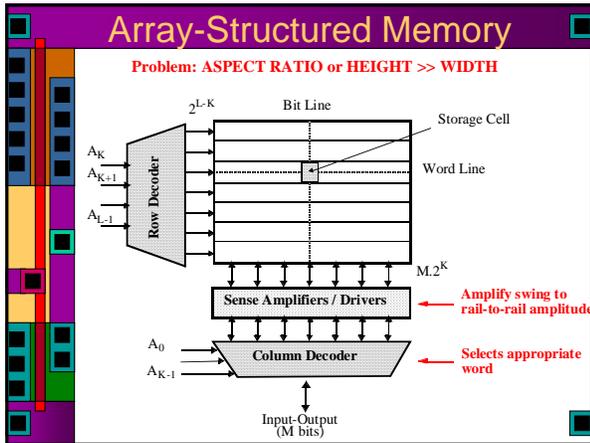
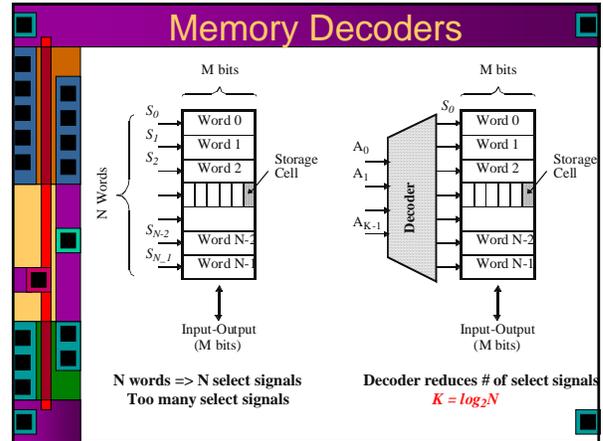
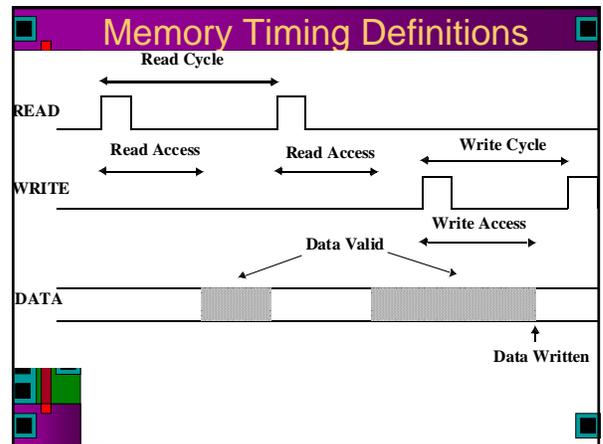
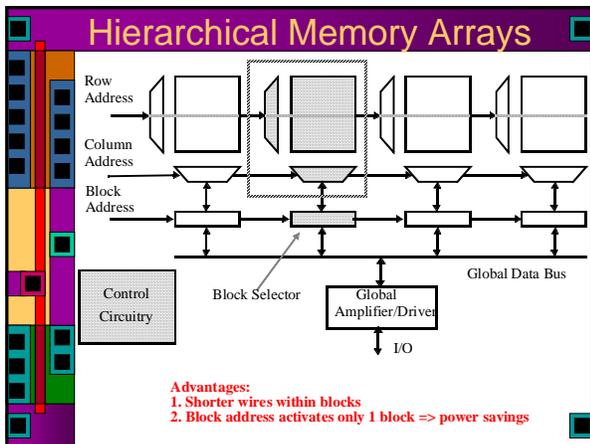


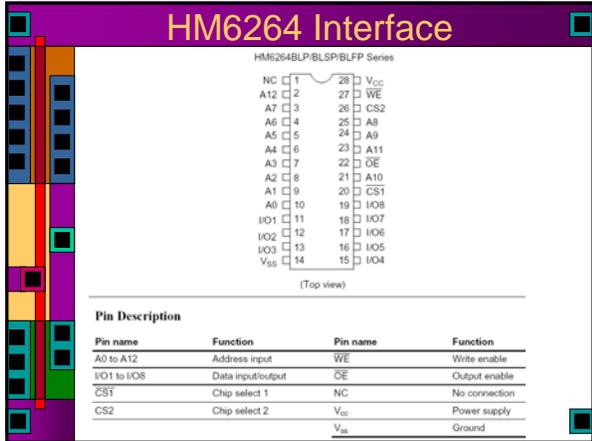
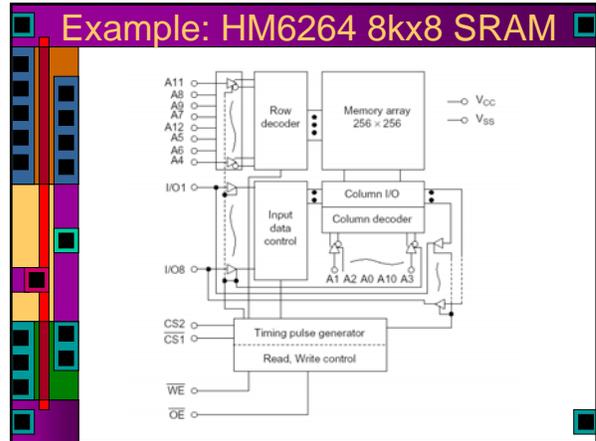
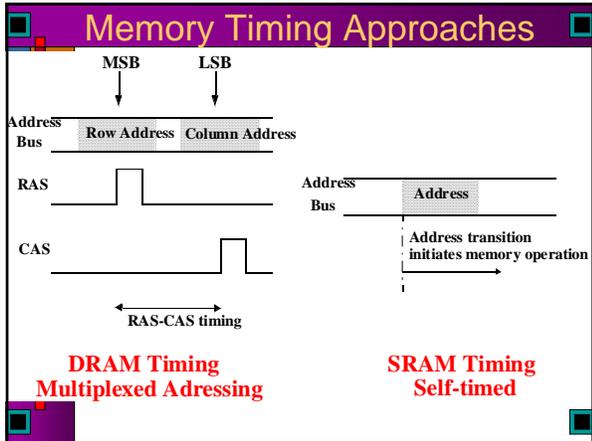
### Memory

| RWM           |                                       | NVRWM                        | ROM                                 |
|---------------|---------------------------------------|------------------------------|-------------------------------------|
| Random Access | Non-Random Access                     | EPROM<br>E <sup>2</sup> PROM | Mask-Programmed Programmable (PROM) |
| SRAM<br>DRAM  | FIFO<br>LIFO<br>Shift Register<br>CAM | FLASH                        |                                     |



- ### Array Decoding
- Typically want an aspect ratio that is not too far from square
  - How to divide up the row, column address decoding?
- Use an 8K x 32 SRAM = 256 Kb =  $2^{18}$
- $2^{18} = 2^9 \text{ rows} \times 2^9 \text{ columns}$
- Row decoder is 9 to 512 decoder
- Every 32 ( $2^5$ ) columns is a 'word', and we only need to decode words. So, column decoder needs to decode  $2^4$  words, so need a 4 to 16 column decoder.





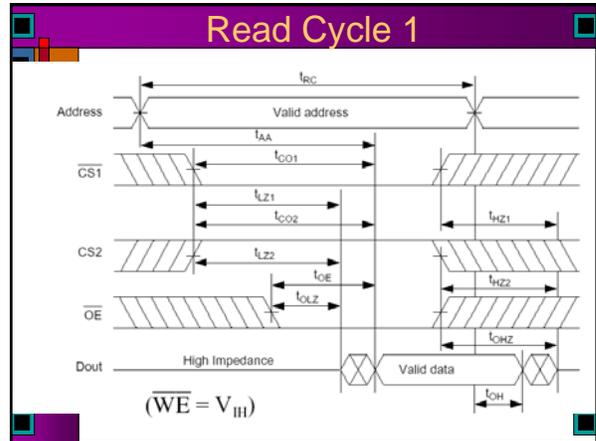
### Function Table

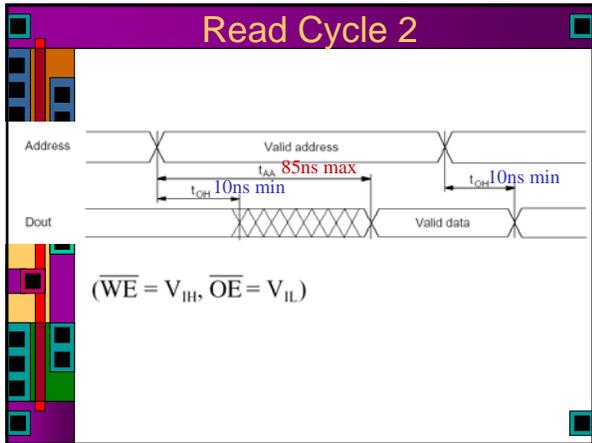
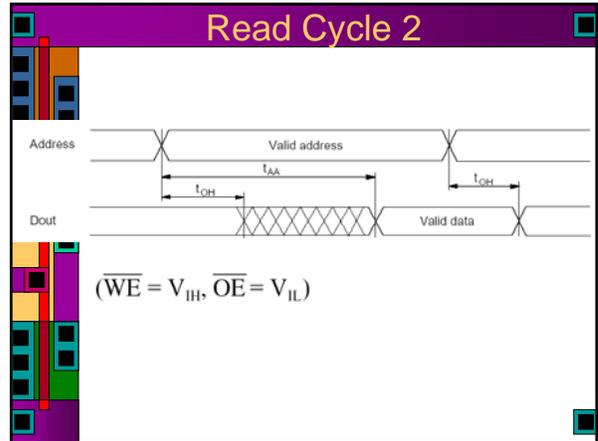
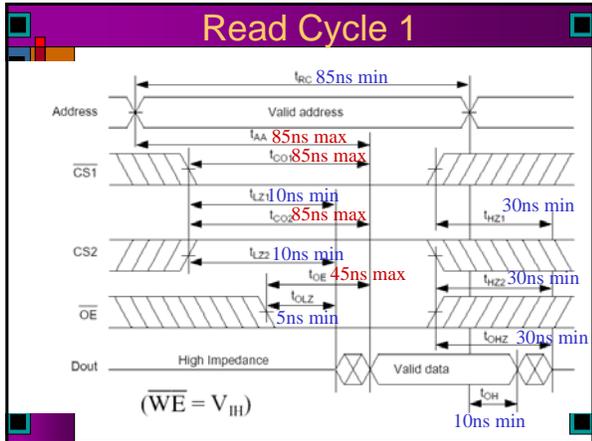
| WE | CS1 | CS2 | OE | Mode                      | V <sub>CC</sub> current | I/O pin | Ref. cycle         |
|----|-----|-----|----|---------------------------|-------------------------|---------|--------------------|
| x  | H   | x   | x  | Not selected (power down) | I <sub>BB-15m</sub>     | High-Z  | —                  |
| x  | x   | L   | x  | Not selected (power down) | I <sub>BB-15m</sub>     | High-Z  | —                  |
| H  | L   | H   | H  | Output disable            | I <sub>CC</sub>         | High-Z  | —                  |
| H  | L   | H   | L  | Read                      | I <sub>CC</sub>         | Dout    | Read cycle (1)-(3) |
| L  | L   | H   | H  | Write                     | I <sub>CC</sub>         | Din     | Write cycle (1)    |
| L  | L   | H   | L  | Write                     | I <sub>CC</sub>         | Din     | Write cycle (2)    |

Note: x = H or L

### Timing

| Parameter                               | Symbol           | HM6264B-8L       |     | HM6264B-10L |     | Unit | Notes |      |
|---|------------------|------------------|-----|-------------|-----|------|-------|------|
|   |                  | Min              | Max | Min         | Max |      |       |      |
| Read cycle time                         | t <sub>RC</sub>  | 85               | —   | 100         | —   | ns   |       |      |
| Address access time                     | t <sub>AA</sub>  | —                | 85  | —           | 100 | ns   |       |      |
| Chip select access time                 | CS1              | t <sub>CO1</sub> | —   | 85          | —   | 100  | ns    |      |
|   |                  | t <sub>COE</sub> | —   | 85          | —   | 100  | ns    |      |
| Output enable to output valid           | t <sub>OE</sub>  | —                | 45  | —           | 50  | ns   |       |      |
| Chip selection to output in low-Z       | CS1              | t <sub>LZ1</sub> | 10  | —           | 10  | —    | ns    | 2    |
|   |                  | t <sub>LZ2</sub> | 10  | —           | 10  | —    | ns    | 2    |
| Output enable to output in low-Z        | t <sub>OLZ</sub> | 5                | —   | 5           | —   | ns   | 2     |      |
| Chip deselection in to output in high-Z | CS1              | t <sub>HZ1</sub> | 0   | 30          | 0   | 35   | ns    | 1, 2 |
|   |                  | t <sub>HZ2</sub> | 0   | 30          | 0   | 35   | ns    | 1, 2 |
| Output disable to output in high-Z      | t <sub>OHZ</sub> | 0                | 30  | 0           | 35  | ns   | 1, 2  |      |
| Output hold from address change         | t <sub>OH</sub>  | 10               | —   | 10          | —   | ns   |       |      |

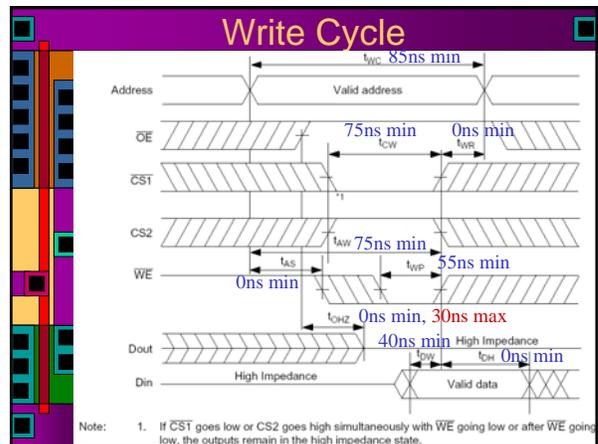
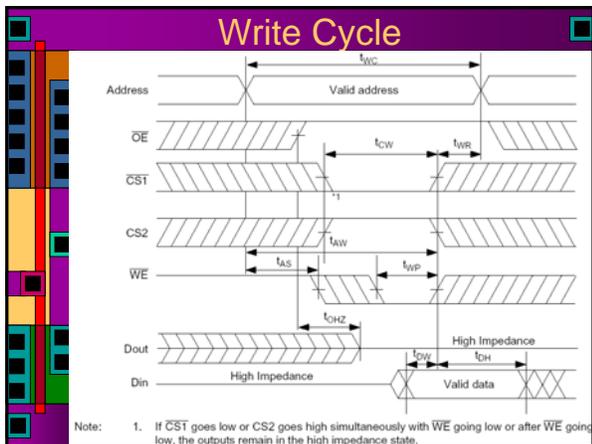




### Write Timing

| Parameter                          | Symbol    | HM6264B-8L |     | HM6264B-10L |     | Unit | Notes |
|------------------------------------|-----------|------------|-----|-------------|-----|------|-------|
|                                    |           | Min        | Max | Min         | Max |      |       |
| Write cycle time                   | $t_{WC}$  | 85         | —   | 100         | —   | ns   |       |
| Chip selection to end of write     | $t_{CW}$  | 75         | —   | 80          | —   | ns   | 2     |
| Address setup time                 | $t_{AS}$  | 0          | —   | 0           | —   | ns   | 3     |
| Address valid to end of write      | $t_{AV}$  | 75         | —   | 80          | —   | ns   |       |
| Write pulse width                  | $t_{WP}$  | 55         | —   | 60          | —   | ns   | 1, 6  |
| Write recovery time                | $t_{WR}$  | 0          | —   | 0           | —   | ns   | 4     |
| WE to output in high-Z             | $t_{WHZ}$ | 0          | 30  | 0           | 35  | ns   | 5     |
| Data to write time overlap         | $t_{DOW}$ | 40         | —   | 40          | —   | ns   |       |
| Data hold from write time          | $t_{DH}$  | 0          | —   | 0           | —   | ns   |       |
| Output active from end of write    | $t_{LOW}$ | 5          | —   | 5           | —   | ns   |       |
| Output disable to output in high-Z | $t_{OHZ}$ | 0          | 30  | 0           | 35  | ns   | 5     |

Notes: 1. A write occurs during the overlap of a low CS1, and high CS2, and a high WE. A write begins at the latest transition among CS1 going low, CS2 going high and WE going low. A write ends at the earliest transition among CS1 going high, CS2 going low and WE going high. Time  $t_{WP}$  is measured from the beginning of write to the end of write.

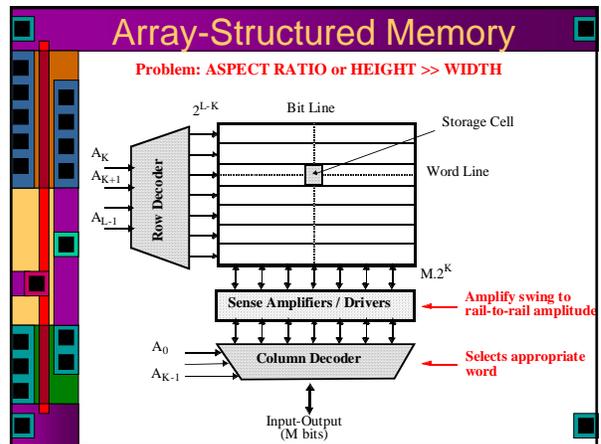
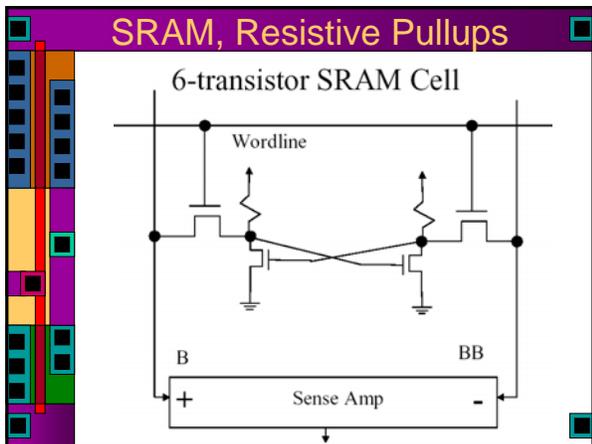
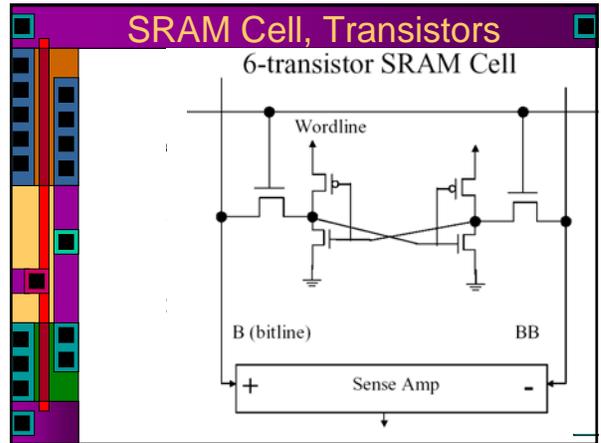
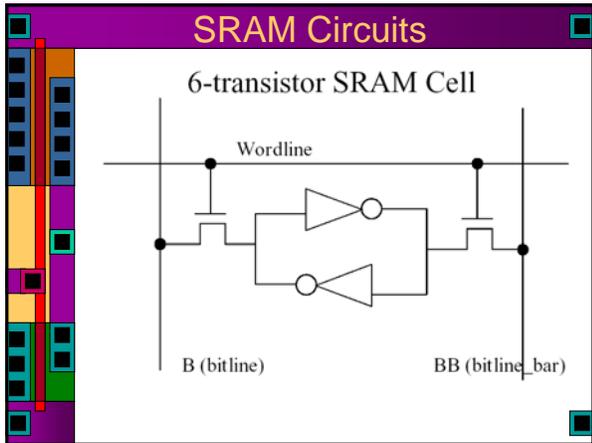


### What Does All This Mean

- ▶ For a read:
  - ▶ If you assert CS1, CS2, address, and OE all at the same time, it will be max 85ns before valid data are available at chip outputs
- ▶ For a write:
  - ▶ You can assert CS1, CS2, address, data, and WE all at the same time if you want to
  - ▶ You need to wait 55ns from WE edge, or 75ns from CS1/CS2 edge for write to have happened

### R/W Memories In General

- **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



### Memory Column

▶ Each column has all the support circuits

The diagram shows an address bus with two segments:  $n-1:k$  and  $k-1:0$ . The  $n-1:k$  segment goes through a Flow Decoder to a RAM cell. The  $k-1:0$  segment goes through a Column Decoder to a Sense Amp Column Mux Write Buffers. The RAM cell is connected to Bit Line Conditioning (clocked) and the Sense Amp. The Sense Amp outputs write-data and read-data, and is also connected to write clocks and read clocks.

### Reading the Bit

The circuit diagram shows a RAM cell connected to a bit line through an inverter. The bit line is precharged and then pulled up by a P-type transistor. The sense amp is connected to the bit line and the inverter output. Waveforms show precharge, bit, -bit, word, and data signals.

- ▶ Single-ended read using an inverter
- ▶ Dynamic pre-charge on the bit lines
- ▶ P-types pull bit lines high

### Reading the Bit 2

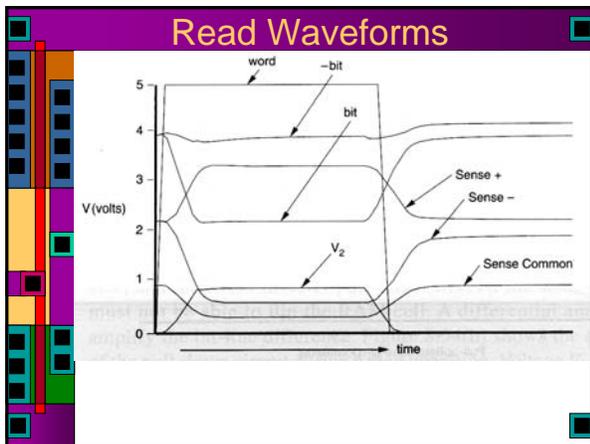
The circuit diagram shows a RAM cell connected to a bit line through an inverter. The bit line is precharged and then pulled up by N-type transistors. The sense amp is connected to the bit line and the inverter output. Waveforms show precharge, bit, -bit, word, and data signals.

- ▶ Single-ended read using an inverter
- ▶ Dynamic pre-charge on the bit lines
- ▶ Note the N-types used as pull-ups

### Reading the Bit 3

The circuit diagram shows a RAM cell connected to a bit line through an inverter. The bit line is precharged and then pulled up by N-type transistors. The sense amp is connected to the bit line and the inverter output. Waveforms show bit, -bit, word, and data signals. A load is connected to the bit line, and the sense amp has a pull-down transistor.

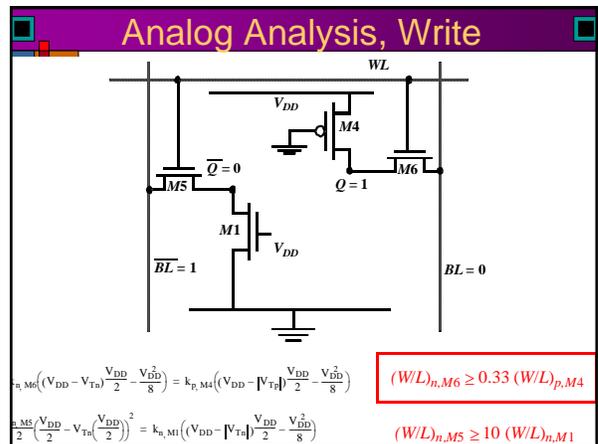
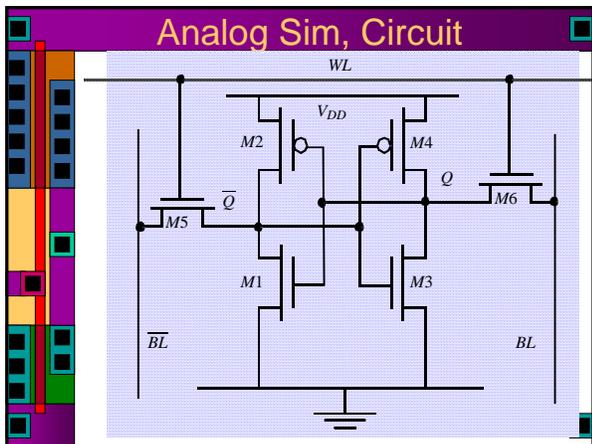
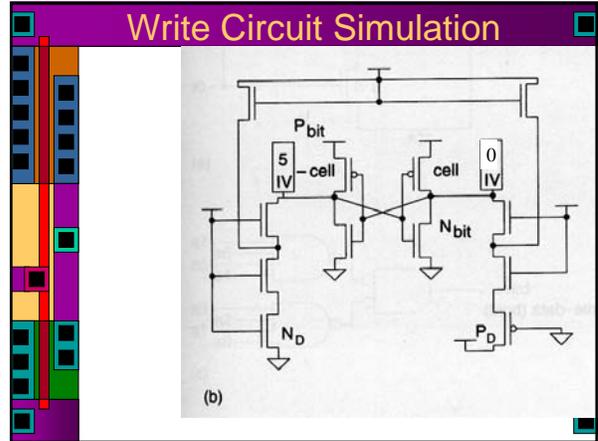
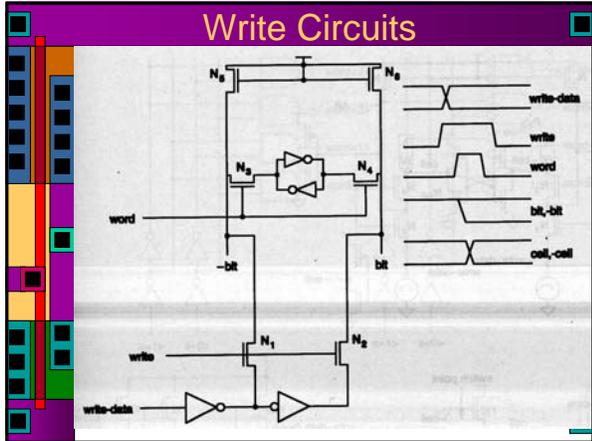
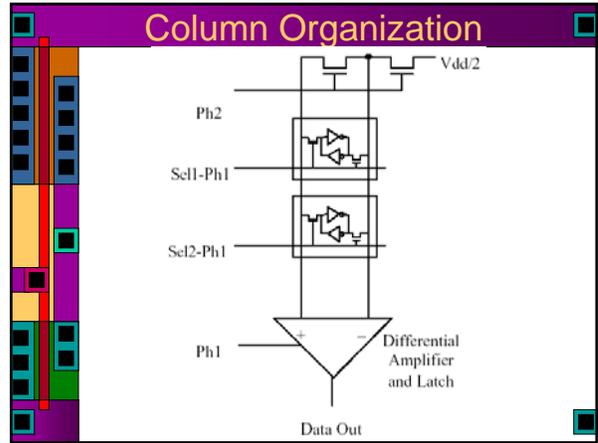
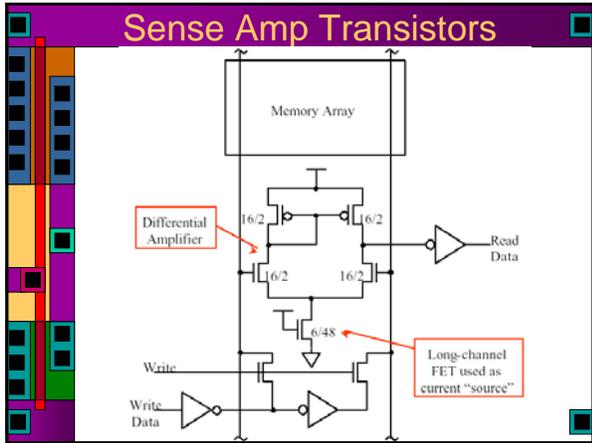
- ▶ Differential read using sense amp
- ▶ Static N-type pullup on the bit lines

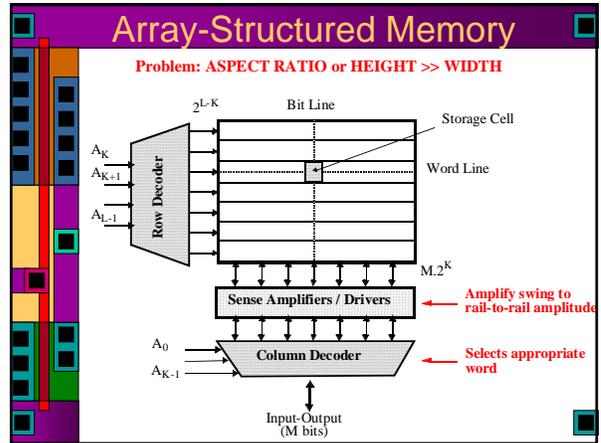
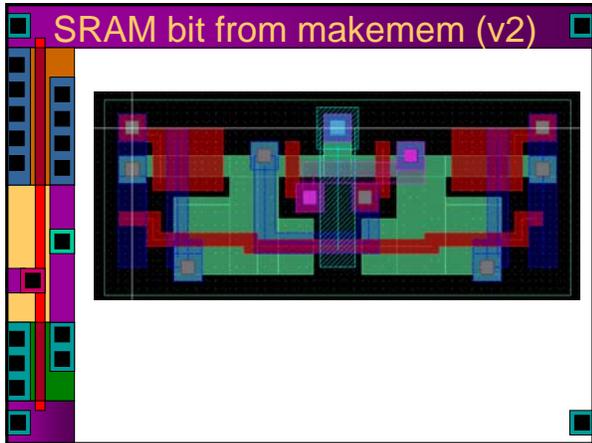
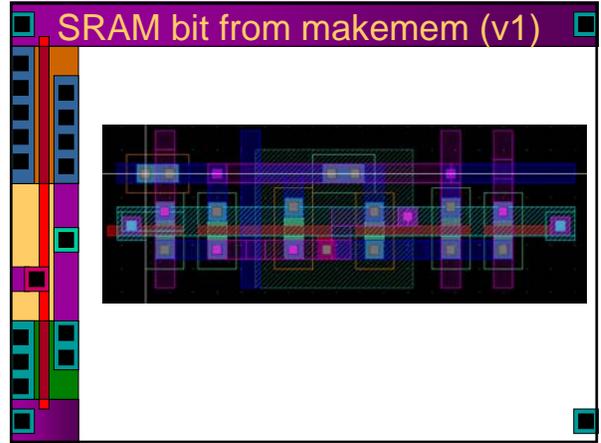
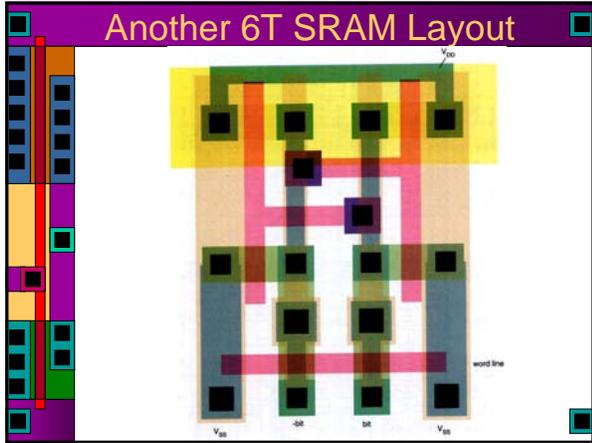
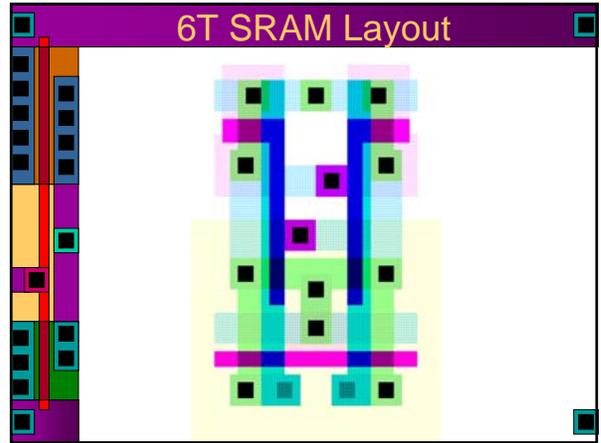
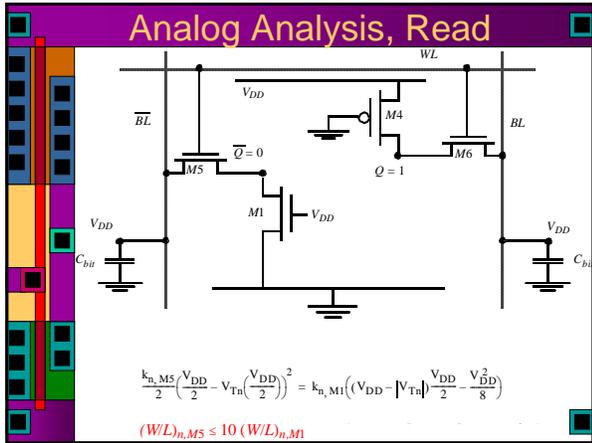


### Sense Amp

The circuit diagram shows a current-mirror sense amplifier. It consists of two PMOS transistors and two NMOS transistors. The PMOS transistors are connected to a common source (SE) and the NMOS transistors are connected to a common drain (BB). The sense amp is used to sense the bit line swing and amplify it to a full swing output.

Only one bit line will swing.  
Want Sense Amplifier turned on for short amount of time in order to save power.  
Job of SA is to sense bit line swing, amplify to full swing output.





### Row Decoders

▶ Select exactly one of the memory rows  
 ▶ Simple versions are just gates

### Row Decoder Gates

- ▶ Standard gates
- ▶ Or, pseudo-nmos gates with static pull up
  - ▶ Easier to make large fan-in NOR

### Pre-decode Row Decoder

▶ Multiple levels of decoding can be more efficient layout

### Pre-decode Row Decoder

▶ Other circuit tricks for building row decoders...

### Array-Structured Memory

**Problem: ASPECT RATIO or HEIGHT >> WIDTH**

Sense Amplifiers / Drivers → Amplify swing to rail-to-rail amplitude  
 Column Decoder → Selects appropriate word

### Array-Structured Memory



### Multi-Port Register

WE  
Re1  
Re0  
Write Data  
Read Data  
D0  
D1

- ▶ Slightly larger cell, but with single-ended read – makes a great register file

### Register File

wr-b-add-c3-0  
rd-a-add-c3-0  
rd-b-add-c3-0  
write-data  
read-data0  
read-data1

- ▶ Slightly larger cell, but with single-ended read – makes a great register file

### Dynamic RAM

word  
- bit  
bit  
(a)

- ▶ Get rid of the pull-ups!
  - ▶ Store info on capacitors
  - ▶ Means that stored information leaks away

### Dynamic RAM...

write  
read  
write-data  
read-data  
(b)  
(c)  
word  
 $V_{DD}$  or  $V_{DD}/2$   
bit  
(d)

- ▶ Once you agree to use a capacitor for charge storage there are other ways to build this...

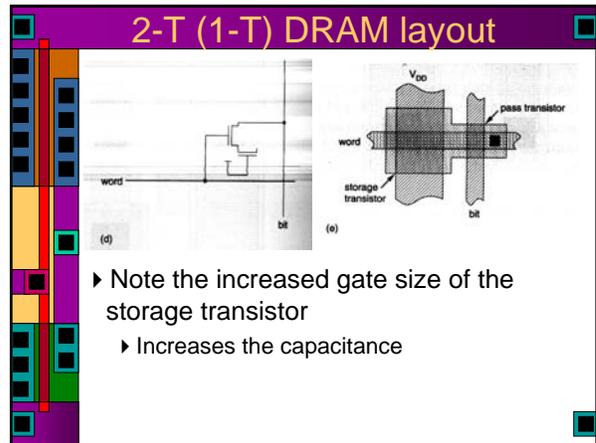
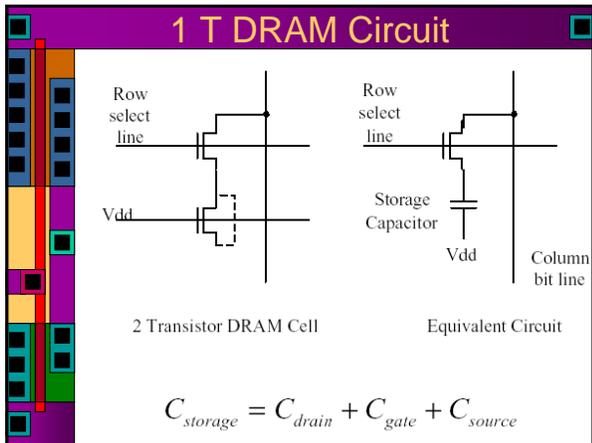
### 3T DRAM Circuit

BL1  
BL2  
WWL  
RWL  
X  
M1  
M2  
M3  
Cs  
V<sub>DD</sub>-V<sub>T</sub>  
V<sub>DD</sub>  
V<sub>DD</sub>-V<sub>T</sub> ΔV

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

BL2  
BL1  
GND  
RWL  
M3  
M2  
WWL  
M1



### 1T DRAM Observations

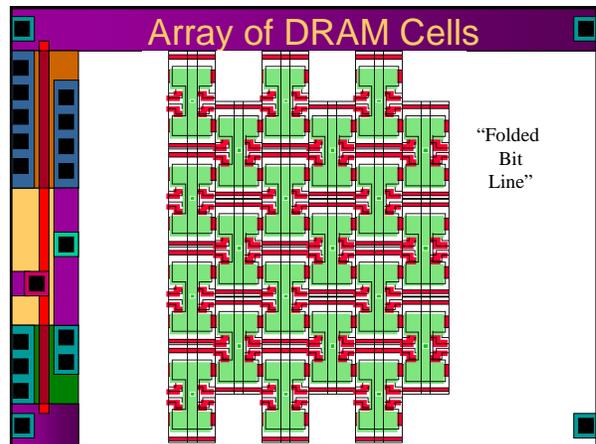
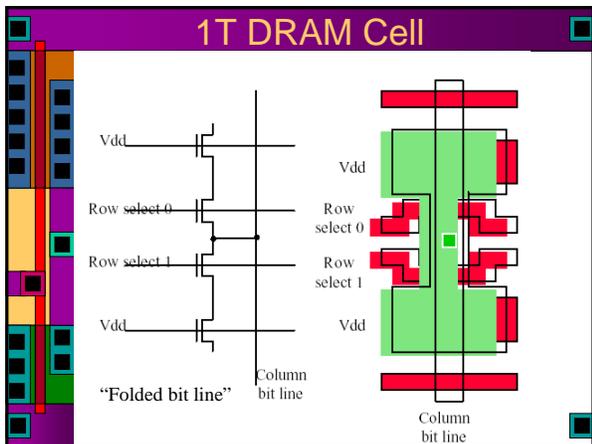
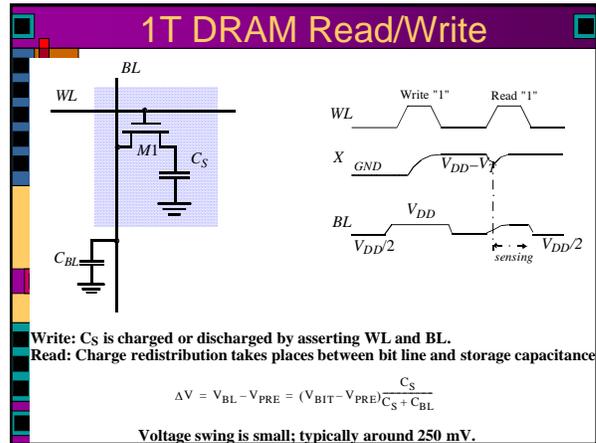
1T DRAM requires a sense amplifier for each bit line, due to charge redistribution read-out.

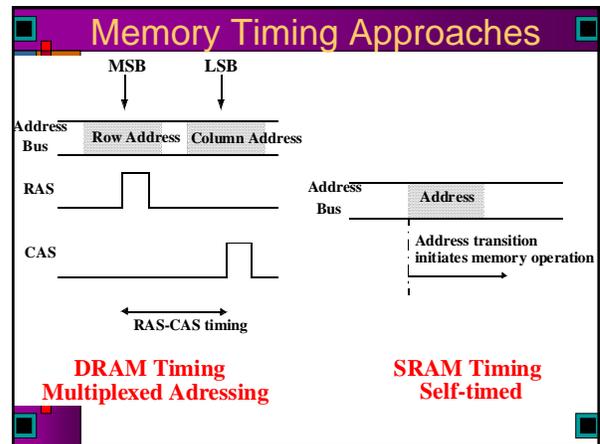
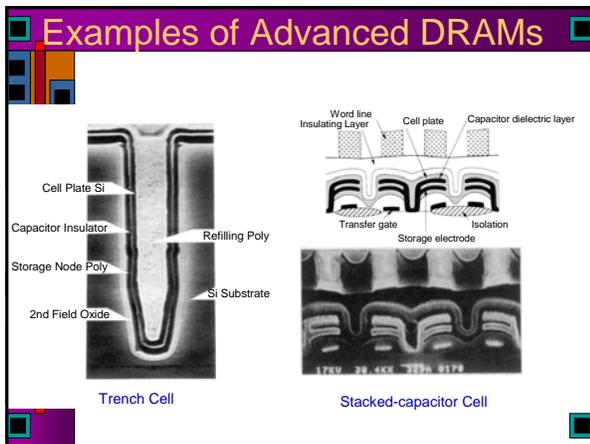
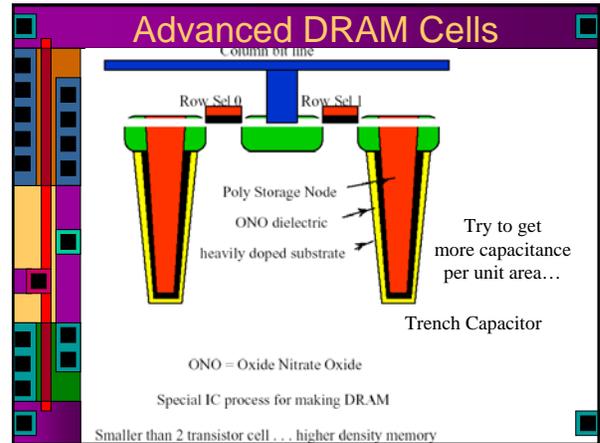
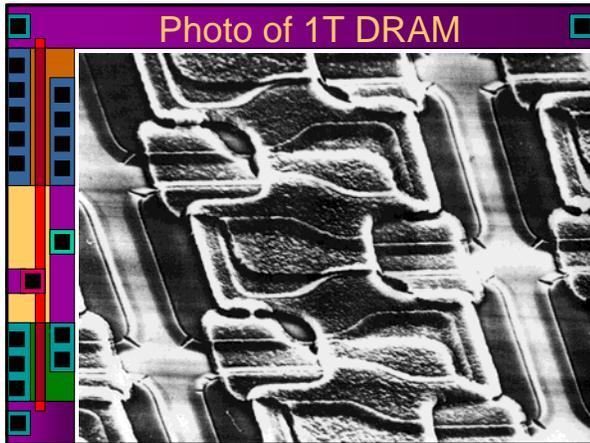
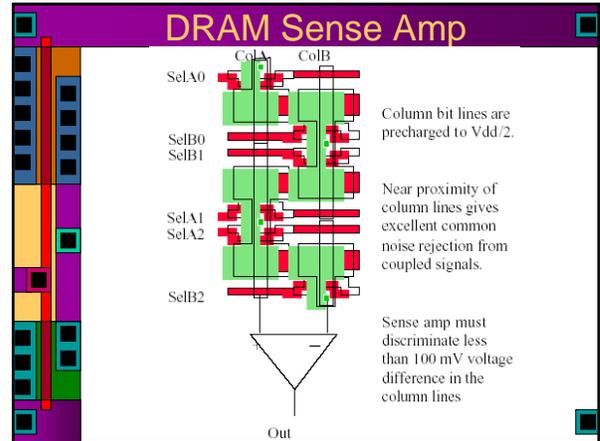
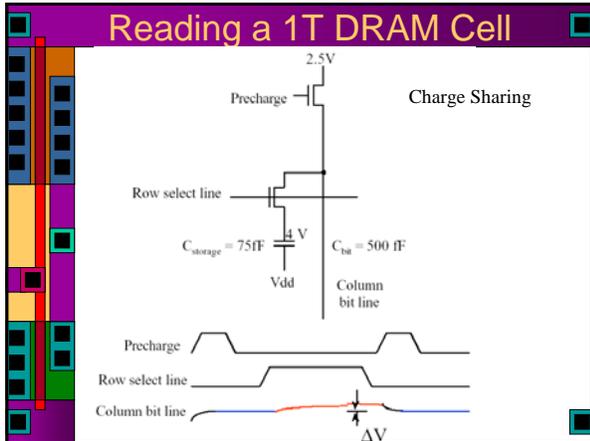
DRAM memory cells are single ended in contrast to SRAM cells.

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  $V_{DD}$ .







### DRAM Timing

- Clock Frequency – 133 Mhz, 100 Mhz
- Two clock latency to first data (20 ns for 100 Mhz clock)
  - SDRAM - 10 ns per location afterwards. For byte-wide, 100 MB/sec transfer rate. 400 MB/sec on 32-bit bus
  - DDR-SDRAM - 5 ns per location afterwards. For byte-wide, 200 MB/sec transfer rate. On 32-bit bus, 800 MB/sec transfer rate.

### RAMBUS DRAM (RDRAM)

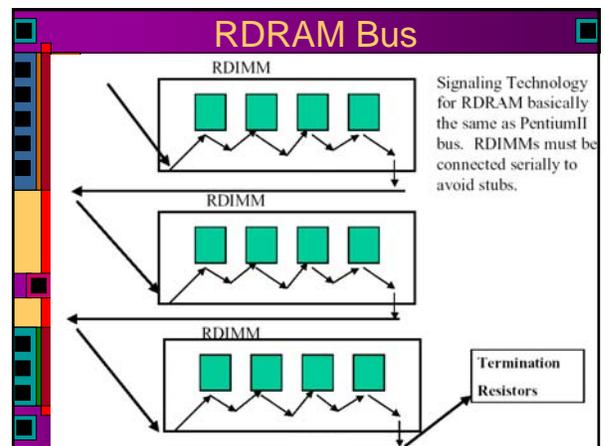
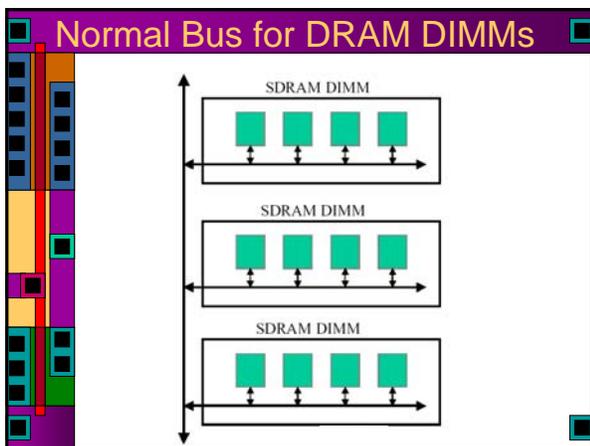
- DRAM with a high speed interface
- 400 Mhz differential clock, data transferred on each edge
- Reduced swing signaling about a reference voltage
  - Termination voltage is 1.5 V
  - Reference Voltage is 1.0 V
  - Signals swing +/- 200 mv about reference voltage
  - All traces are transmission lines

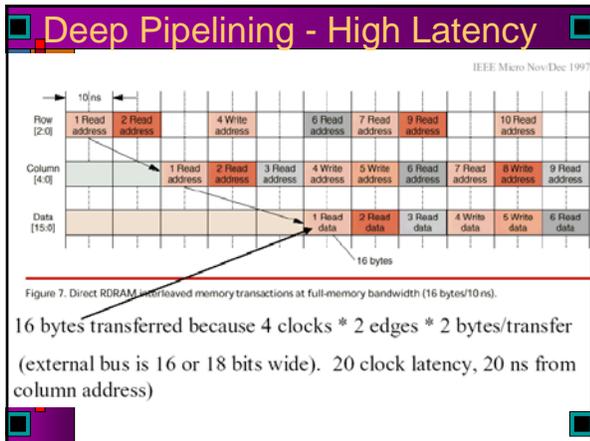
### RDRAM Bandwidth

- External bus is 18 bits wide (2 bytes + 2 parity bits)
- External clock cycle is 400 Mhz, but data is clocked on each edge
  - Actually, external clock is a differential pair and data is sampled at each crossing
- Total Bandwidth is 1.6 GBytes/s
  - $2 \text{ bytes} * 400 \text{ Mhz} * 2 \text{ edges} \Rightarrow 1.6 \text{ Gbytes}$
  - Initial configurations are 4 M x 18 (72 Mbits)

### Maximum Bandwidth

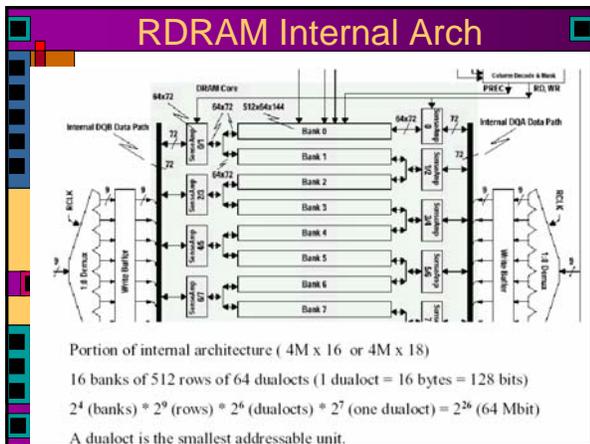
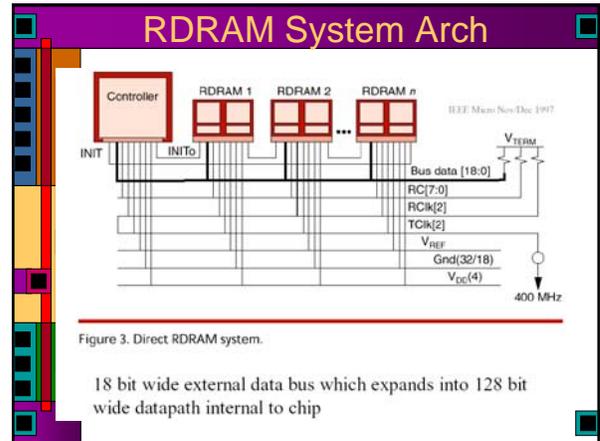
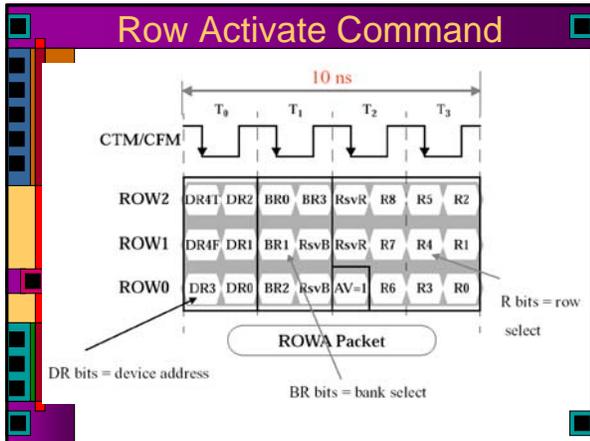
- Note that maximum bandwidth with one RDRAM controller is 1.6GB/s.
  - Only one RDRAM chip can be active at a time on RDRAM bus.
  - More RDRAM chips increase capacity, not bandwidth.
    - With normal DRAM and SDRAM, can increase bandwidth by just adding more DRAM chips in parallel from same DRAM controller
  - To double the bandwidth, would need two separate RDRAM controllers





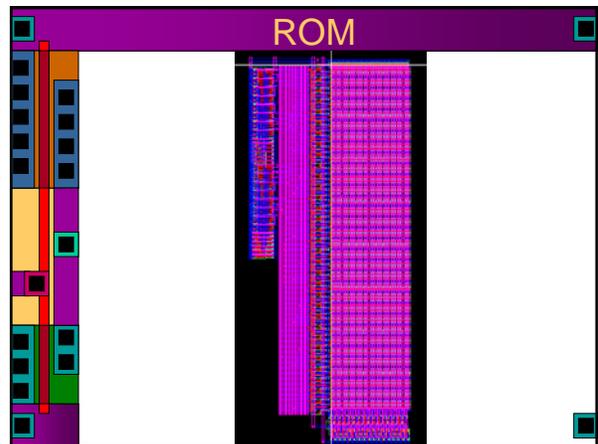
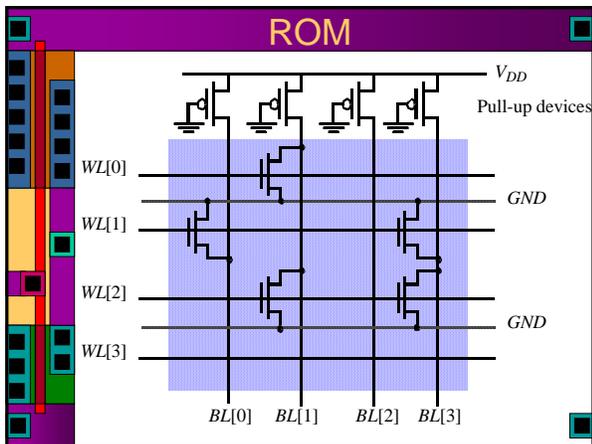
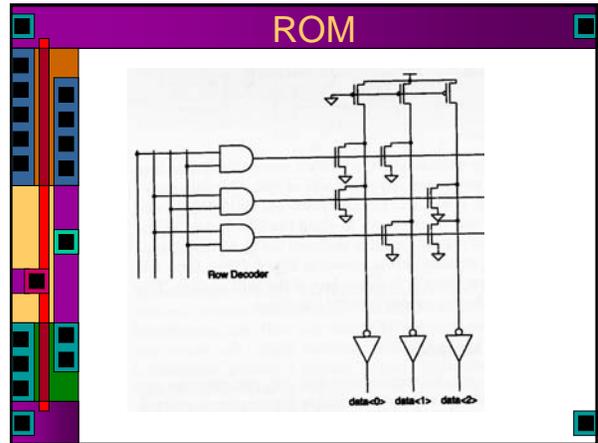
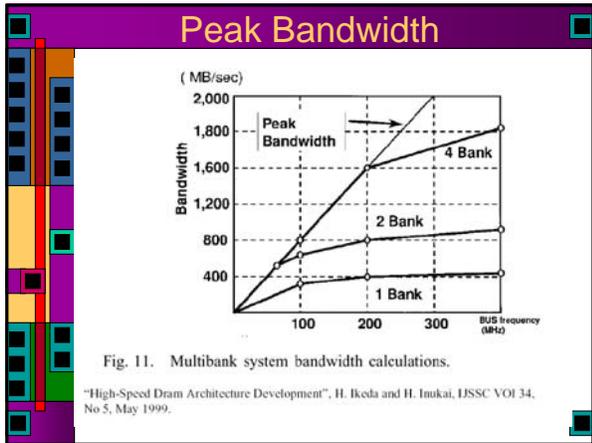
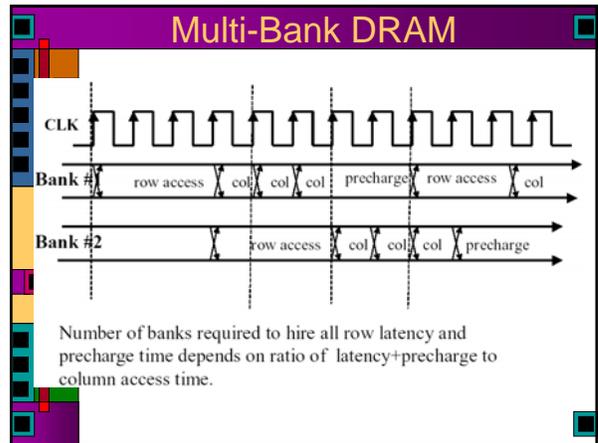
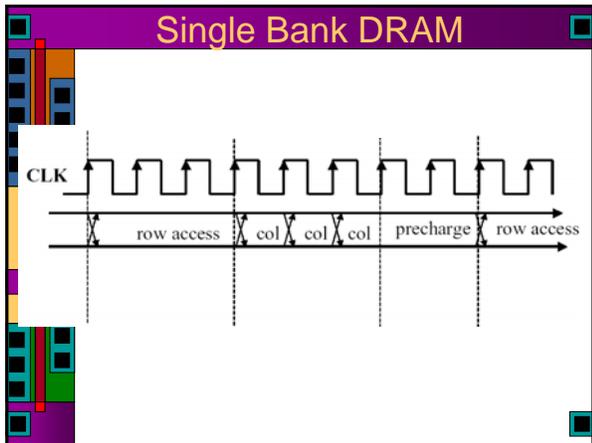
### RDRAM Addressing

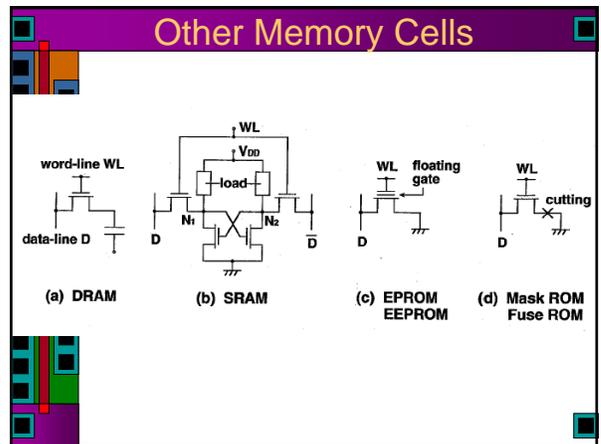
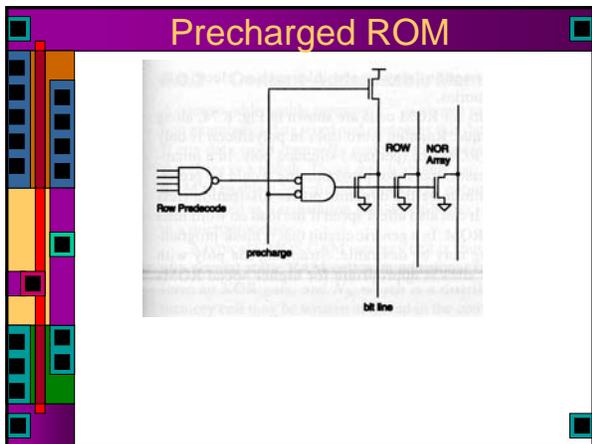
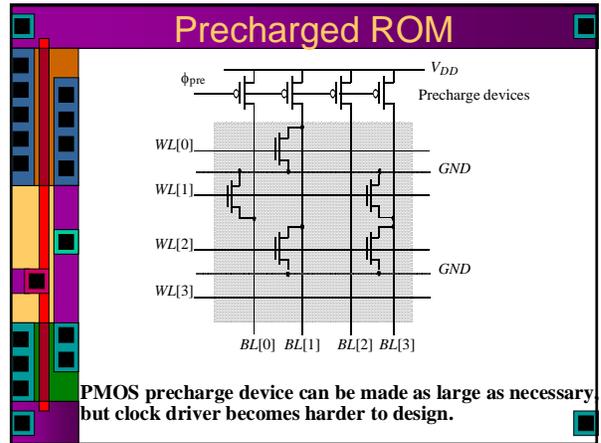
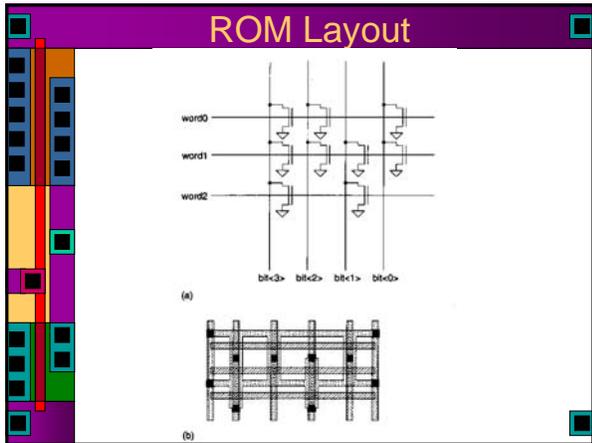
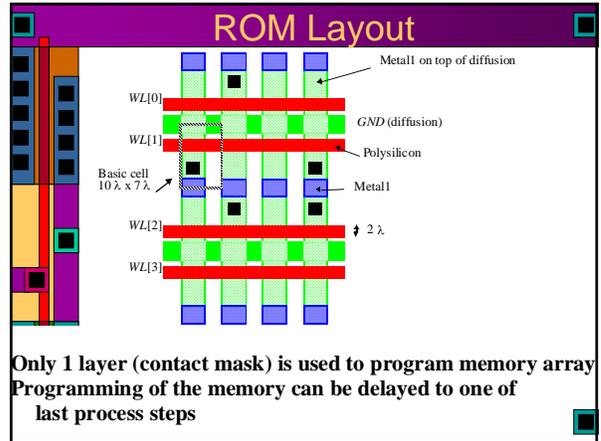
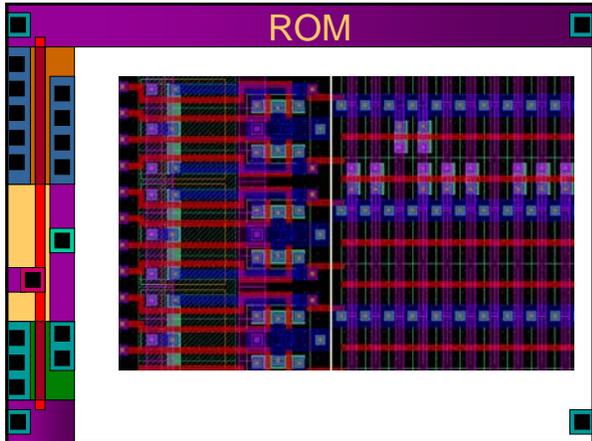
- 3-Bit Row bus used to give commands to RDRAM
- ROW Activate command used for read
  - 4 clocks transfers 8 groups of 3 bits over Row bus due to dual edge clocking (24 bits total)
  - 24 bits in Row Activate command split between device address (6 bits), bank select (4 bits), row select (9 bits), and reserved bits
- There are no chip select lines, internal register holds device address
  - All chips monitor bus - if bus device address matches internal id, then chip is selected.



### Regular DRAM

- Multiple Banks are key to high throughput
- As one DRAM bank is recovering from read operation, next bank is being accessed
- Essentially on-chip memory interleaving
- Goal is to hide latency and bitline precharge time (recovery time)
  - Latency is access to first byte, critical path through row-decode and word line assertion
  - Bitline Precharge time (recovery time to next access) depends on number of bits in a column (number of rows)





## Non-Volatile ROM

- ▶ EPROM
  - ▶ Erasable Programmable ROM
- ▶ EEPROM
  - ▶ Electrically Erasable Programmable ROM
- ▶ Flash EEPROM
  - ▶ Electrically Erasable Programmable ROM that is erased in large chunks
- ▶ All these devices rely on trapping charge on a floating gate

## EPROM

(a) Device cross-section      (b) Schematic symbol

## Programming EPROM

Avalanche injection.      Removing programming voltage leaves charge trapped.      Programming results in higher  $V_T$ .

- ▶ Higher  $V_{th}$  (around 7v) means that 5v  $V_g$ s no longer turns on the transistor
- ▶  $SiO_2$  is an excellent insulator
- ▶ Trapped charge can stay for years

## Erasing an EPROM

- ▶ Erase by shining UV light through window in the package
  - ▶ UV radiation makes oxide slightly conductive
  - ▶ Erasure is slow - from seconds to minutes depending on UV intensity
  - ▶ Also the erase/program cycles are limited (around 1000), mainly as a result of the UV erasing
- ▶ But, EPROMs are simple and dense

## EEPROM

(a) Flotox transistor      (b) Fowler-Nordheim I-V characteristic

Floating Gate Tunneling Oxide transistor

(c) EEPROM cell during a read operation

- ▶ Thin oxide allows erasing in-system
- ▶ Fowler-Nordheim Tunneling

## EEPROM

- ▶ Two transistors instead of one
  - ▶ The second keeps you from removing too much charge during erasure
- ▶ Bigger and not as dense as EPROM
- ▶ But, more erase/program cycles
  - ▶ On the order of  $10^5$
  - ▶ Eventually you get permanently trapped charge in the  $SiO_2$

### Flash EEPROM

Control gate

Floating gate

Thin tunneling oxide

erasure

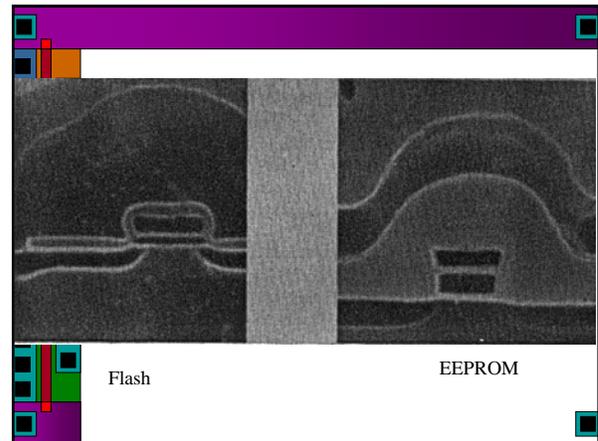
programming

$n^+$  source

$n^+$  drain

$p$ -substrate

- ▶ Essentially the same as EEPROM
- ▶ But, large regions erased at once
- ▶ Means you can monitor the voltages and don't need the extra access transistor



### Realistic PROM Devices

|                                | EPROM [Tomita91]             | EEPROM [Terada89, Pashley89] | Flash EEPROM [Jinbo92]       |
|--------------------------------|------------------------------|------------------------------|------------------------------|
| Memory size                    | 16 Mbit (0.6 $\mu\text{m}$ ) | 1 Mbit (0.8 $\mu\text{m}$ )  | 16 Mbit (0.6 $\mu\text{m}$ ) |
| Chip size                      | 7.18 x 17.39 $\text{mm}^2$   | 11.8 x 7.7 $\text{mm}^2$     | 6.3 x 18.5 $\text{mm}^2$     |
| Cell size                      | 3.8 $\mu\text{m}^2$          | 30 $\mu\text{m}^2$           | 3.4 $\mu\text{m}^2$          |
| Access time                    | 62 nsec                      | 120 nsec                     | 58 nsec                      |
| Erasure time                   | minutes                      | N.A.                         | 4 sec                        |
| Programming time/word          | 5 $\mu\text{sec}$            | 8 msec/word, 4 sec /chip     | 5 $\mu\text{sec}$            |
| Erase/Write cycles [Pashley89] | 100                          | $10^5$                       | $10^3$ - $10^5$              |

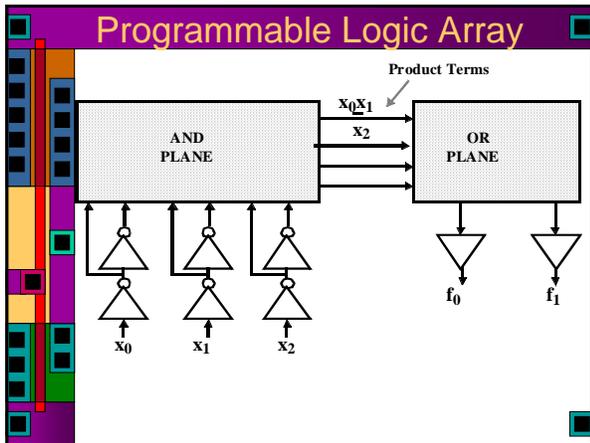
### Content Addressable Mem

- ▶ Asks the question: Are there any locations that hold this value?
- ▶ Used for tag memories in associative caches
- ▶ Or translation lookaside buffers
- ▶ Or other pattern matching applications

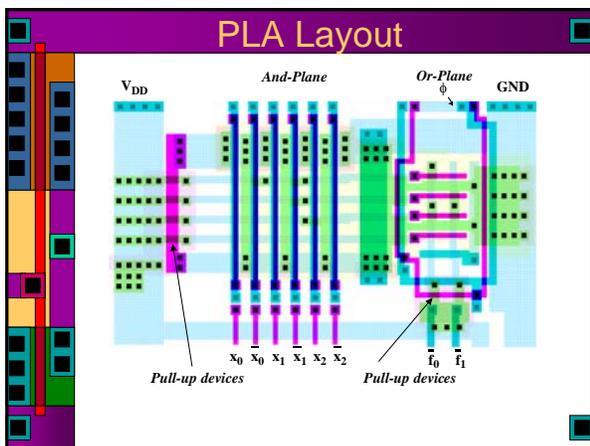
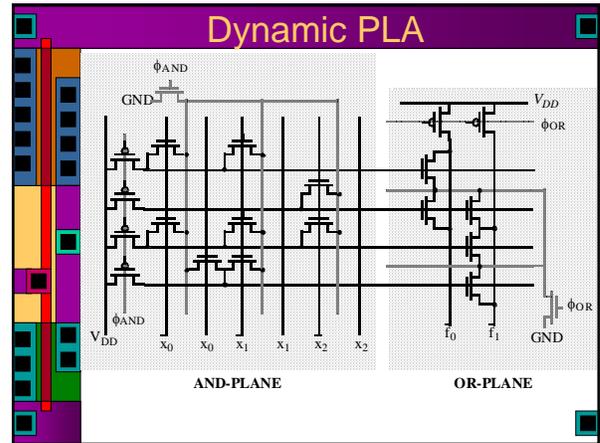
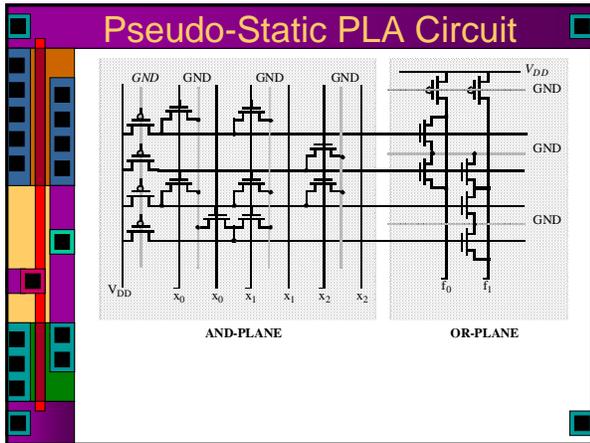
### Content Addressable Mem

- ▶ Add the Match line
- ▶ Essentially a distributed NOR gate

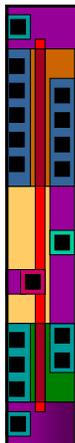
### Content Addressable Mem



- ### PLA
- ▶ Still useful for random combinational logic
    - ▶ Standard cell ASIC tools may be replacing them
  - ▶ They can generate dense AND-OR circuits



- ### PLA vs. 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 (multi-level logic synthesis)

A schematic diagram of an FPGA gate array. It shows a grid of transistors (represented by black squares) connected to a network of metal lines (represented by colored lines: blue, orange, yellow, purple, green). The diagram is enclosed in a purple border with the text 'FPGAs' in yellow at the top center.

## FPGAs

- ▶ Field Programmable Gate Arrays
  - ▶ Array of P-type and N-type transistors
  - ▶ Sources and drains connected to
    - ▶ Power and ground
    - ▶ Metal
  - ▶ Map gate structures to sea of gates
  - ▶ Less expensive – only modify metal masks