**361**
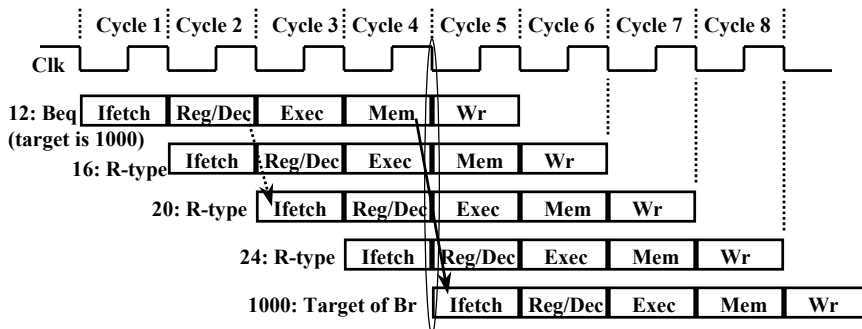**Computer Architecture**
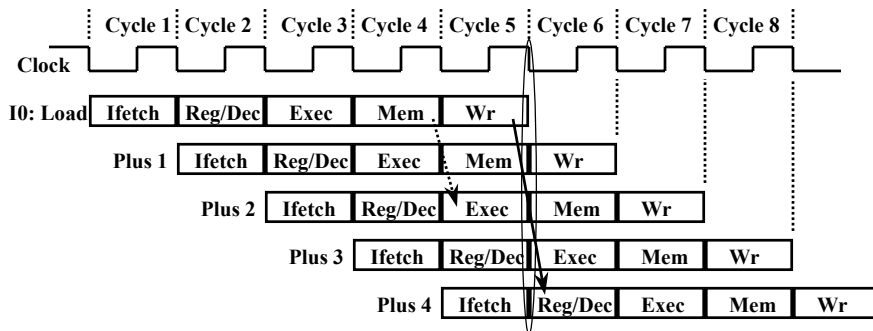**Lecture 16: Memory Systems**

memory.1

---

## 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**

memory.2

## Recap: Solution to Load Hazard



|  | Cycle 1 | Cycle 2 | Cycle 3 | Cycle 4 | Cycle 5 | Cycle 6 | Cycle 7 | Cycle 8 |
|---|---|---|---|---|---|---|---|---|

**Clock**

**I0: Load** Ifetch | Reg/Dec | Exec | Mem | Wr

**Plus 1** Ifetch | Reg/Dec | Exec | Mem | Wr

**Plus 2** Ifetch | Reg/Dec | Exec | Mem | Wr

**Plus 3** Ifetch | Reg/Dec | Exec | Mem | Wr

**Plus 4** Ifetch | Reg/Dec | Exec | Mem | Wr

° **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:**

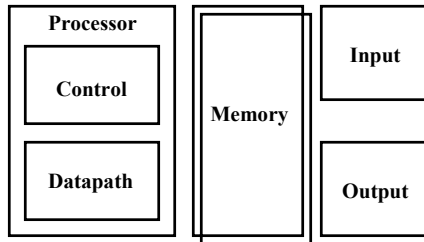- **Forward the data from the pipeline register to the next instruction**

memory.3

---

## Outline of Today's Lecture

° **Recap and Introduction**

° **Memory System: the BIG Picture?**

° **Questions and Administrative Matters**

° **Memory Technology: SRAM**

° **Memory Technology: DRAM**

° **A Real Life Example: SPARCstation 20's Memory System**

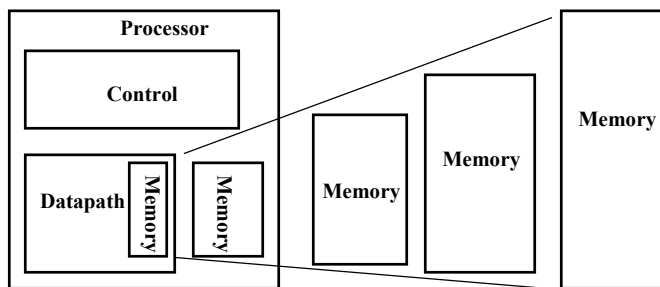° **Summary**

memory.4

# The Big Picture: Where are We Now?

° **The Five Classic Components of a Computer**



° **Today's Topic: Memory System**

# An Expanded View of the Memory System



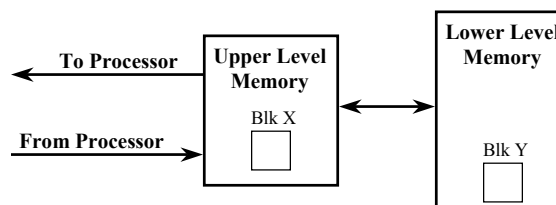| | | |
|---|---|---|
| **Speed:** | Fastest | Slowest |
| **Size:** | Smallest | Biggest |
| **Cost:** | Highest | Lowest |

# 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.**

---

# Memory Hierarchy: Principles of Operation

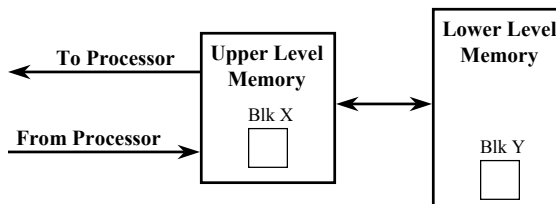° **At any given time, data is copied between only 2 adjacent levels:**
  • **Upper Level: the one closer to the processor**
    - **Smaller, faster, and uses more expensive technology**
  • **Lower Level: the one further away from the processor**
    - **Bigger, slower, and uses less expensive technology**

° **Block:**
  • **The minimum unit of information that can either be present or not present in the two level hierarchy**

To Processor ← | Upper Level Memory — Blk X | ↔ | Lower Level Memory — Blk Y |
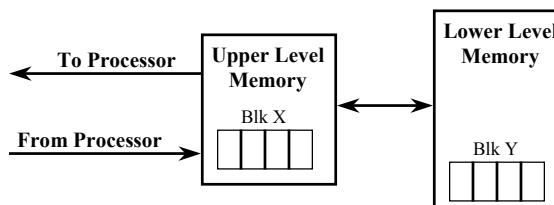
From Processor →

# 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 - (Hit Rate)**
  - **Miss Penalty: Time to replace a block in the upper level +**
    **Time to deliver the block the processor**

° **Hit Time << Miss Penalty**

**To Processor** ← **Upper Level Memory**  Blk X  ↔  **Lower Level Memory**  Blk Y

**From Processor** →

---

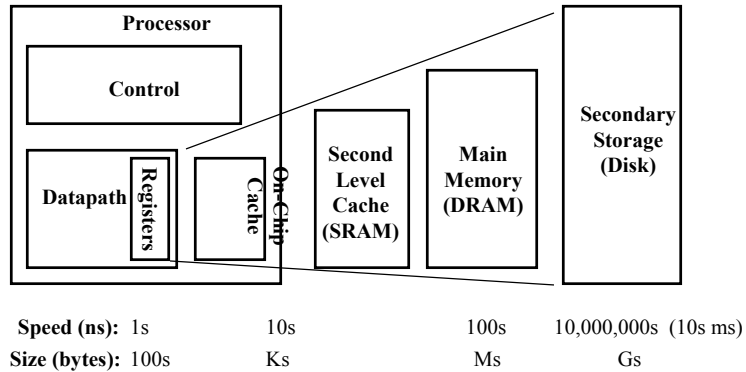# Memory Hierarchy: How Does it Work?

° **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**

**To Processor** ← **Upper Level Memory**  Blk X  ↔  **Lower Level Memory**  Blk Y

**From Processor** →

# 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.**



| | | | |
|---|---|---|---|
| **Speed (ns):** 1s | 10s | 100s | 10,000,000s  (10s ms) |
| **Size (bytes):** 100s | Ks | Ms | Gs |

---

# 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**

# 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**

---

# 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 |

# 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**

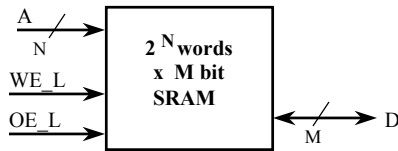    **4. Sense amp on column detects difference**

---

# Typical SRAM Organization: 16-word x 4-bit
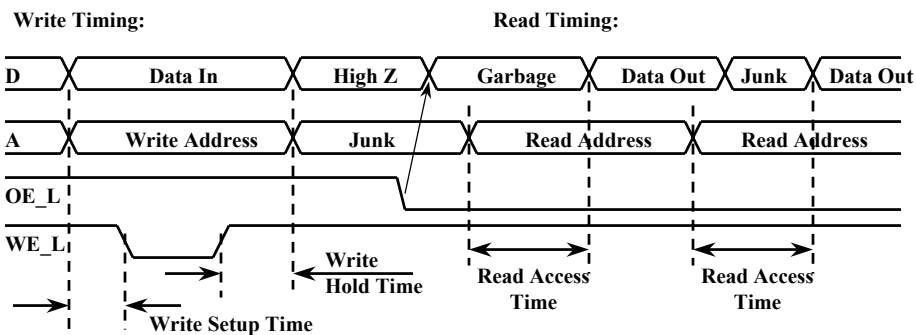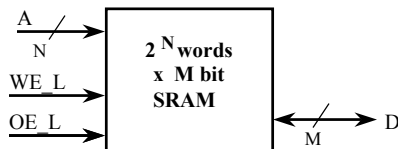
## Logic Diagram of a Typical SRAM

A —/→ N

WE_L →

OE_L →

2 $^N$ words
x M bit
SRAM

←/→ D
M

° **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

A —/→ N

WE_L →

OE_L →

2 $^N$ words
x M bit
SRAM

←/→ D
M

**Write Timing:**                                              **Read Timing:**

| D | Data In | High Z | Garbage | Data Out | Junk | Data Out |

| A | Write Address | Junk | Read Address | | Read Address |

OE_L

WE_L

**Write Hold Time**

**Write Setup Time**

**Read Access Time**

**Read Access Time**

# 1-Transistor Cell

° **Write:**
   • **1. Drive bit line**
   • **2.. Select row**

° **Read:**
   • **1. Precharge bit line to Vdd**
   • **2.. Select row**
   • **3. Sense (fancy sense amp)**
      - **Can detect changes of ~1 million  electrons**
   • **4. Write: restore the value**

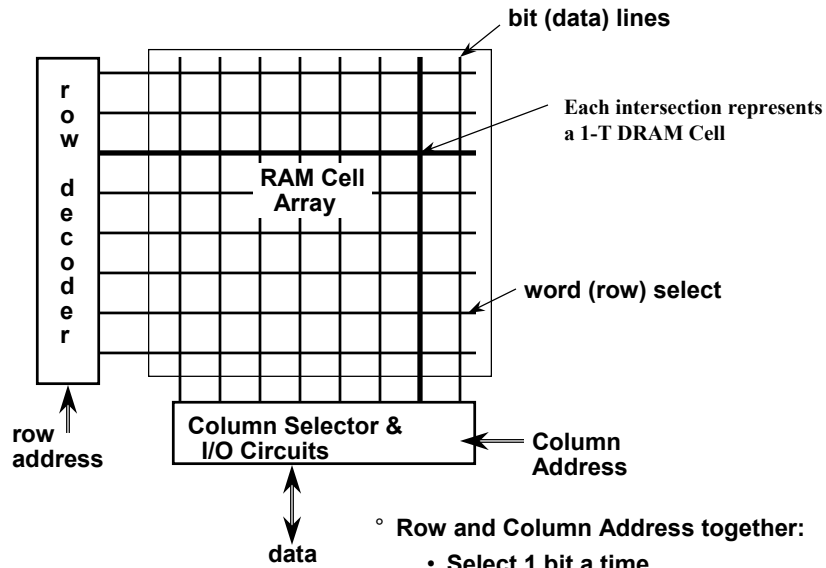° **Refresh**
   • **1. Just do a dummy read to every cell.**



**row select**

**bit**

---

# 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)**
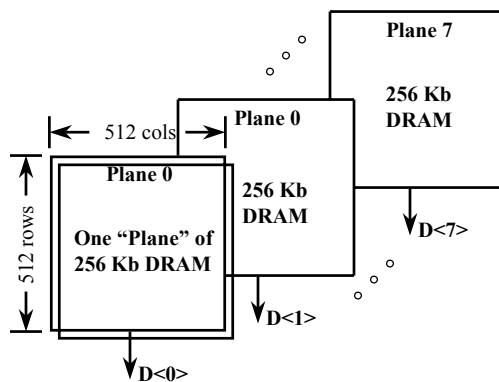      - **Col address strobe (cas)**
   • **Page mode operation**



**N**

**row**

**cell
array
N bits**

**N**

**addr**

**col**

**log N
2**

**sense**

**D**

**one sense amp
less pwr,
less area**

# Classical DRAM Organization



**bit (data) lines**

**Each intersection represents a 1-T DRAM Cell**

**RAM Cell Array**

r o w   d e c o d e r

**word (row) select**

**row address**

**Column Selector & I/O Circuits**

**Column Address**

**data**

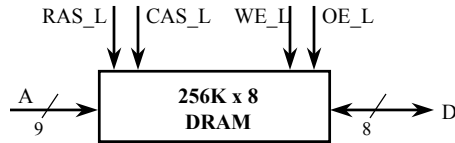° **Row and Column Address together:**
  • **Select 1 bit a time**

---

# 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**



**Plane 7**

**256 Kb DRAM**

**Plane 0**

512 cols

**Plane 0**

**256 Kb DRAM**

512 rows

**One "Plane" of 256 Kb DRAM**

**D<7>**

**D<1>**

**D<0>**

# Logic Diagram of a Typical DRAM

RAS_L    CAS_L    WE_L    OE_L

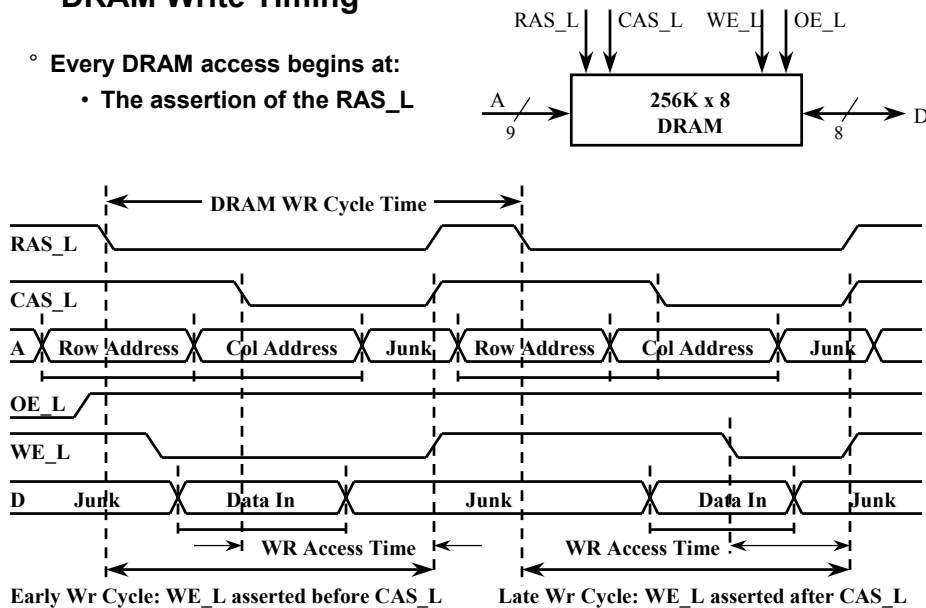A /9 → [ **256K x 8 DRAM** ] ← /8 → D

° **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**

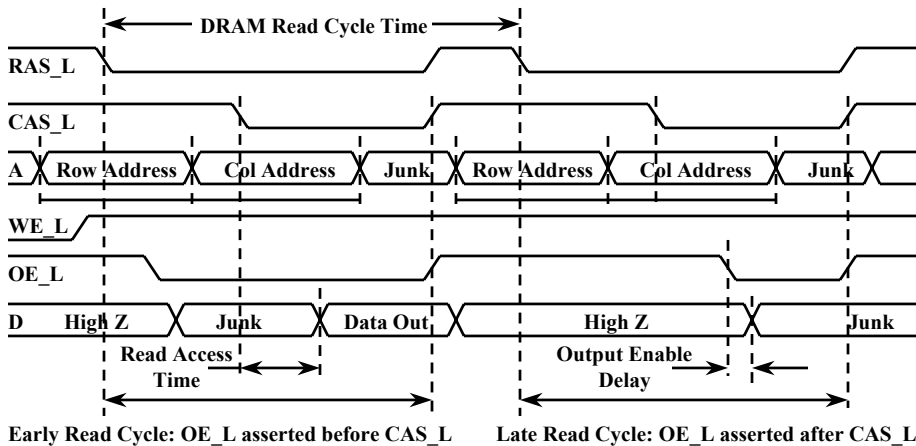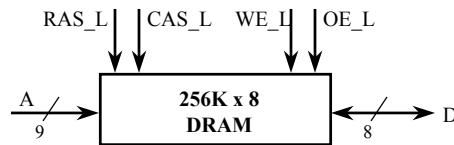memory.23

---

# DRAM Write Timing

° **Every DRAM access begins at:**
  • **The assertion of the RAS_L**

RAS_L    CAS_L    WE_L    OE_L

A /9 → [ **256K x 8 DRAM** ] ← /8 → D

|←——— DRAM WR Cycle Time ———→|

RAS_L

CAS_L

A ⟩⟨ Row Address ⟩⟨ Col Address ⟩⟨ Junk ⟩⟨ Row Address ⟩⟨ Col Address ⟩⟨ Junk ⟩⟨

OE_L

WE_L

D    Junk ⟩⟨ Data In ⟩⟨ Junk ⟩⟨ Data In ⟩⟨ Junk

|→ WR Access Time ←|    WR Access Time

Early Wr Cycle: WE_L asserted before CAS_L          Late Wr Cycle: WE_L asserted after CAS_L

memory.24

## DRAM Read Timing

° **Every DRAM access begins at:**
  • **The assertion of the RAS_L**

RAS_L   CAS_L   WE_L   OE_L

A /9 → | **256K x 8 DRAM** | ↔ /8 → D



RAS_L

CAS_L

A  X Row Address X Col Address X Junk X Row Address X Col Address X Junk X

WE_L

OE_L

D  High Z  X  Junk  X  Data Out  X  High Z  X  Junk

DRAM Read Cycle Time

Read Access Time

Output Enable Delay

**Early Read Cycle: OE_L asserted before CAS_L    Late Read Cycle: OE_L asserted after CAS_L**

---

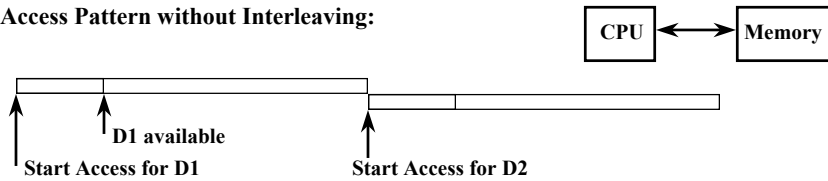## Cycle Time versus Access Time



Cycle Time

Access Time

Time

° **DRAM (Read/Write) Cycle Time  >>  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:**
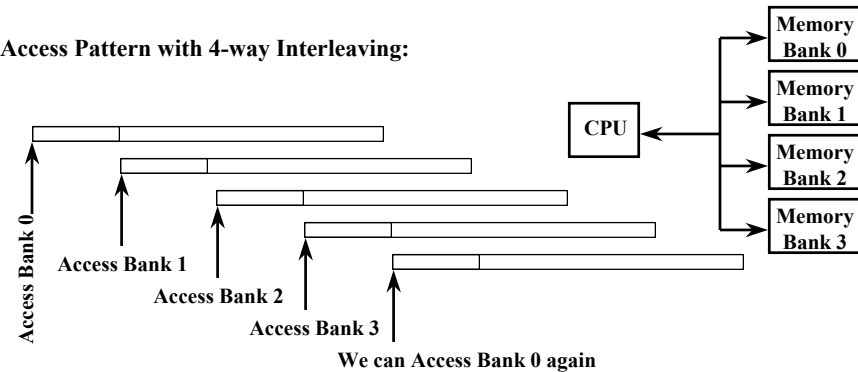  • **What happens if he runs out of money on Wednesday?**

# Increasing Bandwidth - Interleaving

**Access Pattern without Interleaving:**

CPU ◄──► Memory

D1 available

Start Access for D1          Start Access for D2

**Access Pattern with 4-way Interleaving:**

CPU ◄──── Memory Bank 0
          Memory Bank 1
          Memory Bank 2
          Memory Bank 3

Access Bank 0

Access Bank 1

Access Bank 2

Access Bank 3

We can Access Bank 0 again
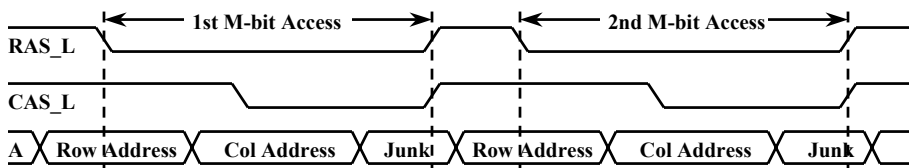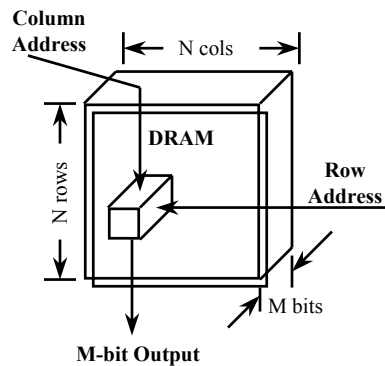
memory.27

---

# Fast Page Mode DRAM

° **Regular DRAM Organization:**
  - **N rows x N column x M-bit**
  - **Read & Write M-bit at a time**
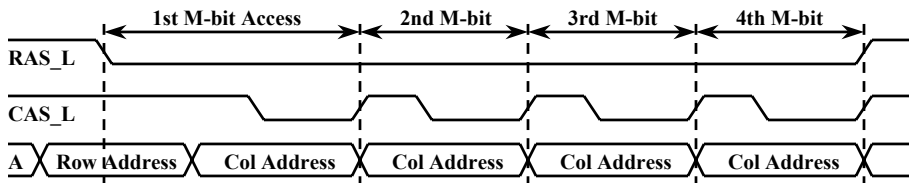  - **Each M-bit access requires a RAS / CAS cycle**

° **Fast Page Mode DRAM**
  - **N x M "register" to save a row**

Column Address

N cols

DRAM

N rows

Row Address

M bits

M-bit Output

| | 1st M-bit Access | | 2nd M-bit Access | |
|---|---|---|---|---|

RAS_L

CAS_L

A ── Row Address ── Col Address ── Junk ── Row Address ── Col Address ── Junk ──
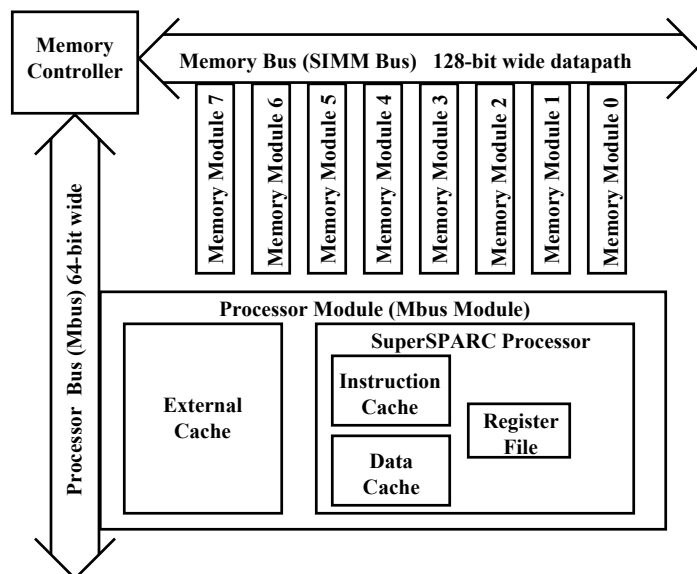
memory.28

# 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**
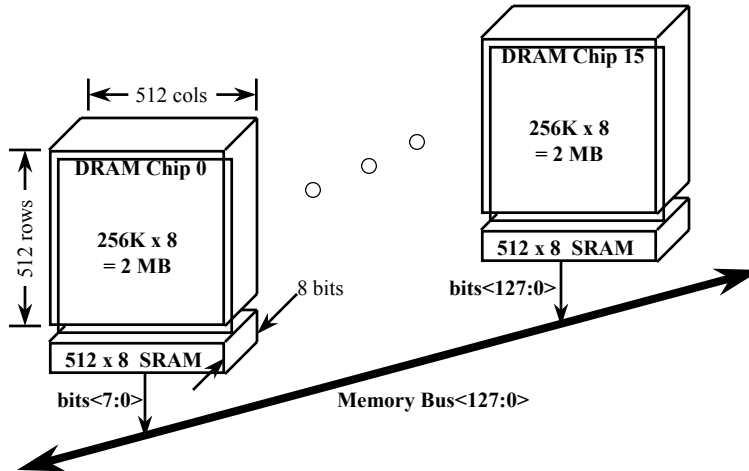
 • **RAS_L remains asserted while CAS_L is toggled**

**Column Address**

N cols

**DRAM**

N rows

**Row Address**

**N x M "SRAM"**

M bits

**M-bit Output**

| 1st M-bit Access | 2nd M-bit | 3rd M-bit | 4th M-bit |
|---|---|---|---|

RAS_L

CAS_L

A   Row Address   Col Address   Col Address   Col Address   Col Address

memory.29

---

# SPARCstation 20's Memory System Overview

**Memory Controller**

**Memory Bus (SIMM Bus)   128-bit wide datapath**

Memory Module 7
Memory Module 6
Memory Module 5
Memory Module 4
Memory Module 3
Memory Module 2
Memory Module 1
Memory Module 0

**Processor Bus (Mbus) 64-bit wide**

**Processor Module (Mbus Module)**

**SuperSPARC Processor**

**External Cache**

**Instruction Cache**

**Register File**

**Data Cache**

memory.30

# SPARCstation 20's Memory Module

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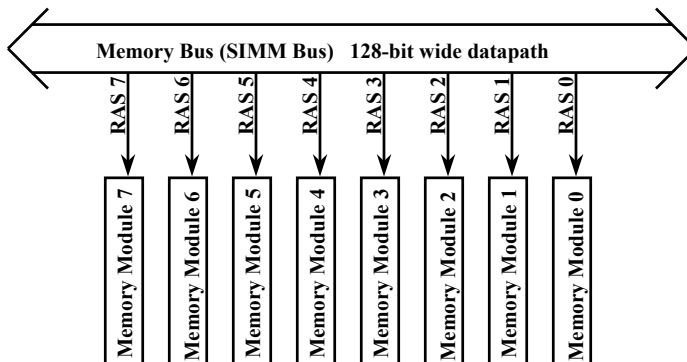  - **Biggest: 64 MB: 32 16Mb chips, 16 KB of Page Mode SRAM**

---

# 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**

# 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.**

# Where to get more information?

° **To be continued ...**