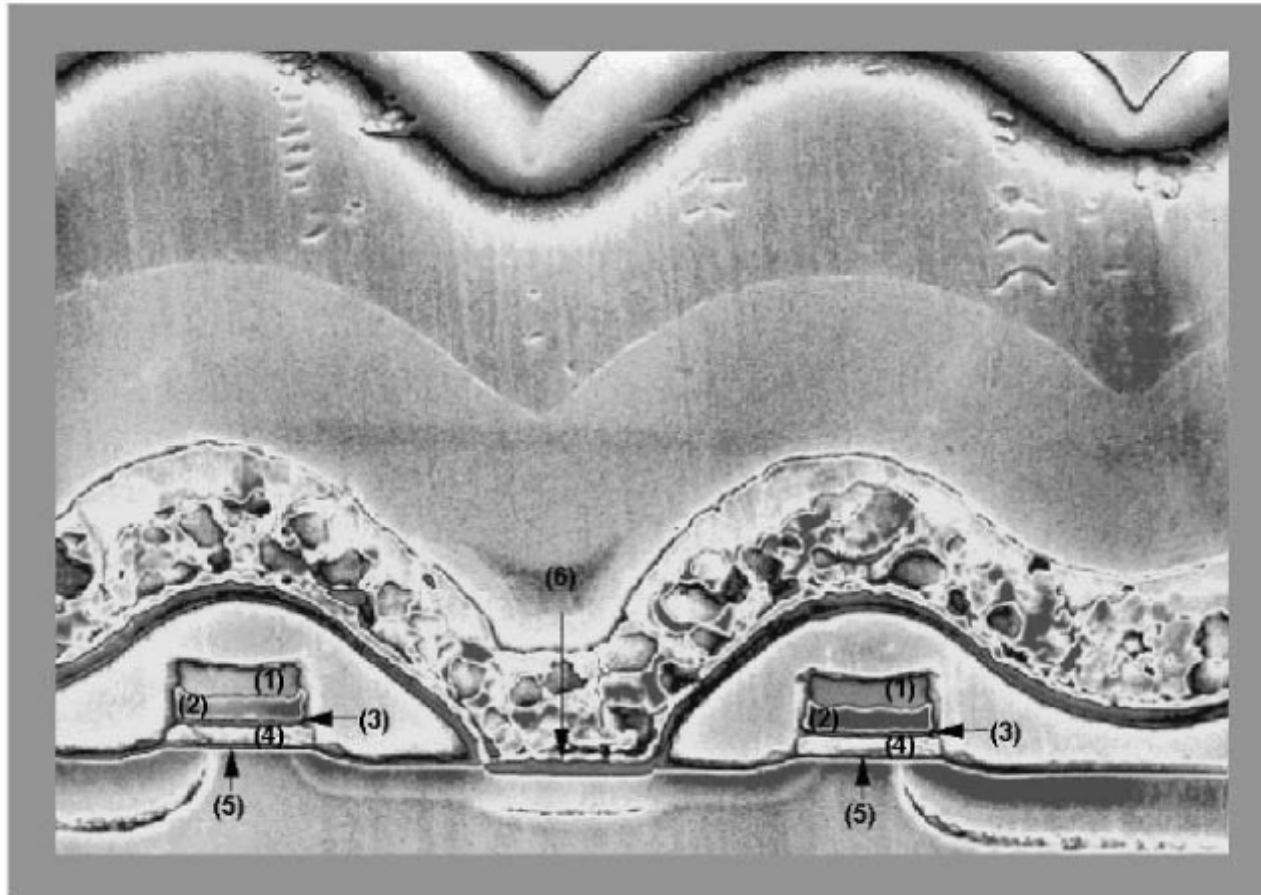
Threshold voltage distribution
for the programmed cells:

$$-V_{TH} > 0$$

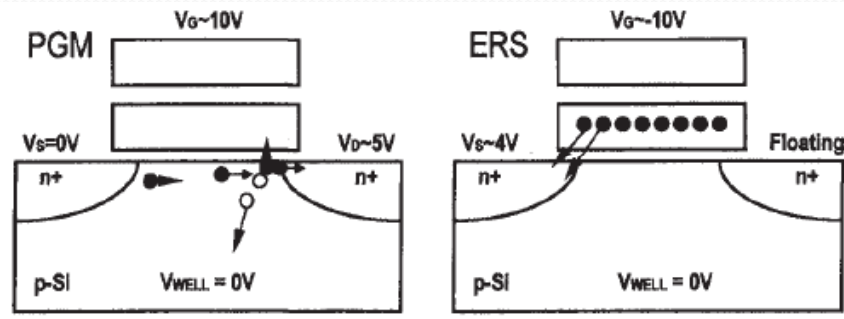
$$-I = 0 @ WL = 0V$$



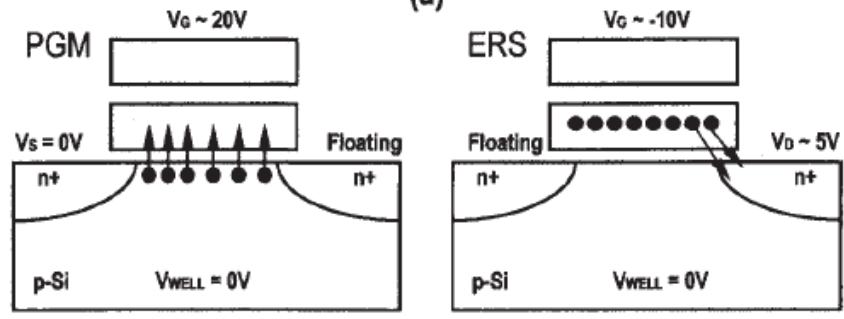
Threshold voltage distributions for erased and programmed cells



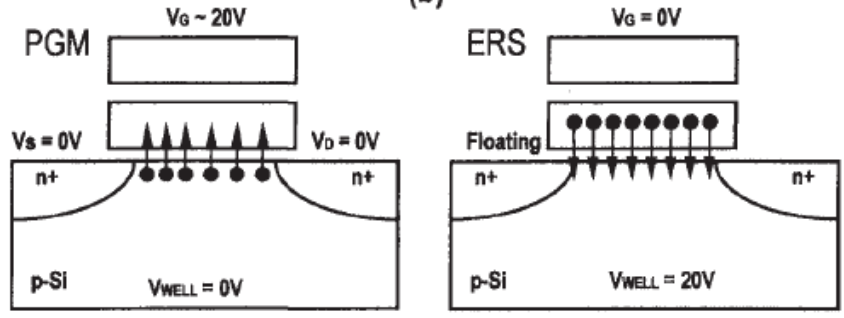
Cross-section of a Flash Memory device; silicide (1) and polysilicon (2) which constitute the control gate, the interpoly oxide layer (3), the floating gate (4), the thin oxide (5), drain contact, shared by the two cells and connected to Metal 1 (the granular wire)(6) and source implant, deeper than the drain implant



(a)



(b)



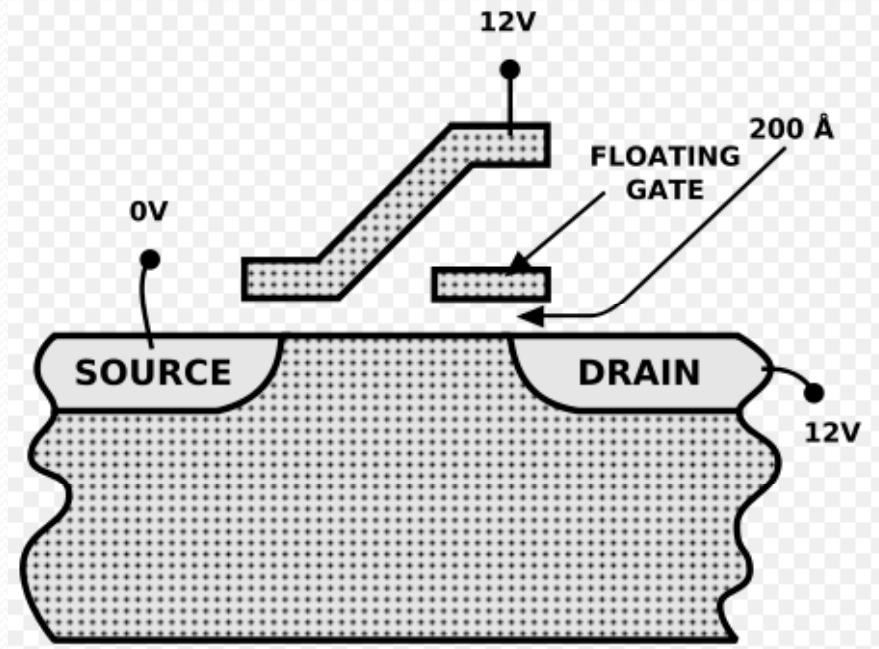
(c)

(b) PGM: DAHC
ERS: FN tunneling (source)

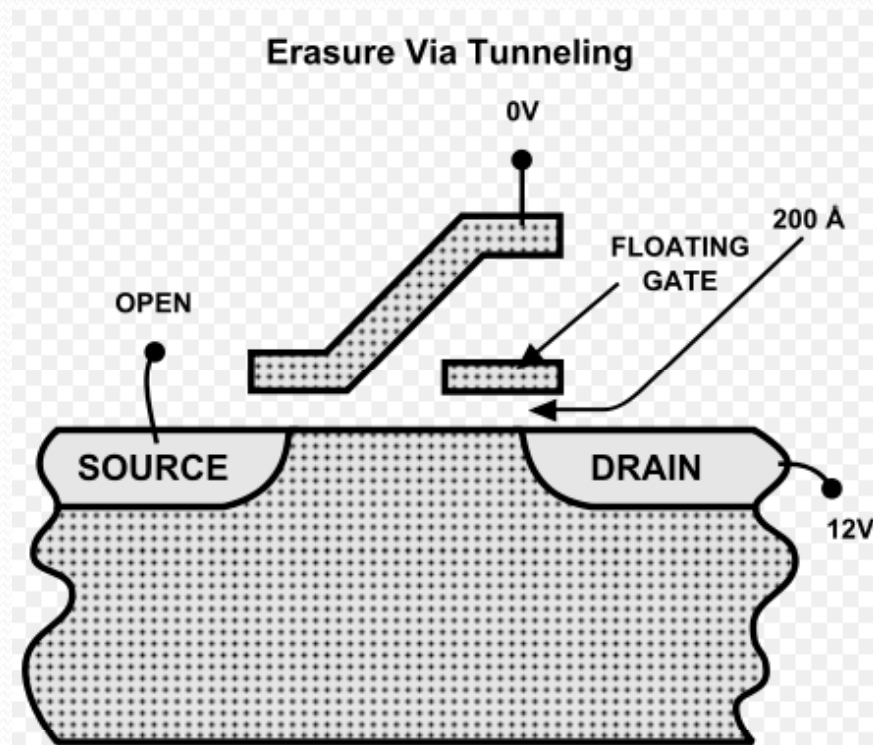
(c) PGM: HCI
ERS: FN tunneling (drain)

(d) PGM: HCI
ERS: FN tunneling

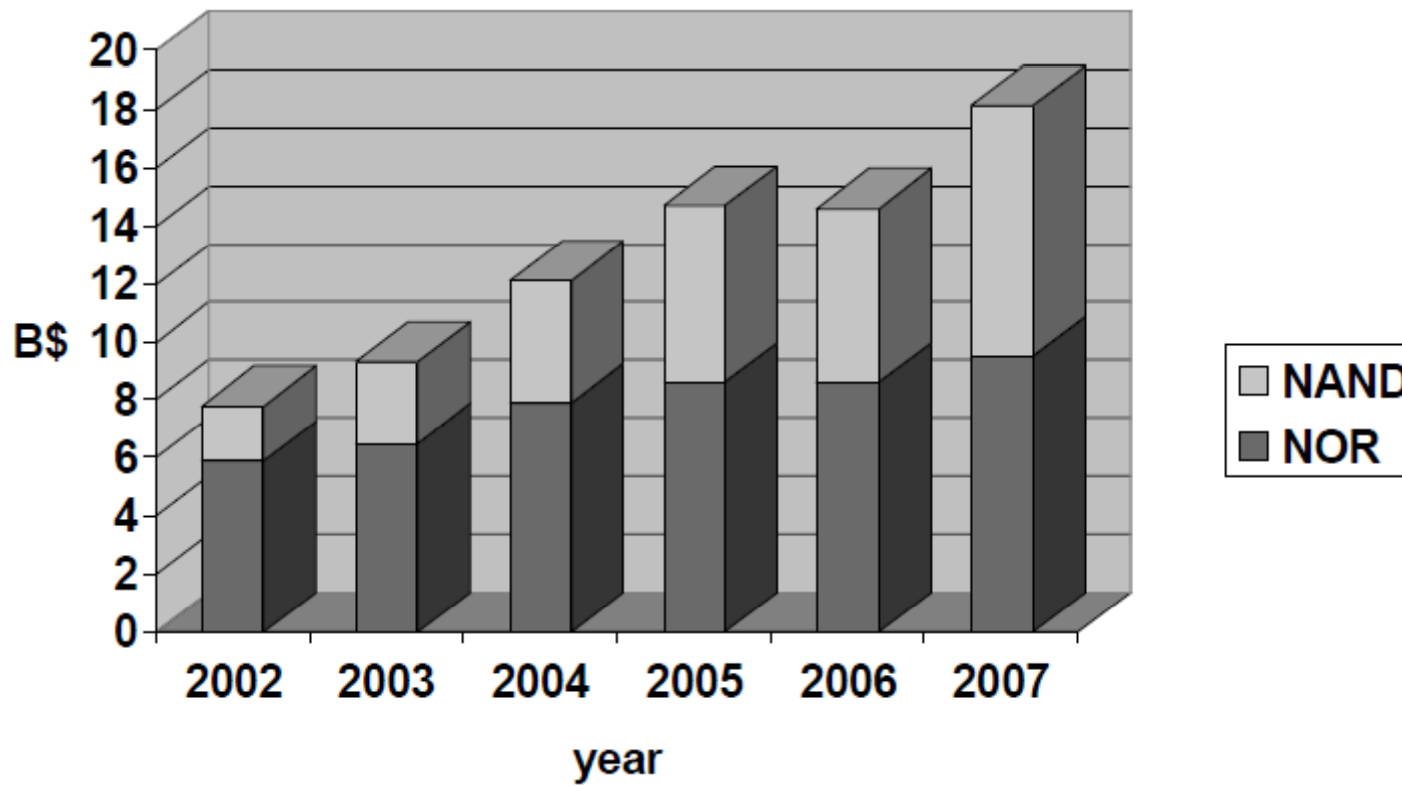
Programming Via Hot Electron Injection




- Default state is “1”, as current will flow through the channel under application of an appropriate voltage to the control gate.
- Set to “0” by applying $> 5V$ (typically) to the CG. Channel is now turned on, so electrons can flow from the source to the drain. Current is then sufficiently large to cause some high energy (hot) electrons to jump through the insulating layer onto the FG.



- To erase a NOR flash cell (resetting it to the "1" state), a large voltage of *the opposite polarity* is applied between the CG and source, pulling the electrons off the FG through quantum tunneling

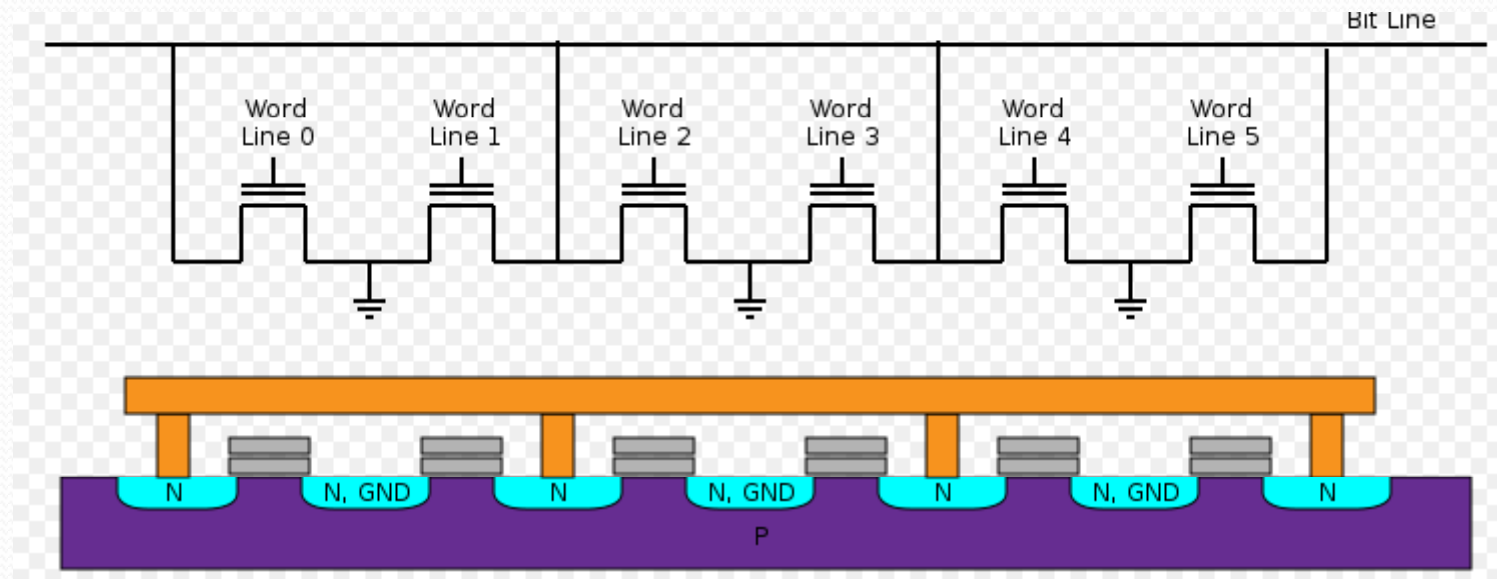


Flash memory market sharing by technology (source: Web Feet Inc.)

- 
- NAND has significantly higher storage capacity than NOR
 - NAND flash has found a market in devices to which large files are frequently uploaded and replaced: MP3 players, digital cameras and USB drives
 - NOR flash is faster, but it's also more expensive. NOR is most often used in mobile phones

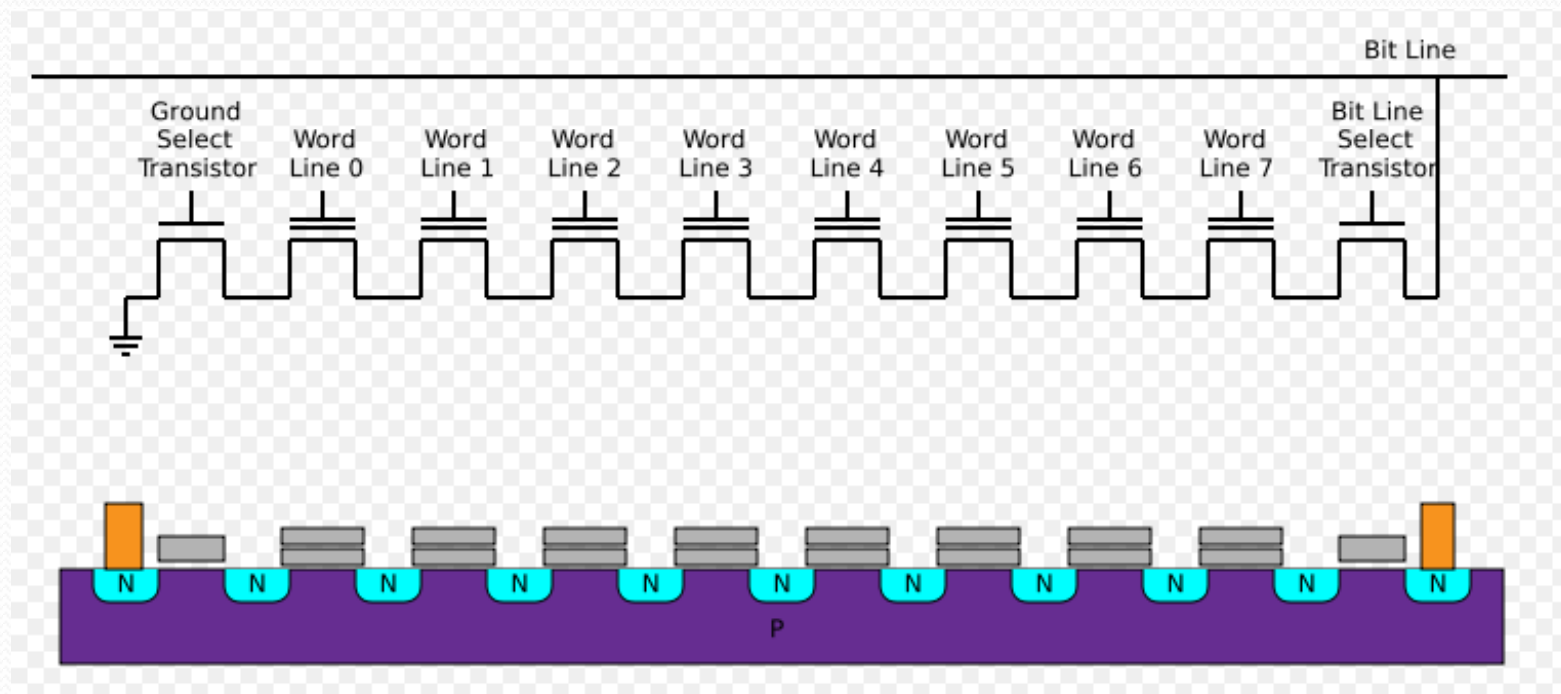
Flash memories

- specific type of EEPROM that is erased and programmed in large blocks



NOR-gate flash (Intel, 1988)

Nor-gate Flash: each cell has one end connected directly to ground, and the other end connected directly to a bit line.



NAND-gate flash (Toshiba 1989)

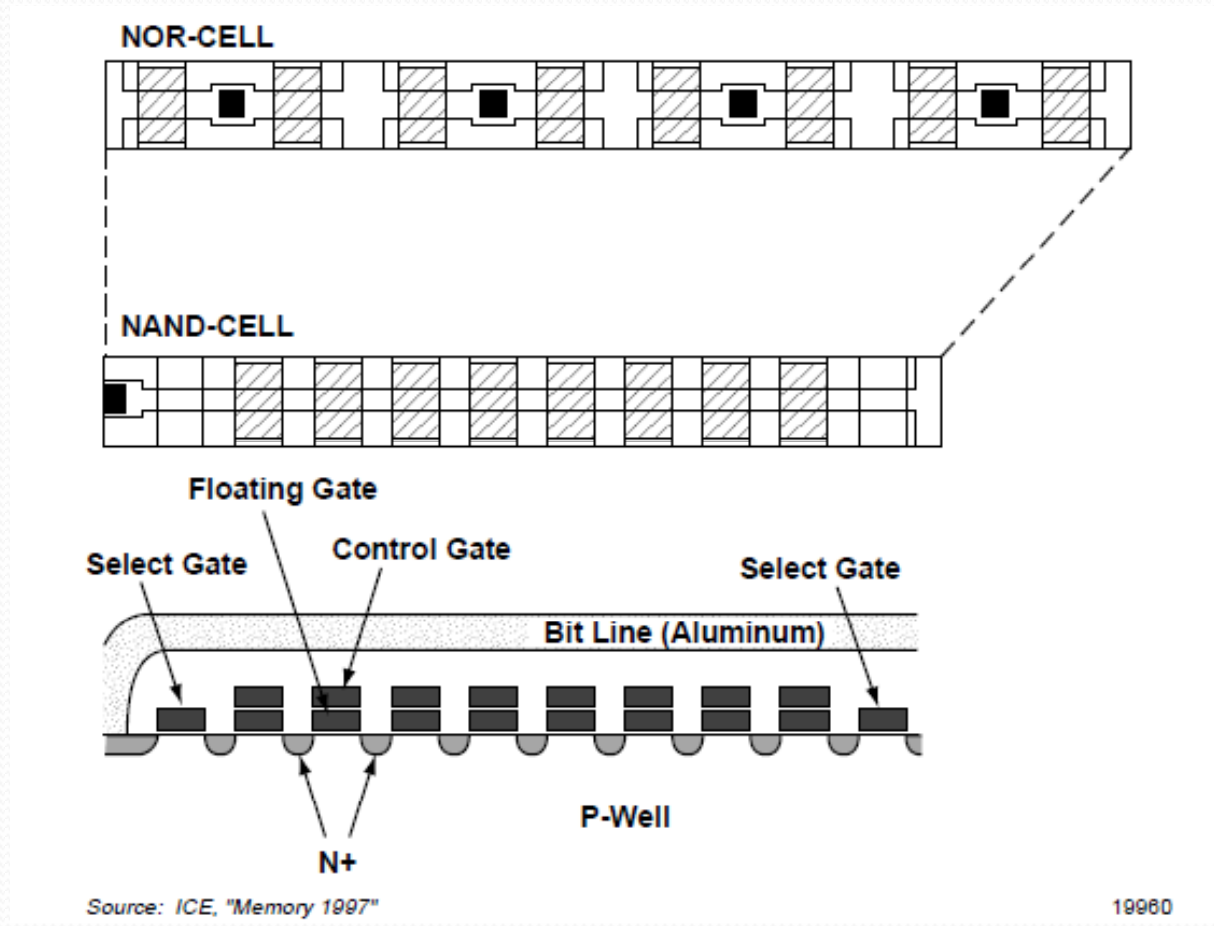
- several transistors (8, 16, 32, ...) are connected in series




When a cell is read, its gate is set to 0V, while the other gates of the stack are biased with a high voltage (typically 4–5 V, say)

➔ the other gates work as pass-transistors, regardless of their threshold voltage.

An erased NAND Flash cell has a negative threshold voltage; on the contrary, a programmed cell has a positive threshold voltage but, in any case, less than 4V.

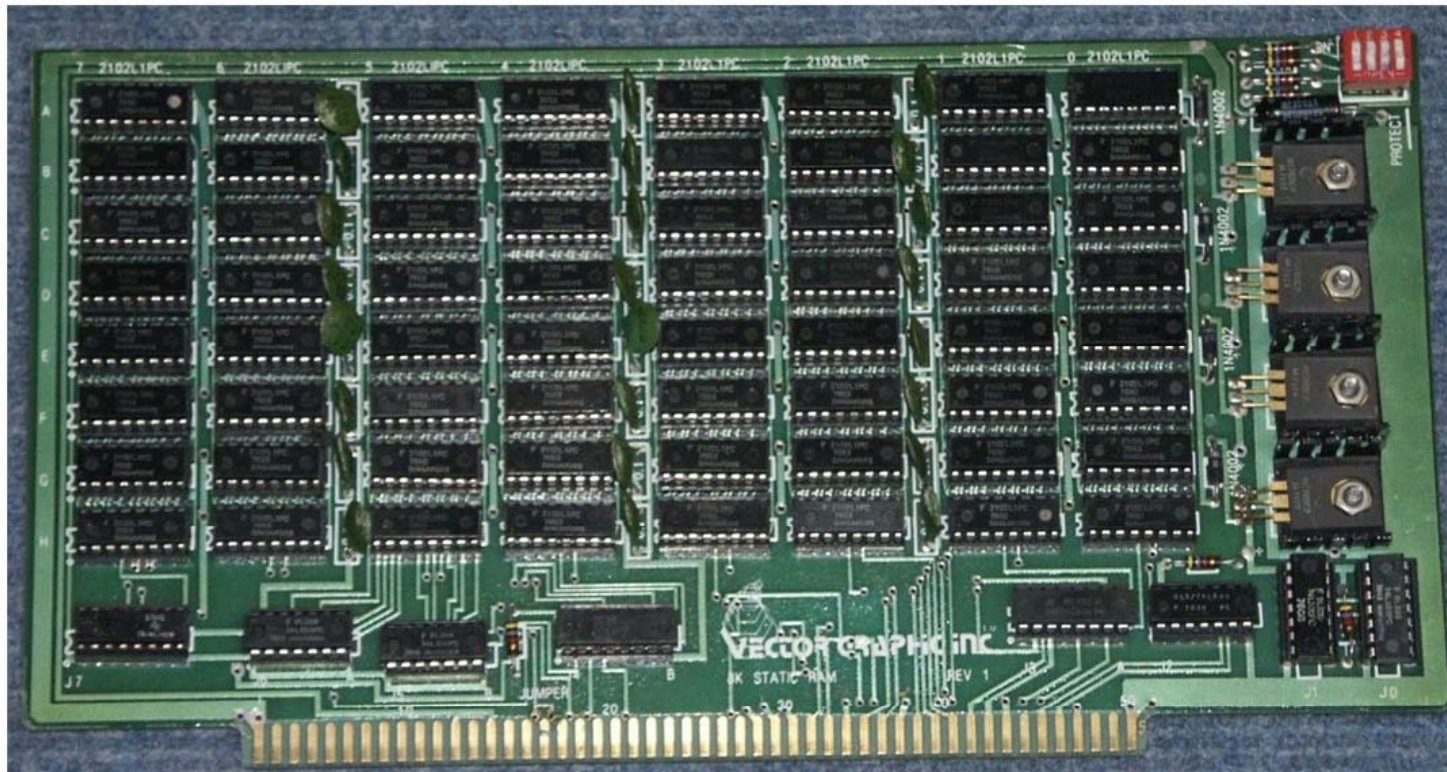


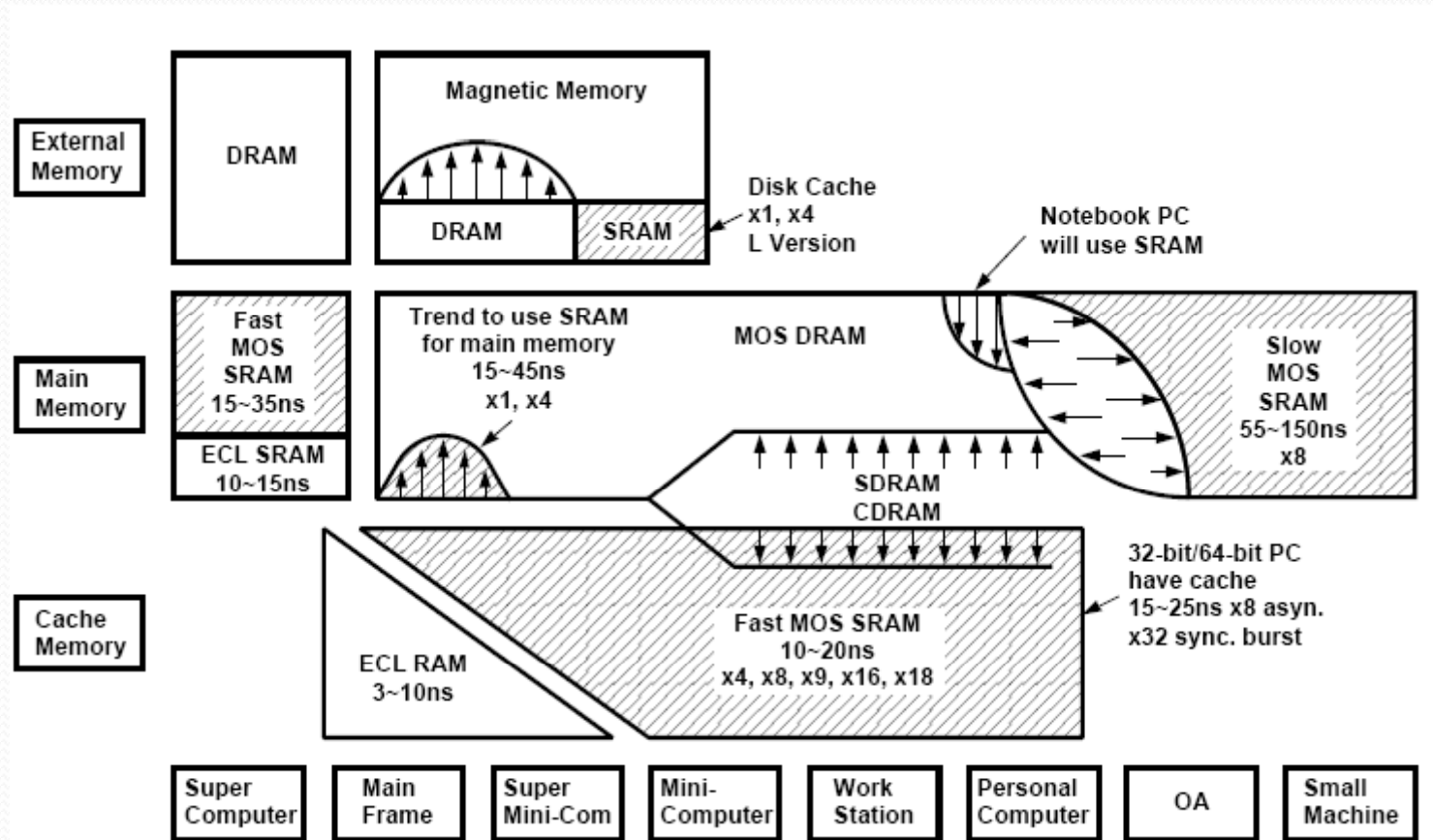


NOR Flash: is a random access device, appropriate for code storage application. It is designed for use in linear program storage for applications such as boot loaders and BIOS, with its key benefit being the ability to satisfy requirements that need to read code wherein each word of data is needed to carry out instructions.

NAND Flash: a sequential access device appropriate for mass storage applications. It is optimized for file structures where each word does not need to be read, but instead provides that sectors of data can be moved to and from media supporting a hard drive like repository structure for data storage to support file systems and allocation tables (FAT)

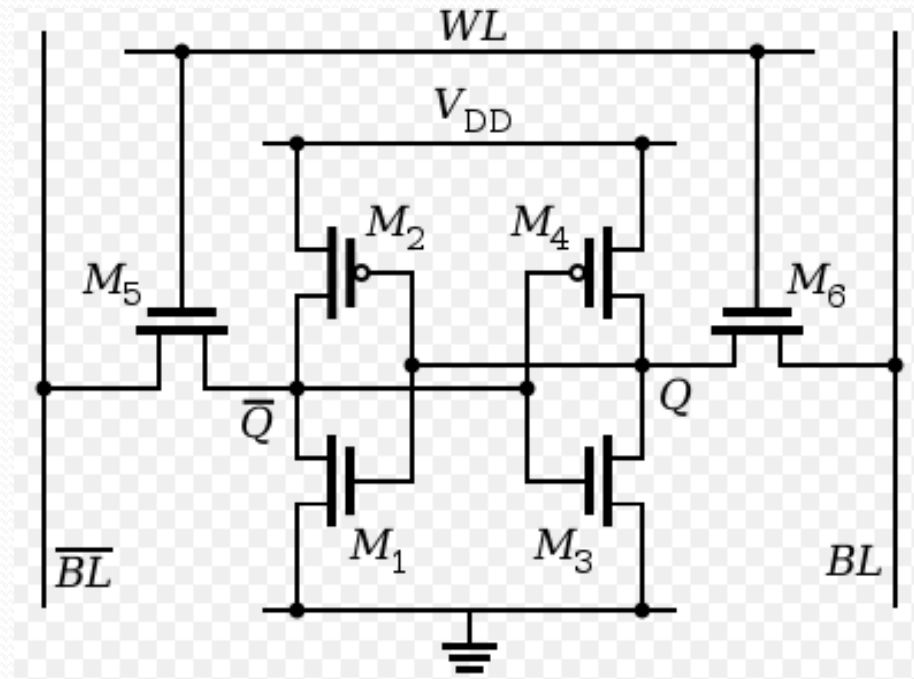
Static RAM (SRAM)





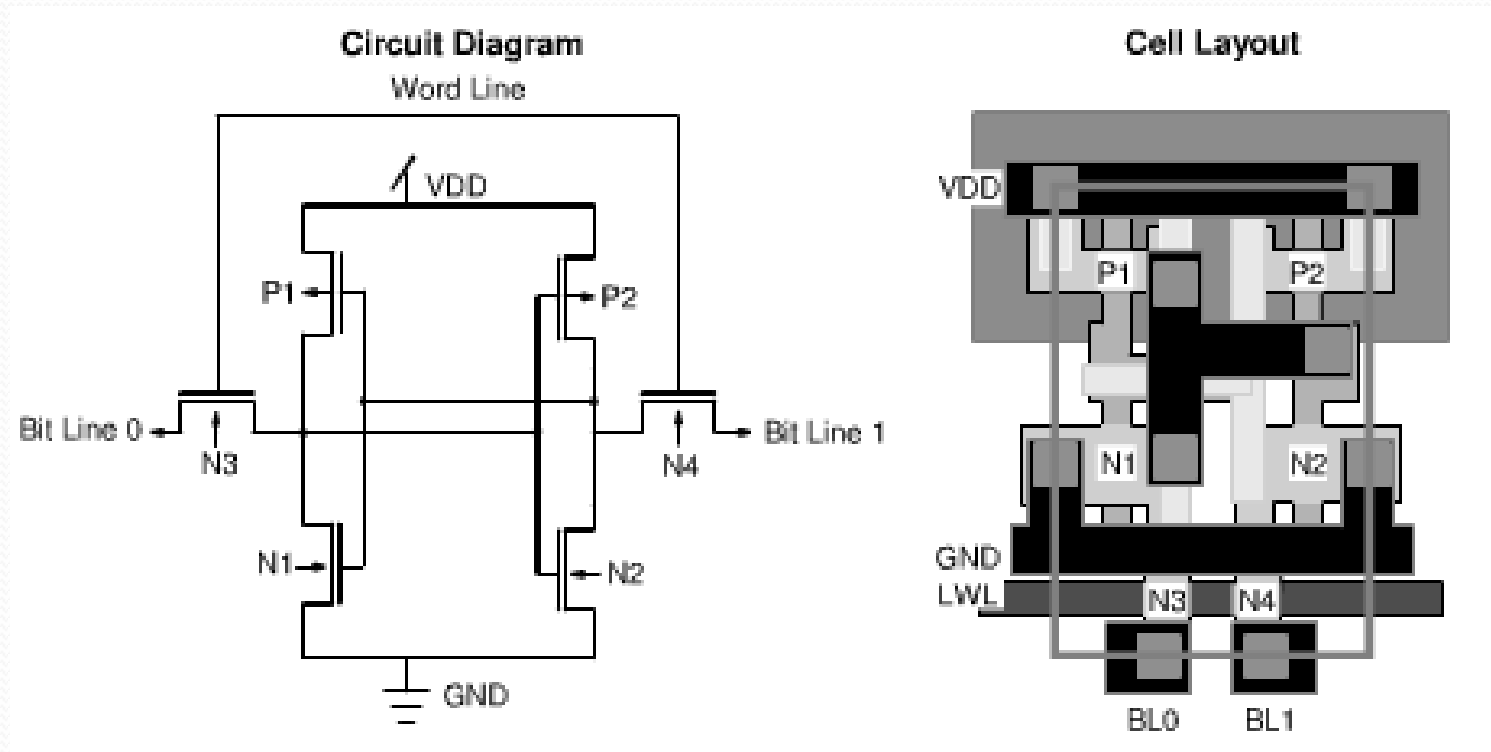
Source: Mitsubishi

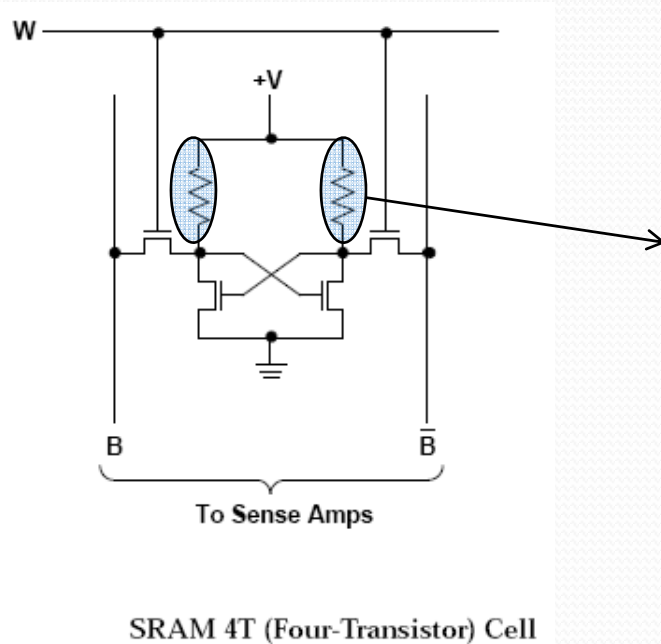
Application of DRAM/SRAM



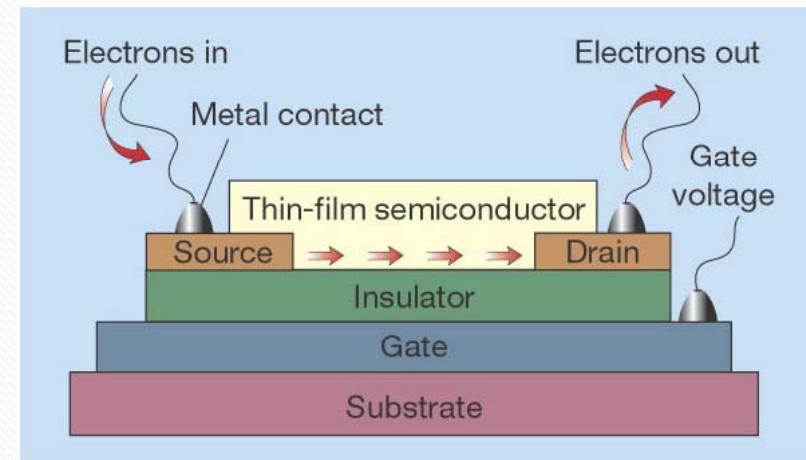
6-transistor SRAM cell

- high speed memory
- data stored as long as power is supplied to the circuit.
- low density \Rightarrow high cost!!

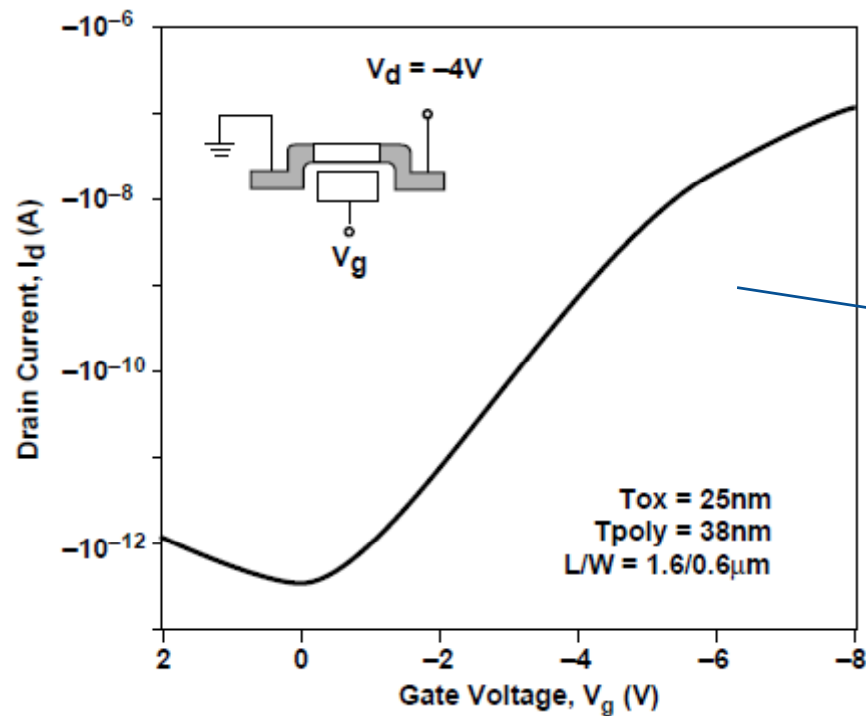




TFT: Thin-Film Transistor



- Load is formed by using polysilicon as a PMOS device. This PMOS transistor is called a Thin Film Transistor (TFT), and it is formed by depositing several layers of polysilicon above the silicon surface.



Source: Hitachi

19953

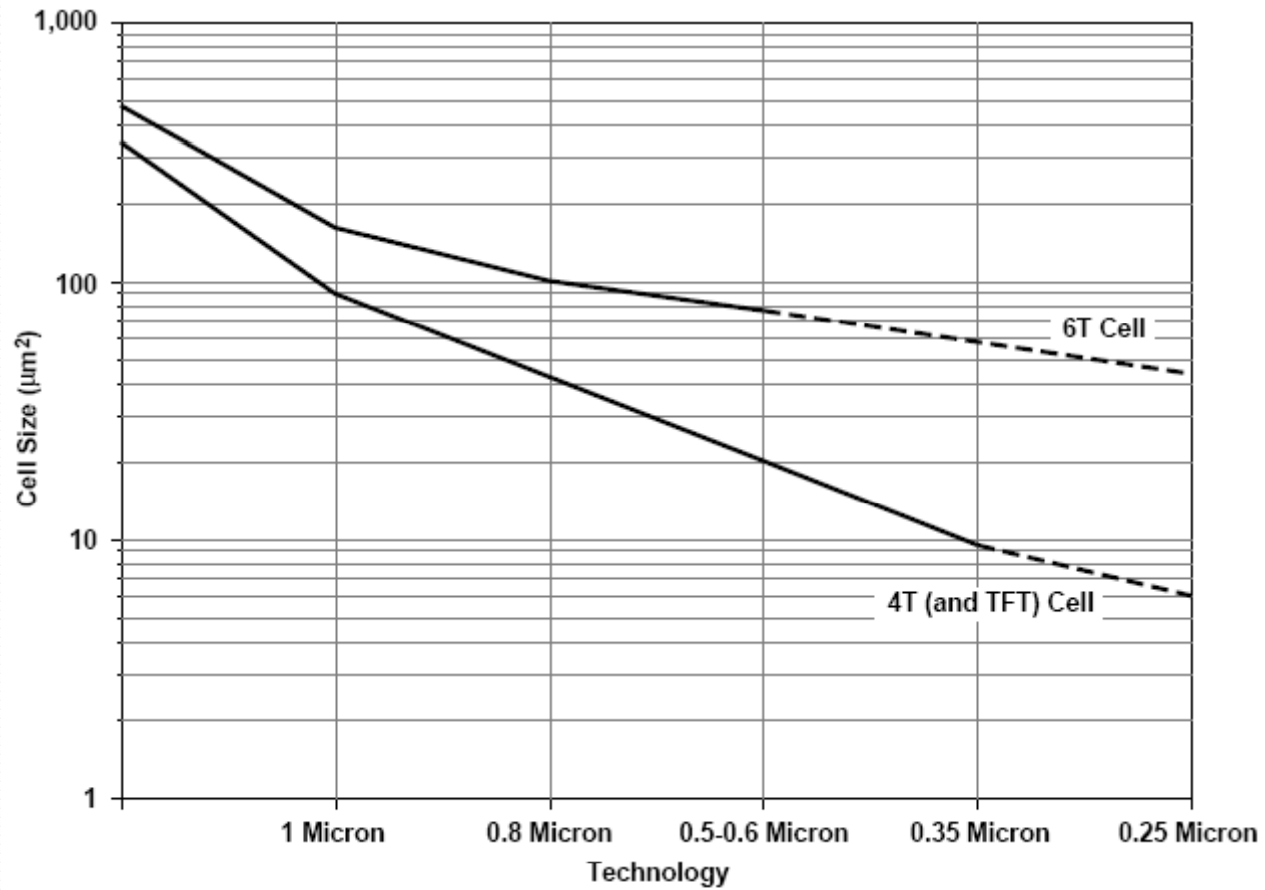
effective resistance ranges from
 $\sim 11 \times 10^{13} \Omega$ to $5 \times 10^9 \Omega$

4T cells have several limitations:

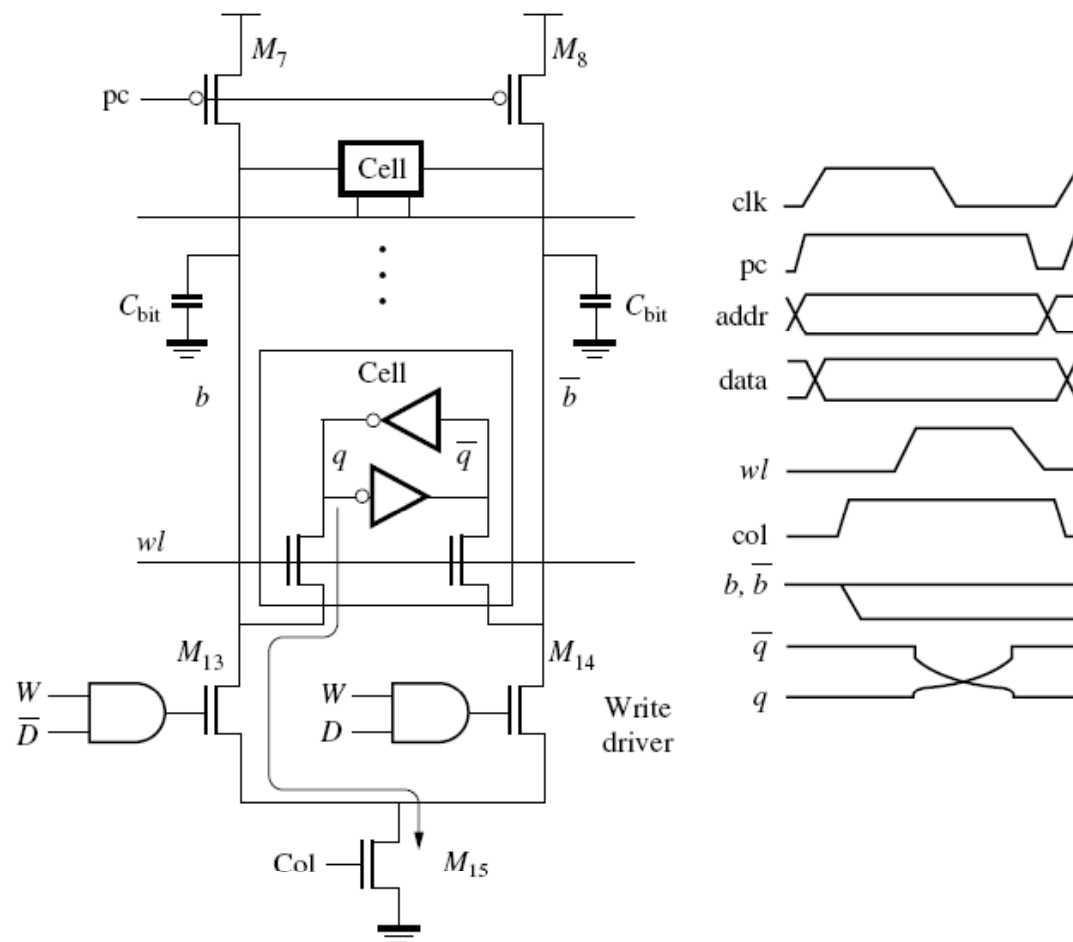
- Each cell has a current flowing in one resistor, i.e., high standby current.
- The cell is sensitive to noise and soft error rate because the resistors are so large.
- The cell is not the fastest cell design available.

	TOSHIBA TC554161FTL-70L 4Mb (x16) 9509	SAMSUNG KM732V588 1Mb (x32) 1995	GALVANTECH GVT7132C32Q7 1Mb (x32) 9524	HITACHI HM67W1664JP-12 1Mb (x16) 9539	NEC D461018LG5-A12 1Mb (x18) 9436	MOTOROLA MCM67C618FN7 1Mb (x18) 9443
Technology	CMOS	BiCMOS	CMOS	BiCMOS	BiCMOS	BiCMOS
Die Size	7.7 x 18.7mm (144mm ²)	5 x 6.6mm (33mm ²)	4.5 x 6.8mm (31mm ²)	6.4 x 10.1mm (64mm ²)	5.7 x 11.7mm (67mm ²)	9.2 x 11.8mm (108mm ²)
Min Gate - (N)	0.65μm	0.5μm	0.4μm	0.45μm	0.6μm	0.6μm
Cell Pitch	3.7 x 6μm	3.0 x 4.75μm	3.4 x 4.9μm	3.3 x 5.7μm	3.5 x 5.5μm	4.9 x 8.2μm
Cell Area	22μm ²	14.25μm ²	16.5μm ²	19μm ²	19μm ²	40μm ²
Cell Type	4T	4T	4T	4T	4T	4T
V _{CC}	5V	3.3V	3.3V	3.3V	3.3V	5V
Access Time	70ns	—	7ns	12ns	12ns	7ns

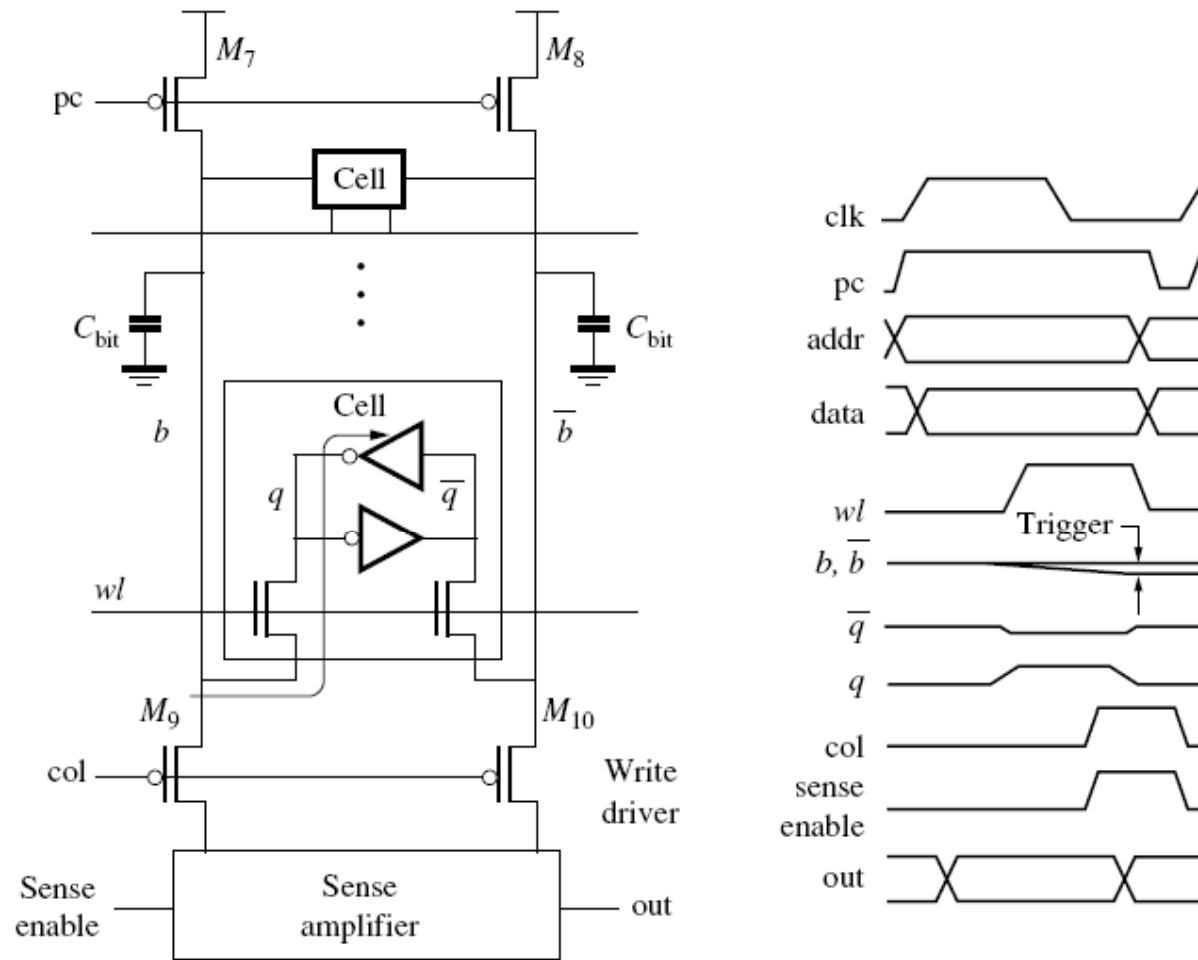
Physical Geometries of SRAMs



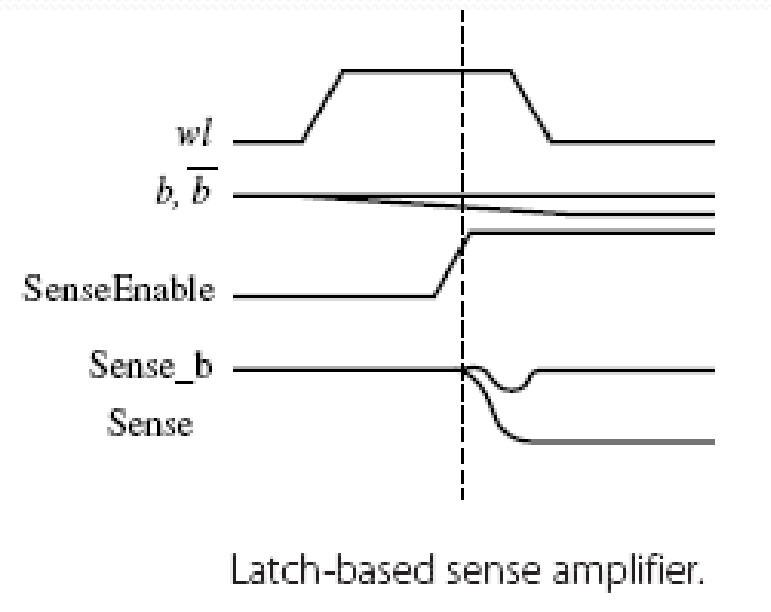
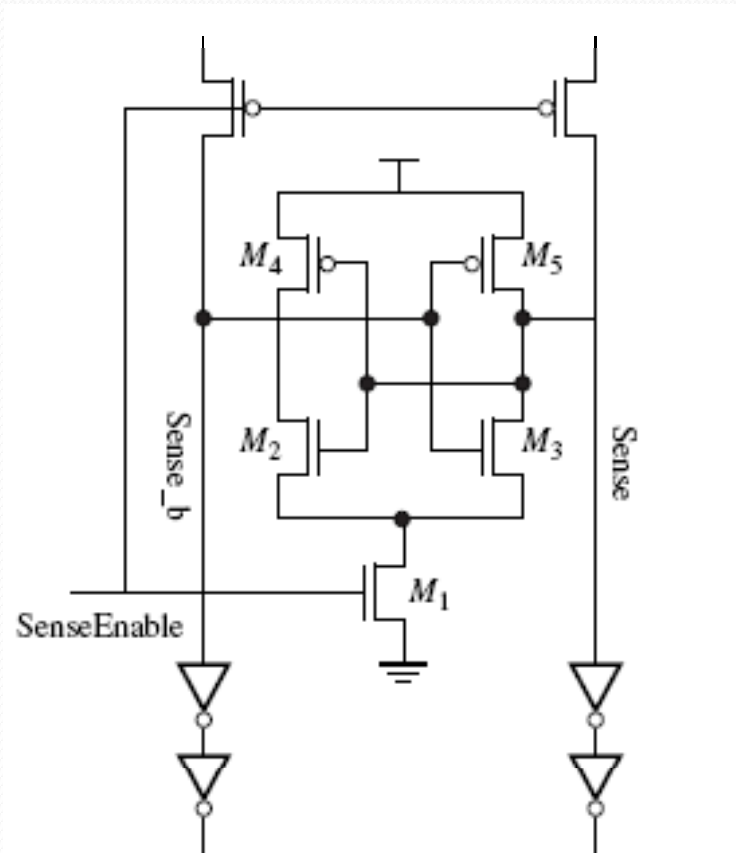
Trend of SRAM Cell Sizes



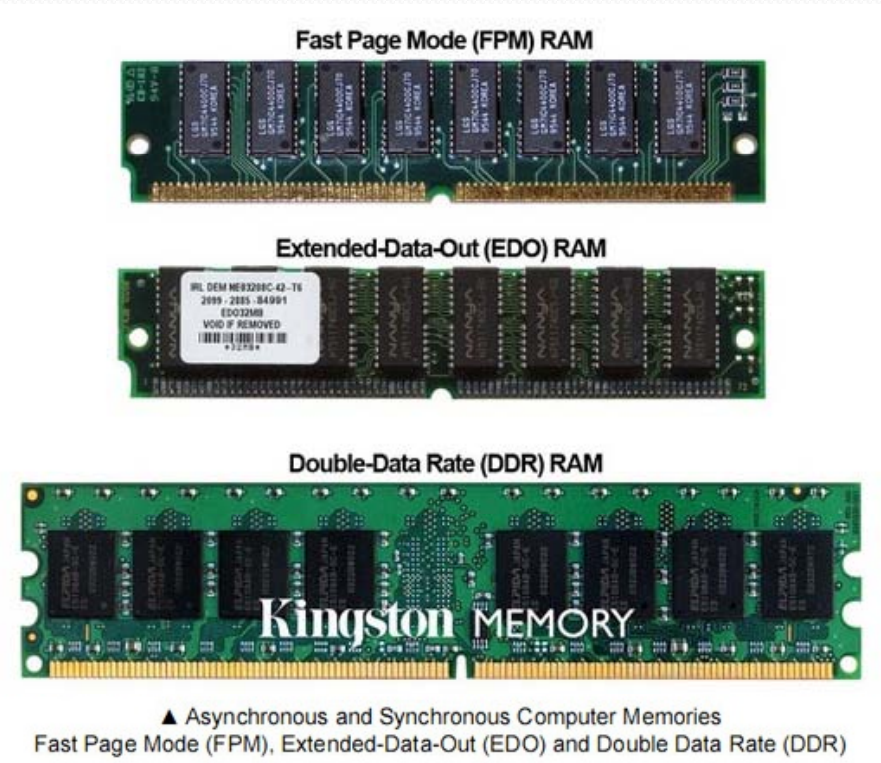
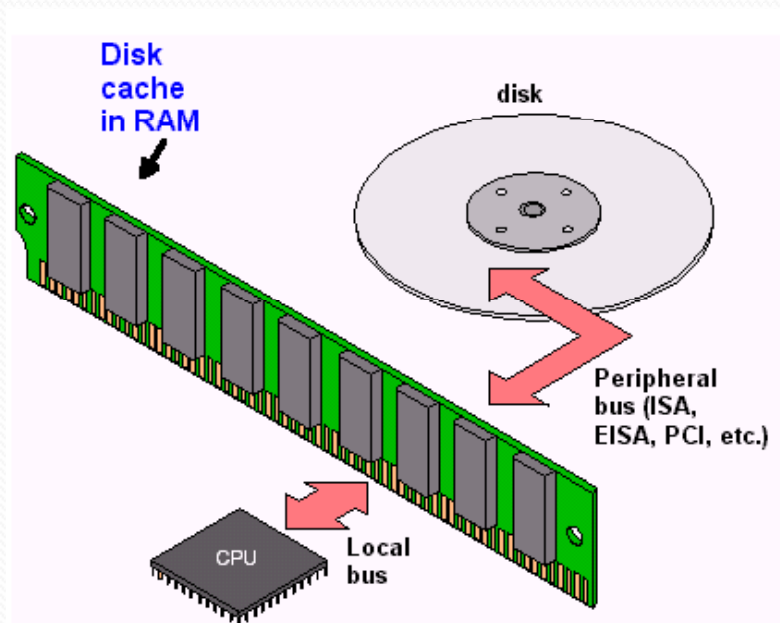
Write driver circuit.



Basic read circuitry.

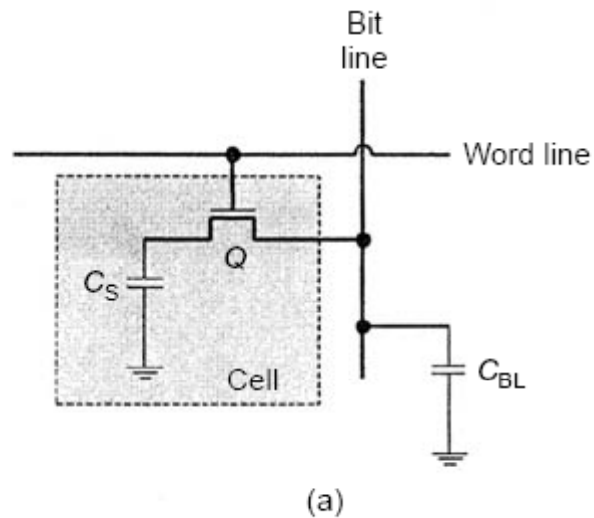


Dynamic RAM (DRAM)

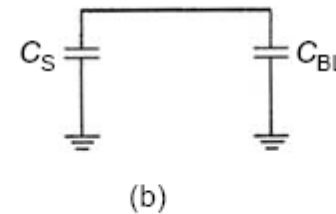


DYNAMIC RAM (DRAM) MEMORY TECHNOLOGIES

Type	Year of Intro.	Maximum Clock Rate	Bus Width	Peak Bandwidth	Volts
FPM	1990	25MHz	64 bits	200 MBps	5v
EDO	1994	40MHz	64 bits	320 MBps	5v
SDRAM	1996	133MHz	64 bits	1.1 GBps	3.3v
RDRAM	1998	400MHz (x2)	16 bits	800 MBps	2.5v
DDR SDRAM	2000	266MHz (x2)	64 bits	4.2 GBps	2.5v
DDR2 SDRAM	2003	533MHz (x2)	64 bits	8.5 GBps	1.8v
DDR3 SDRAM	2007	800MHz (x2)	64 bits	12.8 GBps	1.5v

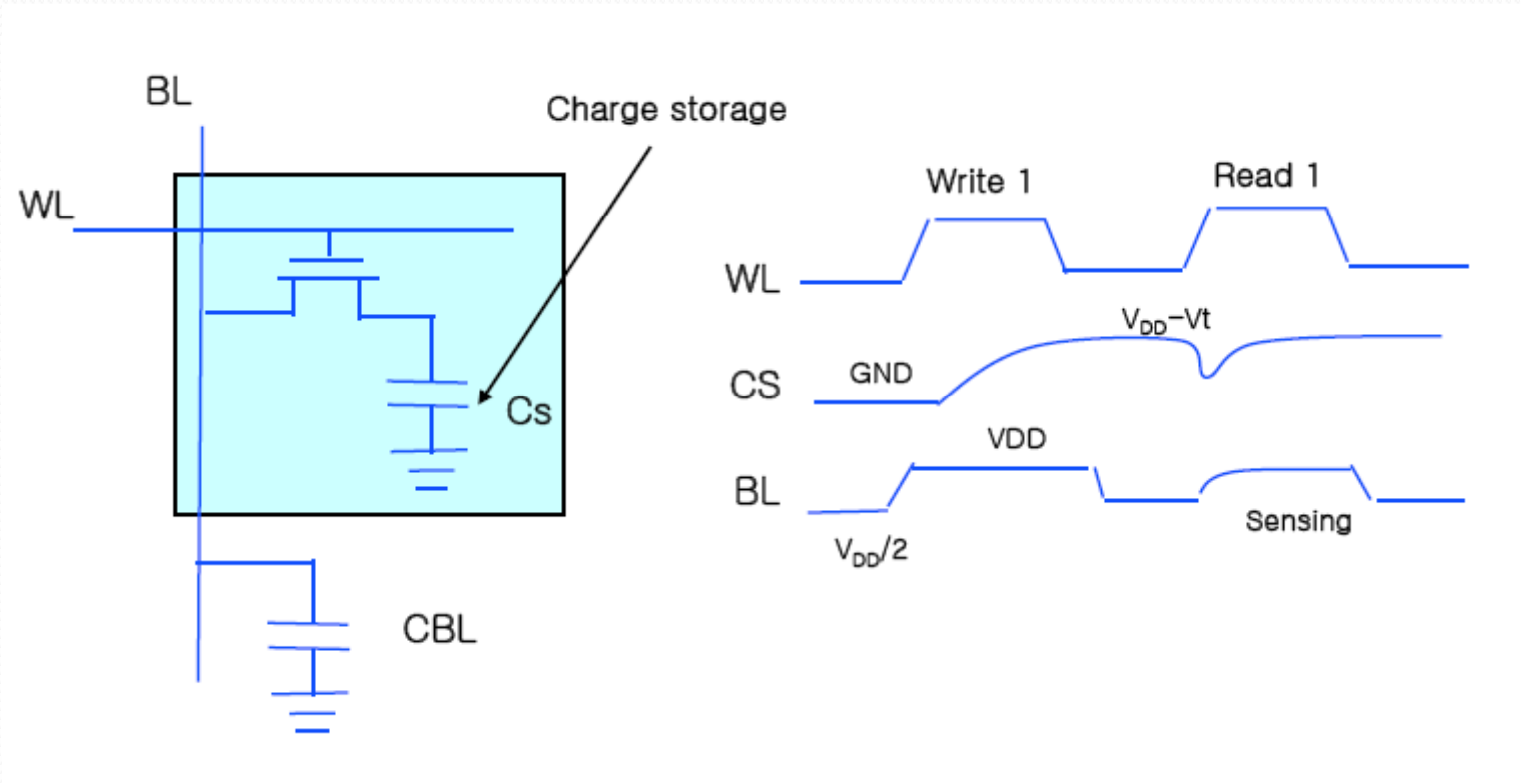


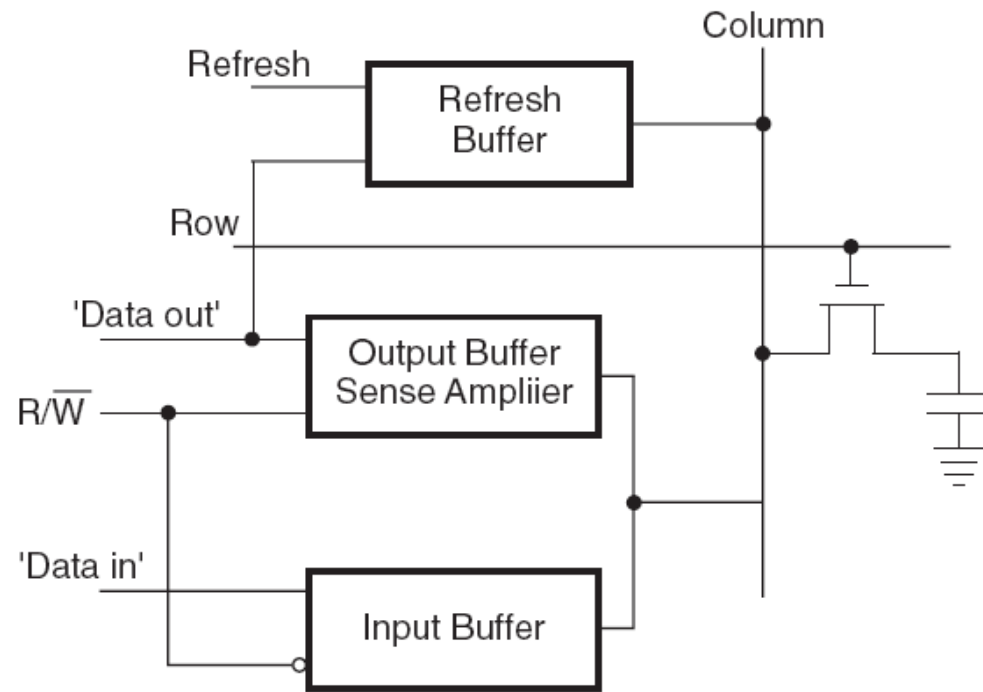
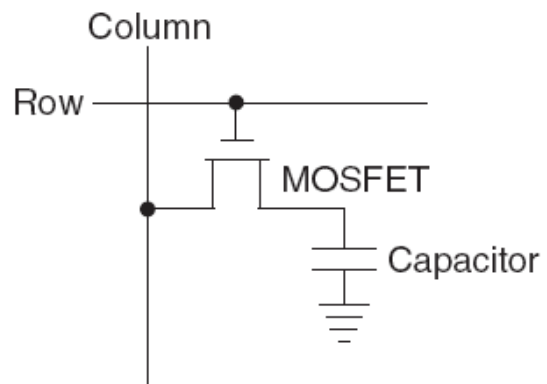
CS: typically, 30-50fF
 "1": $V_{CS} = V_{DD} - V_{TH}$
 "0": $V_{CS} = 0V$



(a) The one-transistor DRAM cell; and (b) during the READ operation, the voltage of the selected word-line is high, thus connecting the storage capacitor C_S to the bit-line capacitance C_{BL} .

READ: the voltage of the selected word-line is high;
readout process is destructive !!!
write process should occur after reading.





Basic memory cell of a DRAM.

- bit-line voltage variation

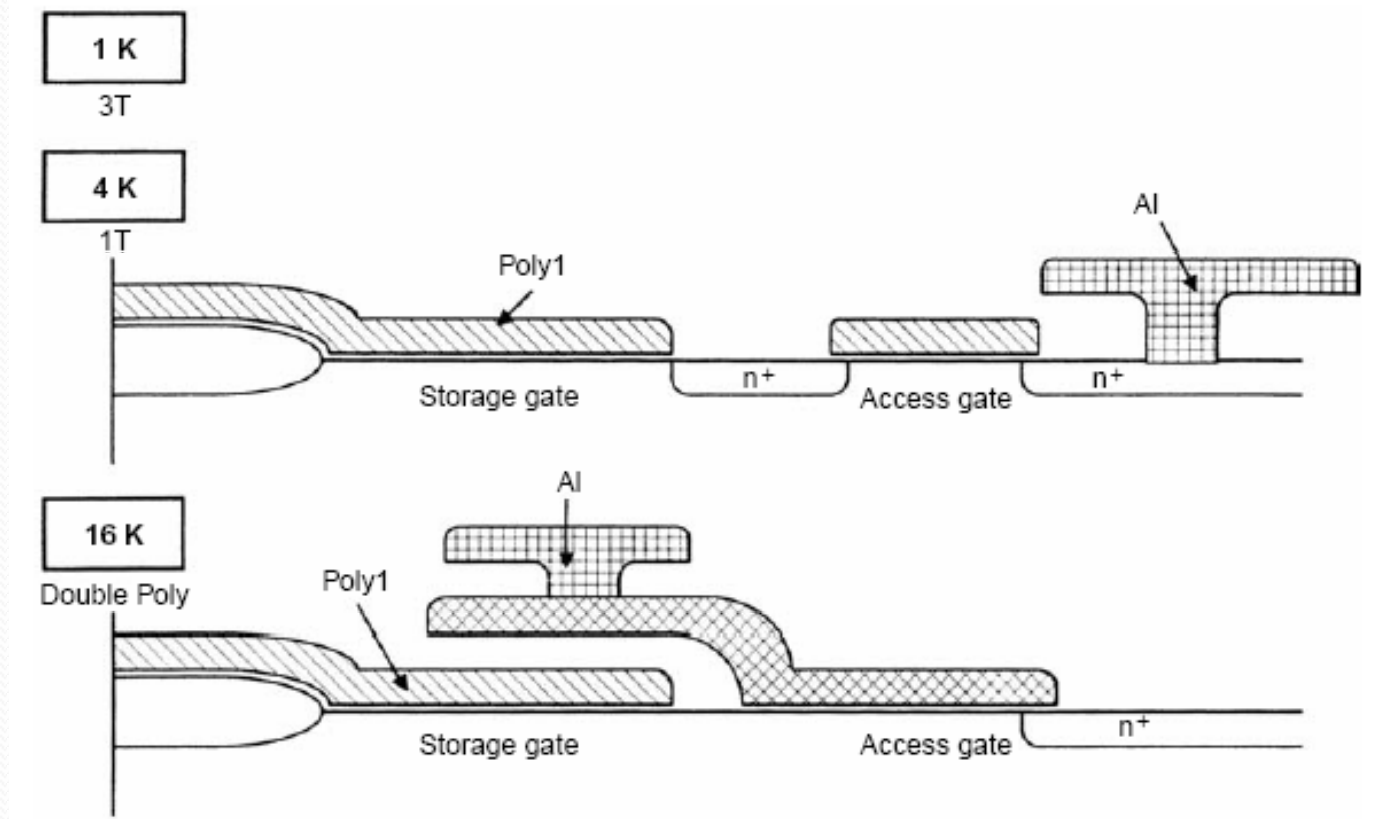
$$V_s = \Delta V_{\text{BL}} = \frac{C_s}{C_{\text{BL}} + C_s} \left(V_{\text{cs}} - \frac{V_{\text{DD}}}{2} \right)$$

$$R = C_{\text{BL}}/C_s$$

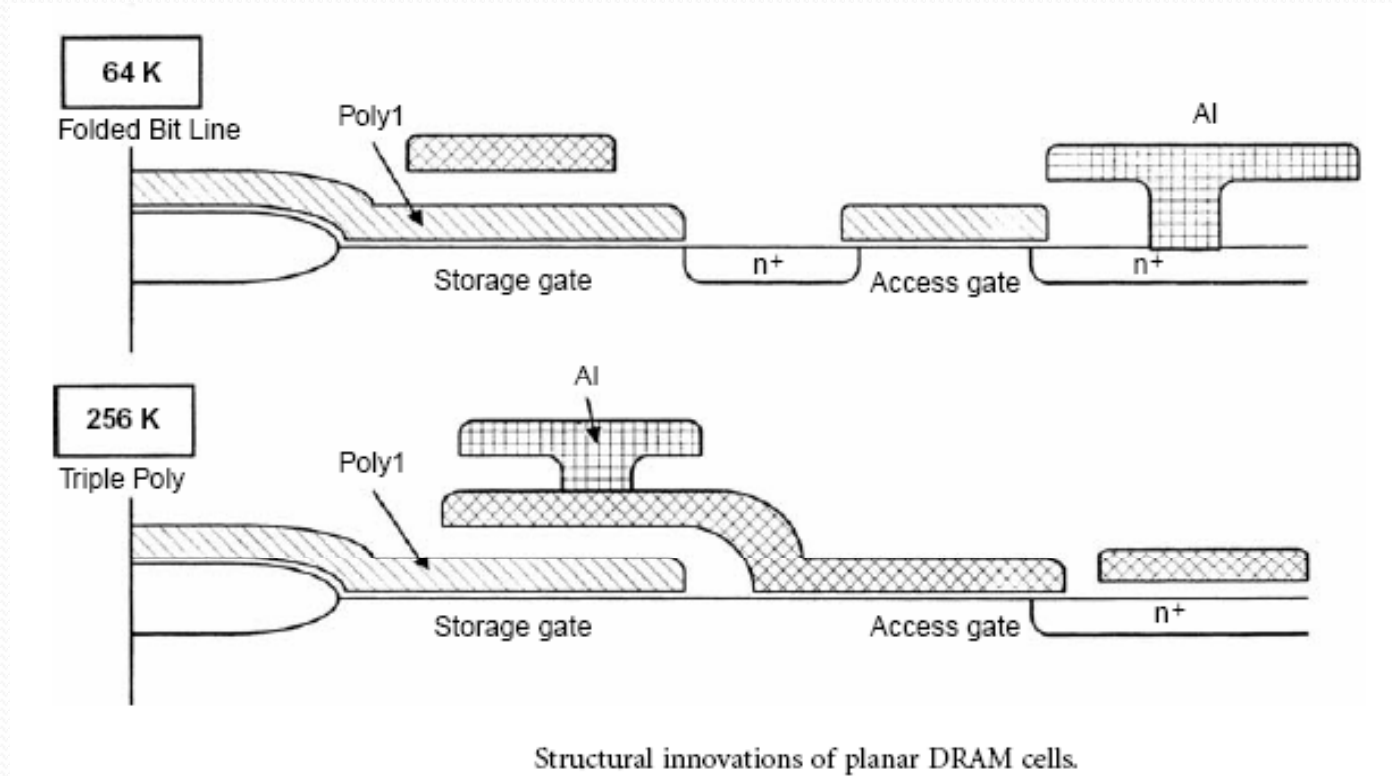
$$\Delta V(1) = \frac{1}{1+R} \left(\frac{V_{\text{DD}}}{2} - V_t \right)$$

$$\Delta V(0) = \frac{1}{1+R} \left(\frac{V_{\text{DD}}}{2} \right)$$

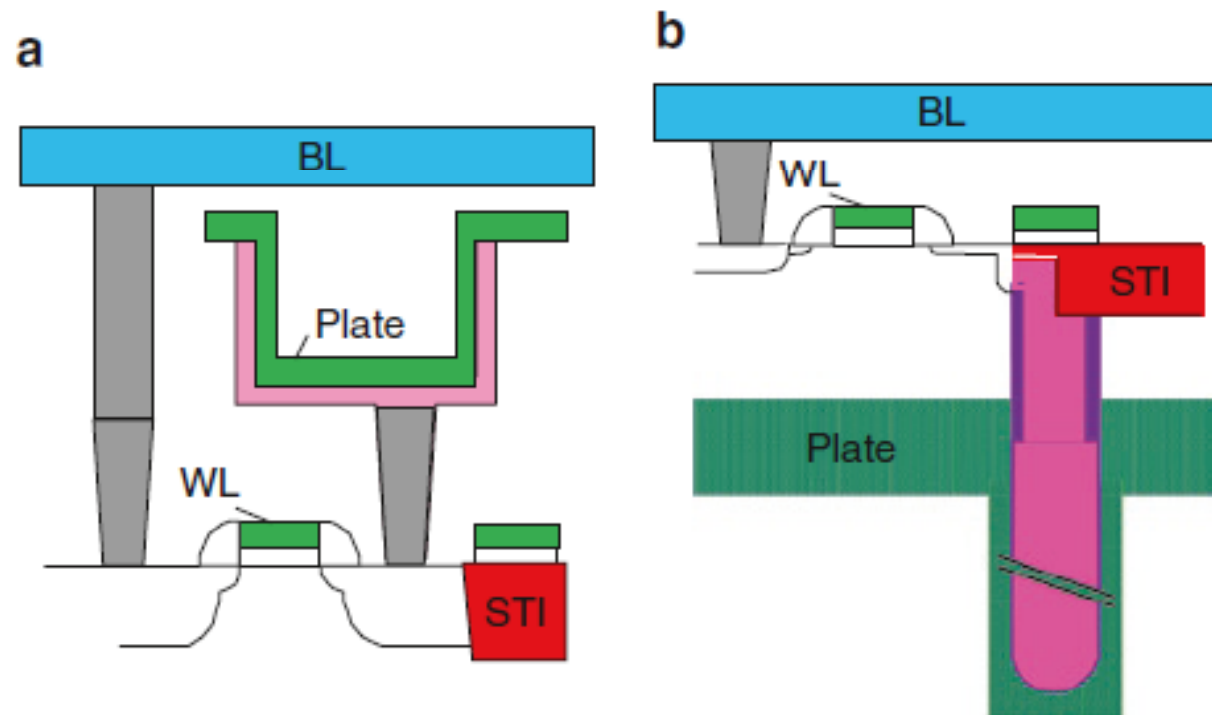
Since ratio $R = C_{\text{BL}}/C_s$ is large, these readout bit-line sense signals $\Delta V(1)$ and $\Delta V(0)$ are very small. Typical values for the sense signal are about 100 mv.



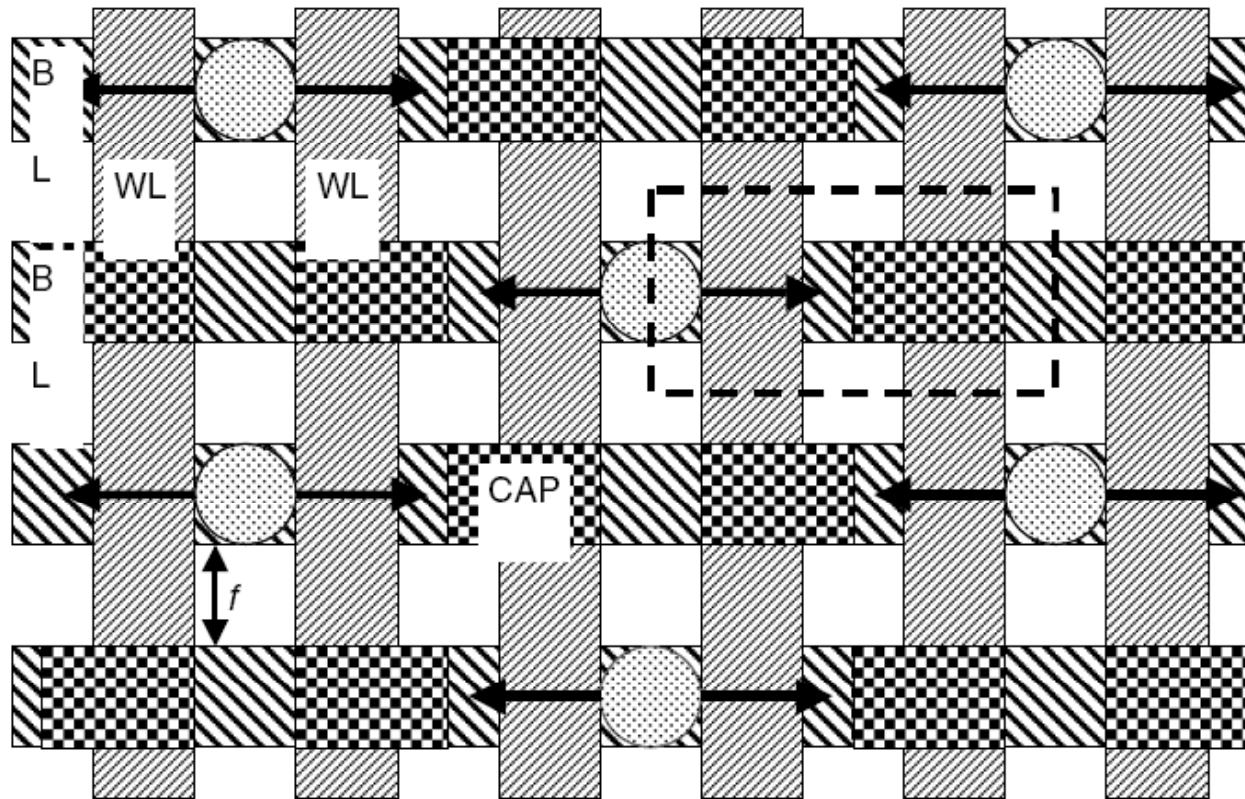
Structural innovations of planar DRAM cells.



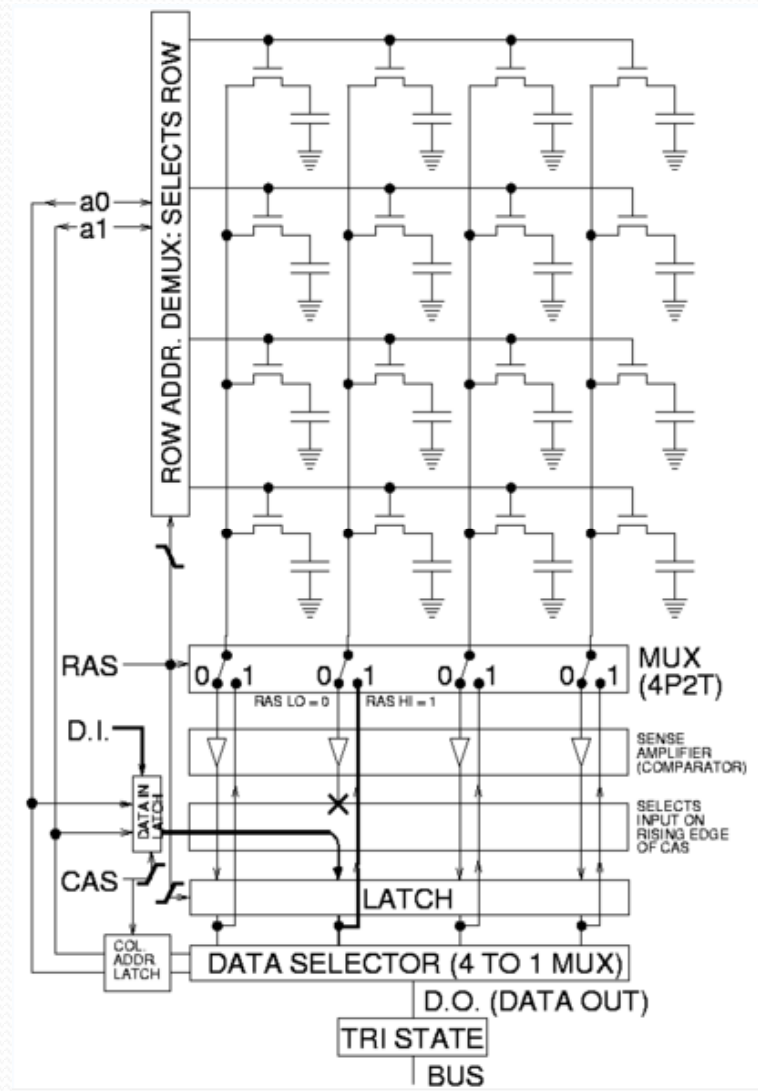
1-Mb DRAM.



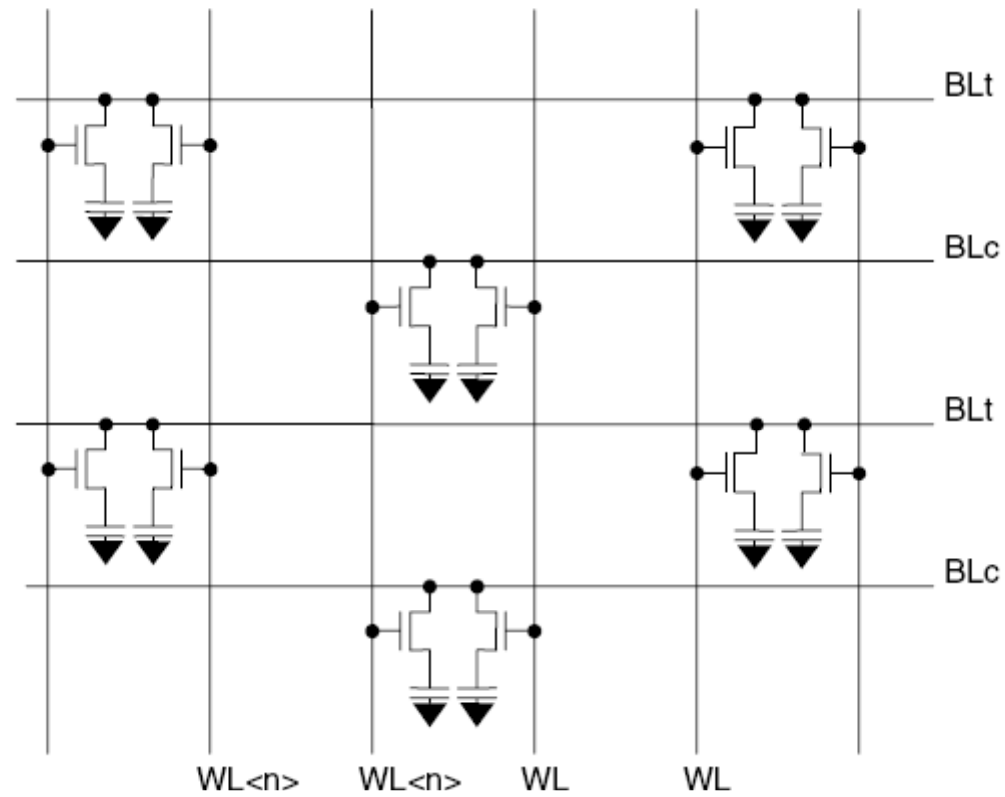
Cells: (a) stacked capacitor and (b) trench capacitor



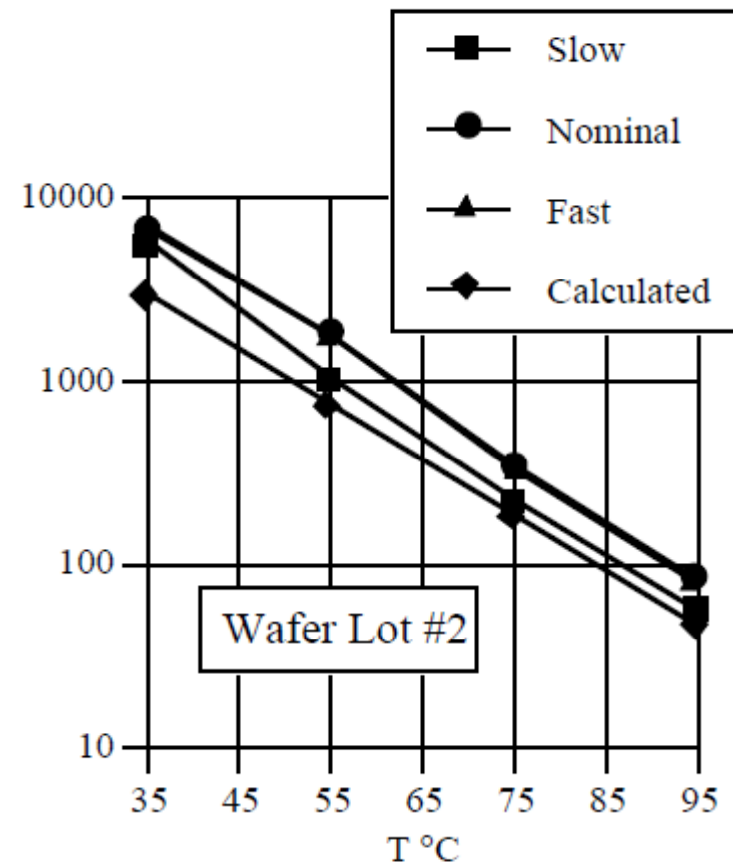
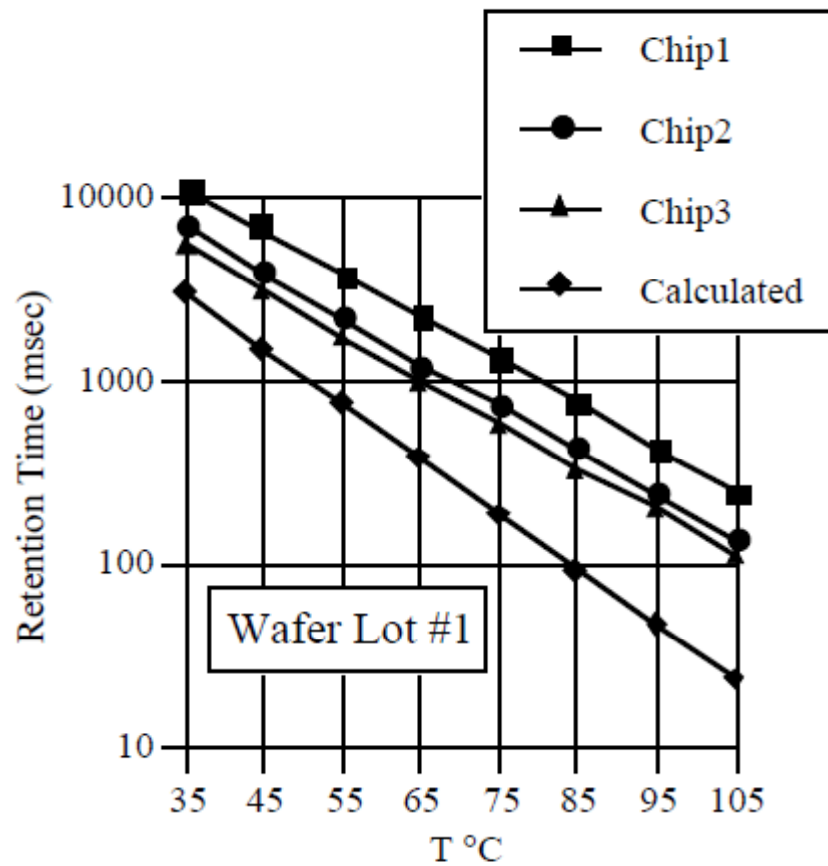
Idealized layout of an $8f^2$ DRAM array. The memory cell outline is given by dashed lines. Arrows indicate the direction of current flow through the transfer device. Circles indicate contact from bit line to diffusion of the transfer device. Word lines (WLs) run vertically; bit lines (BLs) are oriented perpendicularly. A storage capacitor (CAP; checkerboard pattern) may occupy more than $1f^2$ area.



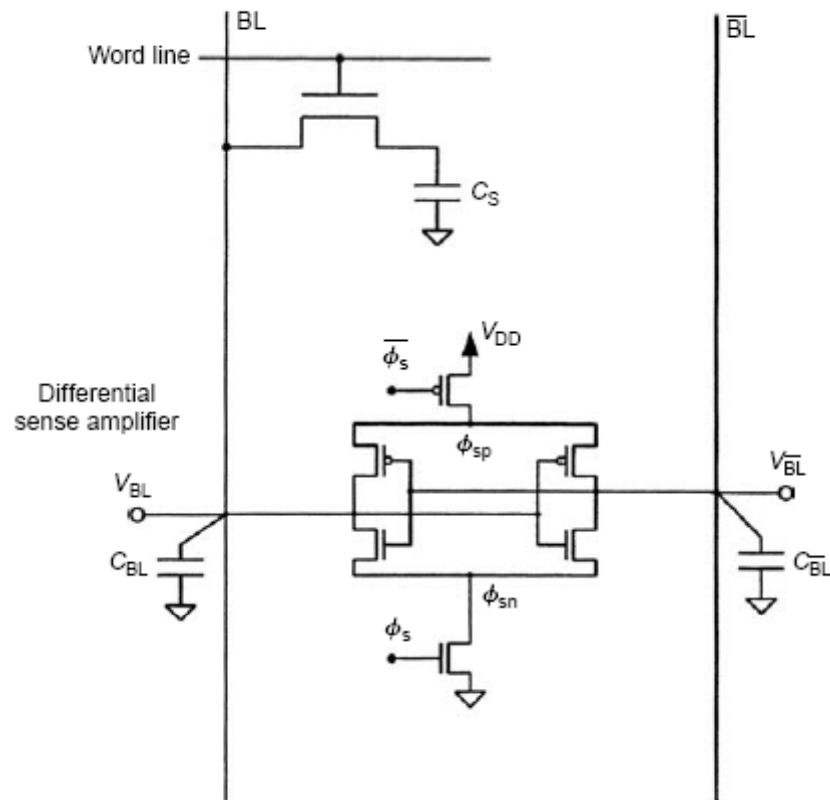
DRAM



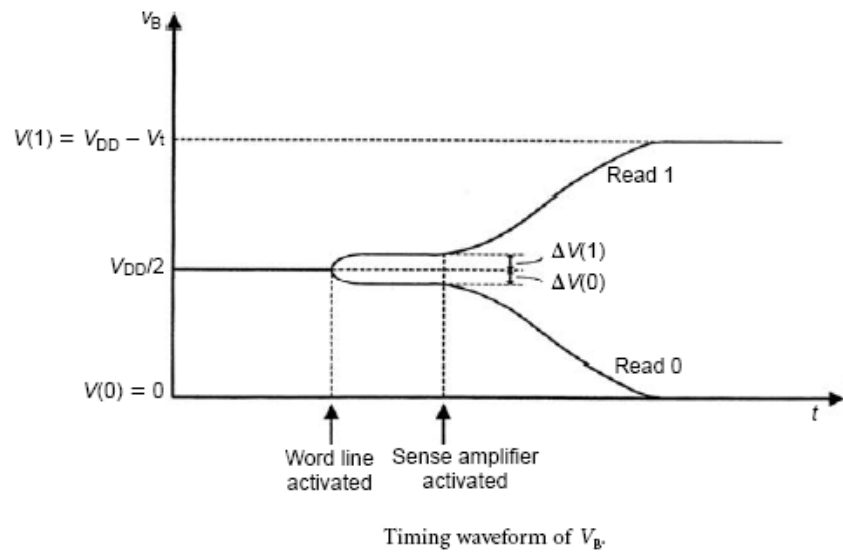
Schematic representation of a DRAM array. Word lines (WLs) are running vertically, bit lines (BLs) horizontally. Logically, bit lines are organized into true/complement pairs.



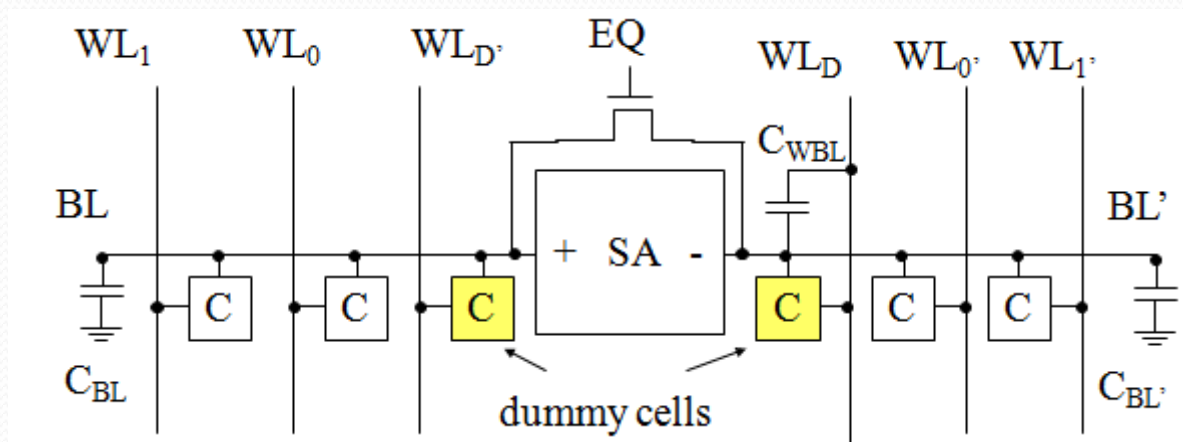
Experimental results for DRAM data retention time.



A differential sense amplifier connected to the bit-line.

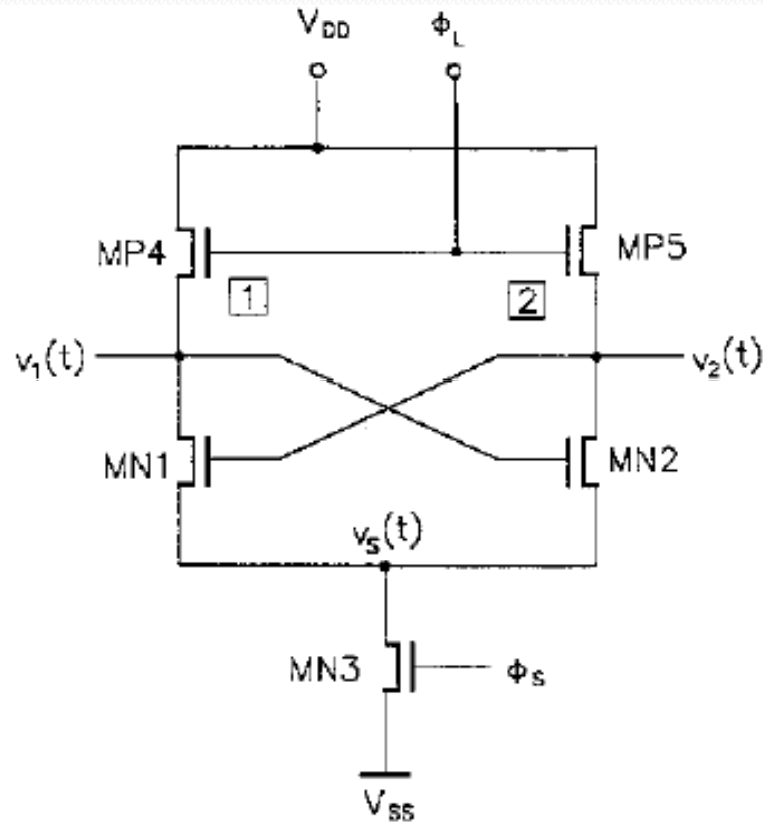


Sense amplifier with dummy-cell sensing structure



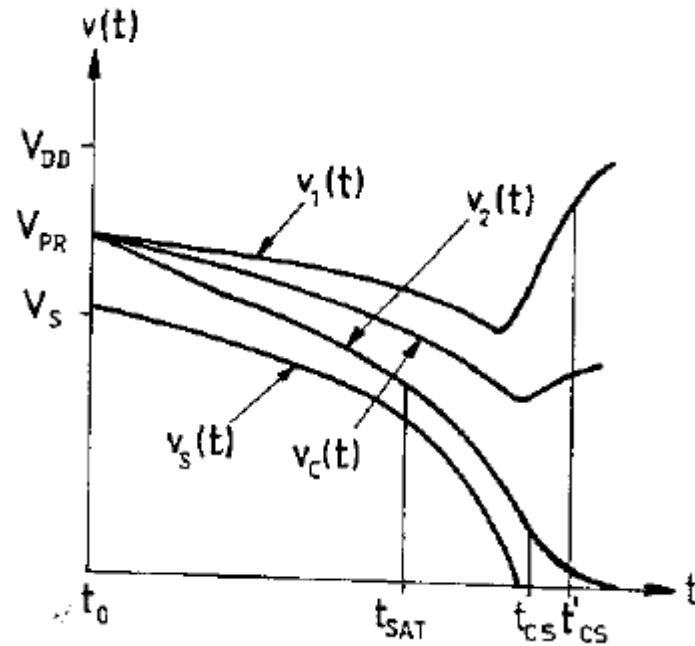
- During reading, if LHS cell is read, the dummy word line on RHS is accessed (and vice-versa).

- sense amplifier

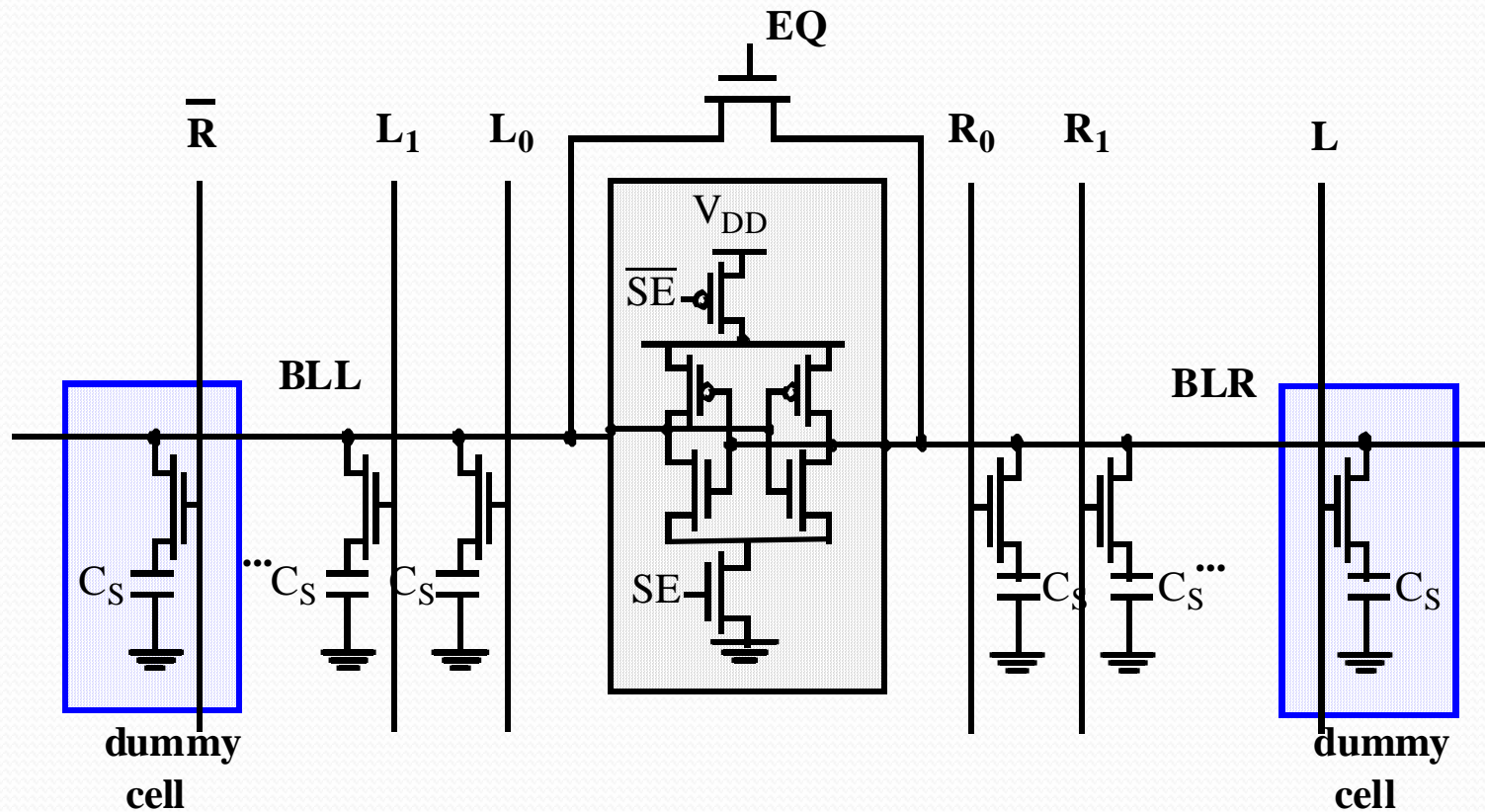


Positive feedback differential voltage sense amplifier circuit.

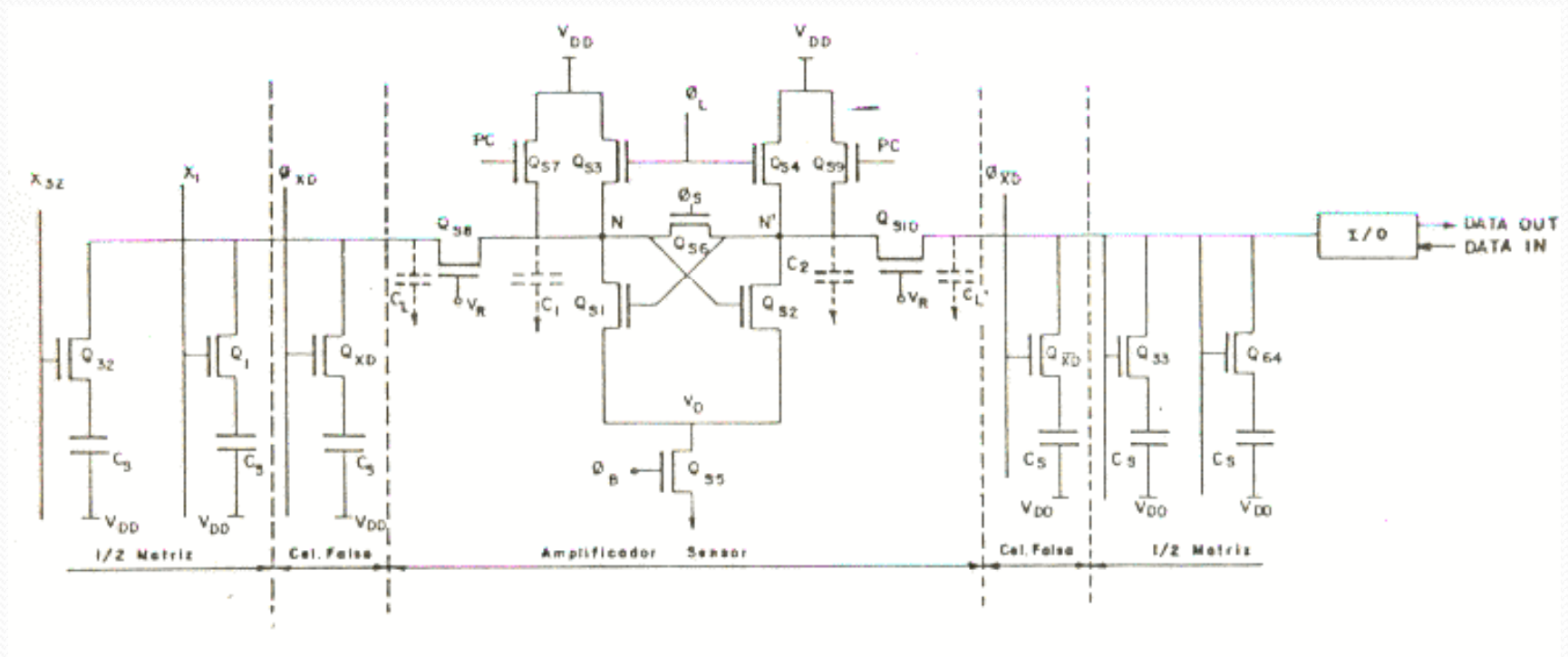
Time [Arbitrary Units]



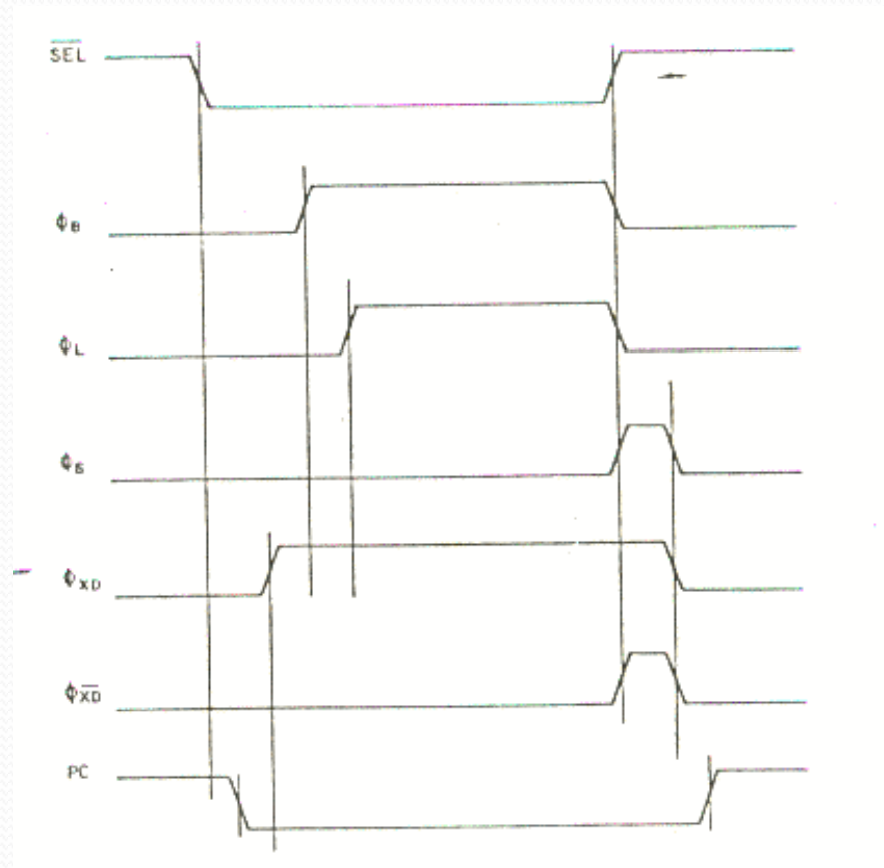
Sense amplifier with dummy-cell sensing structure



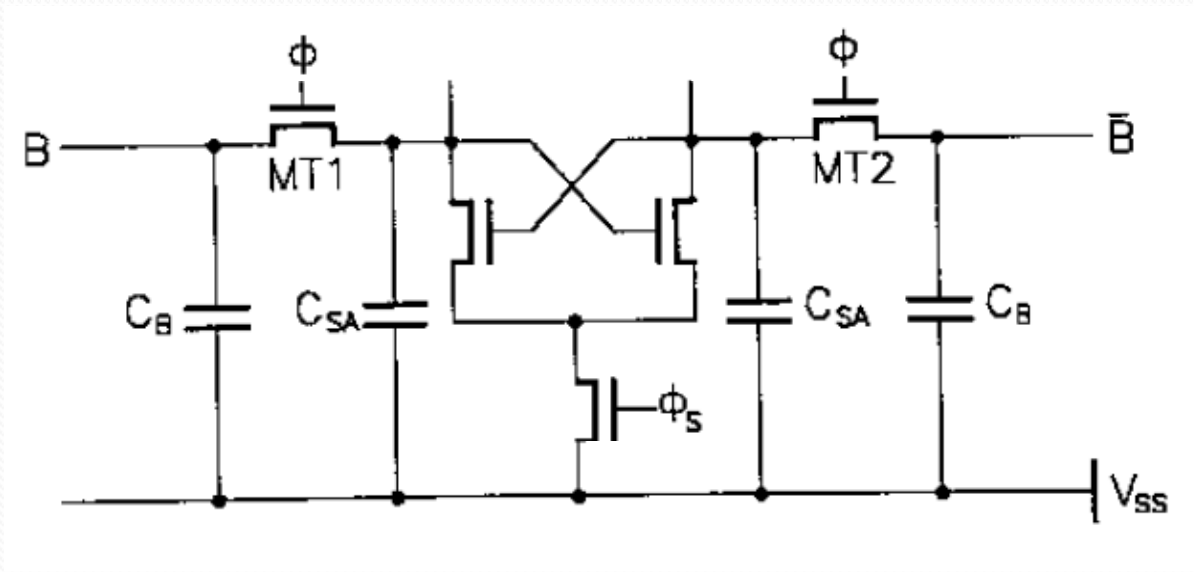
Sense amplifier with dummy-cell sensing structure



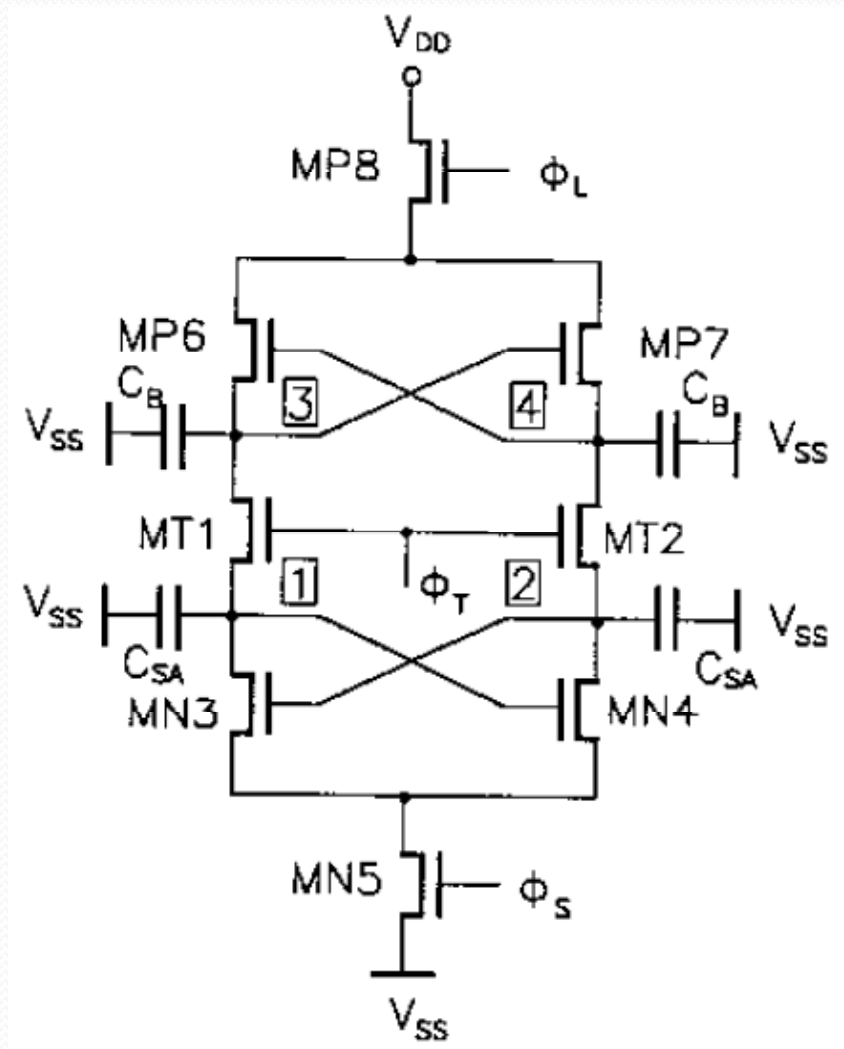
Sense-amplifier timing (dummy-cell approach)



- Decoupling of bit-line capacitance



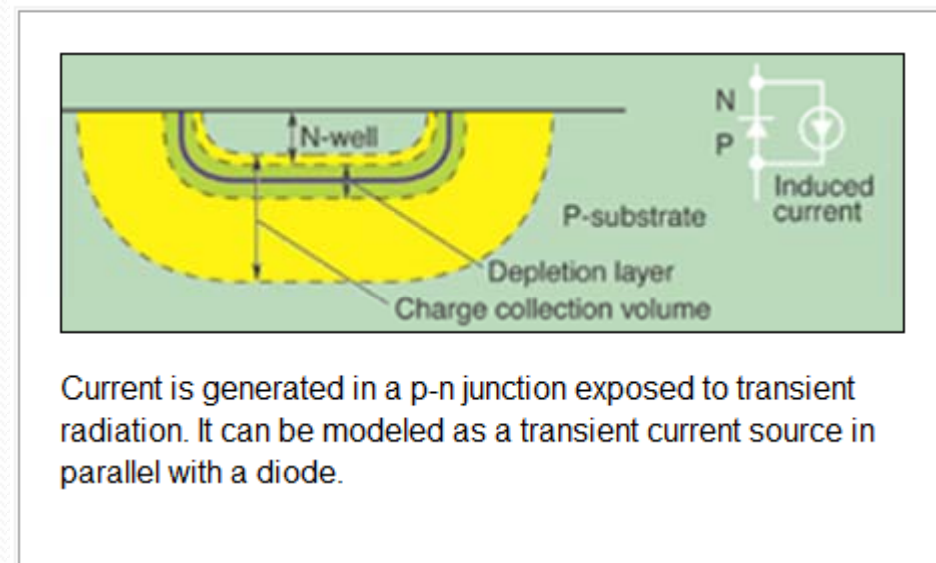
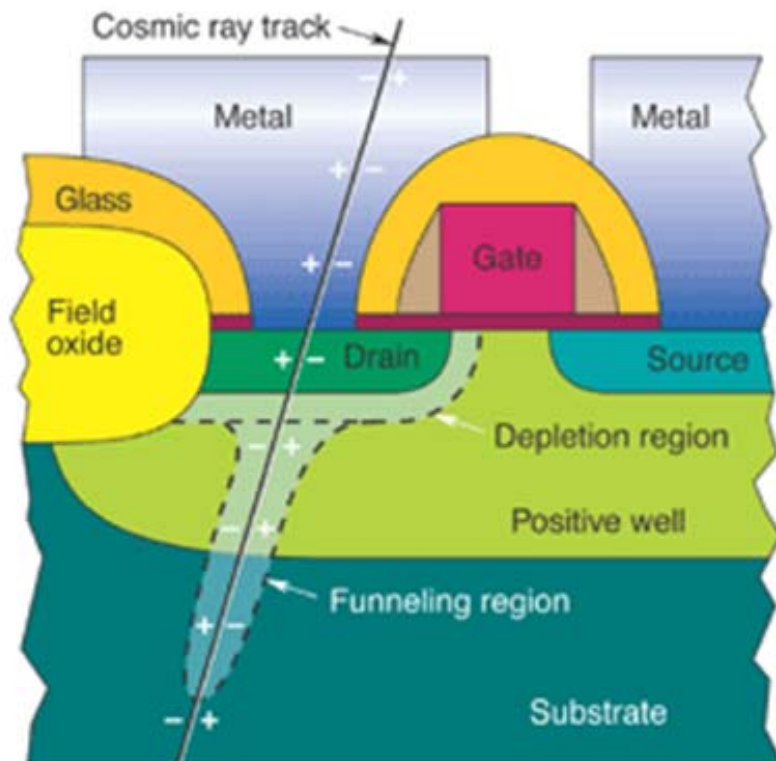
- decoupling devices MT1 and MT2 are initially turned on. SA is not active and load at its input is $C_{BL} + C_{SA}$, with $C_B \gg C_{SA}$.
- AS is activated when differential input-voltage V_{in} is large enough. At this moment, MT1 and MT2 are turned off and load at AS input is only C_{SA} .
 \Rightarrow Faster arbitration.
- during re-writing, MT1 and MT2 are turned on again.



Sense amplifier incorporating decoupler devices.

- initially, Φ_s turns on MN5.
- Φ_t turns on decoupling devices MT1 and MT2. Both are turned off with differential V_{in} is large enough.
- V_{in} is then quickly amplified by MN3, MN4 and MN5, as bit-line capacitances are de-coupled.
- when V_{in} reaches a given value, Φ_t goes high, turning on MT1 and MT2, which allows full V_{in} swing.

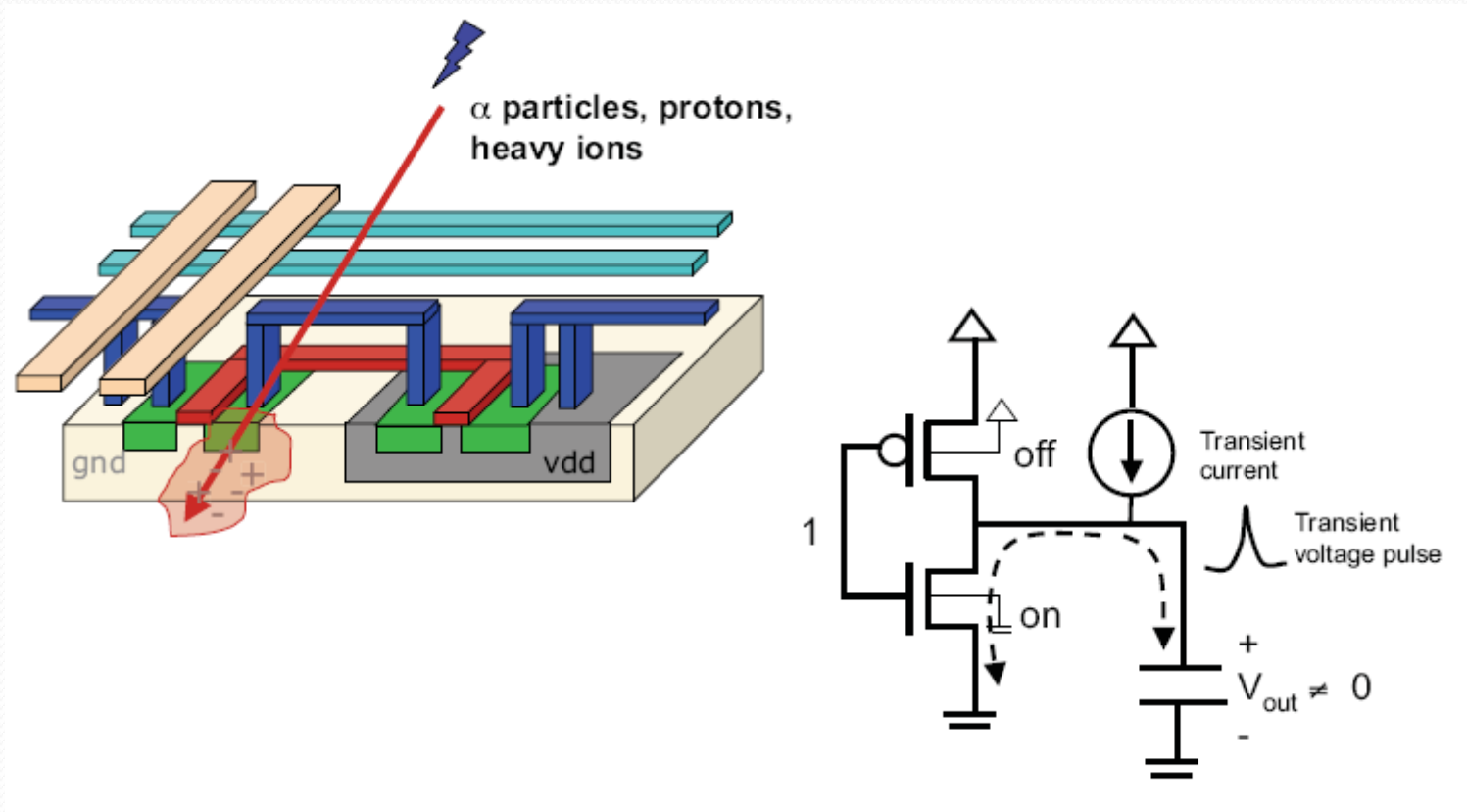
Alpha-particles



Current is generated in a p-n junction exposed to transient radiation. It can be modeled as a transient current source in parallel with a diode.




1 particle ~ 1million carriers

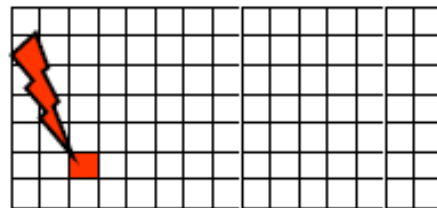
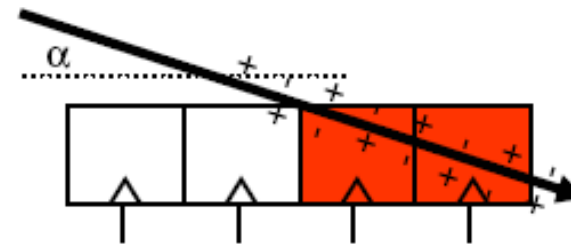
Alpha-particles



Challenges in Sequential Logic

MULTIPLE BIT UPSETS

- Particle incidence angle
- Transistor Dimensions 
- Voltage Supply 
- Memory Array Density 



Single memory cell

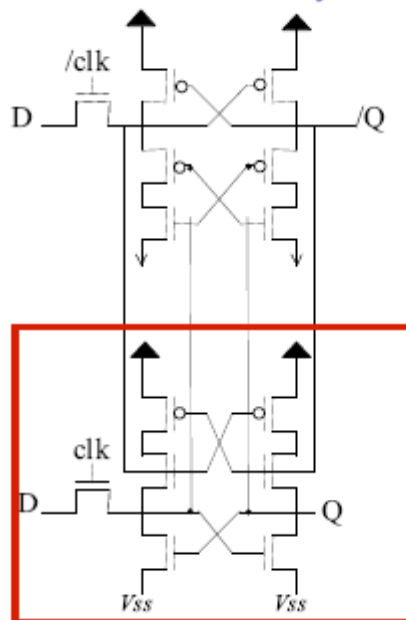


Multiple memory cells

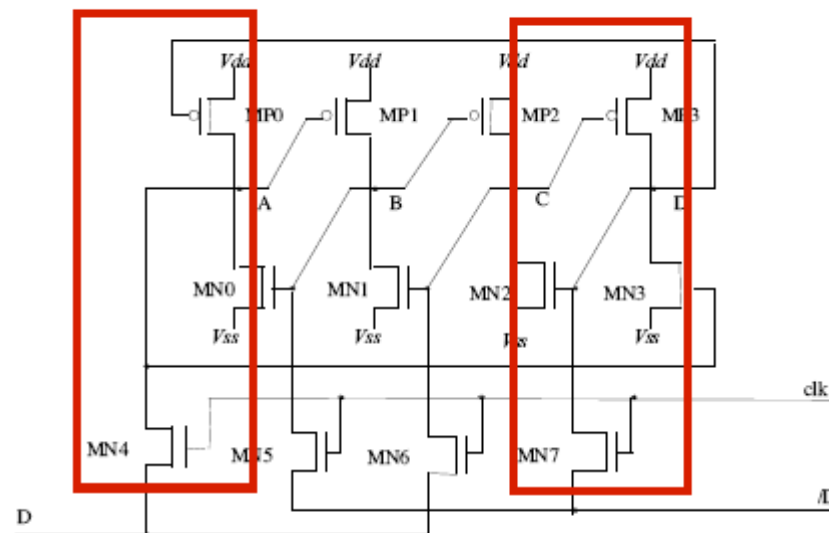
Hardened Memories

- The principle is to store the data in two different locations within the cell in such way that the corrupted part can be restored.

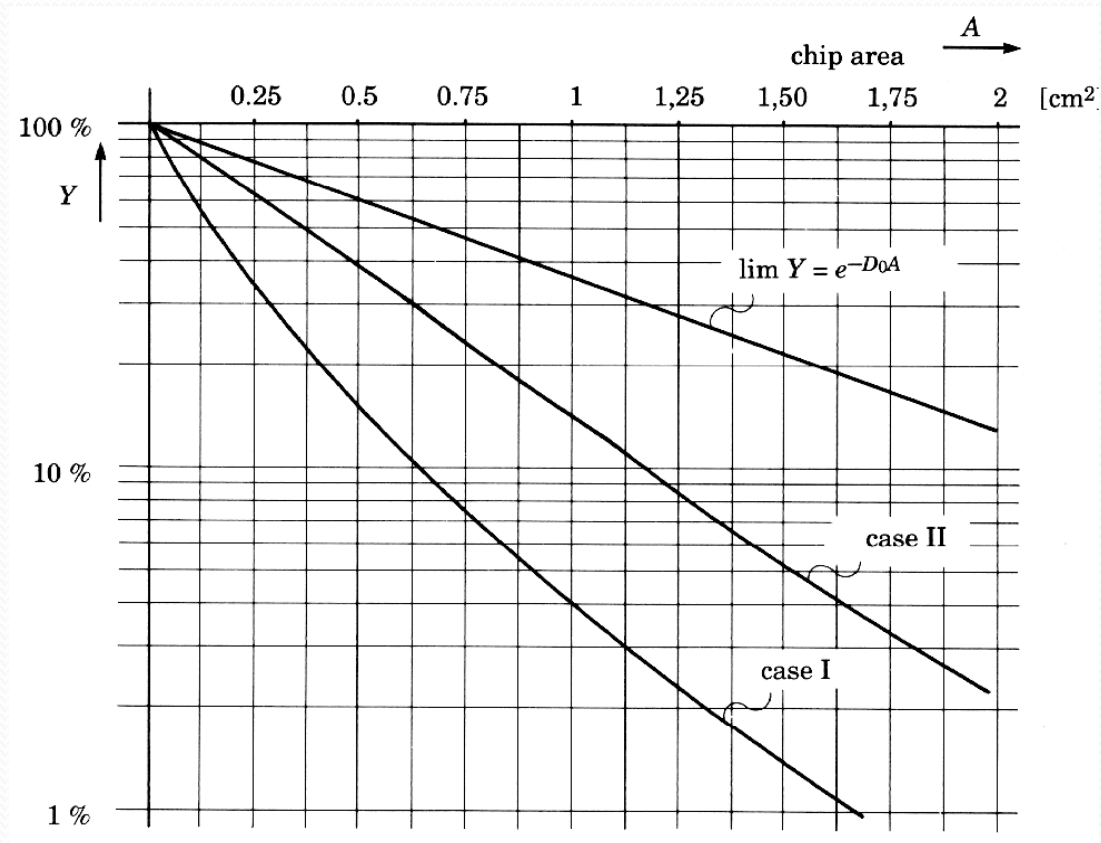
Whitaker/Liu Memory Cell [Liu, 92]



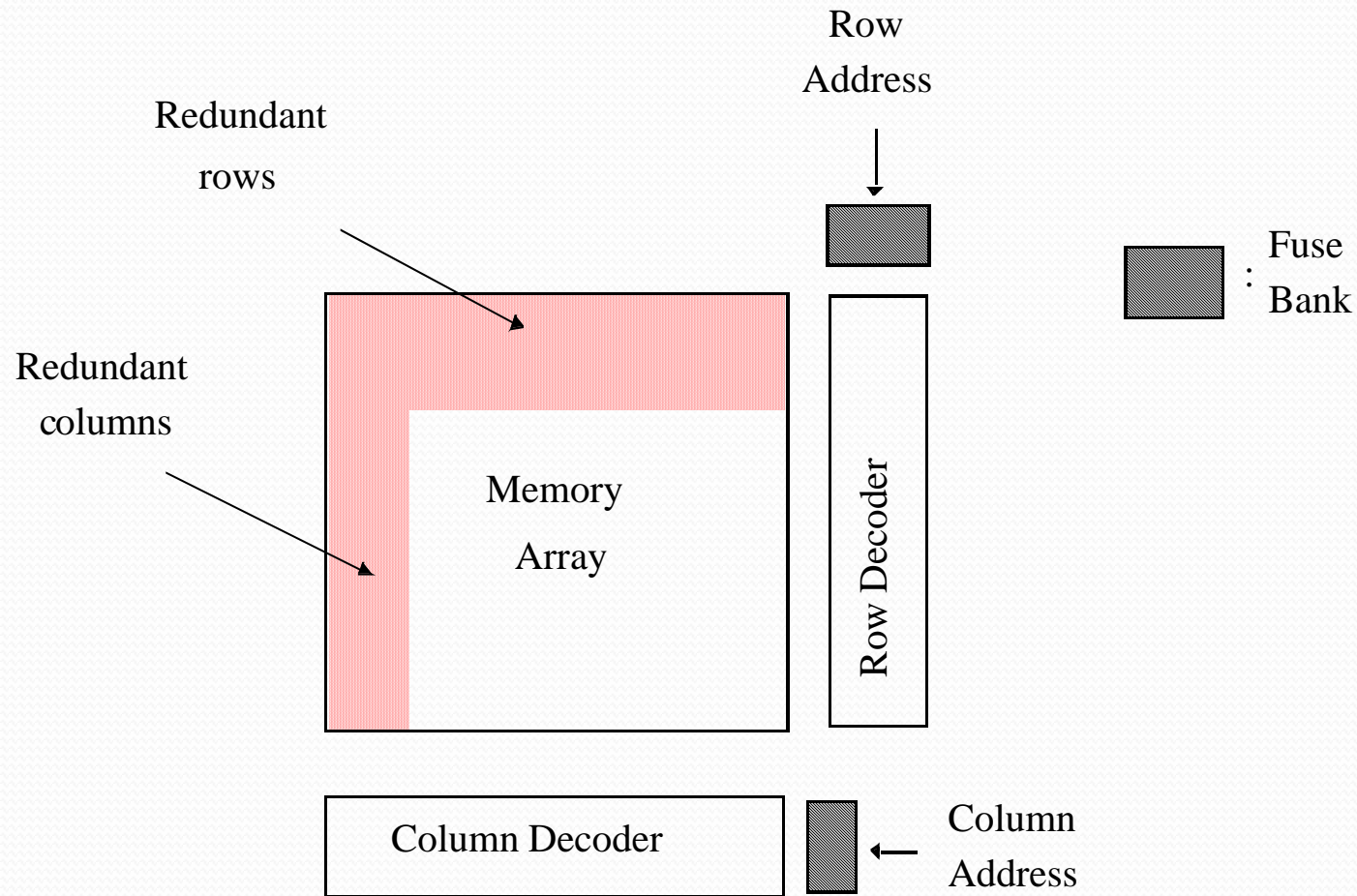
DICE Memory Cell [Calin, 96]



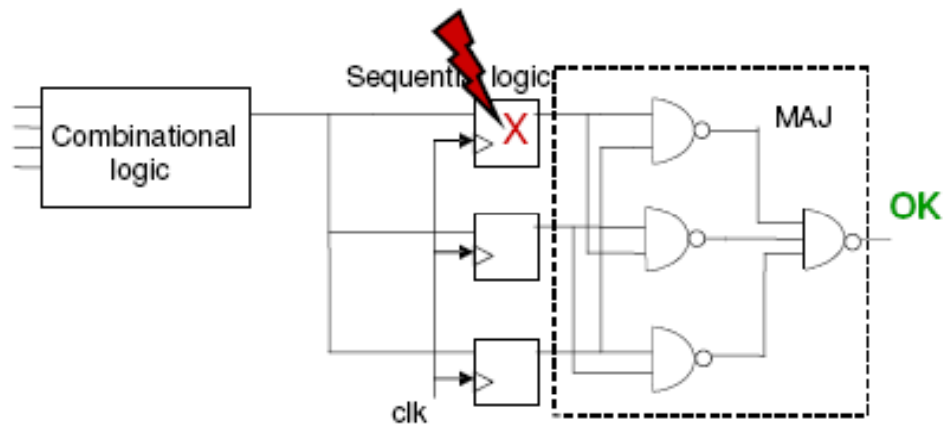
Yield



Redundancy

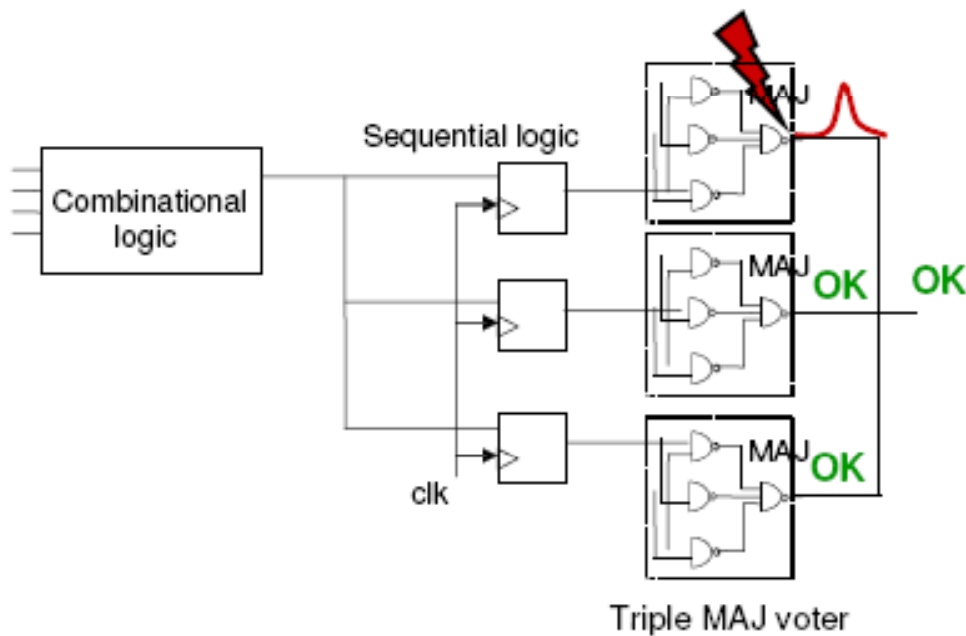


Triple Modular Redundancy



inputs	MAJ
000	0
001	0
010	0
011	1
100	0
101	1
110	1
111	1

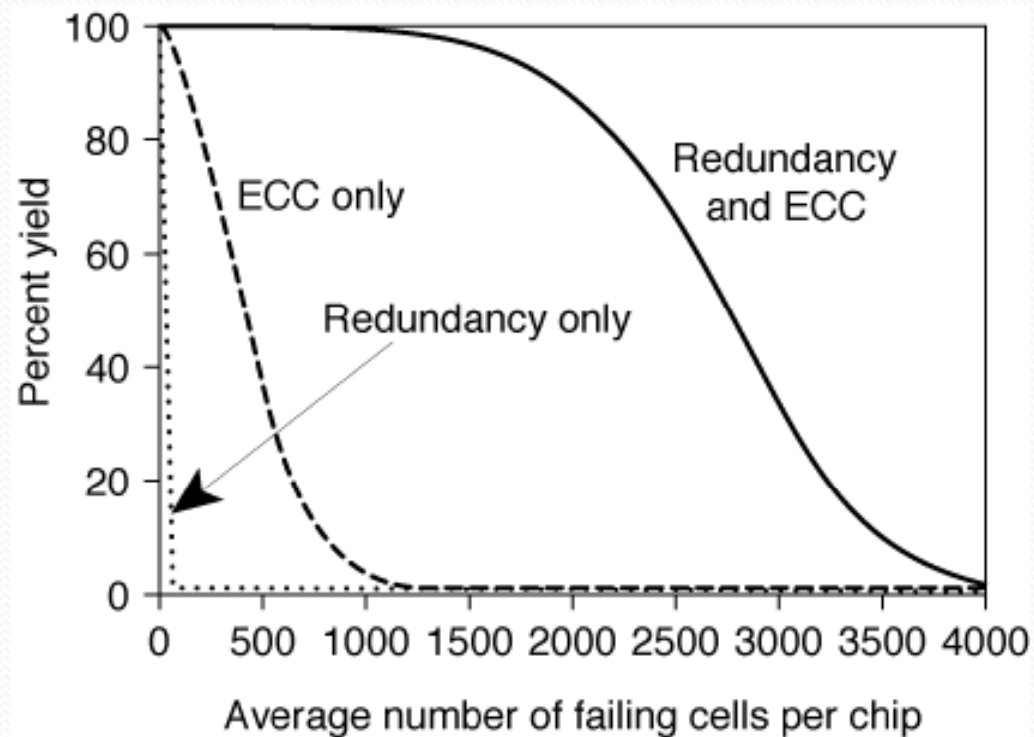
Triple Modular Redundancy



Current strength

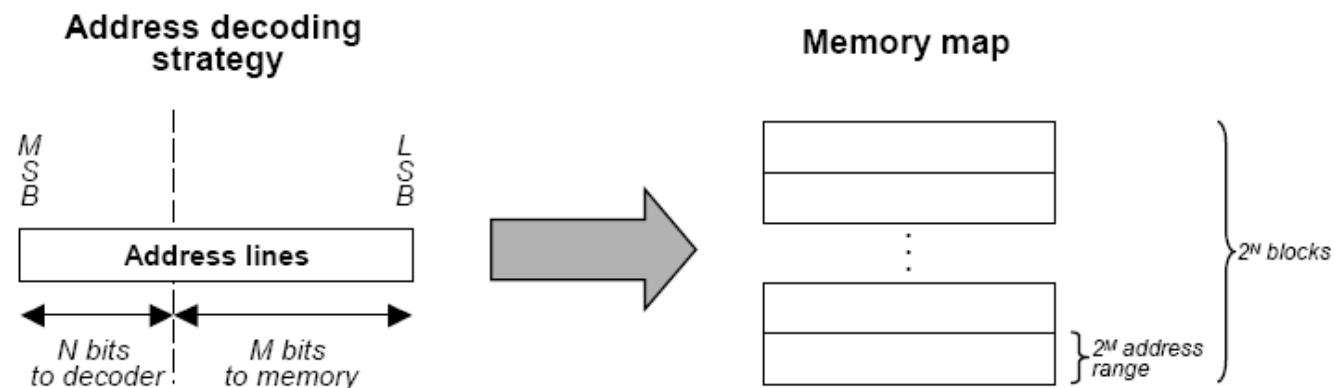
- Increases current drive helping keeping the node in the original value.

Redundancy and Error Correction



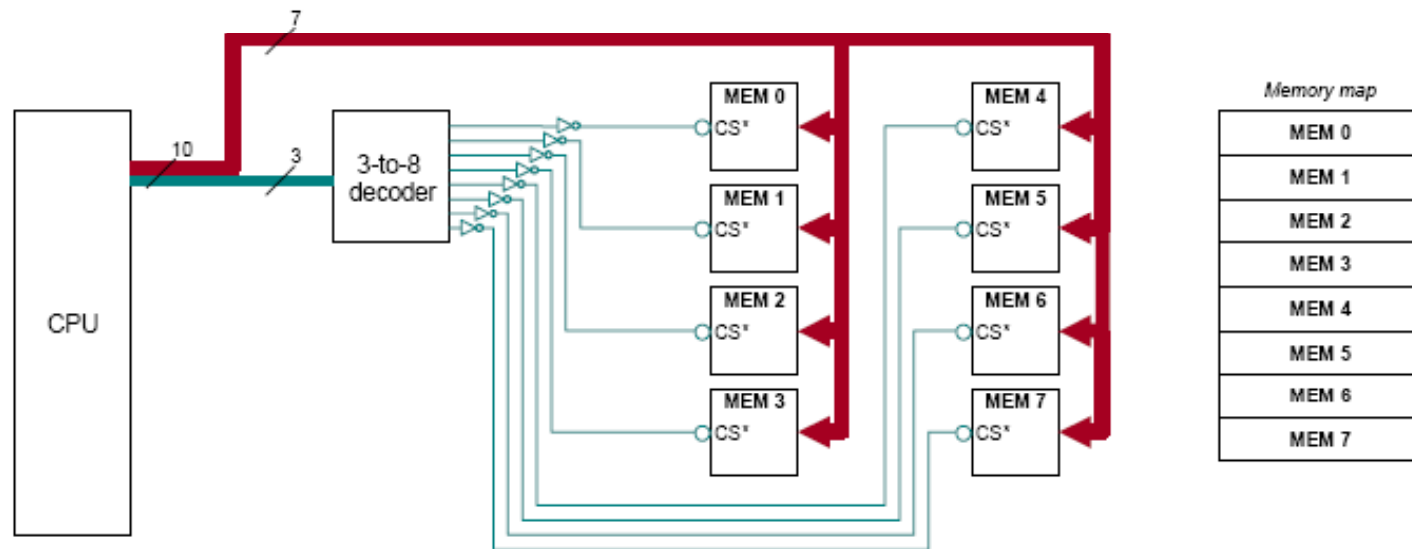
Address Decoding

- Address decoding is the process of generating chip select (CS^*) signals from the address bus for each device in the system
- The address bus lines are split into two sections
 - the N most significant bits are used to generate the CS^* signals for the different devices
 - the M least significant signals are passed to the devices as addresses to the different memory cells or internal registers



Example:

- Let's assume a very simple microprocessor with 10 address lines (1KB memory)
- Let's assume we wish to implement all its memory space and we use 128x8 memory chips
- **SOLUTION**
 - We will need 8 memory chips ($8 \times 128 = 1024$)
 - We will need 3 address lines to select each one of the 8 chips
 - Each chip will need 7 address lines to address its internal memory cells



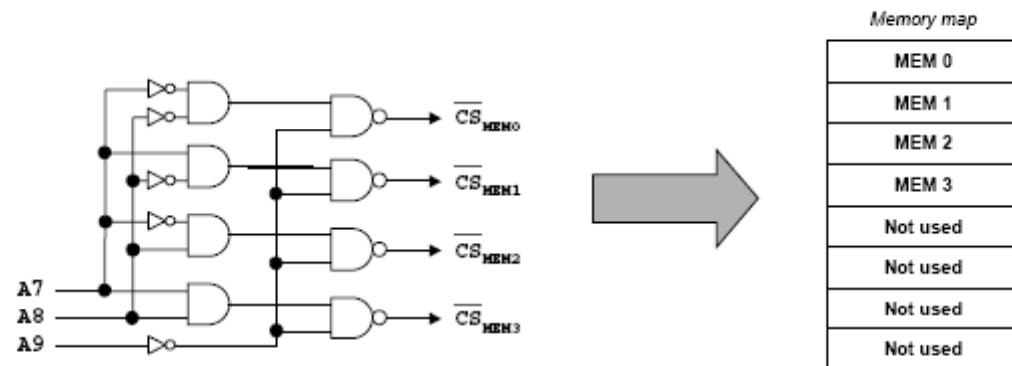
Full address decoding

■ Let's assume the same microprocessor with 10 address lines (1KB memory)

- However, this time we wish to implement only 512 bytes of memory
- We still must use 128-byte memory chips
- Physical memory must be placed on the upper half of the memory map

■ SOLUTION

Device	Used for Address Decoding			Used to reference memory cells on each memory IC						
	A9	A8	A7	A6	A5	A4	A3	A2	A1	A0
MEM 0	0	0	0	X	X	X	X	X	X	X
MEM 1	0	0	1	X	X	X	X	X	X	X
MEM 2	0	1	0	X	X	X	X	X	X	X
MEM 3	0	1	1	X	X	X	X	X	X	X



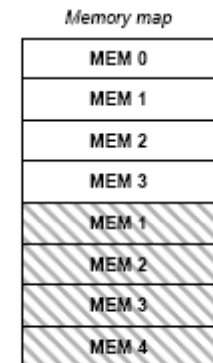
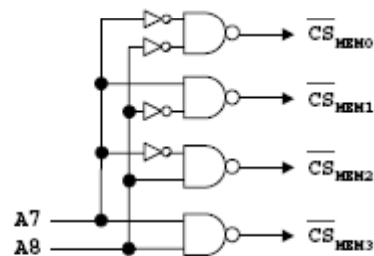
Partial address decoding

- Let's assume the same microprocessor with 10 address lines (1KB memory)

- However, this time we wish to implement only 512 bytes of memory
- We still must use 128-byte memory chips
- Physical memory must be placed on the upper half of the memory map

- SOLUTION**

Device	Not used		Used for Address Decoding		Used to reference memory cells on each memory IC					
	A9	A8	A7	A6	A5	A4	A3	A2	A1	A0
MEM 0	X	0	0	X	X	X	X	X	X	X
MEM 1	X	0	1	X	X	X	X	X	X	X
MEM 2	X	1	0	X	X	X	X	X	X	X
MEM 3	X	1	1	X	X	X	X	X	X	X



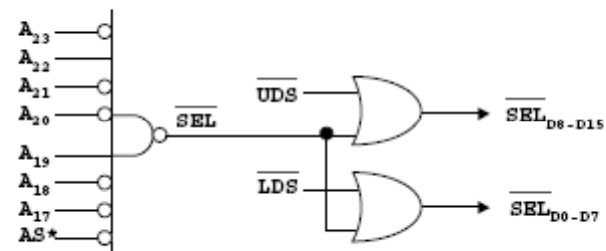
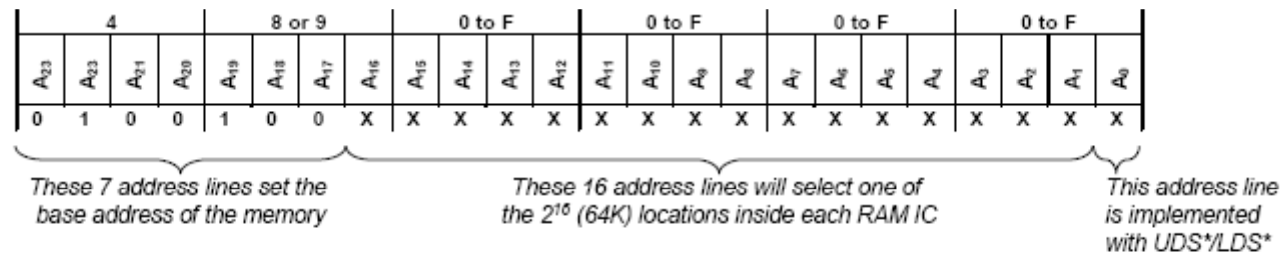
Example 1

- A circuit containing 64K words of RAM is to be interfaced to a 68000-based system, so that the first address of RAM (the base address) is at \$480000

- What is the entire range of RAM addresses?
- Design a FULL address decoder using two 64Kx8 RAM ICs

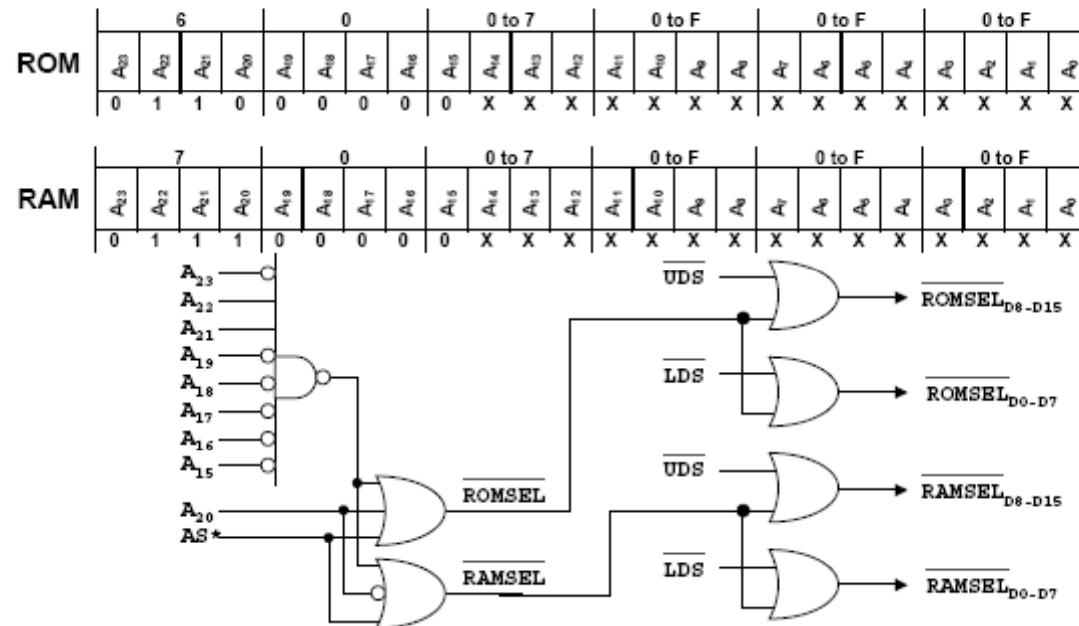
■ Solution

- The address range for the RAM is from \$480000 to \$480000+(128K=\$20000)=\$4A0000-1=\$49FFFF
- The two ICs must be differentiated through UDS*/LDS* (since the 68000 DOES NOT have A₀)



Example 2

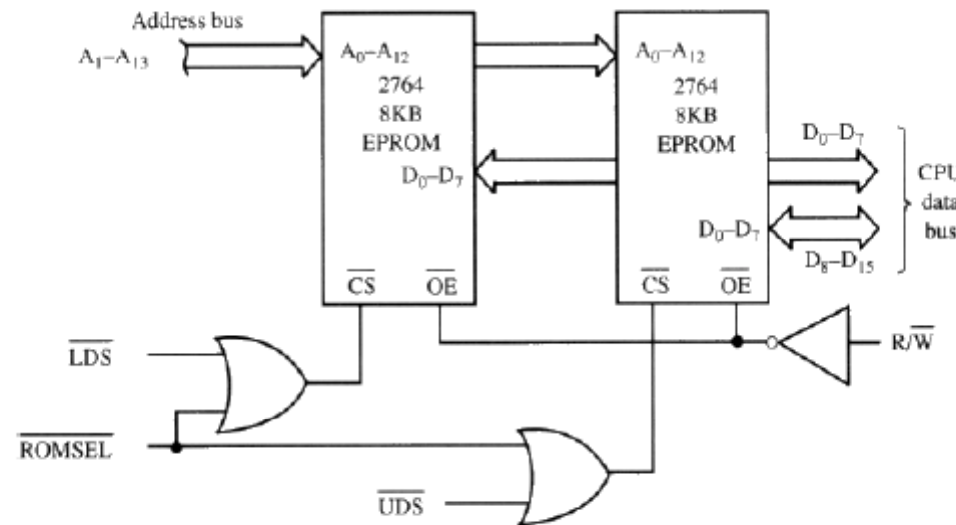
- A 68000-based system is to be built with these memory requirements
 - a 16K word EPROM with a starting address of \$60 0000
 - a 16K word RAM with a starting address of \$70 0000
- Design a FULL address decoder for this application using 16K×8 chips for both EPROM and RAM
 - $\$60\ 0000 + (16\text{KWord}=32\text{KB}=\$8000)-1 = \$60\ 7\text{FFF}$
 - $\$70\ 0000 + (16\text{KWord}=32\text{KB}=\$8000)-1 = \$70\ 7\text{FFF}$



Example 3

- Design a PARTIAL address decoder for a 68000-based system with only 8K words of EPROM space, and a base address at \$4000, using 8Kx8 memory chips

A ₂₃	A ₂₂	A ₂₁	A ₂₀	A ₁₉	A ₁₈	A ₁₇	A ₁₆	A ₁₅	A ₁₄	A ₁₃	A ₁₂	A ₁₁	A ₁₀	A ₉	A ₈	A ₇	A ₆	A ₅	A ₄	A ₃	A ₂	A ₁	A ₀	
0	0	0	0	0	0	0	0	0	1	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X



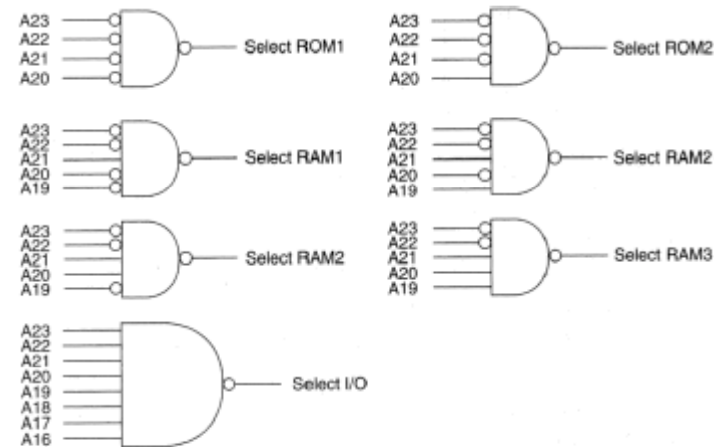
Example 4

Design a partial address decoder for a 68000-based system that contains

- 2MB of EPROM at a starting address \$00 0000 using 512Kx8 chips
- 2MB of RAM at a starting address \$10 0000 using 256Kx8 chips
- 64KB I/O space starting at \$FF0000

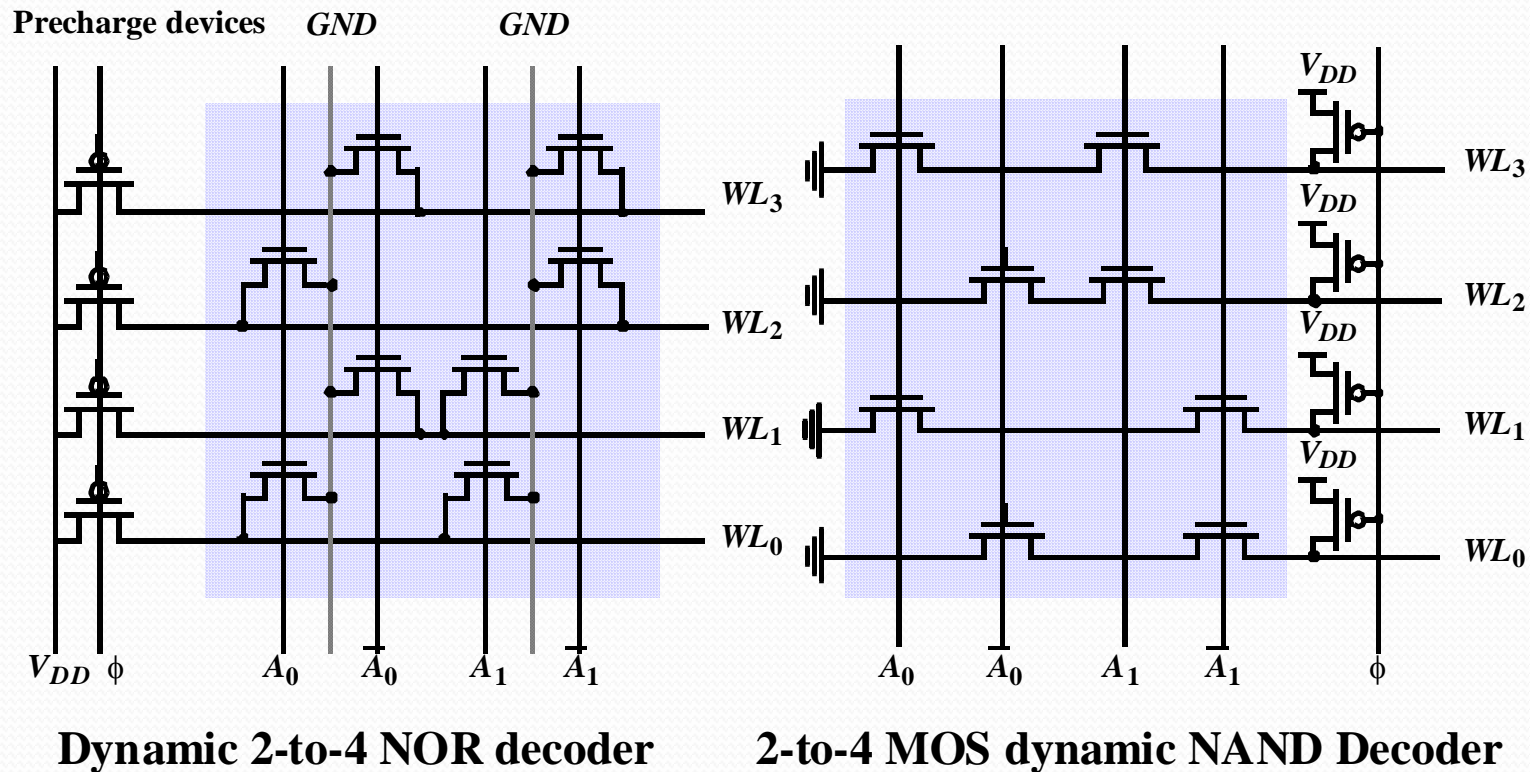
SOLUTION

- For the EPROM we will need 4 512Kx8 chips, organized as 2 pairs of 512x8 chips (in order to use UDS*/LDS*). We will call these pairs ROM1 and ROM2
- For the RAM we will need 8 256Kx8 chips, organized as 4 pairs of 256Kx8: RAM1 to RAM4



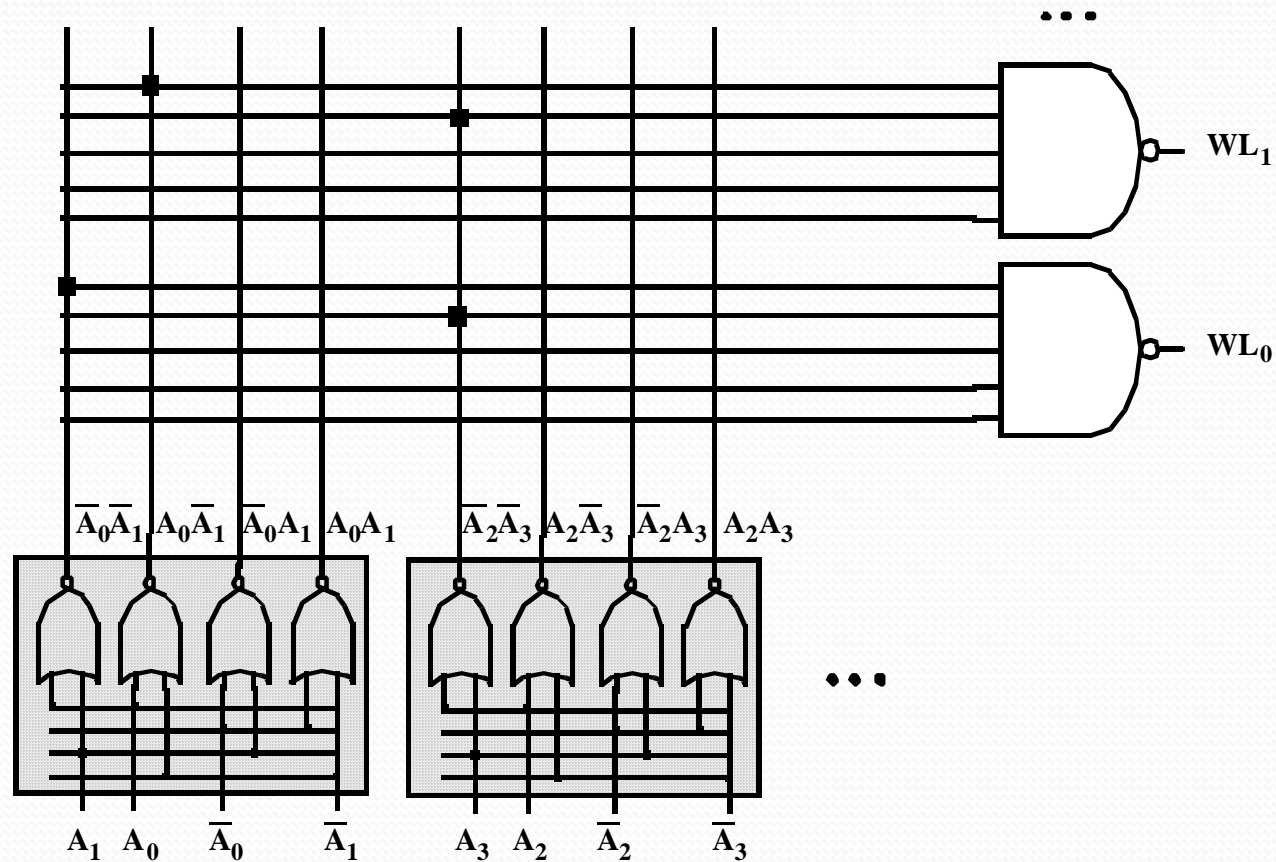
	A ₂₃	A ₂₂	A ₂₁	A ₂₀	A ₁₉	A ₁₈	A ₁₇	A ₁₆	A ₁₅	A ₁₄	A ₁₃	A ₁₂	A ₁₁	A ₁₀	A ₀₉	A ₀₈	A ₀₇	A ₀₆	A ₀₅	A ₀₄	A ₀₃	A ₀₂	A ₀₁	A ₀₀
ROM1	0	0	0	0	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
ROM2	0	0	0	1	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
RAM1	0	0	1	0	0	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
RAM2	0	0	1	0	1	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
RAM3	0	0	1	1	0	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
RAM4	0	0	1	1	1	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
I/O	1	1	1	1	1	1	1	1	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X

Dynamic Decoding



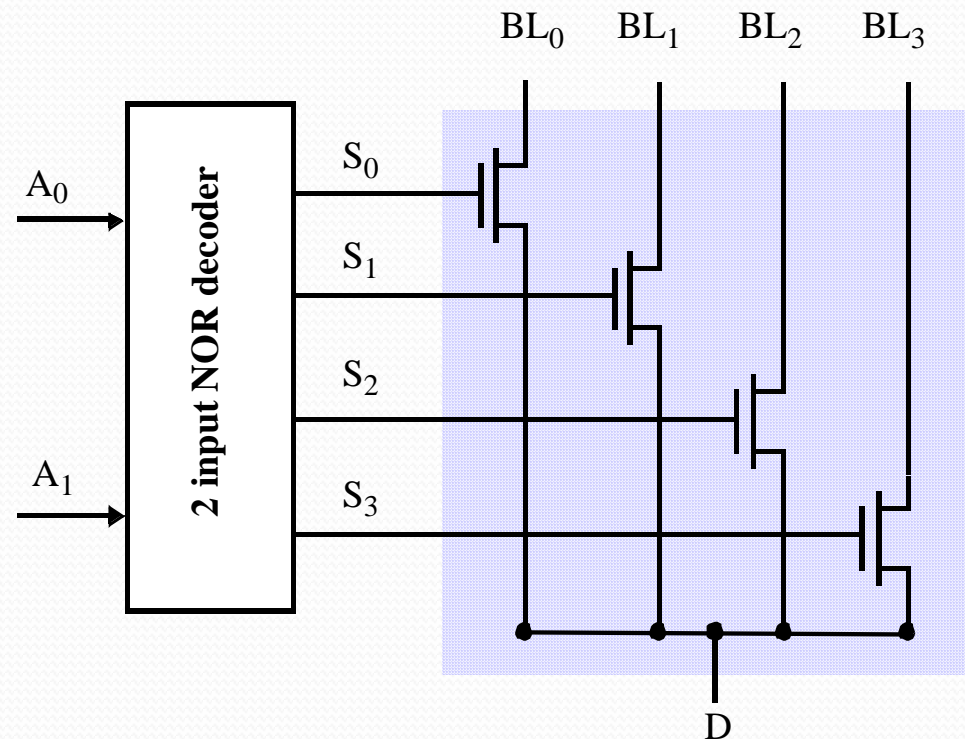
Propagation delay is primary concern

A NAND decoder using 2-input pre-decoders



Splitting decoder into two or more logic layers produces a faster and cheaper implementation

4 input pass-transistor based column decoder

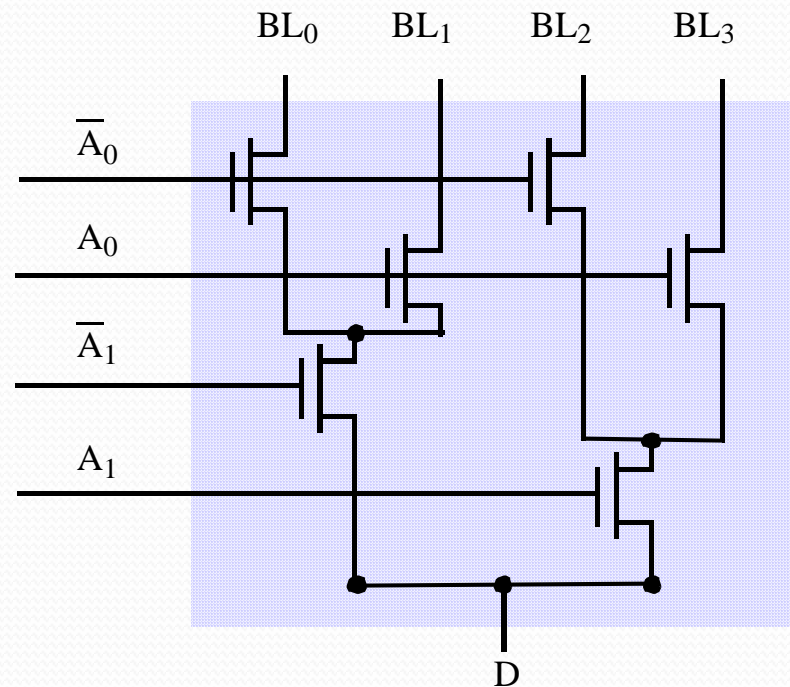


Advantage: speed (t_{pd} does not add to overall memory access time)

only 1 extra transistor in signal path

Disadvantage: large transistor count

4-to-1 tree based column decoder



Number of devices drastically reduced

Delay increases quadratically with # of sections; prohibitive for large decoders

Solutions: buffers

progressive sizing

combination of tree and pass transistor approaches