# EEC 118 Lecture #13: Memories

Rajeevan Amirtharajah University of California, Davis

> Jeff Parkhurst Intel Corporation

# Outline

- Review: Dynamic MOS Logic Circuits: Rabaey 6.3 (Kang & Leblebici, 9.4-9.6)
- Memories: Rabaey 12.1-12.2 (Kang & Leblebici, 10.1-10.6)

# **Review: Dynamic CMOS**

• Operation

- Clk low during Pre-charge
  - Mp is on while Mn is off
  - Output charged to Vdd
- Clk high during evaluate
  - Mn is on while Mp is off
  - Output pulled down according to PDN function

#### • PDN design same as static CMOS



# **Review: Dynamic CMOS Tradeoffs**

#### • Advantages:

- Faster why?
  - Reduced input load
  - No switching contention
- Less layout area



#### • Disadvantages:

- Multiple stage issues
- Charge leakage
- Charge sharing
- Capacitive coupling
- Cannot be cascaded
- Complicated timing/clocking
- Higher power
- Lower noise margins
- Does not scale well with process

# Memory and Performance Trend I

- Memory is becoming a key factor in the performance of a computer
- 1<sup>st</sup> generation computers had just system memory
  - Very slow DRAM
  - As microprocessors got faster, the bottleneck in performance was data access
- 2<sup>nd</sup> generation computers added cache memory
  - This provided faster access to small localized memory that was being read or written
    - Memory placed on front side bus (off-chip)
    - Performance increased
  - As processors got faster memory access again became the bottleneck

### Memory and Performance Trend II

- 3<sup>rd</sup> generation computers added on chip cache
  - These caches started out 16K and were termed level 1 cache
  - Soon we had level 1 and level 2 cache
  - Currently we have three levels of cache on chip that run at processor frequency, then access main memory if data can't be found in any of these caches
  - Moving to stacking memory die on top of processor

- ROM: read-only memory
  - Non-volatile mask programmed
- RWM: read-write memory (RAM, random access memory)
  - SRAM: static memory
    - Data is stored as the state of a bistable circuit
    - State is retained without refresh as long as power is supplied
  - DRAM: dynamic memory
    - Data is stored as a charge on a capacitor
    - State leaks away, refresh is required

# **Types of Memory II**

- NVRWM: non-volatile read-write memory (also called NVRAM, non-volatile random access memory)
  - Flash (EEPROM): ROM at low voltages, writable at high voltages (Electrically Erasable Programmable Read-Only Memory)
  - EPROM: ROM, but erasable with UV light (falling out of common usage)

# Memory Usage in Computers

- DRAM Memory
  - Main memory storage. Used for data and programs
- SRAM Memory
  - Faster than DRAM, however, uses more transistors
    - Used to be used for external cache
    - Variant used in internal cache (on chip cache)
- FLASH Memory and ROM
  - Used for BIOS data storage in PCs
  - Also used to store pictures, MP3 files for digital cameras and MP3 players – eventually for hard disk

### **Basic Memory Array Structure**

 Memory cells () arranged in a rectangular array Rows correspond  $A_0$ to data words Accessed through Memory Array a row decoder 2<sup>N</sup> bits Columns to  $A_{k-1}$ individual bits Selected through a column mux  $2^{k} - 1_{0}$  $2^{N-k}$  - Bit voltage  $A_k \dots A_{N-1}$ amplified by sense (N-k)amplifier

# **Memory Circuit Operation**

- <u>Wordlines</u> (WL) control row (word) access
  - Usually control gates of pass transistors
- <u>Bitlines</u> (BL) route column data (individual bits)
  - Bitlines usually precharged high (like dynamic logic)
  - Memory cells discharge bitline depending on stored data (bitline left high if cell stores 1, bitline discharged if cell stores 0)
  - Bitline swings usually small (10s 100s of mV) and must be amplified by sense amplifiers
  - Synchronous or asynchronous timing can be used
- Memory <u>cells</u> store data value
  - Static vs. dynamic, single or multiple bits, etc.

### DRAM



- Smaller cell size (1 transistor or 1T cell)
  - Reason for inexpensive memory in computers
  - Tradeoff of area (memory density) vs. speed and complexity (refreshing)

# **DRAM Issues**

- Must be periodically refreshed
  - Reads are destructive (modify voltage stored on capacitor)
  - Every read followed by a refresh of the bit (write back of read value)
- No static power dissipation
- Output voltage is charge sharing result of storage capacitor and bitline capacitance
  - More complex sense amplifiers
  - Higher noise susceptibility
- Requires different CMOS process than high performance logic
  - Not compatible with cache in microprocessors

### **Advanced DRAM Process**



#### Vertical transistor, trench capacitor (Beintner, JSSC 04)

# ROM

- Dotted lines refer to either set at '1' or '0'
  - PROM: Replace dotted lines with fuses
- Small cell size (1T cell)
- Not necessary to refresh
- No static power dissipation
- Output voltage is set by WL duration



# **SRAM Cell**



- Cross-coupled inverters: bistable element
- Density is important in memories
  - Single NMOS pass transistor used for reading/writing
  - Transistor sizes should take up minimum area
- Faster than DRAM since typically fewer cells
- No refresh required (nondestructive reads)

# SRAM Design: Read "0"



- Prior to read operation, voltage at node Q = 0V and Q = Vdd, bit lines precharged to Vdd
- Transistors M3 and M4 are turned on by word line (WL) select circuitry

# SRAM Design: Read "0"

- Transistors M3 and M4 are turned on by WL line select circuitry
  - Cc = Vdd to start...capacitance discharges through M1.
  - Need to make sure the ratio between M1 and M3 does not allow Q to go above Vtn.
    - Otherwise node  $\overline{Q}$  accidentally discharged
    - Conservative since there will also be charge sharing at that node as well between small internal node capacitance and large bitline capacitance
  - Sense amp detects that node Q was a stored 0 due to the minor drop of voltage on the bitline

- Data must not be destroyed when bitline voltage different than storage node
  - $V_Q$  must not exceed the threshold of the inverter (assumed to be  $V_{DD}/2$ ), more conservative to keep it below  $V_{TN}$
- Assume V<sub>BL</sub> initially remains at V<sub>DD</sub>: M3 in saturation, M1 in linear  $\frac{k_{n,3}}{2} \left( V_{DD} - V_Q - V_{TN} \right)^2 = \frac{k_{n,1}}{2} \left( 2 \left( V_{DD} - V_{TN} \right) V_Q - V_Q^2 \right)$  $\frac{k_{n,3}}{k_{n,1}} = \frac{\left(\frac{W}{L}\right)_3}{\left(\frac{W}{I}\right)_3} < \frac{2(V_{DD} - 1.5V_{TN})V_{TN}}{(V_{DD} - 2V_{TN})^2} \quad \begin{array}{l} \text{Guarantee V}_{\text{Q}} < \text{V}_{\text{TN}}, \\ \text{plug into I}_{\text{D}} \text{ equations:} \\ \text{I}_{\text{D3}} < \text{I}_{\text{D1}} \text{ at V}_{\text{Q}} = \text{V}_{\text{TN}} \end{array}$

### SRAM Design: Write "0" (1<sup>st</sup> Analysis)



- Data must be forced into the cell
  - $-V_Q$  must fall below the threshold of the <u>inverter</u> to turn M2 off.
  - This allows  $\overline{V_Q}$  to go high enough to go above the Vt of M1
    - This discharges node Q and stores a 0
- Assume V<sub>BL</sub> remains at OV: M3 linear, M5 linear (V<sub>Q</sub>=V<sub>DD</sub>/2) Amirtharajah/Parkhurst, EEC 118 Spring 2010

### SRAM Design: Write "0" (1<sup>st</sup> Analysis)



• Conditions for this to happen: requires M5 to M3 ratio to be relatively small ( $V_{DS} = V_Q = V_{DD}/2$ )

$$\frac{k_{p,5}}{2} \left( \left( V_{DD} - \left| V_{TP} \right| \right) V_{DD} - \frac{V_{DD}^{2}}{4} \right) = \frac{k_{n,3}}{2} \left( \left( V_{DD} - V_{TN} \right) V_{DD} - \frac{V_{DD}^{2}}{4} \right)$$

### SRAM Design: Write "0" (1<sup>st</sup> Analysis)



### SRAM Design: Write "0" (2<sup>nd</sup> Analysis)



- Data must be forced into the cell
  - $-V_Q$  must fall below the threshold of the <u>NMOS</u> (turns M2 off).
  - This allows  $\overline{V_Q}$  to go high enough to go above the Vt of M1
    - This discharges node Q and stores a 0
- Assume V<sub>BL</sub> remains at 0V: M5 sat., M3 linear (V<sub>Q</sub> = V<sub>TN</sub>) Amirtharajah/Parkhurst, EEC 118 Spring 2010

### SRAM Design: Write "0" (2<sup>nd</sup> Analysis)



• Desired current conditions, want  $V_Q < V_{TN}$  (M3 lin., M5 sat.):

$$\frac{k_{n,3}}{2} \left( 2 \left( V_{DD} - V_{TN} \right) V_{TN} - V_{TN}^2 \right) = \frac{k_{p,5}}{2} \left( 0 - V_{DD} - V_{TP} \right)^2$$

### SRAM Design: Write "0" (2<sup>nd</sup> Analysis)



• Required sizing (make M5 relatively weaker than 1<sup>st</sup> case):

$$\frac{\left(\frac{W}{L}\right)_{5}}{\left(\frac{W}{L}\right)_{3}} < \frac{k_{n}'}{k_{p}'} \frac{2(V_{DD} - 1.5V_{TN})V_{TN}}{\left(V_{DD} + V_{TP}\right)^{2}}$$

# Flash Memory Transistor With Floating Gate



# **Flash Memory Operation**

- Two threshold voltages correspond to two states (1 bit)
  - -Bitline is precharged high before a read
  - Low  $V_T$  state, when wordline (control gate) is high the bitline is discharged and a "0" is read
  - High  $V_T$  state, the wordline (control gate) can't go high enough to turn on transistor, bitline stays high and a "1" is read
- Writing a "1": electrons are accelerated by a high field until they accumulate on the floating gate, raising  $V_T$
- Writing a "0": electrons driven off floating gate by a reverse gate-source bias through Fowler-Nordheim tunneling, lowering V<sub>T</sub>

# **Next Topics: Low Power Circuits and Wires**

- Low power design principles and circuit techniques
  - Voltage scaling, activity factor reduction, clock gating, leakage reduction
- On-chip resistance and capacitance
  - Delay estimation
  - Buffering and repeater insertion