


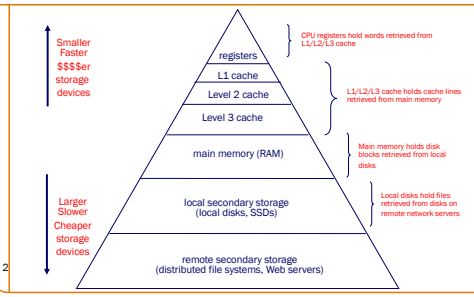
# COS 217: Introduction to Programming Systems

## Storage Hierarchy, Caching, and Locality



1

## Typical Storage Hierarchy



↑ Smaller Faster \$\$\$ storage devices

↓ Larger Slower Cheaper storage devices

- registers: CPU registers hold words retrieved from L1/L2/L3 cache
- L1 cache: L1/L2/L3 cache holds cache lines retrieved from main memory
- Level 2 cache
- Level 3 cache
- main memory (RAM): Main memory holds disk blocks retrieved from local disks
- local secondary storage (local disks, SSDs): Local disks hold files retrieved from disks on remote network servers
- remote secondary storage (distributed file systems, Web servers)

2

## Typical Storage Hierarchy

Factors to consider:


- Capacity
- Latency (how long to do a read)
- Bandwidth (how many bytes/sec can be read)
  - Weakly correlated to latency: reading 1 MB from a hard disk isn't much slower than reading 1 byte
- Volatility
  - Do data persist in the absence of power?

3

## Typical Storage Hierarchy

Registers

- Latency: 0 cycles
- Capacity: 8-256 registers (31 general purpose registers in AArch64)




L1/L2/L3 Cache

- Latency: 1 to 40 cycles
- Capacity: 32KB to 32MB

Main memory (RAM)

- Latency: ~ 50-100 cycles
- 100 times slower than registers
- Capacity: GB





@christianw, @harrisonbroadbent

4

## Typical Storage Hierarchy

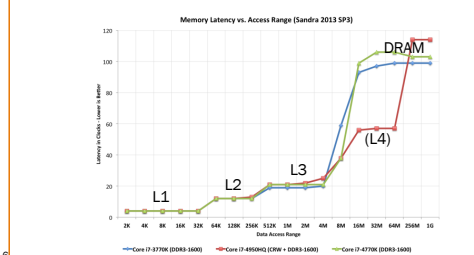
Local secondary storage: disk drives

- Solid-State Disk (SSD):
  - Flash memory (nonvolatile)
  - Latency: 0.1 ms (~ 300k cycles)
  - Capacity: 128 GB - 2 TB
- Hard Disk:
  - Spinning magnetic platters, moving heads
  - Latency: 10 ms (~ 30M cycles)
  - Capacity: 1 - 10 TB

@benjaminiskman, Samsung Belgium

5

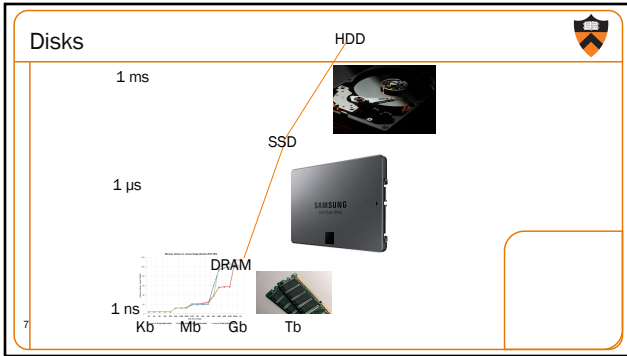
## Cache / RAM Latency



Memory Latency vs. Access Range (Sandra 2013 SP3)

1 clock =  $3 \cdot 10^{-10}$  sec <http://www.anandtech.com/show/9563/real-time-cpu-DNAX-graphics-performance-7-950-the-king-2>

6



7

### Typical Storage Hierarchy

Remote secondary storage (a.k.a. "the cloud")

- **Latency:** tens of milliseconds
- Limited by network bandwidth
- **Capacity:** essentially unlimited

@TheDigitalArtist

8

### Storage Device Speed vs. Size

Facts:

- CPU needs sub-nanosecond access to data to run instructions at full speed
- **Fast** storage (sub-nanosecond) is small (100-1000 bytes)
- **Big** storage (gigabytes) is slow (15 nanoseconds)
- **Huge** storage (terabytes) is *glacially* slow (milliseconds)

Goal:

- Need many *g*igabytes of memory,
- but with fast (sub-nanosecond) average access time

Solution: **locality** allows **caching**

- Most programs exhibit good **locality**
- A program that exhibits good **locality** will benefit from proper **caching**, which enables good **average** performance

9

### Locality

Two kinds of **locality**

- **Temporal** locality
  - If a program references item X now, it probably will reference X again soon
- **Spatial** locality
  - If a program references item X now, it probably will reference item at address  $X \pm 1$  soon

Most programs exhibit good temporal and spatial locality

10

### Locality Example

Locality example

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

Typical code (good locality)

- **Temporal locality**
  - *Data:* Whenever the CPU accesses `sum`, it accesses `sum` again shortly thereafter
  - *Instructions:* Whenever the CPU executes `sum += a[i]`, it executes `sum += a[i]` again shortly thereafter
- **Spatial locality**
  - *Data:* Whenever the CPU accesses `a[i]`, it accesses `a[i+1]` shortly thereafter
  - *Instructions:* Whenever the CPU executes `sum += a[i]`, it executes `i++` shortly thereafter

11

### Caching

**Cache**

- Fast access, small capacity storage device
- Acts as a staging area for a subset of the items in a slow access, large capacity storage device

Good locality + proper caching  
 ⇒ Most storage accesses can be satisfied by cache  
 ⇒ Overall storage performance improved

12

### Caching in a Storage Hierarchy

Level k: [ 4 | 9 | 10 | 3 ]

Smaller, faster device at level k caches a subset of the blocks from level k+1

Blocks copied between levels

Level k+1:

0	1	2	3
4	5	6	7
8	9	10	11
12	13	14	15

Larger, slower device at level k+1 is partitioned into blocks

13

13

### Cache Hits and Misses

Level k: [ 4 | 9 | 10 | 3 ]

Level k is a cache for level k+1

Level k+1:

0	1	2	3
4	5	6	7
8	9	10	11
12	13	14	15

14

14

- Cache hit**
- E.g., request for block 10
  - Access block 10 at level k
  - Fast!
- Cache miss**
- E.g., request for block 8
  - **Evict** some block from level k
  - Load block 8 from level k+1 to level k
  - Access block 8 at level k
  - Slow!
- Caching goal:**
- Maximize cache hits
  - Minimize cache misses



### Cache Eviction Policies

**Best** eviction policy: **"oracle"**

- Always evict a block that is *never* accessed again, or...
- Always evict the block accessed the *furthest in the future*
- Impossible in the general case

**Worst** eviction policy

- Always evict the block that will be accessed next!
- Causes **thrashing**
- Impossible in the general case!

15

15

### Cache Eviction Policies

**Reasonable** eviction policy: **LRU policy**

- Evict the "Least Recently Used" (LRU) block
- With the assumption that it will not be used again (soon)
- Good for straight-line code
- (can be) bad for (large) loops
- Expensive to implement
- Often simpler approximations are used
- See Wikipedia "Page replacement algorithm" topic

16

16

### Do Exam Questions Exhibit Temporal Locality?

Here's a real question from an old exam:  
For caching in a memory hierarchy,  
what is the best motivation for a *larger* cache block size?

A. Temporal Locality                      B. Spatial Locality makes use of subsequent data after a given read, so having more data to keep reading is a win.

C. Both

D. Neither

17

17

### Locality/Caching Example: Matrix Multiplication

**Matrix multiplication**

- Matrix = two-dimensional array
- Multiply n-by-n matrices A and B
- Store product in matrix C

**Performance depends upon**

- Effective use of caching (as implemented by **system**)
- Good locality (as implemented by **you**)

18

18

### Locality/Caching Example: Matrix Mult

Two-dimensional arrays are stored in either **row-major** or **column-major** order

	0	1	2
0	18	19	20
1	21	22	23
2	24	25	26

row-major	col-major
a[0][0] 18	a[0][0] 18
a[0][1] 19	a[1][0] 21
a[0][2] 20	a[2][0] 24
a[1][0] 21	a[0][1] 19
a[1][1] 22	a[1][1] 22
a[1][2] 23	a[2][1] 25
a[2][0] 24	a[0][2] 20
a[2][1] 25	a[1][2] 23
a[2][2] 26	a[2][2] 26

C uses **row-major** order

- Access in row-major order ⇒ good spatial locality
- Access in column-major order ⇒ poor spatial locality

19

19

### Locality/Caching Example: Matrix Mult

```

for (i=0; i<n; i++)
  for (j=0; j<n; j++)
    for (k=0; k<n; k++)
      c[i][j] += a[i][k] * b[k][j];
    
```

**Reasonable cache effects**

- Good locality for A
- Bad locality for B
- Good locality for C

20

20

### Locality/Caching Example: Matrix Mult

```

for (j=0; j<n; j++)
  for (k=0; k<n; k++)
    for (i=0; i<n; i++)
      c[i][j] += a[i][k] * b[k][j];
    
```

**Poor cache effects**

- Bad locality for A
- Bad locality for B
- Bad locality for C

21

21

### Locality/Caching Example: Matrix Mult

```

for (i=0; i<n; i++)
  for (k=0; k<n; k++)
    for (j=0; j<n; j++)
      c[i][j] += a[i][k] * b[k][j];
    
```

**Good cache effects**

- Good locality for A
- Good locality for B
- Good locality for C

22

22

### Another ghost of exams past ...

Suppose that C laid out arrays in memory in column-major order instead of row-major order. What would be the most efficient loop ordering for this matrix multiplication to maximize performance through good locality?

A. i k j (Same as row-major)  
 B. i j k  
 C. j k i  
 D. j i k  
 E. k i j  
 F. k j i

C: j k i

Exactly what makes this bad for all three in row-major makes it ideal for column-major:  
 a and c have good spatial  
 b has good temporal, spatial

23

23

### Storage Hierarchy & Caching Issues

**Issue: Block size?**

Large block size:  
 + do data transfer less often  
 + take advantage of spatial locality  
 - longer time to complete data transfer  
 - less advantage of temporal locality

Small block size: the opposite

Typical: Lower in pyramid ⇒ slower data transfer ⇒ larger block sizes

Device	Block Size
Register	8 bytes
L1/L2/L3 cache line	128 bytes
Main memory page	4KB or 64KB
Disk block	512 bytes to 4KB
Disk transfer block	4KB (4096 bytes) to 64MB (67108864 bytes)

24

24

### Storage Hierarchy & Caching Issues

Issue: Who manages the cache?

Device	Managed by:
Registers (cache of L1/L2/L3 cache and main memory)	Compiler, using complex code-analysis techniques Assembly lang programmer
L1/L2/L3 cache (cache of main memory)	Hardware, using simple algorithms
Main memory (cache of local sec storage)	Hardware and OS, using virtual memory with complex algorithms (since accessing disk is expensive)
Local secondary storage (cache of remote sec storage)	End user, by deciding which files to download

25

25

### Next time ...

Getting started with ARM!

26

[Lobsterhermidor, Rauschen](#)

26