Many of the following slides are taken with permission from

### Complete Powerpoint Lecture Notes for Computer Systems: A Programmer's Perspective (CS:APP)

Randal E. Bryant and David R. O'Hallaron

http://csapp.cs.cmu.edu/public/lectures.html

The book is used explicitly in CS 2505 and CS 3214 and as a reference in CS 2506.

# An Example Memory Hierarchy



**Computer Organization II** 

### Key features

- **RAM** is traditionally packaged as a chip.
- Basic storage unit is normally a cell (one bit per cell).
- Multiple RAM chips form a memory.

### Static RAM (SRAM)

- Each cell stores a bit with a four or six-transistor circuit.
- Retains value indefinitely, as long as it is kept powered.
- Relatively insensitive to electrical noise (EMI), radiation, etc.
- Faster and more expensive than DRAM.

### Dynamic RAM (DRAM)

- Each cell stores bit with a capacitor. One transistor is used for access
- Value must be refreshed every 10-100 ms.
- More sensitive to disturbances (EMI, radiation,...) than SRAM.
- Slower and cheaper than SRAM.

A bus is a collection of parallel wires that carry address, data, and control signals.

Buses are typically shared by multiple devices.



# Memory Read Transaction (1)



CS@VT

**Computer Organization II** 

Main memory reads A from the memory bus, retrieves word x, and places it on the bus.



CS@VT

### Memory Read Transaction (3)



**Computer Organization II** 

# Memory Write Transaction (1)

Memory Hierarchy 8

CPU places address A on bus. Main memory reads it and waits for the corresponding data word to arrive.



CS@VT

**Computer Organization II** 

# Memory Write Transaction (2)



**Computer Organization II** 

## Memory Write Transaction (3)



CS@VT

**Computer Organization II** 

# The Bigger Picture: I/O Bus



# Storage Trends

#### SRAM

Metric	1980	1985	1990	1995	2000	2005	2010	2010:1980
\$/MB	19,200	2,900	320	256	100	75	60	320
access (ns)	300	150	35	15	3	2	1.5	200

#### DRAM

Metric	1980	1985	1990	1995	2000	2005	2010	2010:1980
\$/MB	8,000	880	100	30	1	0.1	0.06	130,000
access (ns)	375	200	100	70	60	50	40	9
typical size (MB)	0.064	0.256	4	16	64	2,000	8,000	125,000

#### Disk

Metric	1980	1985	1990	1995	2000	2005	2010	2010:1980
\$/MB	500	100	8	0.30	0.01	0.005	0.0003	1,600,000
access (ms)	87	75	28	10	8	4	3	29
typical size (MB)	1	10	160	1,000	20,000	160,000	1,500,000	1,500,000

### The CPU-Memory Gap

### The gap widens between DRAM, disk, and CPU speeds.



CS@VT

**Computer Organization II** 

### Locality

Principle of Locality: Programs tend to use data and instructions with addresses near or equal to those they have used recently

### **Temporal locality:**

 Recently referenced items are likely to be referenced again in the near future



### **Spatial locality:**

 Items with nearby addresses tend to be referenced close together in time



### Locality Example

```
sum = 0;
for (i = 0; i < n; i++)
    sum += a[i];
return sum;
```

Data references

- Reference array elements in succession (stride-1 reference pattern).
- Reference variable sum each iteration.

Spatial locality Temporal locality

#### Instruction references

- Reference instructions in sequence.
- Cycle through loop repeatedly.

Spatial locality Temporal locality

Use the characteristics of the memory hierarchy

large and slow vs small and fast

Store everything on disk or SSD (virtual memory)

Copy recently accessed (and nearby) items from disk to smaller DRAM memory

Copy more recently accessed (and nearby) items from DRAM to even smaller SRAM memory

# Is Locality There to Be Exploited?

Locality depends on a number of factors:

Developer's choice of data organization

Developer's choice of algorithms

Developer's other coding decisions

Developer's choice of programming language

Organization of memory hardware

And... the locality exhibited by a process tends to change during execution...

# An Example Memory Hierarchy



**Computer Organization II** 

### Caches

*Cache:* a smaller, faster storage device that acts as a staging area for a subset of the data in a larger, slower device.

Fundamental idea of a memory hierarchy:

 For each k, the faster, smaller device at level k serves as a cache for the larger, slower device at level k+1.

Why do memory hierarchies work?

- Because of locality, programs tend to access the data at level k more often than they access the data at level k+1.
- Thus, the storage at level k+1 can be slower, and thus larger and cheaper per bit.

*Big Idea:* The memory hierarchy creates a large pool of storage that costs as much as the cheap storage near the bottom, but that serves data to programs at the rate of the fast storage near the top.

### **General Cache Concepts**



**Computer Organization II** 

### General Cache Concepts: Hit



**Computer Organization II** 

### **General Cache Concepts: Miss**



### Cold (compulsory) miss

- Cold misses occur because the cache is empty.
- ... or because we are accessing something for the first time.

### **Conflict miss**

- Most caches limit blocks at level k+1 to a small subset (sometimes a singleton) of the block positions at level k.
  - E.g. Block i at level k+1 must be placed in block (i mod 4) at level k.
- Conflict misses occur when the level k cache is large enough, but multiple data objects all map to the same level k block.
  - E.g. Referencing blocks 0, 8, 0, 8, 0, 8, ... would miss every time.

#### Capacity miss

Occurs when the set of active cache blocks (working set) is larger than the cache.

Cache Type	What is Cached?	Where is it Cached?	Latency (cycles)	Managed By
Registers	4-8 bytes words	CPU core	0	Compiler
TLB	Address translations	On-Chip TLB	0	Hardware
L1 cache	64-bytes block	On-Chip L1	1	Hardware
L2 cache	64-bytes block	On/Off-Chip L2	10	Hardware
Virtual Memory	4-KB page	Main memory	100	Hardware + OS
Buffer cache	Parts of files	Main memory	100	OS
Disk cache	Disk sectors	Disk controller	100,000	Disk firmware
Network buffer cache	Parts of files	Local disk	10,000,000	AFS/NFS client
Browser cache	Web pages	Local disk	10,000,000	Web browser
Web cache	Web pages	Remote server disks	1,000,000,000	Web proxy server

## **Cache Memories**

Cache memories are small, fast SRAM-based memories managed automatically in hardware.

- Hold frequently accessed blocks of main memory

CPU looks first for data in caches (e.g., L1, L2, and L3), then in main memory.

Typical system structure:

