# *Digital Integrated Circuits*
## *A Design Perspective*

Jan M. Rabaey
Anantha Chandrakasan
Borivoje Nikolic

# Semiconductor Memories

*December 20, 2002*
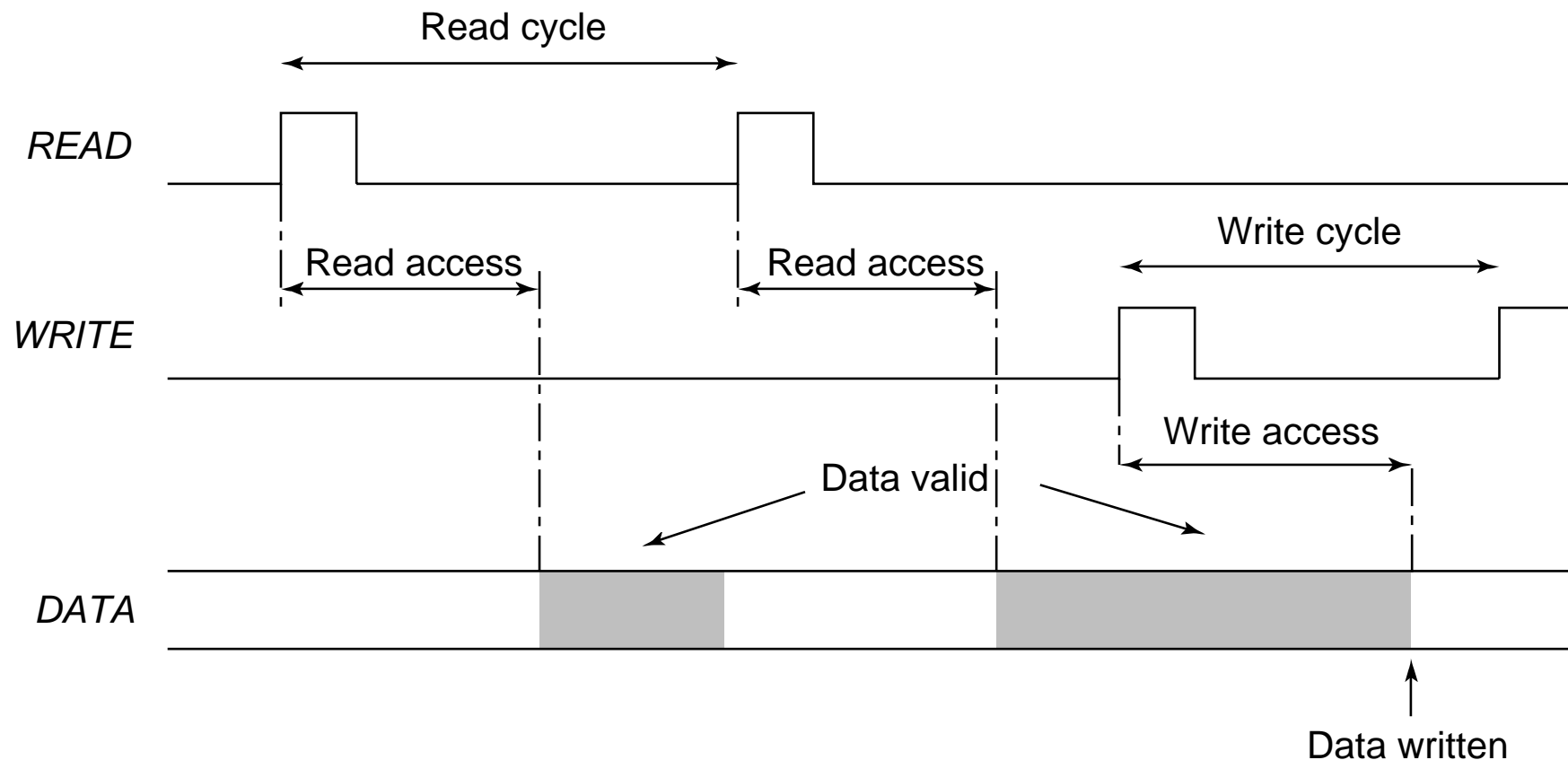
# *Chapter Overview*

- ❑ **Memory Classification**
- ❑ **Memory Architectures**
- ❑ **The Memory Core**
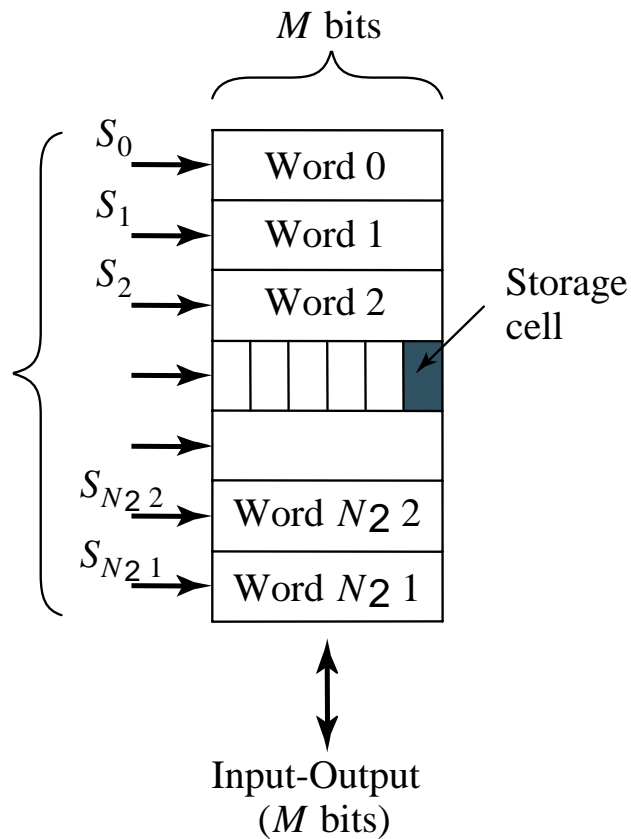- ❑ **Periphery**
- ❑ **Reliability**
- ❑ **Case Studies**

# *Semiconductor Memory Classification*

| Read-Write Memory | | Non-Volatile Read-Write Memory | Read-Only Memory |
|---|---|---|---|
| **Random Access** | **Non-Random Access** | EPROM $E^2PROM$ FLASH | Mask-Programmed Programmable (PROM) |
| SRAM DRAM | FIFO LIFO Shift Register CAM | | |

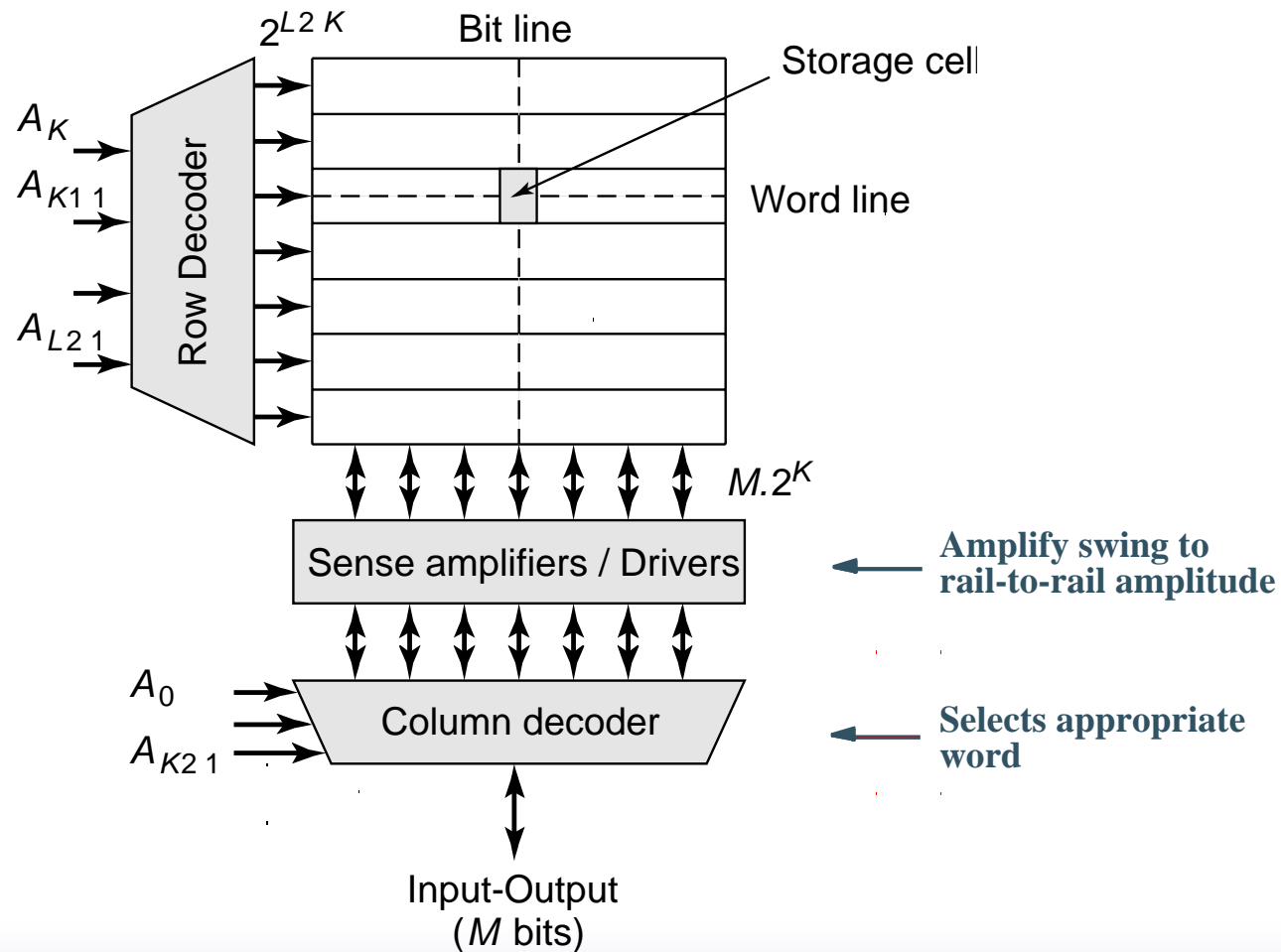# *Memory Timing: Definitions*

# *Memory Architecture: Decoders*

$M$ bits

$S_0$ → Word 0
$S_1$ → Word 1
$S_2$ → Word 2

→ Storage cell

→

$S_{N2\,2}$ → Word $N_2$ 2
$S_{N2\,1}$ → Word $N_2$ 1

Input-Output
($M$ bits)

**Intuitive architecture for N x M memory
Too many select signals:
N words == N select signals**

$M$ bits

$S_0$ → Word 0
→ Word 1

$A_0$ →
$A_1$ → Word 2

→ Storage cell

→

$A_{K2\,1}$ →
→ Word $N_2$ 2
→ Word $N_2$ 1

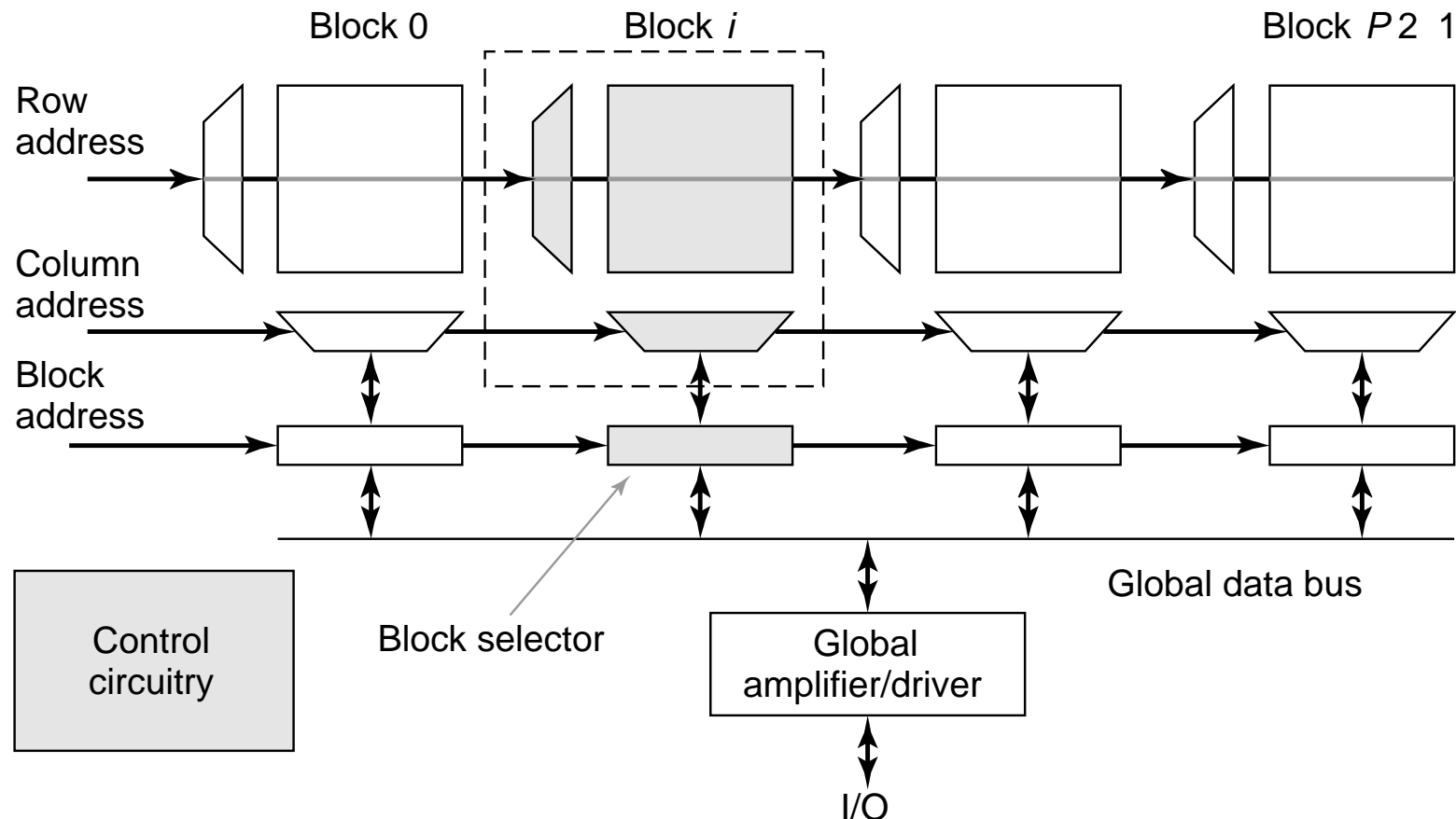$K\ 5\ \log_2 N$

Input-Output
($M$ bits)

**Decoder reduces the number of select signals**
$$K = log_2 N$$

# Array-Structured Memory Architecture
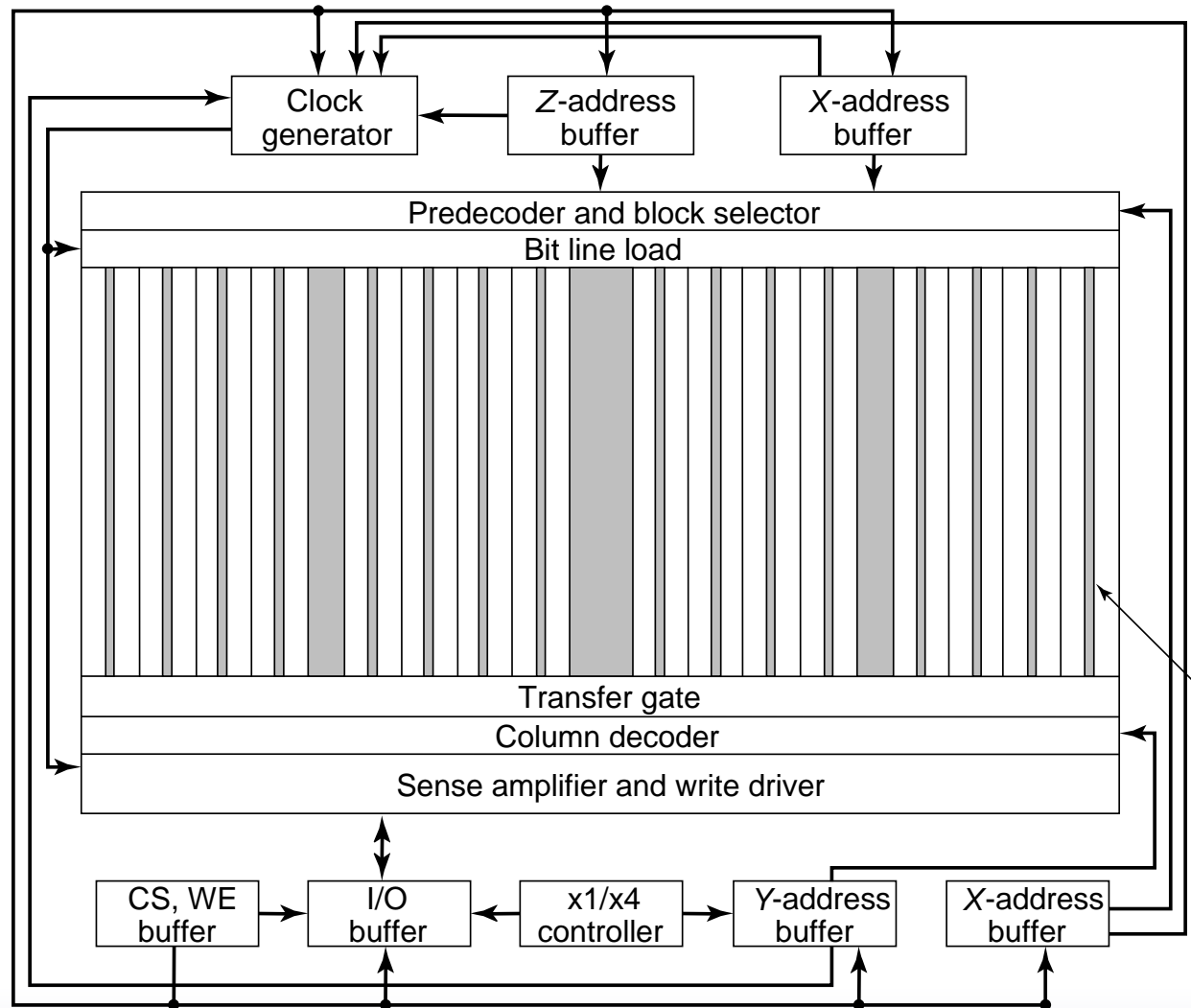
## Problem: ASPECT RATIO or HEIGHT >> WIDTH
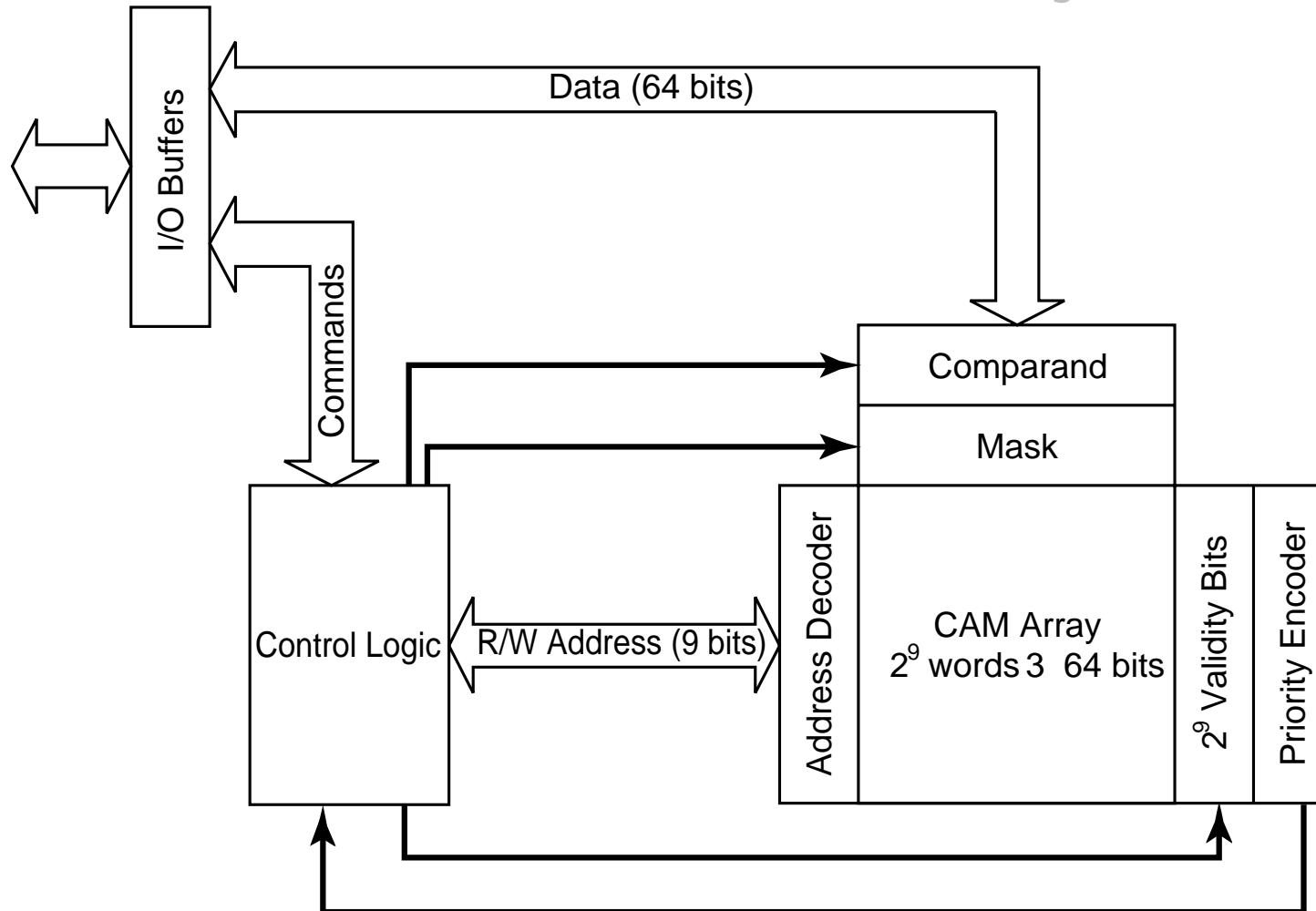
# Hierarchical Memory Architecture



**Advantages:**
**1. Shorter wires within blocks**
**2. Block address activates only 1 block => power savings**
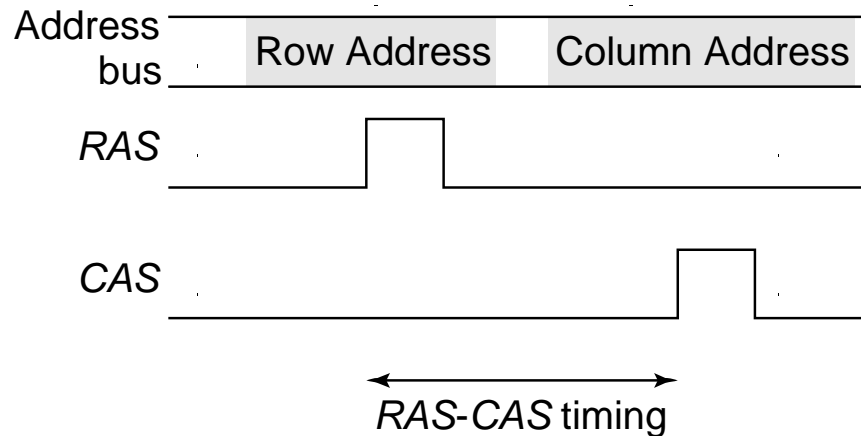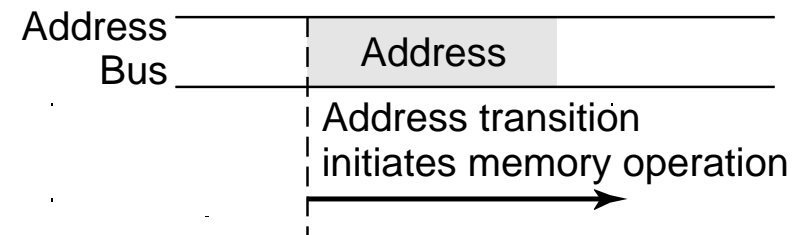
# *Block Diagram of 4 Mbit SRAM*

# *Contents-Addressable Memory*



Data (64 bits)
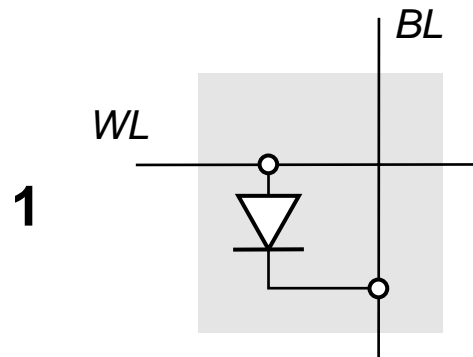
I/O Buffers

Commands

Comparand

Mask

Control Logic

R/W Address (9 bits)

Address Decoder

CAM Array
$2^9$ words $\times$ 64 bits

$2^9$ Validity Bits

Priority Encoder

# *Memory Timing: Approaches*

Address bus: Row Address | Column Address

*RAS*

*CAS*

*RAS-CAS* timing

**DRAM Timing
Multiplexed Adressing**

Address Bus: Address

Address transition
initiates memory operation

**SRAM Timing
Self-timed**

# Read-Only Memory Cells

**1**

BL

WL
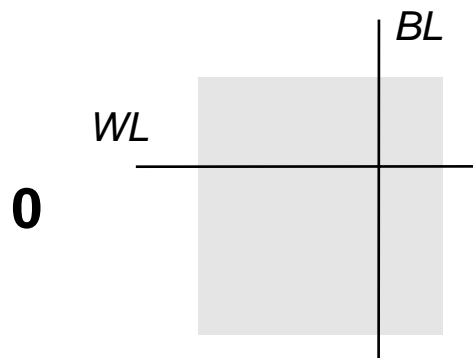
BL

$V_{DD}$

WL

BL

WL

**0**

BL

WL

BL

WL

BL

WL

GND

**Diode ROM**

**MOS ROM 1**

**MOS ROM 2**

# MOS OR ROM

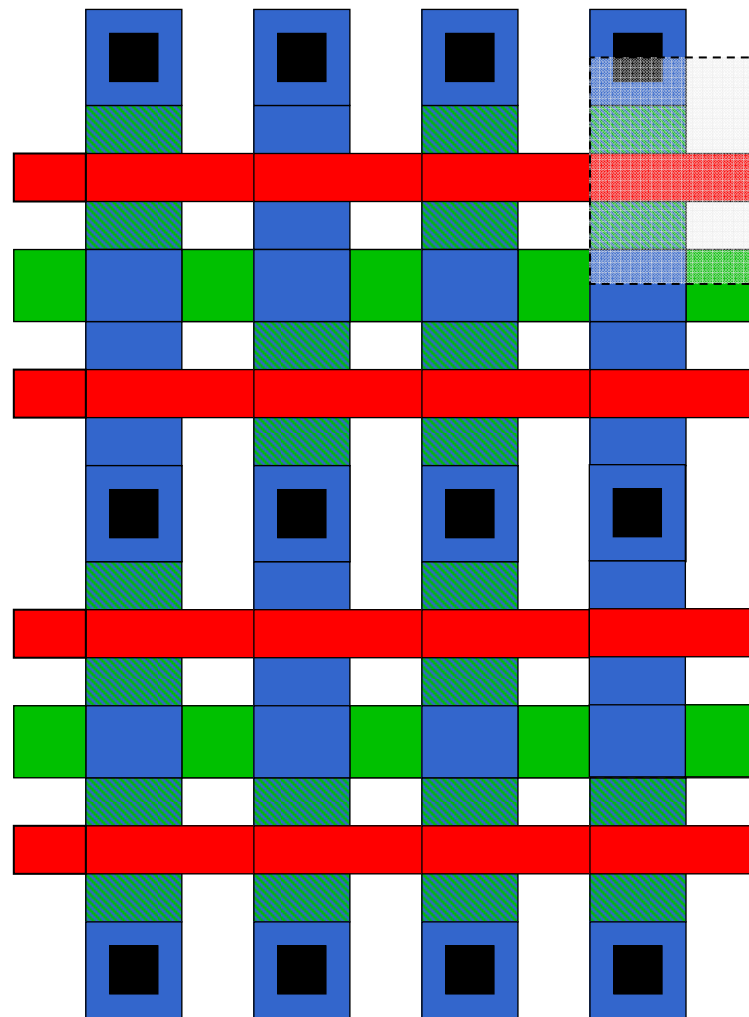BL[0]  BL[1]  BL[2]  BL[3]

WL[0]

$V_{DD}$

WL[1]

WL[2]

$V_{DD}$

WL[3]

$V_{bias}$

**Pull-down loads**

# MOS NOR ROM

# MOS NOR ROM Layout

Cell (9.5λ x 7λ)

**Programming using the Active Layer Only**

Polysilicon

Metal1

Diffusion

Metal1 on Diffusion

# MOS NOR ROM Layout

Cell (11λ x 7λ)

**Programmming using the Contact Layer Only**

| | |
|---|---|
| 🟥 | Polysilicon |
| 🟦 | Metal1 |
| 🟩 | Diffusion |
| 🟩 | Metal1 on Diffusion |

# MOS NAND ROM



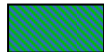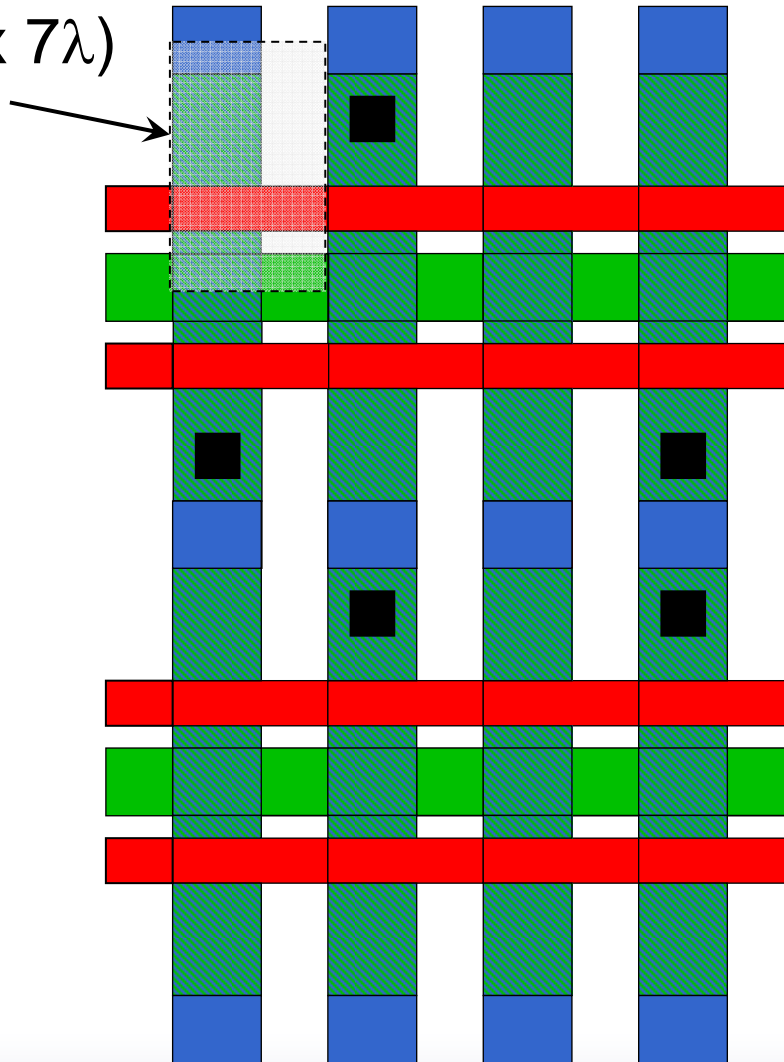**All word lines high by default with exception of selected row**

# MOS NAND ROM Layout

Cell ($8\lambda \times 7\lambda$)

**Programmming using the Metal-1 Layer Only**

**No contact to VDD or GND necessary; drastically reduced cell size**
**Loss in performance compared to NOR ROM**

■ Polysilicon

■ Diffusion

■ Metal1 on Diffusion

# NAND ROM Layout

Cell (5λ x 6λ)

**Programmming using Implants Only**

Polysilicon

Threshold-altering implant

Metal1 on Diffusion

# Decreasing Word Line Delay

Driver

WL             Polysilicon word line

Metal word line

## (a) Driving the word line from both sides

Metal bypass

WL     K cells                      Polysilicon word line

## (b) Using a metal bypass

## (c) Use silicides

# Precharged MOS NOR ROM



**PMOS precharge device can be made as large as necessary, but clock driver becomes harder to design.**

# Non-Volatile Memories
# The Floating-gate transistor (FAMOS)



**Device cross-section**

**Schematic symbol**

# Floating-Gate Transistor Programming



**Avalanche injection**

**Removing programming voltage leaves charge trapped**

**Programming results in higher $V_T$.**

# A "Programmable-Threshold" Transistor

# FLOTOX EEPROM



Floating gate      Gate

Source      Drain

20–30 nm

$n^1$      $n^1$

Substrate
p

10 nm

**FLOTOX transistor**

$I$

-10 V

$V_{GD}$

10 V

**Fowler-Nordheim
I-V characteristic**

# EEPROM Cell



BL

WL

$V_{DD}$

**Absolute threshold control is hard**
**Unprogrammed transistor might be depletion**
⇨ **2 transistor cell**

# Flash EEPROM

Control gate

Floating gate

*erasure*

Thin tunneling oxide

$n^1$ source

$n^1$ drain

*programming*

*p*-substrate

**Many other options …**

# Cross-sections of NVM cells



**Flash**

**EPROM**

*Courtesy Intel*

# Basic Operations in a NOR Flash Memory – Erase

cell

array

# Basic Operations in a NOR Flash Memory – Write

# Basic Operations in a NOR Flash Memory– Read

# NAND Flash Memory

Word line(poly)

**Unit Cell**

Source line
(Diff. Layer)

Gate

ONO

Gate
Oxide

FG

# NAND Flash Memory

Select transistor          Word lines

Active area

STI

Bit line contact          Source line contact

# *Characteristics of State-of-the-art NVM*

**Table 12-1** Comparison between nonvolatile memories ([Itoh01]).
$V_{DD}$ = 3.3 or 5 V; $V_{PP}$ = 12 or 12.5 V.

| | Cell— Nr. of Transistors | Cell Area (ratio wrt EPROM) | Mechanism | | External Power Supply | | Program/ Erase Cycles |
|---|---|---|---|---|---|---|---|
| | | | Erase | Write | Write | Read | |
| MASK ROM | 1 T (NAND) | 0.35–5 | — | — | — | $V_{DD}$ | 0 |
| EPROM | 1 T | 1 | UV Exposure | Hot electrons | $V_{PP}$ | $V_{DD}$ | ~100 |
| EEPROM | 2 T | 3–5 | FN Tunneling | FN Tunneling | $V_{PP}$ (int) | $V_{DD}$ | $10^4$–$10^5$ |
| Flash Memory | 1 T | 1–2 | FN Tunneling | Hot electrons | $V_{PP}$ | $V_{DD}$ | $10^4$–$10^5$ |
| | | | FN Tunneling | FN Tunneling | $V_{PP}$ (int) | $V_{DD}$ | $10^4$–$10^5$ |

# Read-Write Memories (RAM)

❑ **STATIC (SRAM)**

Data stored as long as supply is applied
Large (6 transistors/cell)
Fast
Differential

❑ **DYNAMIC (DRAM)**

Periodic refresh required
Small (1-3 transistors/cell)
Slower
Single Ended

# 6-transistor CMOS SRAM Cell

# CMOS SRAM Analysis (Read)



$$k_{n, M5}\left((V_{DD} - \Delta V - V_{Tn})V_{DSATn} - \frac{V_{DSATn}^2}{2}\right) = k_{n, M1}\left((V_{DD} - V_{Tn})\Delta V - \frac{\Delta V^2}{2}\right)$$

$$\Delta V = \frac{V_{DSATn} + CR(V_{DD} - V_{Tn}) - \sqrt{V_{DSATn}^2(1 + CR) + CR^2(V_{DD} - V_{Tn})^2}}{CR}$$

# CMOS SRAM Analysis (Read)



$$CR = \frac{W_1/L_1}{W_5/L_5}$$

# CMOS SRAM Analysis (Write)



$$k_{n,M6}\left((V_{DD} - V_{Tn})V_Q - \frac{V_Q^2}{2}\right) = k_{p,M4}\left((V_{DD} - |V_{Tp}|)V_{DSATp} - \frac{V_{DSATp}^2}{2}\right)$$

$$V_Q = V_{DD} - V_{Tn} - \sqrt{(V_{DD} - V_{Tn})^2 - 2\frac{\mu_p}{\mu_n}PR\left((V_{DD} - |V_{Tp}|)V_{DSATp} - \frac{V_{DSATp}^2}{2}\right)},$$

# CMOS SRAM Analysis (Write)

# 6T-SRAM — Layout

# *Resistance-load SRAM Cell*



Static power dissipation -- Want $R_L$ large
Bit lines precharged to $V_{DD}$ to address $t_p$ problem

# SRAM Characteristics

**Table 12-2** Comparison of CMOS SRAM cells used in 1-Mbit memory (from [Takada91])

| | Complementary CMOS | Resistive Load | TFT Cell |
|---|---|---|---|
| Number of transistors | 6 | 4 | 4 (+2 TFT) |
| Cell size | 58.2 $\mu m^2$ (0.7-$\mu$m rule) | 40.8 $\mu m^2$ (0.7-$\mu$m rule) | 41.1 $\mu m^2$ (0.8-$\mu$m rule) |
| Standby current (per cell) | $10^{-15}$ A | $10^{-12}$ A | $10^{-13}$ A |

# 3-Transistor DRAM Cell



**No constraints on device ratios**

**Reads are non-destructive**

**Value stored at node X when writing a "1" = $V_{WWL} - V_{Tn}$**

# 3T-DRAM — Layout

# 1-Transistor DRAM Cell



**Write: C$_S$ is charged or discharged by asserting WL and BL.**
**Read: Charge redistribution takes places between bit line and storage capacitance**

$$\Delta V = V_{BL} - V_{PRE} = V_{BIT} - V_{PRE} \frac{C_S}{C_S + C_{BL}}$$

**Voltage swing is small; typically around 250 mV.**

# DRAM Cell Observations

❏ 1T DRAM requires a sense amplifier for each bit line, due to charge redistribution read-out.

❏ DRAM memory cells are single ended in contrast to SRAM cells.

❏The read-out of the 1T DRAM cell is destructive; read and  refresh operations are necessary for correct operation.

❏ Unlike 3T cell, 1T cell requires presence of an extra capacitance that must be explicitly included in the design.

❏ When writing a "1" into a DRAM cell, a threshold voltage is lost. This charge loss can be circumvented by bootstrapping the word lines to a higher value than $V_{DD}$

# Sense Amp Operation

# 1-T DRAM Cell



Metal word line

SiO$_2$

Poly

Field Oxide

$n^+$    $n^+$

Poly

Inversion layer induced by plate bias

**Cross-section**

Capacitor

M$_1$ word line

Diffused bit line

**Polysilicon gate**

**Polysilicon plate**

**Layout**

**Uses Polysilicon-Diffusion Capacitance**

**Expensive in Area**

# SEM of poly-diffusion capacitor 1T-DRAM

# Advanced 1T DRAM Cells

Cell Plate Si

Capacitor Insulator

Storage Node Poly

2nd Field Oxide

Refilling Poly

Si Substrate

Word line
Insulating Layer

Cell plate

Capacitor dielectric layer

Transfer gate

Isolation

Storage electrode

17KV  38.4KX  325n 0170

**Trench Cell**

**Stacked-capacitor Cell**

# Static CAM Memory Cell

Bit  $\overline{\text{Bit}}$  Bit  $\overline{\text{Bit}}$

Word

**CAM** ••• **CAM**

Word

**CAM** ••• **CAM**

**Wired-NOR Match Line**

Bit  $\overline{\text{Bit}}$

M8  M9  M5

M4

M6  M7

Word

$\overline{S}$  S

int

M3  M2

Match

M1

# CAM in Cache Memory

# *Periphery*

- ❑ **Decoders**
- ❑ **Sense Amplifiers**
- ❑ **Input/Output Buffers**
- ❑ **Control / Timing Circuitry**

# *Row Decoders*

**Collection of $2^M$ complex logic gates
Organized in regular and dense fashion**

**(N)AND Decoder**

$$WL_0 = A_0 A_1 A_2 A_3 A_4 A_5 A_6 A_7 A_8 A_9$$

$$WL_{511} = \overline{A_0} A_1 A_2 A_3 A_4 A_5 A_6 A_7 A_8 A_9$$

**NOR Decoder**

$$WL_0 = \overline{A_0 + A_1 + A_2 + A_3 + A_4 + A_5 + A_6 + A_7 + A_8 + A_9}$$

$$WL_{511} = \overline{A_0 + \overline{A_1} + \overline{A_2} + \overline{A_3} + \overline{A_4} + \overline{A_5} + \overline{A_6} + \overline{A_7} + \overline{A_8} + \overline{A_9}}$$

# Hierarchical Decoders

## Multi-stage implementation improves performance



**NAND decoder using 2-input pre-decoders**

# Dynamic Decoders



**2-input NOR decoder**

**2-input NAND decoder**

# 4-input pass-transistor based column decoder

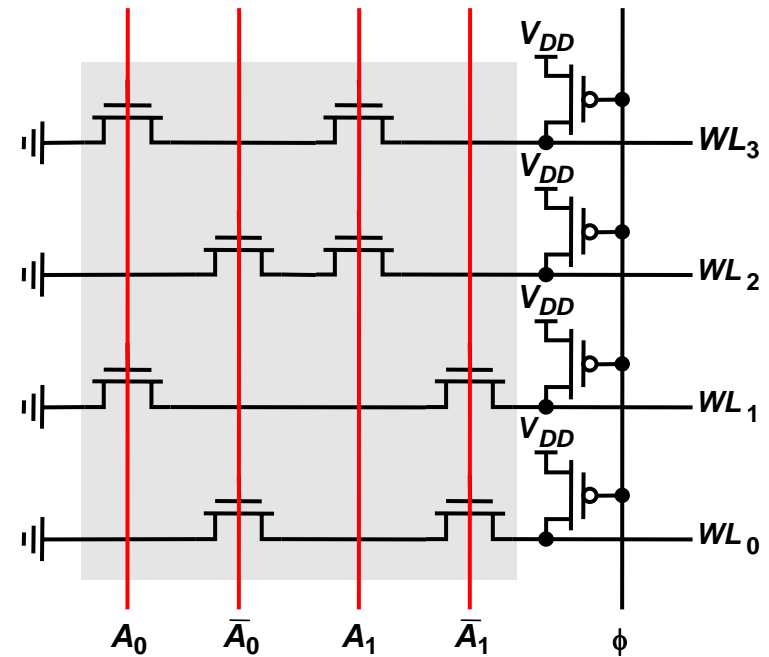$BL_0$  $BL_1$  $BL_2$  $BL_3$

$A_0$

$S_0$
$S_1$
$S_2$
$A_1$
$S_3$

$D$

**Advantages: speed ($t_{pd}$ does not add to overall memory access time)**
    **Only one 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**

# Decoder for circular shift-register

# *Sense Amplifiers*

make $\Delta V$ as small
as possible

$$t_p = \frac{C \cdot \Delta V}{I_{av}}$$

large          small

**Idea: Use Sense Amplifer**

**small
transition**  **s.a.**

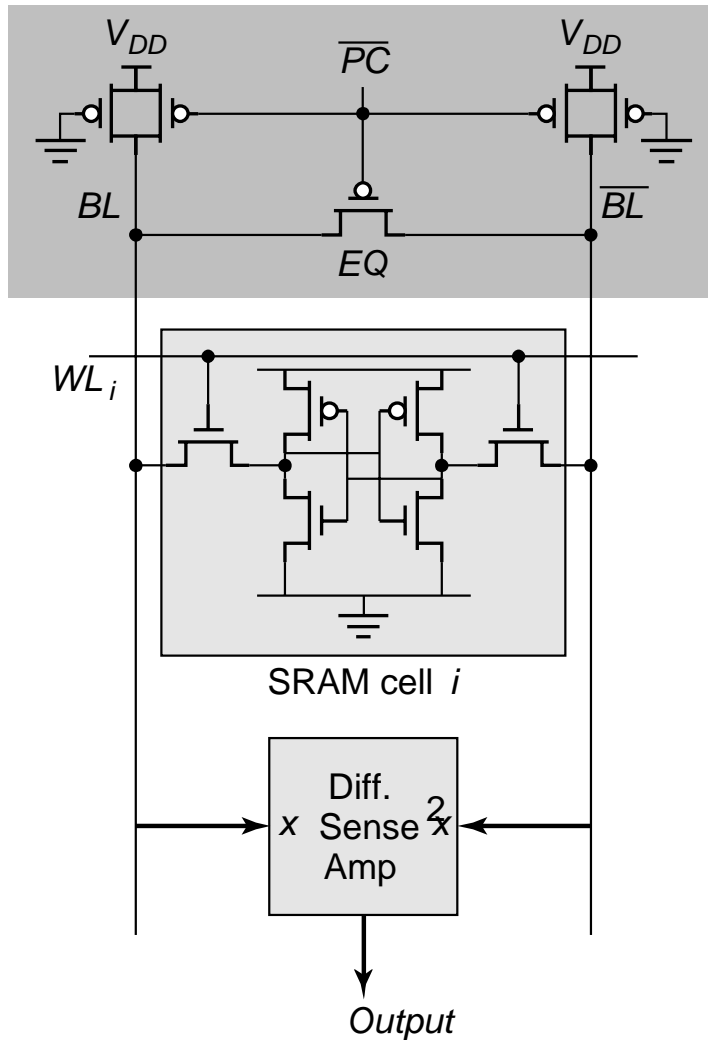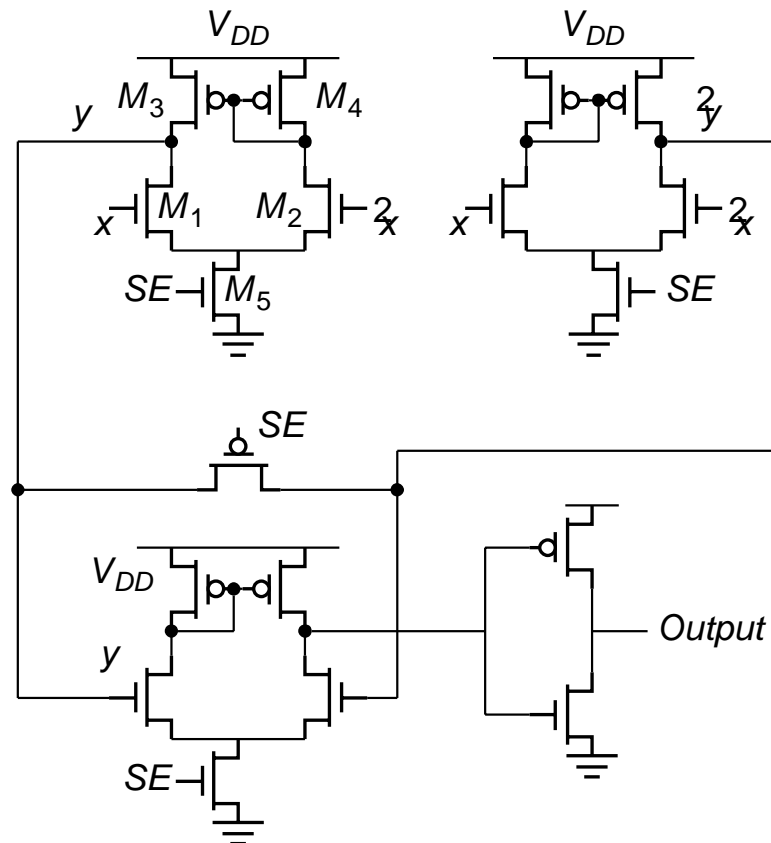input          output

# Differential Sense Amplifier



**Directly applicable to SRAMs**

# Differential Sensing – SRAM



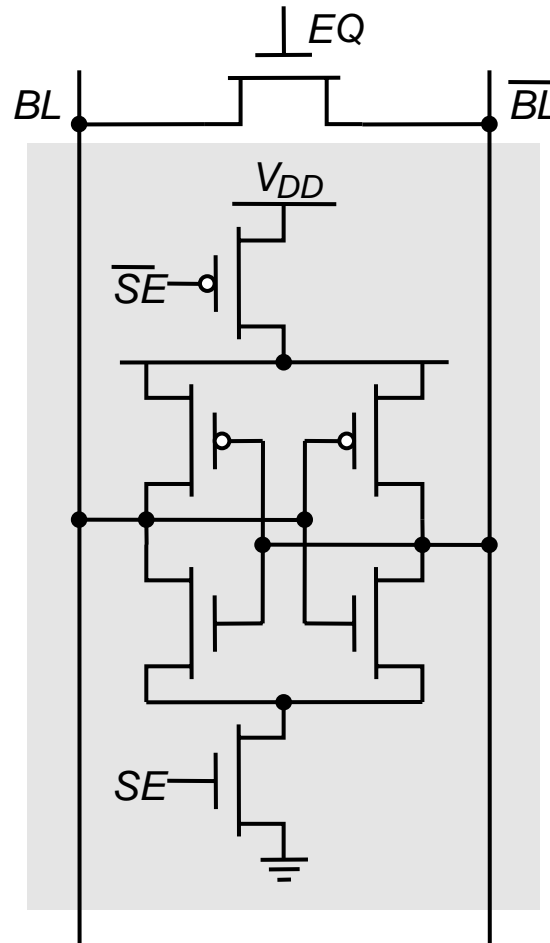(a) SRAM sensing scheme

(b) two stage differential amplifier
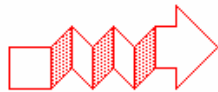
# Latch-Based Sense Amplifier (DRAM)



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.

# *Reliability and Yield*

- Semiconductor memories trade off noise-margin for density and performance

⟹ Highly Sensitive to Noise (Crosstalk, Supply Noise)

- High Density and Large Die size cause Yield Problems

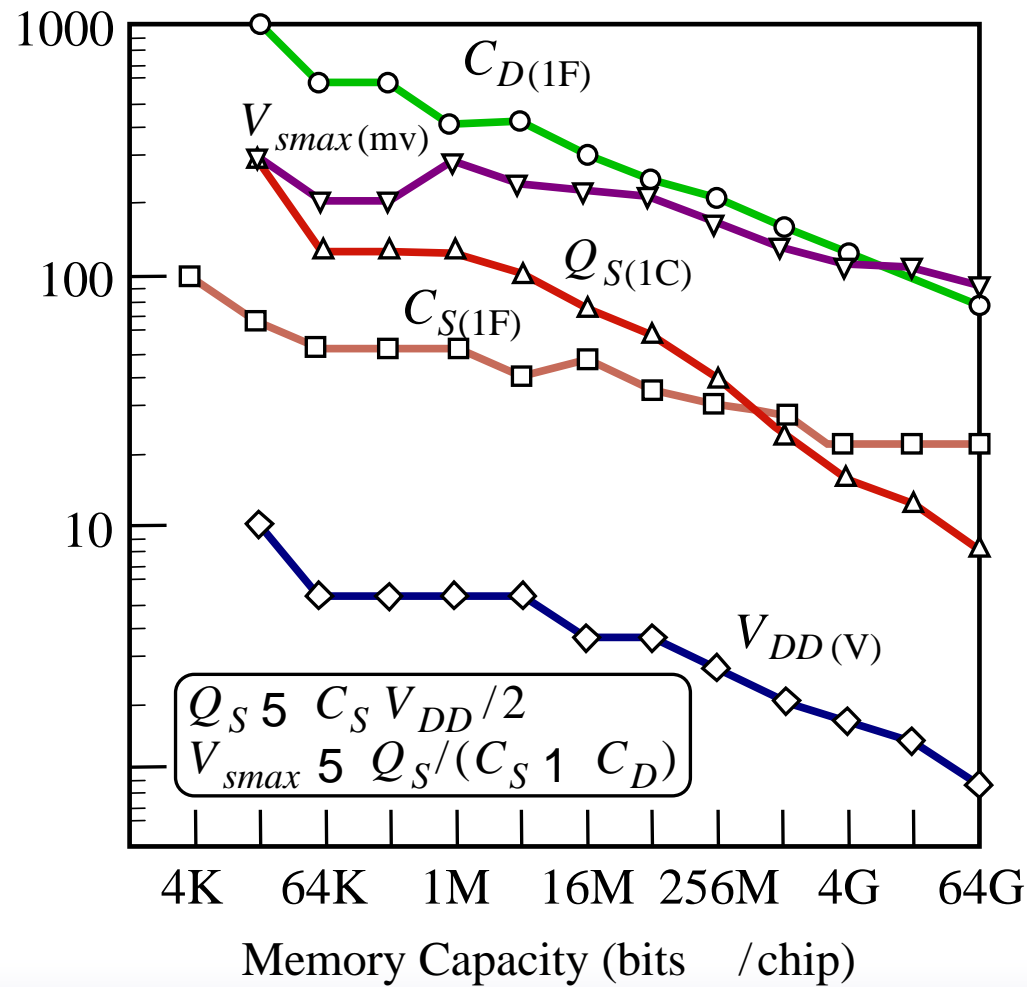$$Y = 100 \frac{Number\ of\ Good\ Chips\ on\ Wafer}{Number\ of\ Chips\ on\ Wafer}$$

$$Y = \left[ \frac{1 - e^{-AD}}{AD} \right]^2$$

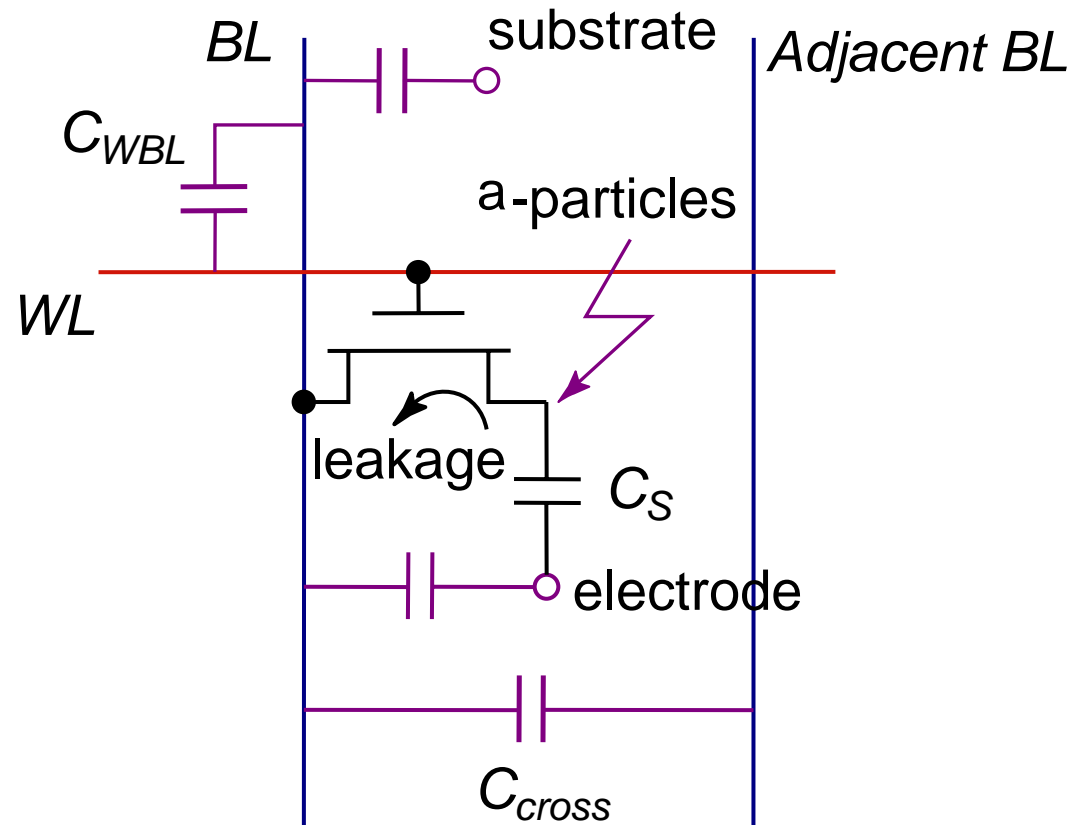**Increase Yield using Error Correction and Redundancy**
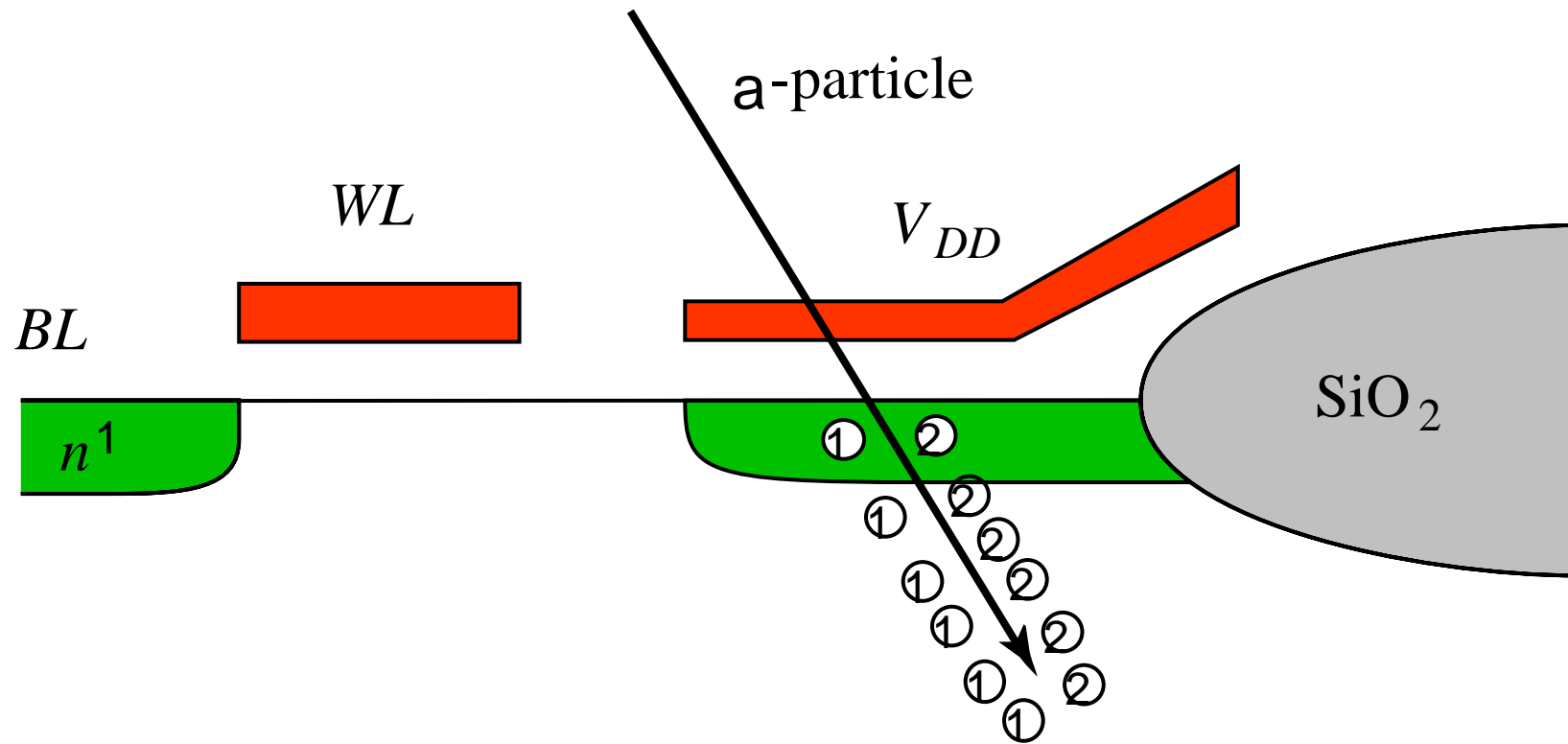
# Sensing Parameters in DRAM

# *Noise Sources in 1T DRam*

# Alpha-particles (or Neutrons)



**1 Particle ~ 1 Million Carriers**

# Yield



Yield curves at different stages of process maturity
(from [Veendrick92])

# Redundancy

Redundant
rows

Redundant
columns

Row
Address

Memory
Array

Column Decoder

Column
Address

. Fuse
. Bank

# Error-Correcting Codes

## Example: Hamming Codes

$$P_1 \, P_2 \, B_3 \, P_4 \, B_5 \, B_6 \, B_7$$

**with**

$$P_1 \oplus B_3 \oplus B_5 \oplus B_7 = 0$$

$$P_2 \oplus B_3 \oplus B_6 \oplus B_7 = 0$$

$$P_4 \oplus B_5 \oplus B_6 \oplus B_7 = 0$$

*e.g. B3 Wrong*

1

1    *= 3*

0

# *Redundancy and Error Correction*

# *Case Studies*

❑ Programmable Logic Array

❑ SRAM

❑ Flash Memory

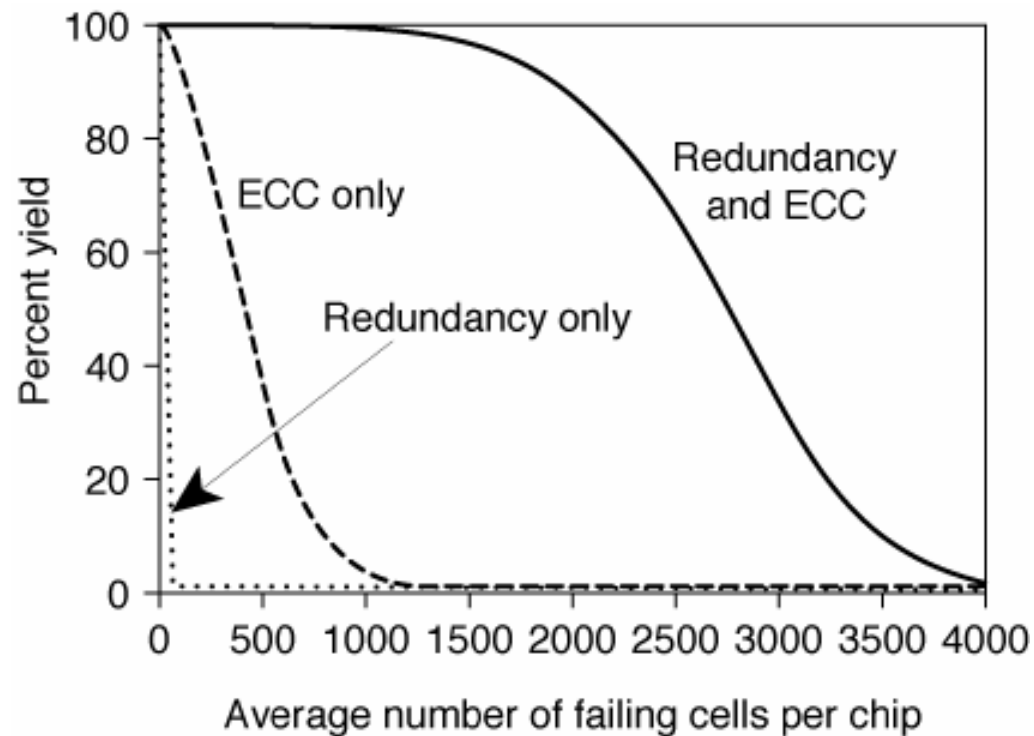# PLA versus ROM

❑ **Programmable Logic Array**

structured approach to random logic
"two level logic implementation"
NOR-NOR (product of sums)
NAND-NAND (sum of products)

IDENTICAL TO ROM!

❑ **Main difference**

ROM: fully populated
PLA: one element per minterm
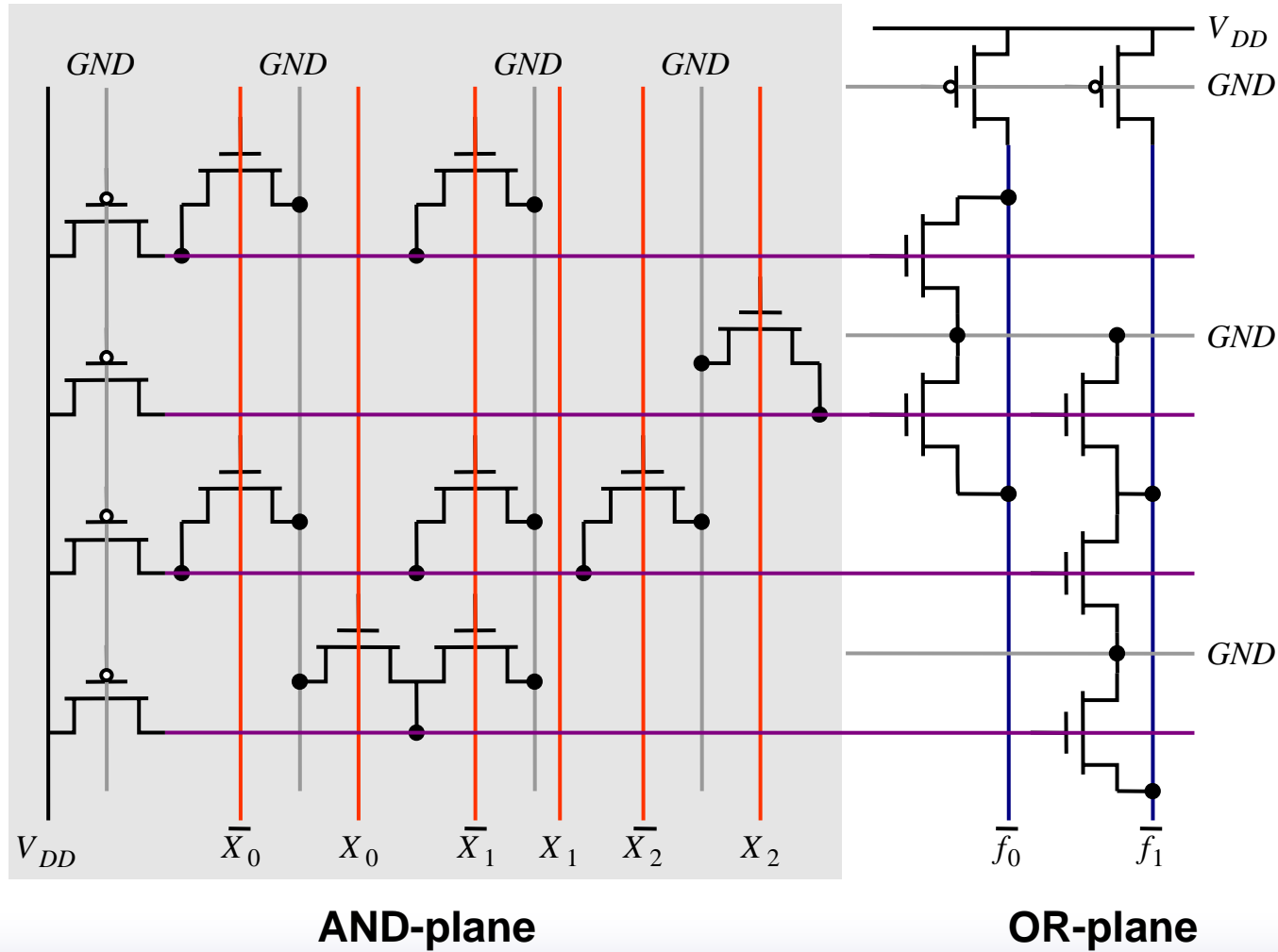
Note: Importance of PLA's has drastically reduced
1. slow
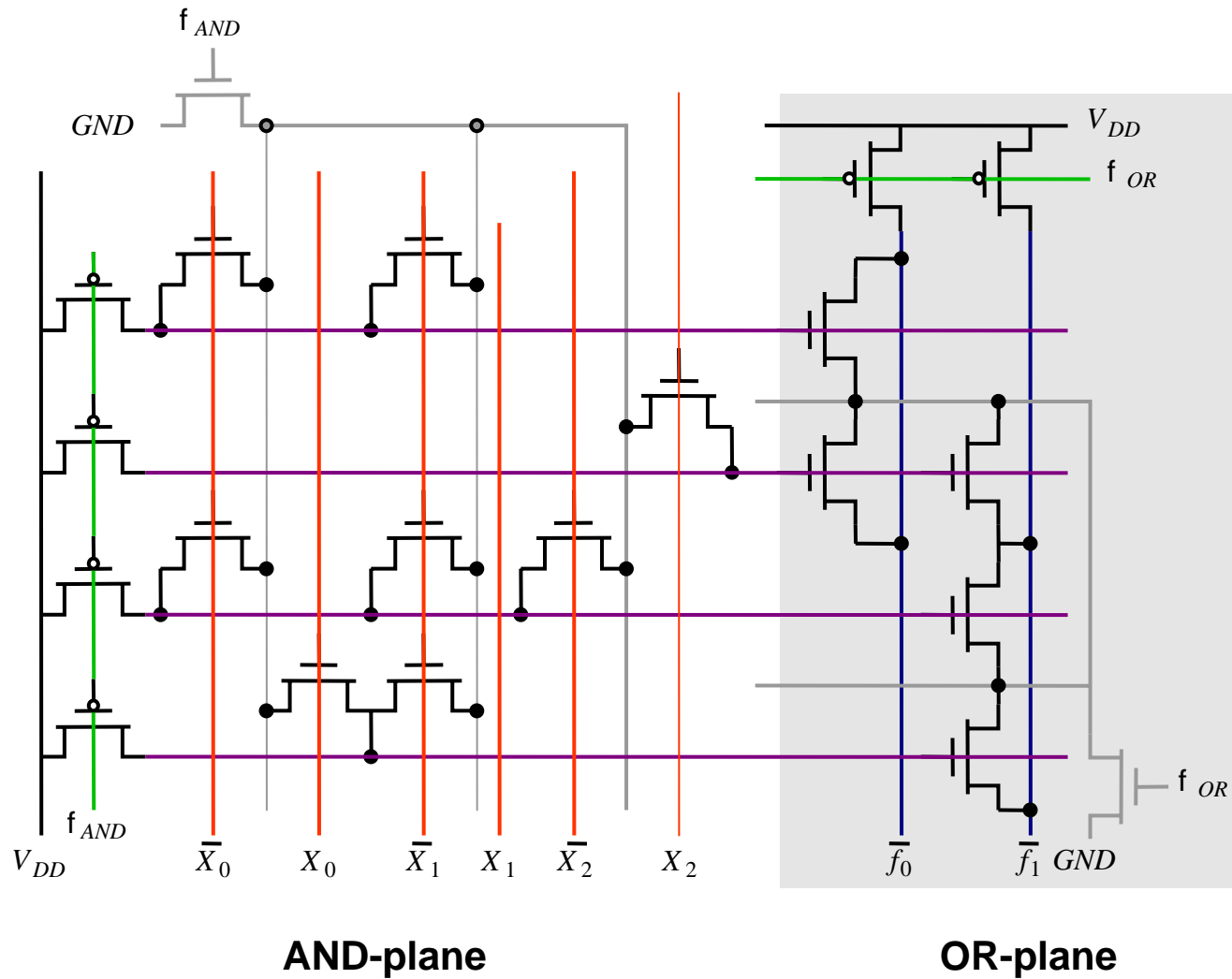2. better software techniques (mutli-level logic synthesis)

*But …*

# Programmable Logic Array

## Pseudo-NMOS PLA



**AND-plane**

**OR-plane**

# *Dynamic PLA*



**AND-plane**                    **OR-plane**

# Clock Signal Generation
# for self-timed dynamic PLA



f

f AND

$t_{pre}$  $t_{eval}$

f OR

(a) Clock signals

f

Dummy AND row

f AND

f AND

Dummy AND row

f OR

(b) Timing generation circuitry

# PLA Layout



$V_{DD}$       *And-Plane*       *Or-Plane* $\phi$    GND

$x_0$   $\bar{x}_0$   $x_1$   $\bar{x}_1$   $x_2$   $\bar{x}_2$      $\bar{f}_0$   $\bar{f}_1$

*Pull-up devices*      *Pull-up devices*

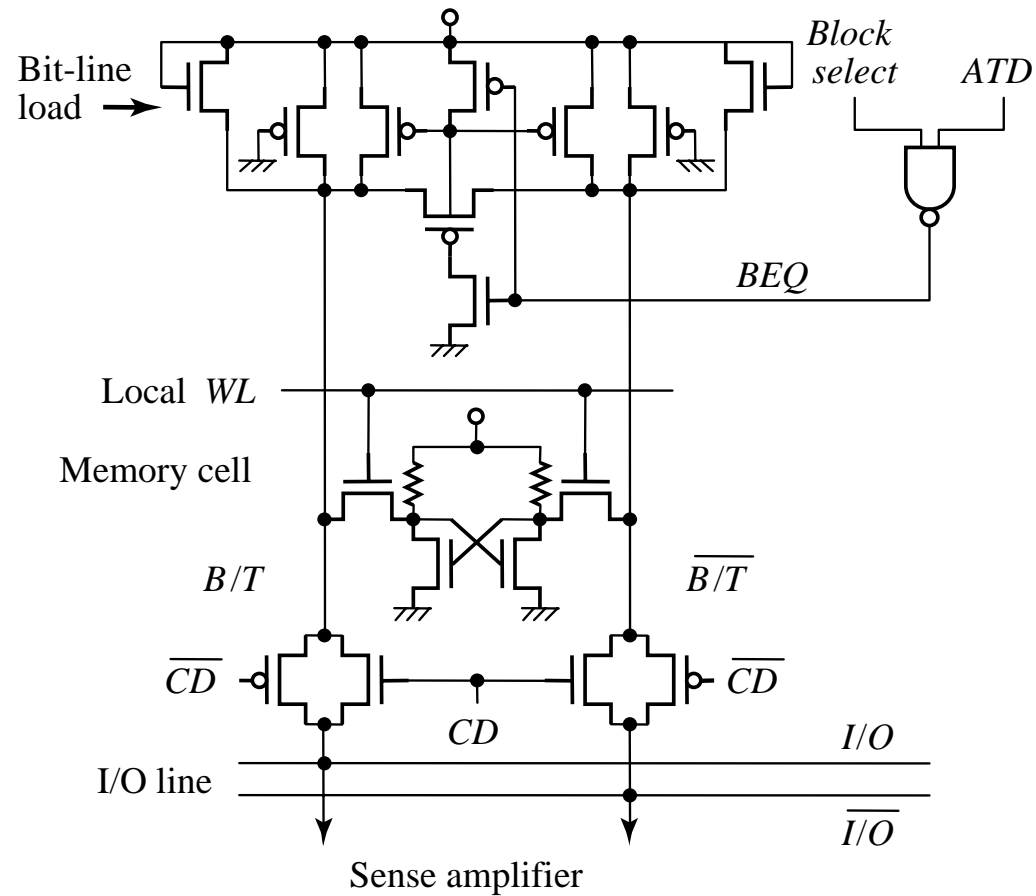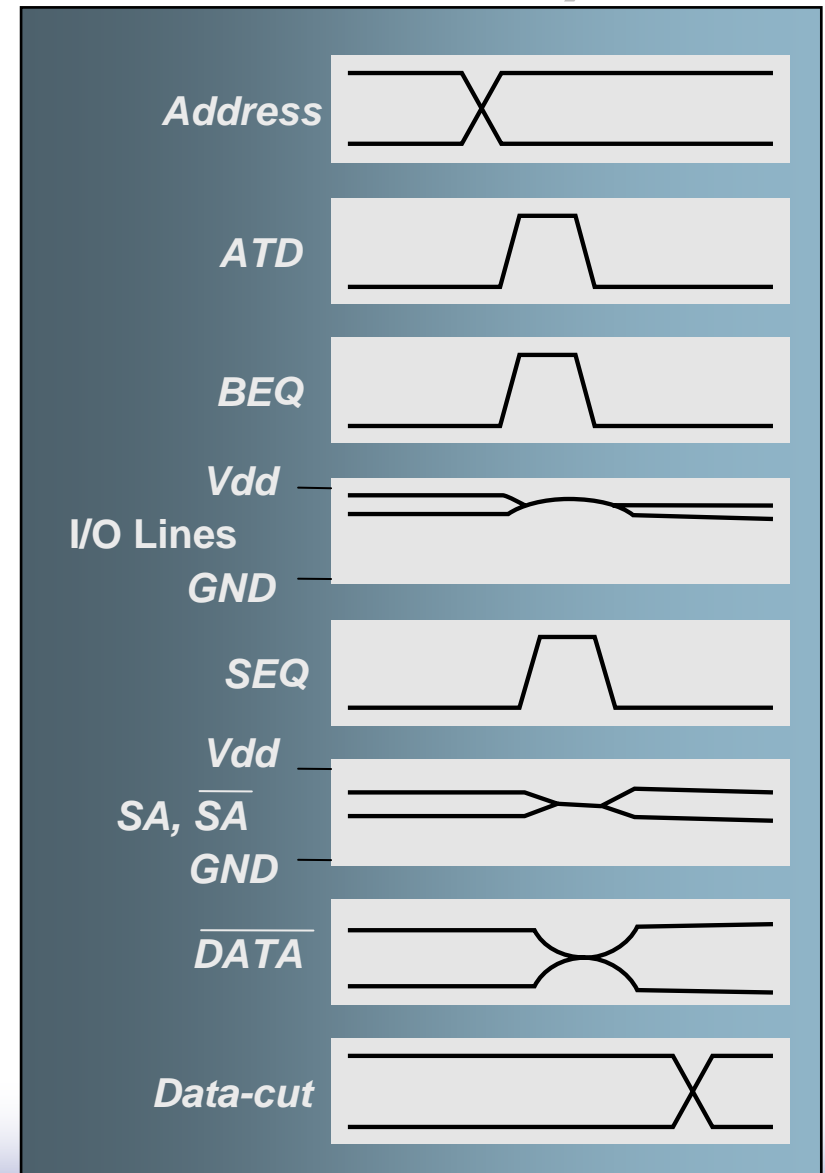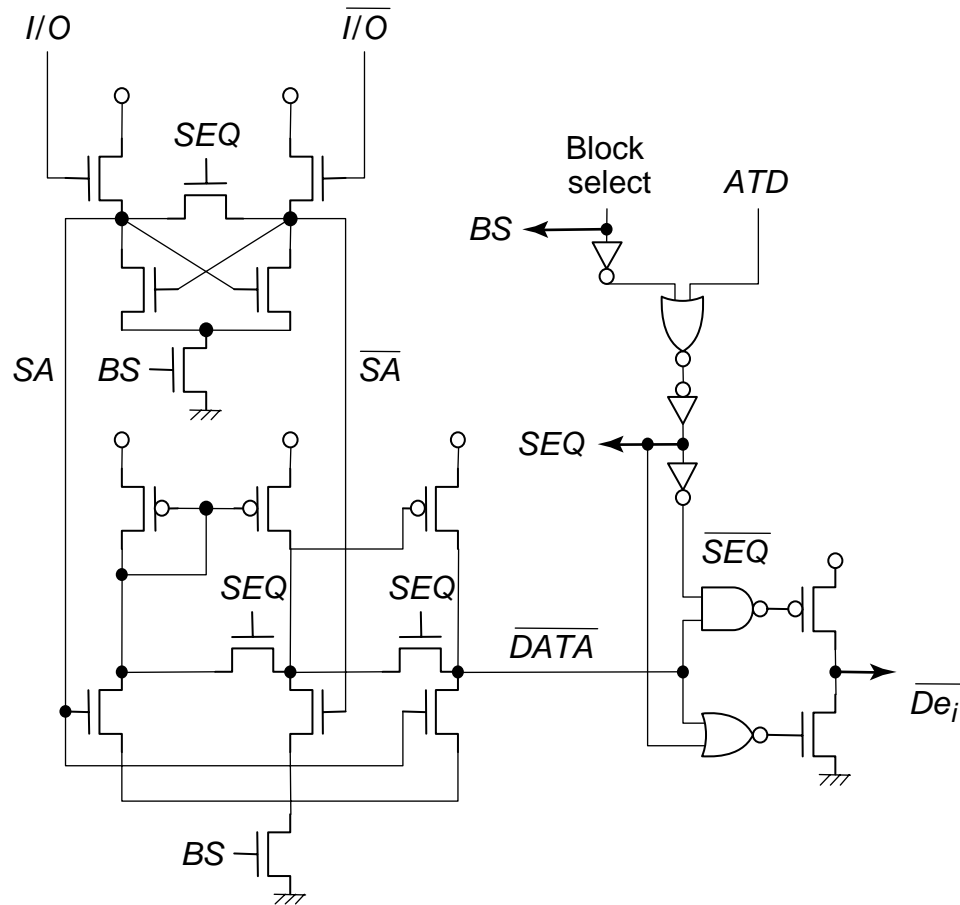# 4 Mbit SRAM
# Hierarchical Word-line Architecture

# Bit-line Circuitry
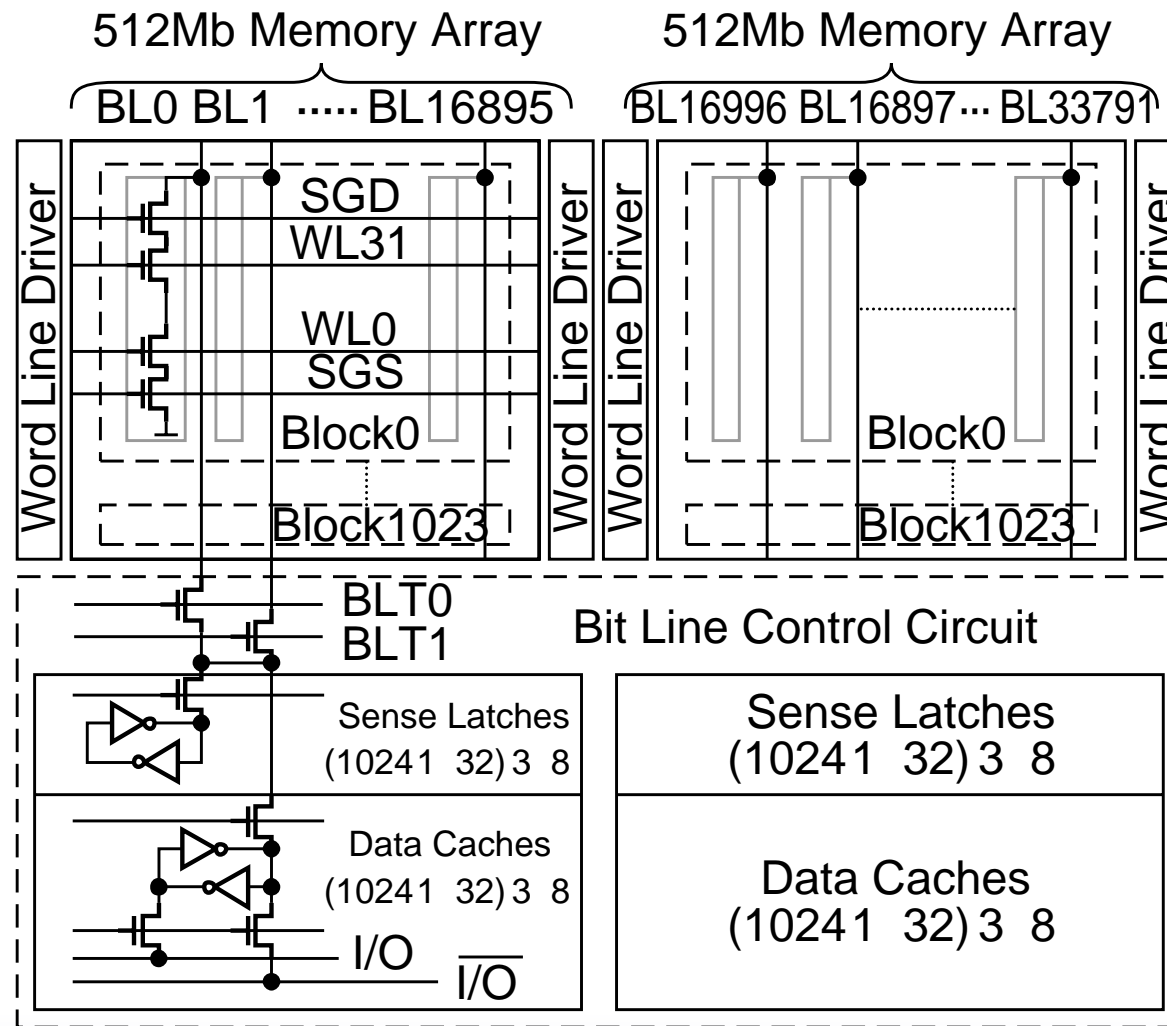


Sense amplifier

# Sense Amplifier (and Waveforms)

# 1 Gbit Flash Memory



512Mb Memory Array     512Mb Memory Array

BL0 BL1 ····· BL16895     BL16996 BL16897 ··· BL33791

Word Line Driver

SGD
WL31
WL0
SGS
Block0
Block1023

Word Line Driver

Word Line Driver

Block0
Block1023

Word Line Driver

BLT0
BLT1

Bit Line Control Circuit

Sense Latches
(1024 1 32) 3 8

Data Caches
(1024 1 32) 3 8

I/O
I/O

Sense Latches
(1024 1 32) 3 8

Data Caches
(1024 1 32) 3 8

**From [Nakamura02]**     Memories

# Writing Flash Memory

Verify level 5  0.8 V   Word-line level 5  4.5 V

Number of memory cells

Result of 4 times program

0V  1V  2V  3V  4V

Vt of memory cells

**Evolution of thresholds**

$10^8$

$10^6$

$10^4$

$10^2$

$10^0$

0V  1V  2V  3V  4V

Vt of memory cells

**Final Distribution**

# 125mm² 1Gbit NAND Flash Memory



32 word lines
x 1024 blocks

16896 bit lines

Charge pump

2kB Page buffer & cache

10.7mm
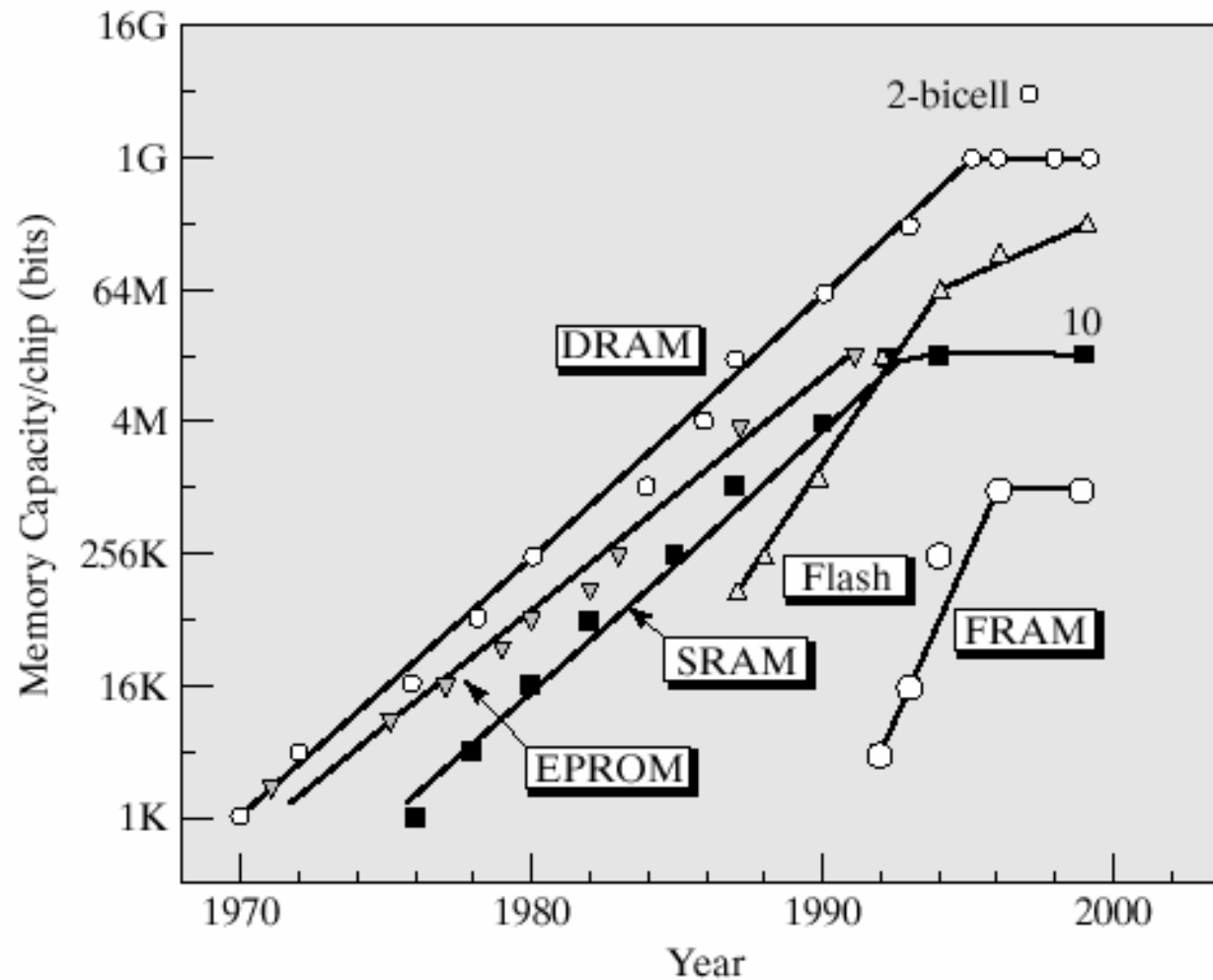
11.7mm

# 125mm² 1Gbit NAND Flash Memory

- ❑ **Technology** **0.13μm p-sub CMOS triple-well**
  **1poly, 1polycide, 1W, 2Al**
- ❑ **Cell size** **0.077μm2**
- ❑ **Chip size** **125.2mm2**
- ❑ **Organization** **2112 x 8b x 64 page x 1k block**
- ❑ **Power supply** **2.7V-3.6V**
- ❑ **Cycle time** **50ns**
- ❑ **Read time** **25μs**
- ❑ **Program time** **200μs / page**
- ❑ **Erase time** **2ms / block**
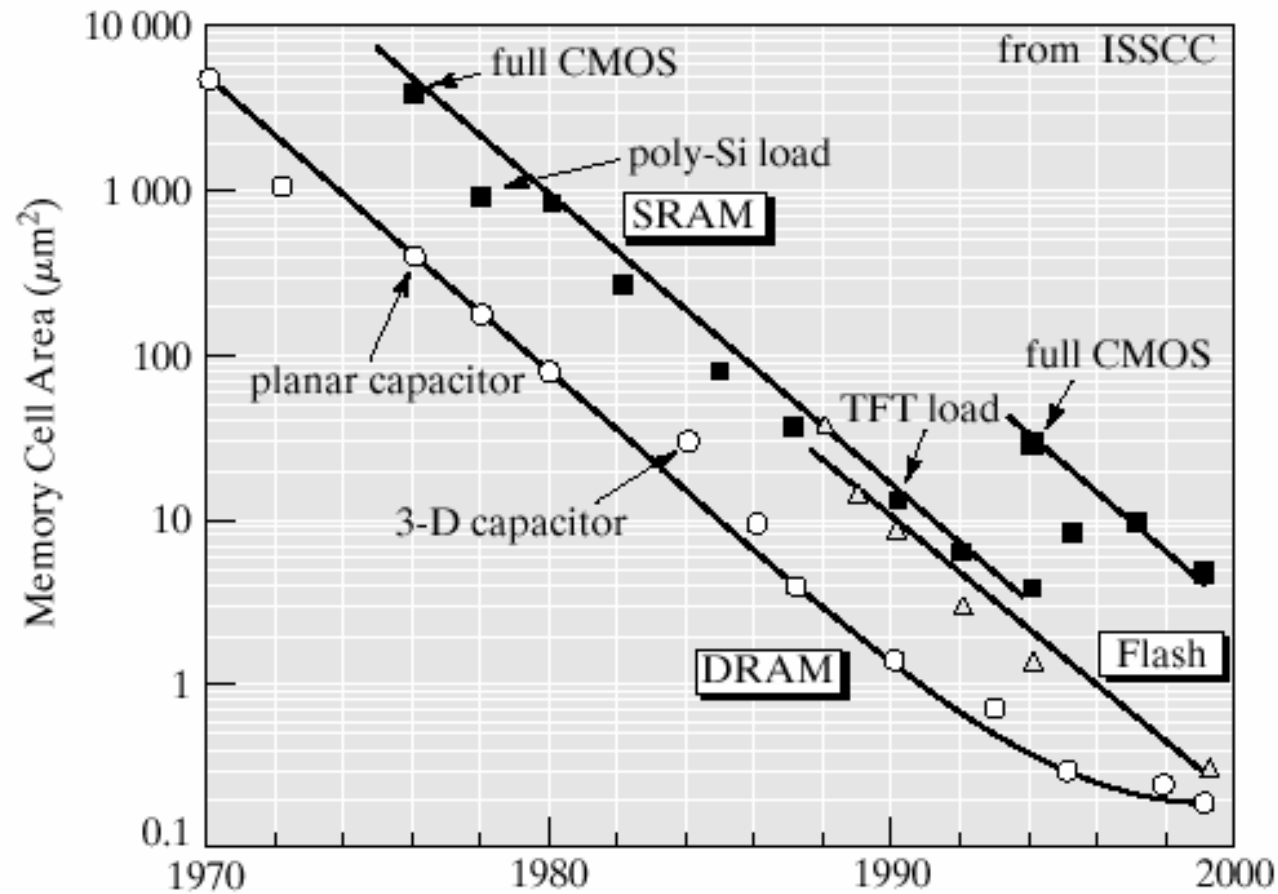
# Semiconductor Memory Trends
# (up to the 90's)



Memory Size as a function of time: x 4 every three years

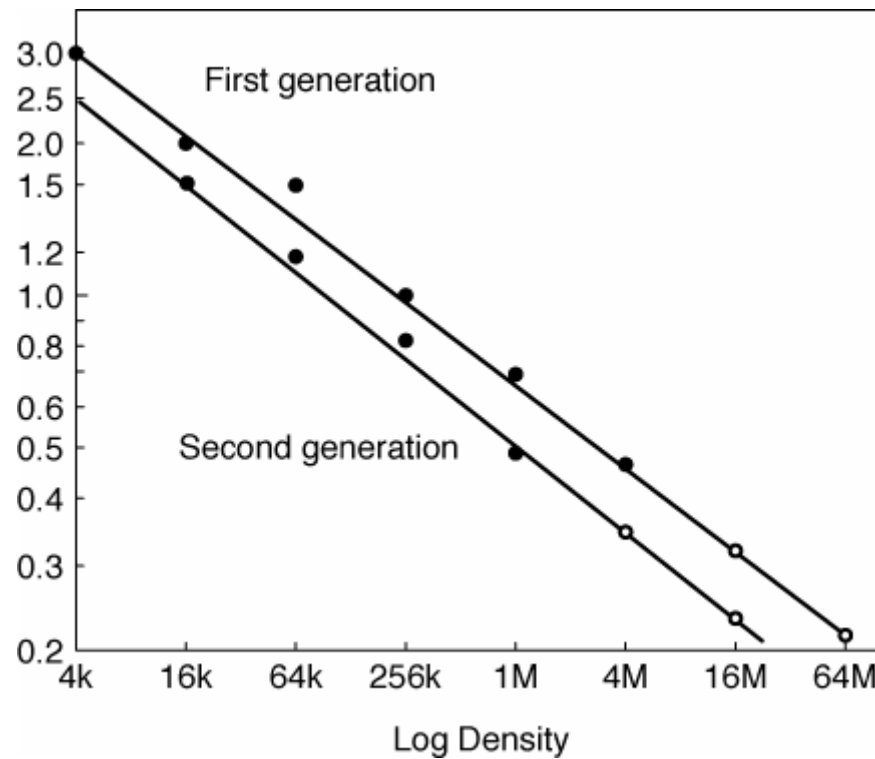# Semiconductor Memory Trends (updated)

# Trends in Memory Cell Area

# Semiconductor Memory Trends



Technology feature size for different SRAM generations