

Memory Hierarchy & Data Locality

CSE 373
Data Structures & Algorithms
Ruth Anderson

11/22/2013

1

Why do we need to know about the memory hierarchy/locality?

- One of the assumptions that Big-Oh makes is that *all operations take the same amount of time.*
- Is that really true?

11/22/2013

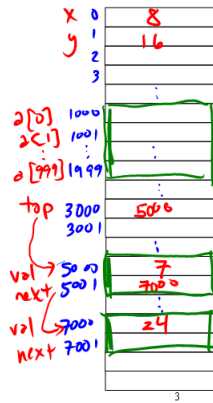
2

Where are these values in memory?

```
int x = 8;
int y = 2 * x;

int[] a = new int[1000];
z = a[0] + a[1] + a[999];

ListNode top = new ListNode(7);
top.next = new ListNode(24);
ListNode temp = top.next;
```



11/22/2013

3

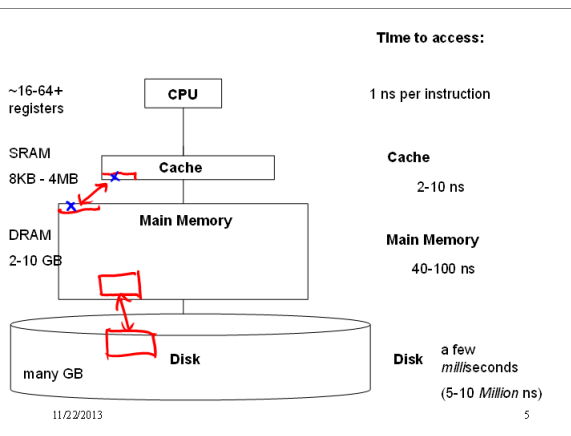
Definitions

Cycle – (for our purposes) the time it takes to execute a single simple instruction. (ex. Add 2 registers together)

Memory Latency – time it takes to access memory

11/22/2013

4



11/22/2013

5

Morals

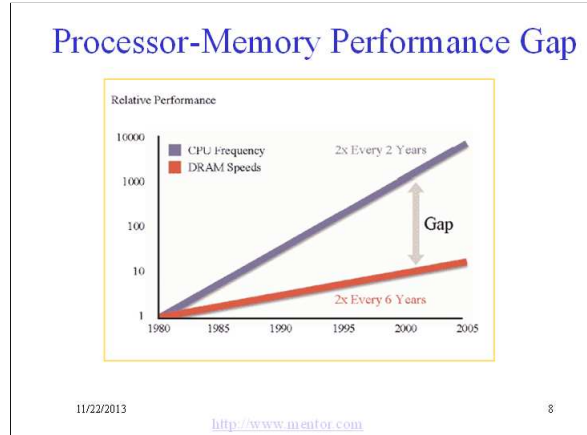
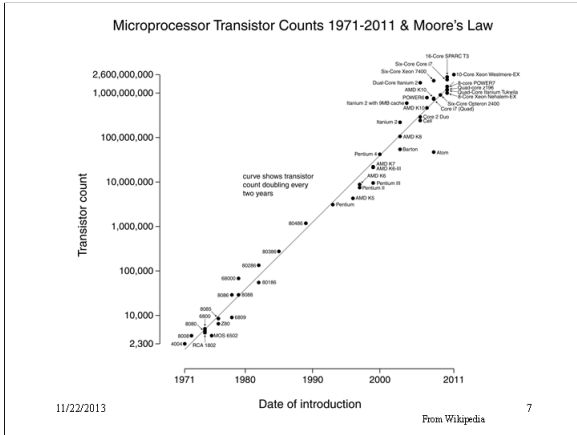
It is much faster to do:	Than:
5 million arithmetic ops	1 disk access
2500 L2 cache accesses	1 disk access
400 main memory accesses	1 disk access

Why are computers built this way?

- Physical realities (speed of light, closeness to CPU)
- Cost (price per byte of different technologies)
- Disks get much bigger not much faster
 - Spinning at 7200 RPM accounts for much of the slowness and unlikely to spin faster in the future
- Speedup at higher levels (e.g. a faster processor) makes lower levels *relatively slower*. Argh!

11/22/2013

6



What can be done?

- **Goal:** Attempt to reduce the number of accesses to the slower levels.
- How?

11/22/2013 9

So, what can we do?

The hardware automatically moves data into the caches from main memory for you

- Replacing items already there
- Algorithms are much faster if "data fits in cache" (often does)

Disk accesses are done by software (e.g., ask operating system to open a file or database to access some data)

So most code "just runs" but sometimes it's worth designing algorithms / data structures with knowledge of memory hierarchy

- And when you do, you often need to know one more thing...

11/22/2013 10

Locality

Temporal Locality (locality in time) – If an item (a location in memory) is referenced, that same location will tend to be referenced again soon.

Spatial Locality (locality in space) – If an item is referenced, items whose addresses are close by will tend to be referenced soon.

11/22/2013 11

How does data move up the hierarchy?

- Moving data up the memory hierarchy is slow because of *latency* (think distance-to-travel)
 - Since we're making the trip anyway, may as well carpool
 - Get a block of data in the same time it would take to get a byte
 - Sends *nearby memory* because: Spatial Locality
 - It's easy
 - Nearby memory is likely to be asked for soon (think fields/arrays)
- Side note: Once a value is in cache, may as well keep it around for awhile; accessed once, a value is more likely to be accessed again in the near future (more likely than some random other value) Temporal locality

11/22/2013 12

Cache Facts

- Each level is a **sub-set** of the level below.

Definitions:

- Cache Hit** – address requested is in cache
- Cache Miss** – address requested is NOT in cache
- Block or Page size** - the number of contiguous bytes moved from **disk** into **memory**
- Cache line size** - the number of contiguous bytes moved from **memory** into **cache**

11/22/2013

13

Examples

```
x = a + 6; miss    x = a[0] + 6; miss
y = a + 5; hit    y = a[1] + 5; hit
z = 8 * a; hit    z = 8 * a[2]; hit
```

temporal locality

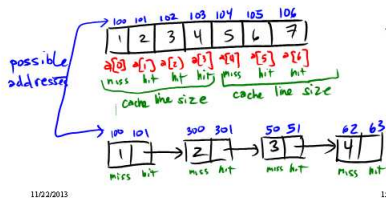
spatial locality

11/22/2013

14

Locality and Data Structures

- Which has (at least the potential for) better spatial locality, arrays or linked lists?



11/22/2013

e.g. traversing elements in a list/array in order

miss on first item in a cache line

miss on every item (unless more than one item randomly happened to be on the same cache line)

15

Where is the Locality?

```
for (i = 1; i < 100; i++) {
    a = a * 7;
    b = b + x[i];
    c = y[5] + d;
}
```

temporal locality

spatial on locations in array x

11/22/2013

16