**Memory**

Can be categorized into:

- **Read Write Memory** (RWM)
  - Random Access Memory (RAM): static SRAM (faster) verses dynamic DRAM (smaller) structures possible. Access time independent of physical location of data.
  - Non-RAM: Serial Access Memory (FIFO, LIFO, Shift register) and Content Access Memory (CAM). Non-uniform access time.
- **Non-volatile Read Write Memory** (NVRWM): write time much larger than read time.

  - EPROM, E$^2$PROM, FLASH
- **Read Only Memory** (ROM)

A second classification for RAMs and ROMs:
- Static-load: no clock required.
- Synchronous: require a clock edge to enable memory operation.
- Asynchronous: recognize address changes and output new data. More difficult to build.

## Memory Architecture

In order to build an *N-word* memory where each word is *M bits* wide (typically 1, 4 or 8 bits), a straightforward approach is to stack memory:

A word is selected by setting exactly **one** of the select bits, $S_x$, high.

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

Storage cell

N words

$S_{N-2}$ → Word N-2
$S_{N-1}$ → Word N-1

Input-Output
(M bits)

This approach works well for small memories but has problems for large memories.

For example, to build a 1Mword (where word = 8 bits) memory, requires 1M select lines, provided by some off-chip device.

This approach is not practical.
What can we do?

## Memory Architecture

Add a decoder to solve the package problem:



$K = \log_2 N$

one-hot    Input-Output (M bits)

Storage cell

This reduces the number of external address pins from 1M to 20.

This does not address the **memory aspect ratio** problem:

The memory is 128,000 time higher than wide ($2^{20}/2^3$) !

Besides the bizarre shape factor, the design is *extremely slow* since the vertical wires are VERY long (delay is at least linear to length).

**Memory Architecture**

The vertical and horizontal dimensions are usually very similar, for an aspect ratio of *unity*.

Multiple words are stored in each row and selected simultaneously:

Row address =
$A_K$ to $A_{L-1}$

$A_K$

$A_{K+1}$

$A_{K+2}$

$A_{L-1}$

**Row Decoder**

$S_0$
$S_1$
$S_2$

$S_{N-2}$
$S_{N-1}$

Bit line

Storage cell

Word line

Column address =
$A_0$ to $A_{K-1}$

$A_0$

$A_{K-1}$

**Column decoder**

Sense amps
and drivers
not shown

Input-Output
(M bits)

A column decoder is added to select the desired **word** from a row.

**Memory Architecture**

This strategy works well for memories up to 64 Kbits to 256 Kbits.

Larger memories start to suffer excess delay along bit and word lines.

A **third dimension** is added to the address space to solve this problem:



Global Data bus

Address: [**Row**][**Block**][**Col**]

Global amplifier/driver

I/O

4 Mbit: P = 32 blocks with 128Kbits/block.

128Kbit block: 1024 rows and 128 columns.

**Memory: Architecture**

An example:

$2^{m+k}$ bits

$A_k$

Row decoder

$A_k+1$

Row decoder

Row decoder

$A_{n-1}$

Row decoder

$2^{n-k}$ bits

$A_0$

Column decoder

column mux, sense amp, write buffers

$A_{k-1}$

$[A_{n-1}..A_k][A_{k-1}..A_0]$

For example: Let N = 1,048,576 and M = 8 bits for a 1 million byte memory.
n = $\log_2 N$ = 20, k = 8 and m = $\log_2 M$ = 3.

Then there are $2^{n-k}$ rows = $2^{12}$ = 4096 and $2^{k+m}$ columns/$2^3$ bits per word = $2^8$
= 256 words.

# ROM

ROM cells are permanently fixed: Several possibilities:



Diode supplies current to raise BL (bitline) for all cells on the row.

BJT supplies current to raise BL for each cell on the row. Requires $V_{DD}$ to be routed.

psuedo n-MOS NOR gate.

Resistance of n/p should be at least 4.

p-MOS used to hold BL high. n-MOS provides pull-down path.

## Non-volatile Read-Write Memories

Virtually identical in structure to ROMs.

Selective enabling/disabling of transistors is accomplished through modifications to **threshold voltage**. This is accomplished through a floating gate.



Applying a high voltage (15 to 20 V) between source and gate-drain create high electric field and causes avalanche injection to occur.

Hot electrons traverse first oxide and get trapped on floating gate, leaving it negatively charged.

This increases the threshold voltage to ~7V. Applying 5V to the gate does not permit the device to turn on.

**Non-volatile Read-Write Memories**

The method of erasing is the main differentiating factor between the various
classes of reprogrammable nonvolatile memories.

- **EPROM**:

    UV light renders oxide slightly conductive.

    Erase is slow (seconds to several minutes).

    Programming is slow (5-10 microsecs per word).

    Limited number of programming cycles - about 1000.

    Very dense - single transistor functions as both the programming and
    access device.

**Non-volatile Read-Write Memories**

- **EEPROM** or **E$^2$PROM**:

    Very thin oxide allows electrons to fl ow to and from the gate via Fowler-
    Nordheim tunneling with $V_{GD}$ applied.

    Erasure is achieved by reversing the voltage applied during writing.

Source                                    Drain

Gate

Floating Gate                    $t_{ox}$      10V

$-$ $-$ $-$ $-$ $-$ $-$ $-$ $-$

$t_{ox}$

n+          thin tunneling ox          n+

Substrate

WL

BL

Threshold control becomes a problem:

Removing too much charge results in a
depletion device that cannot be turned off.

$V_{DD}$

Remedy: Add an access transistor.

**Non-volatile Read-Write Memories**
  • **Flash EEPROM**:
     Combines density adv. of EPROM with versatility of EEPROM.
     Uses avalanche hot-electron-injection approach to program.
     Erasure performed using Fowler-Nordheim tunneling.
     Monitoring control hardware checks the value of the threshold during
       erasure - making sure the unprogrammed transistor remains an
       enhancement device.

Source                               12V    Drain

Gate

Floating Gate                    $t_{ox}$

12V                                                    12V

erasure  ⊖ ⊖ ⊖ ⊖ ⊖ ⊖ ⊖ ⊖

thin tunneling ox

12V

n+          ⊖    ⊖                         n+
                    programming

Substrate

Programming performed by applying 12V to gate and drain.
Erasure performed with gate grounded and source at 12V.

**Read-Write Memories (RAM)**

SRAM:



word line

$\overline{bit}$                                                                bit

$V_{DD}$

**Read-Write Memories (RAM)**

Generic RAM circuit:



Bit Line Conditioning

clocks

bit

$\overline{\text{bit}}$

RAM cell

row decoder

n-1;k

word line

column decoder

k-1;0

Sense Amp
Column Mux
Write Buffers

Address

read-data    write-data

**Read-Write Memories (RAM)**

   SRAM: Read Operation

   Precharging bit and bit_bar to 5V before enabling the word line
   improves performance.

To optimize speed,
use n-channels as
precharge devices.

$\overline{bit}$         precharge         bit

word

data

precharge

$V_{DD}$ —         bit, $\overline{bit}$

word

data

**Read-Write Memories (RAM)**

SRAM: Write Operation:



$N_5$  $N_6$

0

$\overline{cell}$  $N_3$  $P_{bit}$  cell  $N_4$

$1{\rightarrow}0$  $0{\rightarrow}1$

word

$\overline{bit}$  bit

write  $N_1$  $N_2$

write-data

1  $N_d$

write-data

write

word

bit, $\overline{bit}$

cell, $\overline{cell}$

Zero stored in cell originally.
$N_d$, $N_1$, and $N_3$ have to pull $P_{bit}$ below
the inverter threshold.

# Read-Write Memories (RAM)

Register files:

Single-write-port, double-read-port

Overpowers weak feedback inverter

2/2

2/3

4/1

4/1

4/1

4/1

4/1

4/1

2/1

8/1  Biased toward $V_{SS}$ to help write.

Adv: No matter what the load, cell cannot be fl ipped.

decode addr<3:0>

write-data

read-data0

read-data1

**Read-Write Memories (RAM)**

   DRAM:

      Refresh: Compensate for charge loss by periodically rewriting the cell
       contents.

      Read followed by a write operation.

      Typical refresh cycles occur every 1 to 4 milliseconds.

   4 transistor DRAM created by simply eliminating the p tree in an SRAM cell.



word line

bit                                                    $\overline{bit}$

      Logic 1 values are, of course, a threshold below $V_{DD}$.

## Read-Write Memories (RAM)

3T DRAM:



Note that this cell is inverting

bit2 is either clamped to $V_{DD}$ or is precharged to either $V_{DD}$ or $V_{DD}-V_T$.

No device ratioing necessary here !

Most common method of refresh is to read *bit2*, place its inverse on *bit1* and assert *write*.

Precharge method of 'setting' *bit2* is preferred (no steady-state current).

Memory structure of choice in ASICs because of its relative simplicity in both design and operation.

**Read-Write Memories (RAM)**

   1T DRAM



$$\Delta V = V_{bit} - V_{pre} = (V_x - V_{pre})\frac{C_x}{(C_x + C_{bit})}$$

   During read operation, charge redistribution occurs between node X and
      node bit.

   $C_x$ is typically 1 or 2 orders of magnitude smaller than $C_{bit}$ so the delta-V
      value is typically 250 mV.

   Most pervasive DRAM cell in commercial memory design.

**Read-Write Memories (RAM)**

    1T DRAM observations:

- Amplification of delta-V (through a sense amplifier) is necessary in order for the cell to be functional.

      Other cell designs used sense amps only to speed up the read operation.

- The read-out operation is destructive ! Output of sense amp is imposed onto the bit line with word-line high during read-out.

Sense amp activated

$V_{pre}$

V(1)

V(0)

Word-line activated

- 1T transistor requires an explicit capacitor (3T used gate capacitance).

    Capacitance must be large (~30fF) but area small - key challenge in design.

- Bootstrapping word-line to a value larger than $V_{DD}$ circumvents $V_T$ loss on storage capacitor.

**Read-Write Memories (RAM)**

Content Access Memory (CAM):

Determines if a match exists between a data word with a stored word.

Used in Translation-look-aside buffers.



$\overline{bit}$

$V_{DD}$

bit

Match is 0 if ANY SRAM cell has $bit / \overline{cell}$ or $\overline{bit} / cell$ equal to 1.

word line

SRAM with extra n-channels to implement **XOR** function.

cell          $\overline{cell}$

Each bit of the word is tied to the match line.

match

Dynamic or Pseudo n-MOS implementations possible.