

CS152

Computer Architecture and Engineering

Lecture 16: Memory System

March 15, 1995

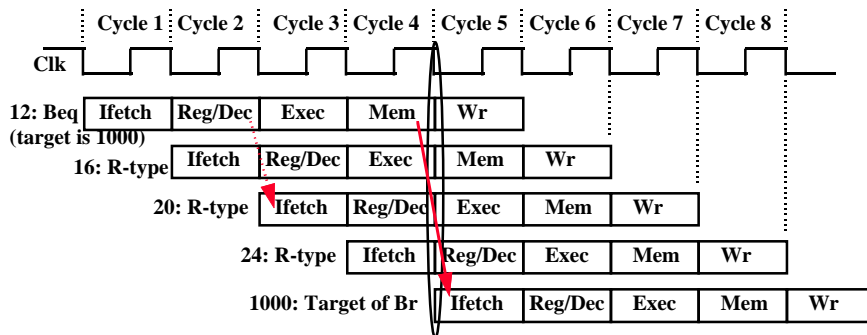
Dave Patterson (patterson@cs) and
Shing Kong (shing.kong@eng.sun.com)

Slides available on <http://http.cs.berkeley.edu/~patterson>

cs 152 memory.1

©DAP & SIK 1995

Recap: Solution to Branch Hazard

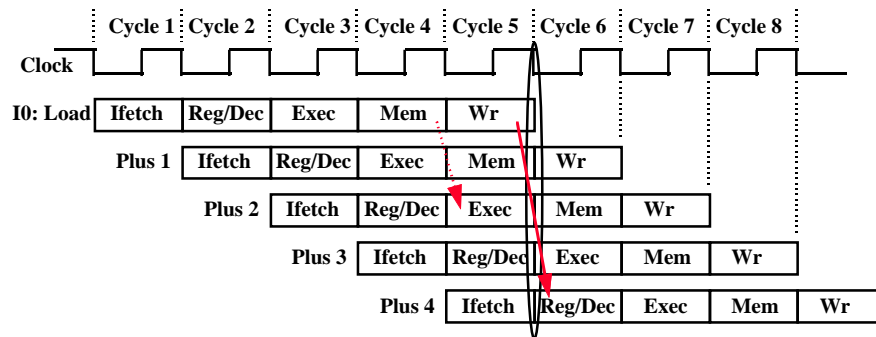


- In the Simple Pipeline Processor if a Beq is fetched during Cycle 1:
 - Target address is NOT written into the PC until the end of Cycle 4
 - Branch's target is NOT fetched until Cycle 5
 - 3-instruction delay before the branch take effect
- This Branch Hazard can be reduced to 1 instruction if in Beq's Reg/Dec:
 - Calculate the target address
 - Compare the registers using some "quick compare" logic

cs 152 memory.2

©DAP & SIK 1995

Recap: Solution to Load Hazard



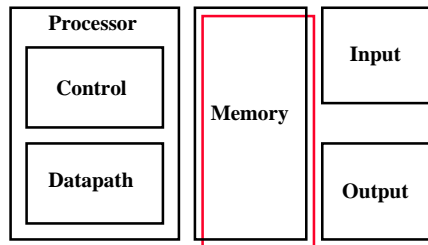
- In the Simple Pipeline Processor if a Load is fetched during Cycle 1:
 - The data is NOT written into the Reg File until the end of Cycle 5
 - We cannot read this value from the Reg File until Cycle 6
 - 3-instruction delay before the load take effect
- This Data Hazard can be reduced to 1 instruction if we:
 - Forward the data from the pipeline register to the next instruction

Outline of Today's Lecture

- Recap and Introduction (5 minutes)
- Memory System: the BIG Picture? (15 minutes)
- Questions and Administrative Matters (5 minutes)
- Memory Technology: SRAM and Register File (25 minutes)
- Break (5 minutes)
- Memory Technology: DRAM (15 minutes)
- A Real Life Example: SPARCstation 20's Memory System (5 minutes)
- Summary (5 minutes)

The Big Picture: Where are We Now?

- The Five Classic Components of a Computer

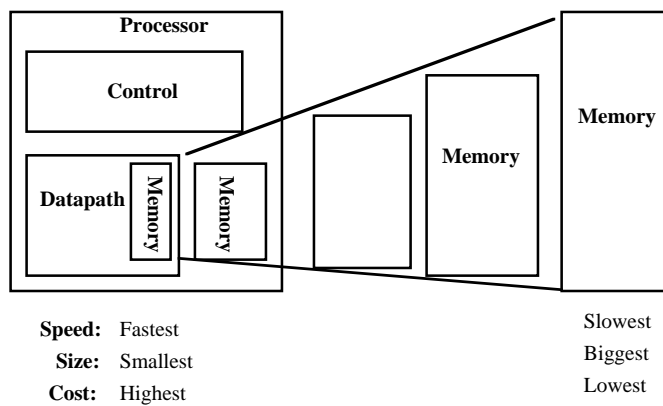


- Today's Topic: Memory System

cs 152 memory.5

©DAP & SIK 1995

An Expanded View of the Memory System



cs 152 memory.6

©DAP & SIK 1995

The Principle of Locality

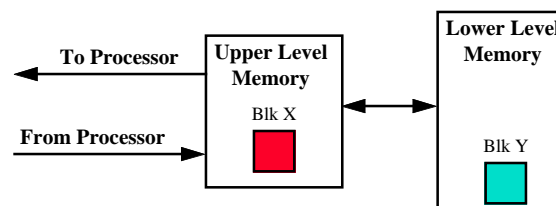
- **The Principle of Locality:**
 - Program access a relatively small portion of the address space at any instant of time.
- **Two Different Types of Locality:**
 - **Temporal Locality (Locality in Time):** If an item is referenced, it will tend to be referenced again soon.
 - **Spatial Locality (Locality in Space):** If an item is referenced, items whose addresses are close by tend to be referenced soon.

cs 152 memory.7

©DAP & SIK 1995

Memory Hierarchy: Principles of Operation

- **At any given time, data is copied between only 2 adjacent levels:**
 - **Upper Level:** the one closer to the processor
 - Smaller, faster, and uses more expensive technology
 - **Lower Level:** the one further away from the processor
 - Bigger, slower, and uses less expensive technology
- **Block:**
 - The minimum unit of information that can either be present or not present in the two level hierarchy

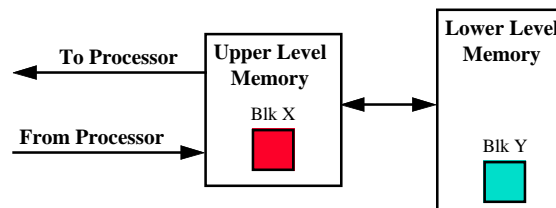


cs 152 memory.8

©DAP & SIK 1995

Memory Hierarchy: Terminology

- **Hit:** data appears in some block in the upper level (example: Block X)
 - **Hit Rate:** the fraction of memory access found in the upper level
 - **Hit Time:** Time to access the upper level which consists of
RAM access time + Time to determine hit/miss
- **Miss:** data needs to be retrieve from a block in the lower level (Block Y)
 - **Miss Rate** = $1 - (\text{Hit Rate})$
 - **Miss Penalty:** Time to replace a block in the upper level +
Time to deliver the block the processor
- **Hit Time** \ll **Miss Penalty**

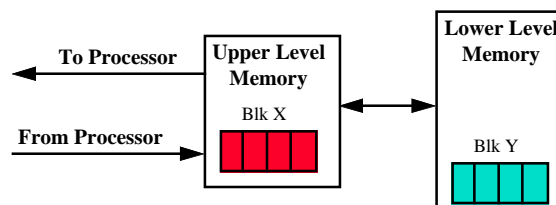


cs 152 memory.9

©DAP & SIK 1995

Memory Hierarchy: How Does it Work?

- **Temporal Locality (Locality in Time):** If an item is referenced, it will tend to be referenced again soon.
 - Keep more recently accessed data items closer to the processor
- **Spatial Locality (Locality in Space):** If an item is referenced, items whose addresses are close by tend to be referenced soon.
 - Move blocks consists of contiguous words to the upper levels

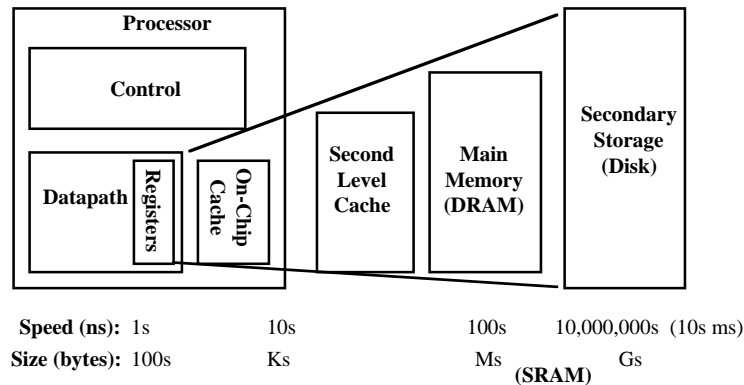


cs 152 memory.10

©DAP & SIK 1995

Memory Hierarchy of a Modern Computer System

- By taking advantage of the principle of locality:
 - Present the user with as much memory as is available in the cheapest technology.
 - Provide access at the speed offered by the fastest technology.



cs 152 memory.11

©DAP & SIK 1995

Memory Hierarchy Technology

- Random Access:
 - “Random” is good: access time is the same for all locations
 - DRAM: Dynamic Random Access Memory
 - High density, low power, cheap, slow
 - Dynamic: need to be “refreshed” regularly
 - SRAM: Static Random Access Memory
 - Low density, high power, expensive, fast
 - Static: content will last “forever”
- “Non-so-random” Access Technology:
 - Access time varies from location to location and from time to time
 - Examples: Disk, tape drive, CDROM
- The next two lectures will concentrate on random access technology
 - The Main Memory: DRAMs
 - Caches: SRAMs

cs 152 memory.12

©DAP & SIK 1995

Questions and Administrative Matters (5 Minutes)

cs 152 memory.13

©DAP & SIK 1995

Random Access Memory (RAM) Technology

- Why do computer designers need to know about RAM technology?
 - Processor performance is usually limited by memory bandwidth
 - As IC densities increase, lots of memory will fit on processor chip
 - Tailor on-chip memory to specific needs
 - Instruction cache
 - Data cache
 - Write buffer
- What makes RAM different from a bunch of flip-flops?
 - Density: RAM is much more denser

cs 152 memory.14

©DAP & SIK 1995

Technology Trends

	Capacity	Speed
Logic:	2x in 3 years	2x in 3 years
DRAM:	4x in 3 years	1.4x in 10 years
Disk:	2x in 3 years	1.4x in 10 years

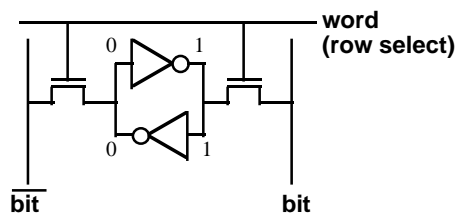
DRAM		
Year	Size	Cycle Time
1980	64 Kb	250 ns
1983	256 Kb	220 ns
1986	1 Mb	190 ns
1989	4 Mb	165 ns
1992	16 Mb	145 ns
1995	64 Mb	120 ns

cs 152 memory.15

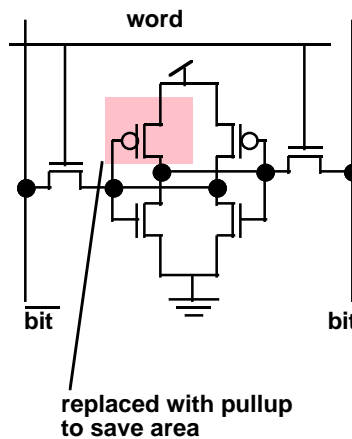
©DAP & SIK 1995

Static RAM Cell

6-Transistor SRAM Cell



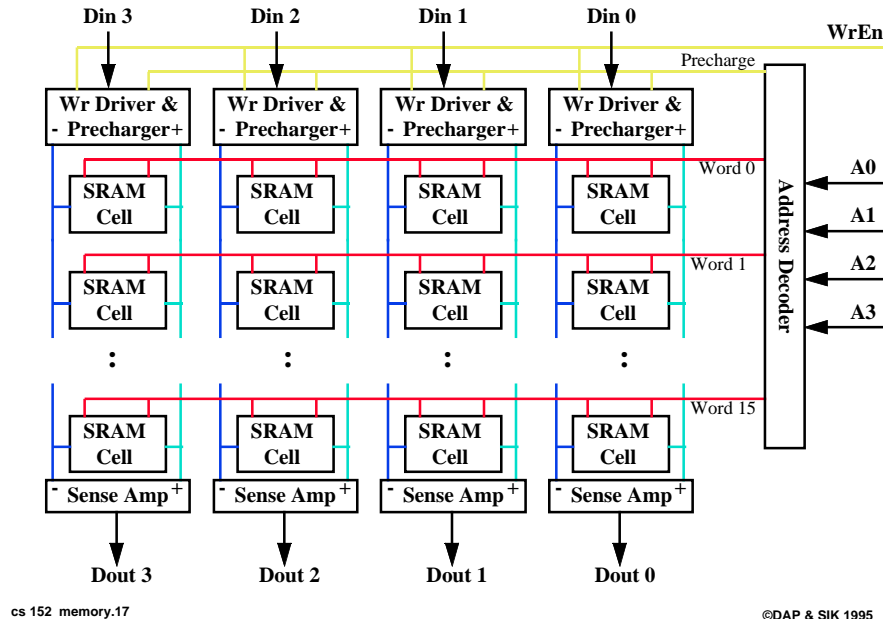
- Write:
 1. Drive bit lines
 - 2.. Select row
- Read:
 1. Precharge bit and bit to Vdd
 - 2.. Select row
 3. Cell pulls one line low
 4. Sense amp on column detects difference



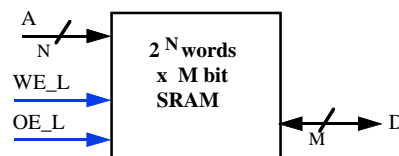
cs 152 memory.16

©DAP & SIK 1995

Typical SRAM Organization: 16-word x 4-bit

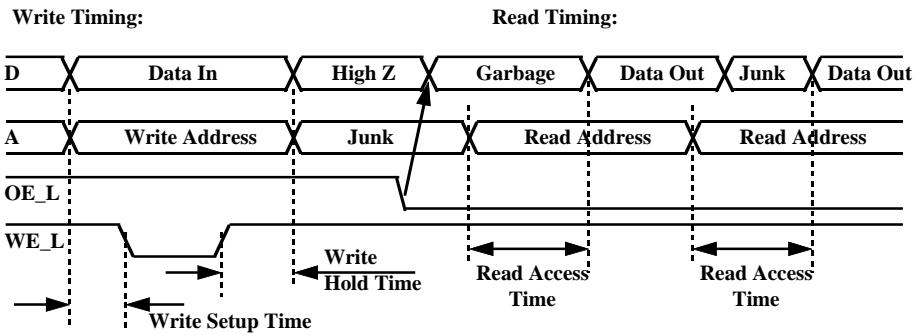
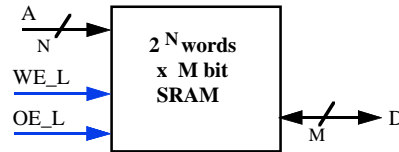


Logic Diagram of a Typical SRAM



- Write Enable is usually active low (WE_L)
- Din and Dout are combined:
 - A new control signal, output enable (OE_L) is needed
 - WE_L is asserted (Low), OE_L is disasserted (High)
 - D serves as the data input pin
 - WE_L is disasserted (High), OE_L is asserted (Low)
 - D is the data output pin
 - Both WE_L and OE_L are asserted:
 - Result is unknown. Don't do that!!!

Typical SRAM Timing

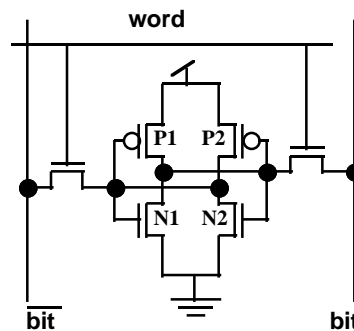
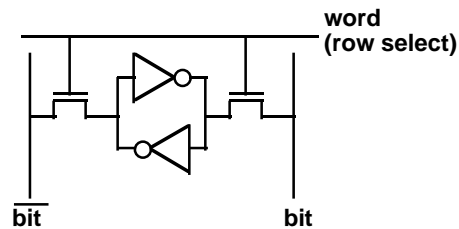


cs 152 memory.19

©DAP & SIK 1995

A Closer Look at the SRAM Cell

6-Transistor SRAM Cell



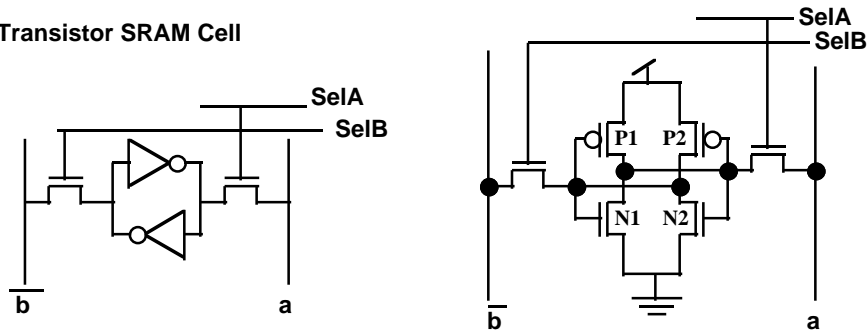
- Typical SRAM has a lot of words (rows)
 - The bit lines are very long and have a lot of capacitance
 - Transistors N1, N2, P1, and P2 must be very small
- Transistors N1 P1 not strong enough to drive “bit” quickly:
 - Need to build a sense amplifier to compare “bit” and “not bit”

cs 152 memory.20

©DAP & SIK 1995

Dual-ported (Read) SRAM Cell for Register File

6-Transistor SRAM Cell

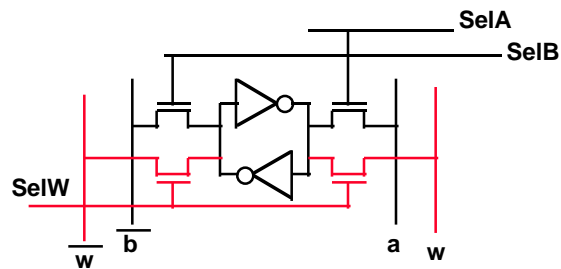


- With short bit lines and larger transistors:
 - N1 and P1 can drive bit line “a” quickly
 - N2 and P2 can drive bit line “not b” quickly
- We can read two words simultaneously

cs 152 memory.21

©DAP & SIK 1995

Single-ported (Write) SRAM Cell for Register File

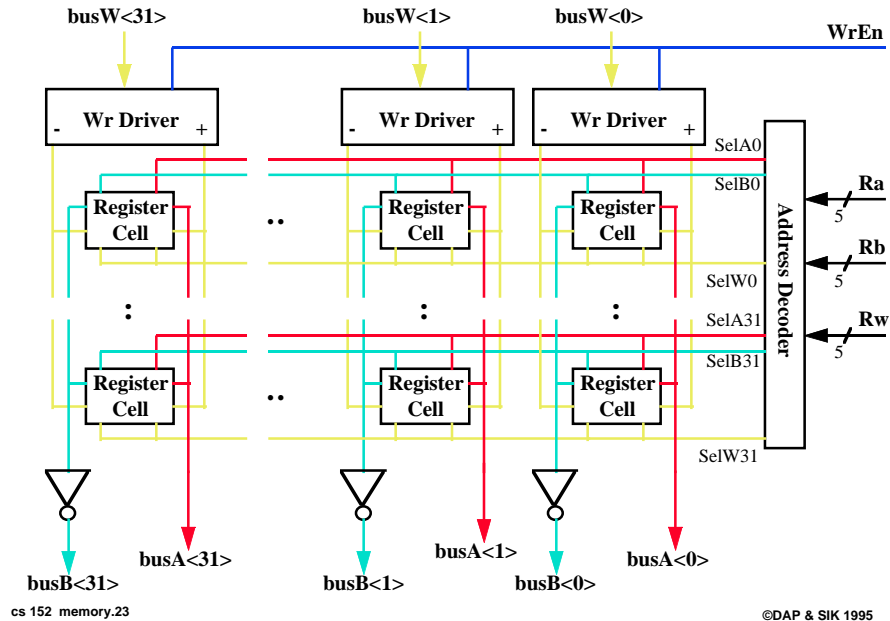


- In order to write a new value into the cell:
 - We need to drive both sides simultaneously
 - We can only write one word at a time
- Extra pair of bit lines (“w” and “not w”)
 - Read and write can occur simultaneously

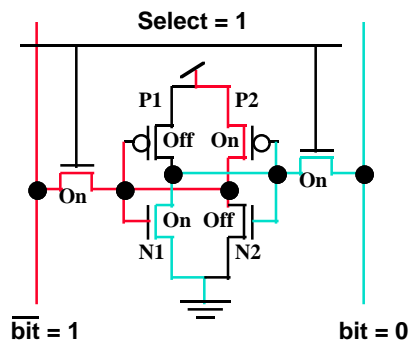
cs 152 memory.22

©DAP & SIK 1995

Dual-ported Read Single-ported Write Register File

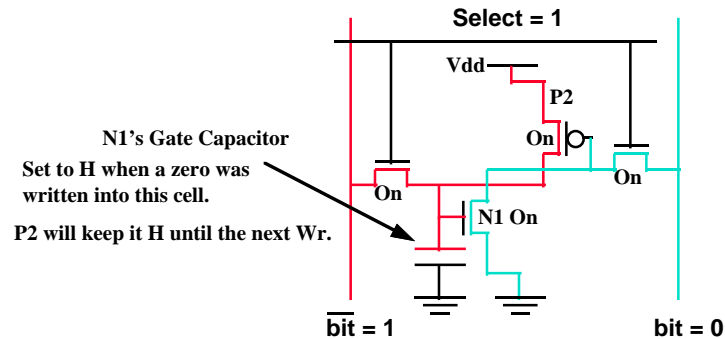


Problems with SRAM



- Six transistors use up a lot of area
- Consider a “Zero” is stored in the cell:
 - Transistor N1 will try to pull “bit” to 0
 - Transistor P2 will try to pull “bit bar” to 1
- But bit lines are precharged to high: Are P1 and P2 necessary?

Problems with SRAM (Continue)



- The P-type transistor (P2) has three functions:
 - Drive the “bit bar” line to HI during read (Select = 1)
 - Keep N1’s gate at HI until the next write
 - Prevent N1’s gate capacitor from “leaking” all its charges to “bit bar” during read

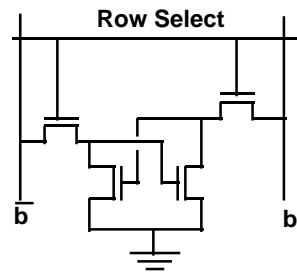
cs 152 memory.25

©DAP & SIK 1995

4-Transistor RAM cell

- Read:
 1. Precharge b and \bar{b} to Vdd
 - 2.. Select row
 3. Sense
 4. Amplify data
 5. Write

} Voltages on the gates converge during read and must be restored.
- Refresh:
 - Dummy read cycle
- Write:
 1. Drive bit lines
 2. Select row



Advantage:
Smaller: Eliminates 2 load devices and 1 power rail

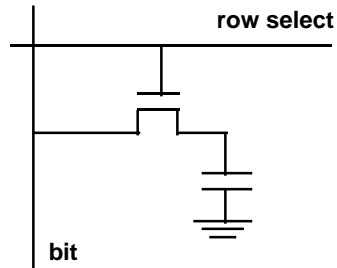
Disadvantages:
Additional refresh cycle
Lower noise margin

cs 152 memory.26

©DAP & SIK 1995

1-Transistor Cell

- Write:
 - 1. Drive bit line
 - 2.. Select row
- Read:
 - 1. Precharge bit line to Vdd
 - 2.. Select row
 - 3. Cell and bit line share charges
 - Very small voltage changes on the bit line
 - 4. Sense (fancy sense amp)
 - Can detect changes of ~1 million electrons
 - 5. Write: restore the value
- Refresh
 - 1. Just do a dummy read to every cell.



cs 152 memory.27

©DAP & SIK 1995

Break (5 Minutes)

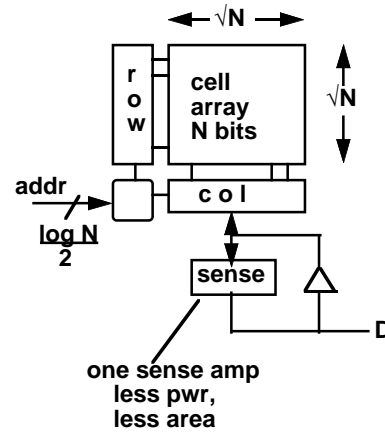
cs 152 memory.28

©DAP & SIK 1995

Introduction to DRAM

- Dynamic RAM (DRAM):

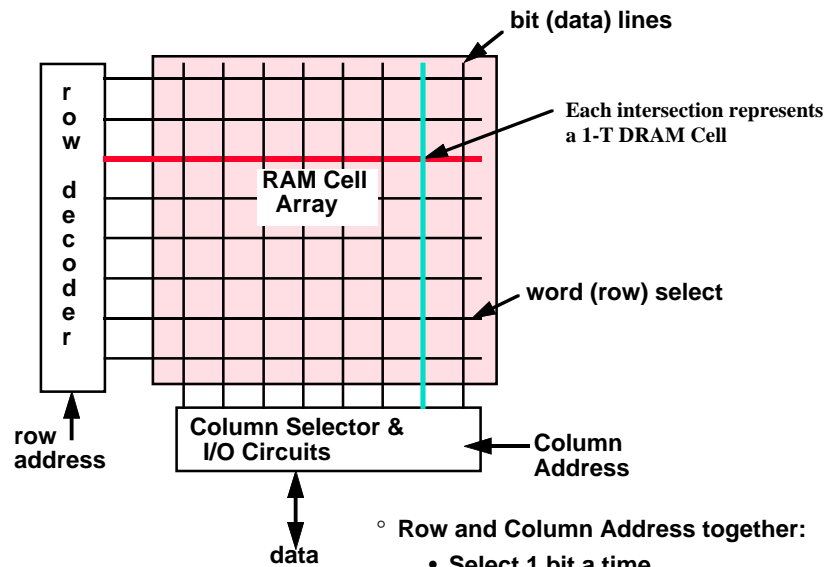
- Refresh required
- Very high density
- Low power (.1 - .5 W active, .25 - 10 mW standby)
- Low cost per bit
- Pin sensitive:
 - Output Enable (OE_L)
 - Write Enable (WE_L)
 - Row address strobe (ras)
 - Col address strobe (cas)
- Page mode operation



cs 152 memory.29

©DAP & SIK 1995

Classical DRAM Organization

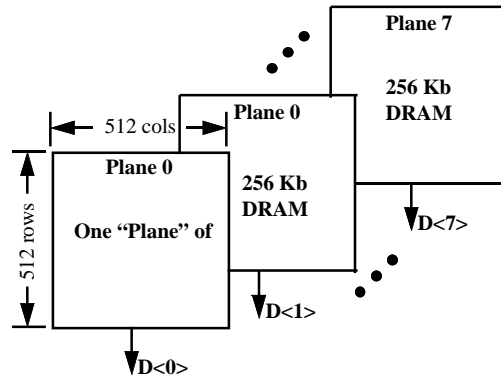


cs 152 memory.30

©DAP & SIK 1995

Typical DRAM Organization

- Typical DRAMs: access multiple bits in parallel
 - Example: 2 Mb DRAM = 256K x 8 = 512 rows x 512 cols x 8 bits
 - Row and column addresses are applied to all 8 planes in parallel

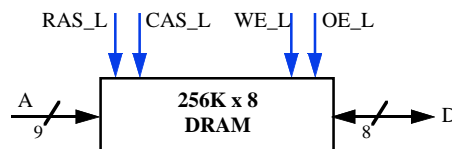


cs 152 memory.31

256 Kb DRAM

©DAP & SIK 1995

Logic Diagram of a Typical DRAM



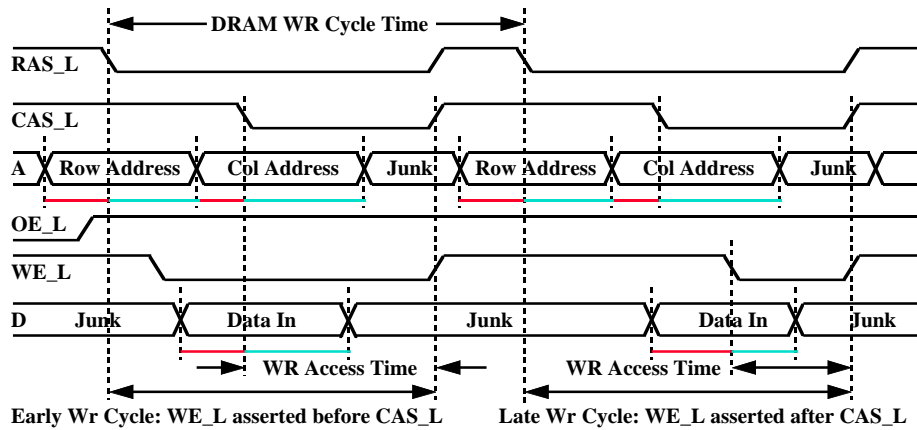
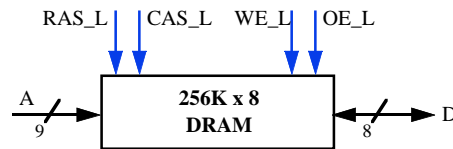
- Control Signals (RAS_L, CAS_L, WE_L, OE_L) are all active low
- Din and Dout are combined (D):
 - WE_L is asserted (Low), OE_L is disasserted (High)
 - D serves as the data input pin
 - WE_L is disasserted (High), OE_L is asserted (Low)
 - D is the data output pin
- Row and column addresses share the same pins (A)
 - RAS_L goes low: Pins A are latched in as row address
 - CAS_L goes low: Pins A are latched in as column address

cs 152 memory.32

©DAP & SIK 1995

DRAM Write Timing

- Every DRAM access begins at:
 - The assertion of the RAS_L

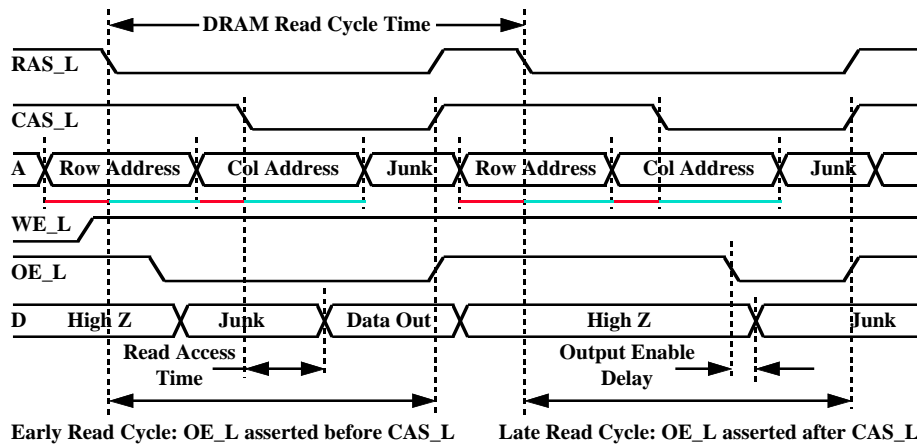
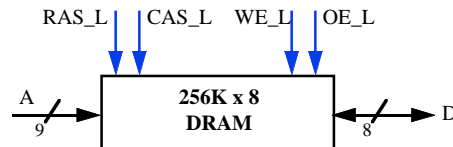


cs 152 memory.33

©DAP & SIK 1995

DRAM Read Timing

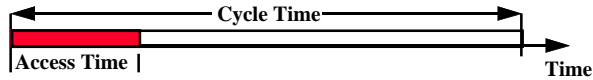
- Every DRAM access begins at:
 - The assertion of the RAS_L



cs 152 memory.34

©DAP & SIK 1995

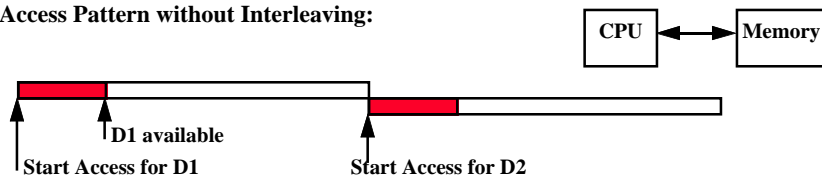
Cycle Time versus Access Time



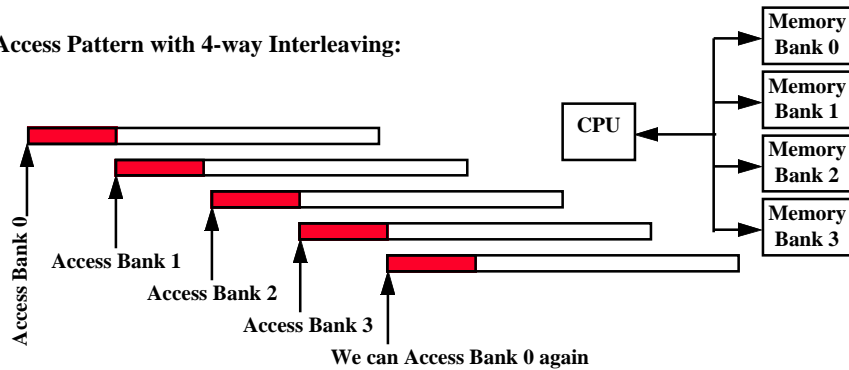
- DRAM (Read/Write) Cycle Time \gg DRAM (Read/Write) Access Time
- DRAM (Read/Write) Cycle Time :
 - How frequent can you initiate an access?
 - Analogy: A little kid can only ask his father for money on Saturday
- DRAM (Read/Write) Access Time:
 - How quickly will you get what you want once you initiate an access?
 - Analogy: As soon as he asks, his father will give him the money
- DRAM Bandwidth Limitation analogy:
 - What happens if he runs out of money on Wednesday?

Increasing Bandwidth - Interleaving

Access Pattern without Interleaving:

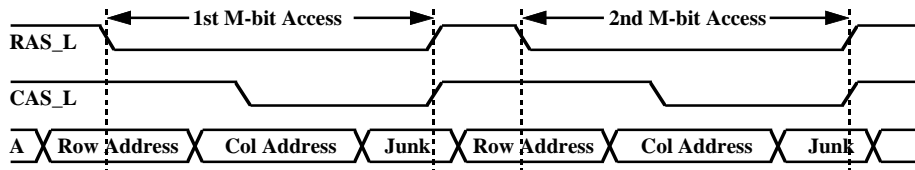
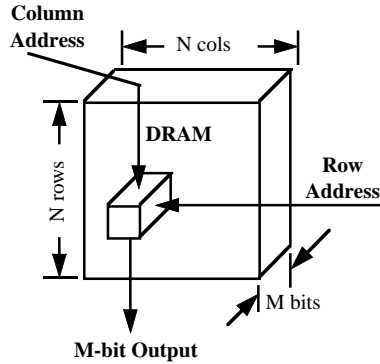


Access Pattern with 4-way Interleaving:



Fast Page Mode DRAM

- Regular DRAM Organization:
 - N rows x N column x M-bit
 - Read & Write M-bit at a time
 - Each M-bit access requires a RAS / CAS cycle
- Fast Page Mode DRAM
 - N x M "register" to save a row

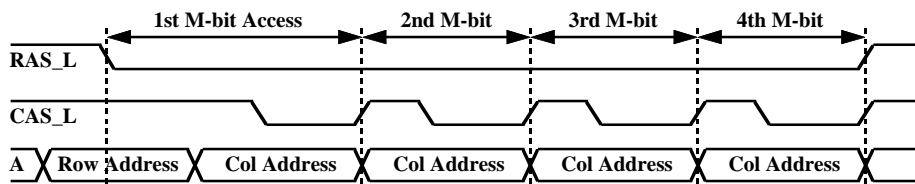
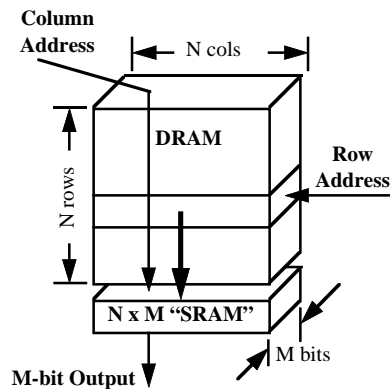


cs 152 memory.37

©DAP & SIK 1995

Fast Page Mode Operation

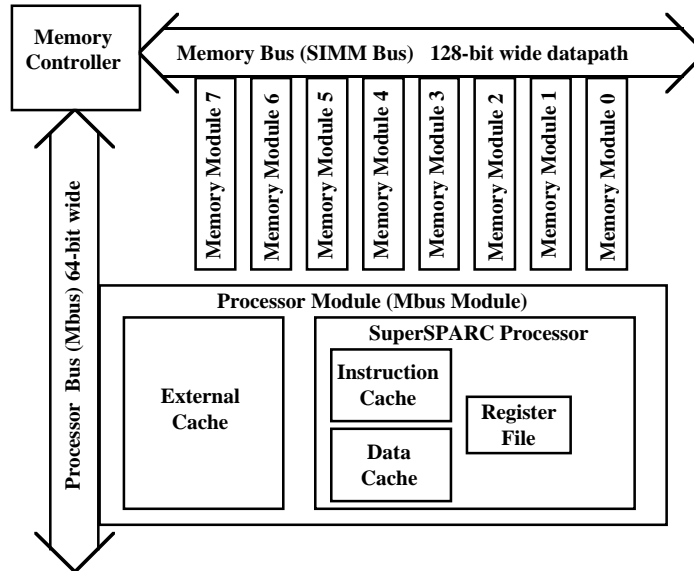
- Fast Page Mode DRAM
 - N x M "SRAM" to save a row
- After a row is read into the register
 - Only CAS is needed to access other M-bit blocks on that row
 - RAS_L remains asserted while CAS_L is toggled



cs 152 memory.38

©DAP & SIK 1995

SPARCstation 20's Memory System Overview

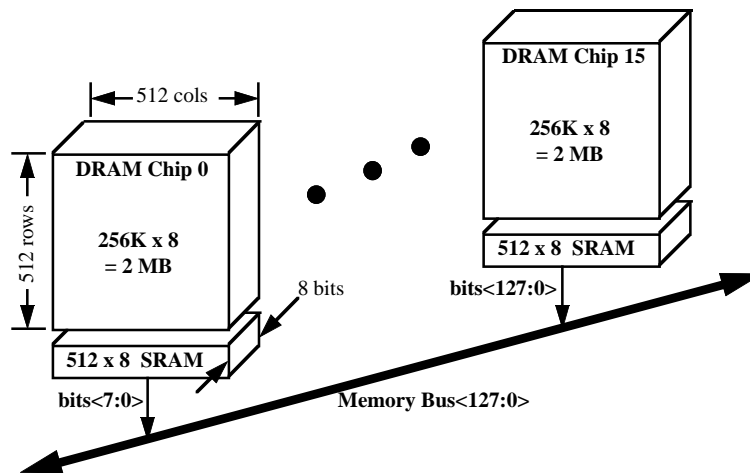


cs 152 memory.39

©DAP & SIK 1995

SPARCstation 20's Memory Module

- Supports a wide range of sizes:
 - Smallest 4 MB: 16 2Mb DRAM chips, 8 KB of Page Mode SRAM
 - Biggest: 64 MB: 32 16Mb chips, 16 KB of Page Mode SRAM

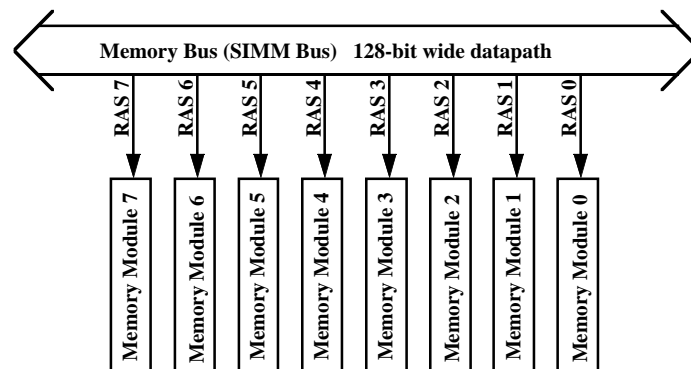


cs 152 memory.40

©DAP & SIK 1995

SPARCstation 20's Main Memory

- **Biggest Possible Main Memory :**
 - 8 64MB Modules: 8 x 64 MB DRAM 8 x 16 KB of Page Mode SRAM
- **How do we select 1 out of the 8 memory modules?**
Remember: every DRAM operation start with the assertion of RAS
 - SS20's Memory Bus has 8 separate RAS lines



cs 152 memory.41

©DAP & SIK 1995

Summary:

- **Two Different Types of Locality:**
 - **Temporal Locality (Locality in Time):** If an item is referenced, it will tend to be referenced again soon.
 - **Spatial Locality (Locality in Space):** If an item is referenced, items whose addresses are close by tend to be referenced soon.
- **By taking advantage of the principle of locality:**
 - Present the user with as much memory as is available in the cheapest technology.
 - Provide access at the speed offered by the fastest technology.
- **DRAM is slow but cheap and dense:**
 - Good choice for presenting the user with a **BIG** memory system
- **SRAM is fast but expensive and not very dense:**
 - Good choice for providing the user **FAST** access time.

cs 152 memory.42

©DAP & SIK 1995

Where to get more information?

- To be continued ...