# Section 7: Memory and Caches

- ~~Cache basics~~
- ~~Principle of locality~~
- Memory hierarchies
- Cache organization
- Program optimizations that consider caches

# Cost of Cache Misses

■ **Huge difference between a hit and a miss**

   ▪ Could be 100x, if just L1 and main memory

■ **Would you believe 99% hits is twice as good as 97%?**

   ▪ Consider:
   
   Cache hit time of 1 cycle
   Miss penalty of 100 cycles

   ▪ Average access time:

      ▪ 97% hits:  1 cycle + 0.03 * 100 cycles = 4 cycles
      ▪ 99% hits:  1 cycle + 0.01 * 100 cycles = 2 cycles

■ **This is why "miss rate" is used instead of "hit rate"**

# Cache Performance Metrics

- **Miss Rate**
  - Fraction of memory references not found in cache (misses / accesses) = 1 - hit rate
  - Typical numbers (in percentages):
    - 3% - 10% for L1

- **Hit Time**
  - Time to deliver a line in the cache to the processor
    - Includes time to determine whether the line is in the cache
  - Typical hit times: 1 - 2 clock cycles for L1

- **Miss Penalty**
  - Additional time required because of a miss
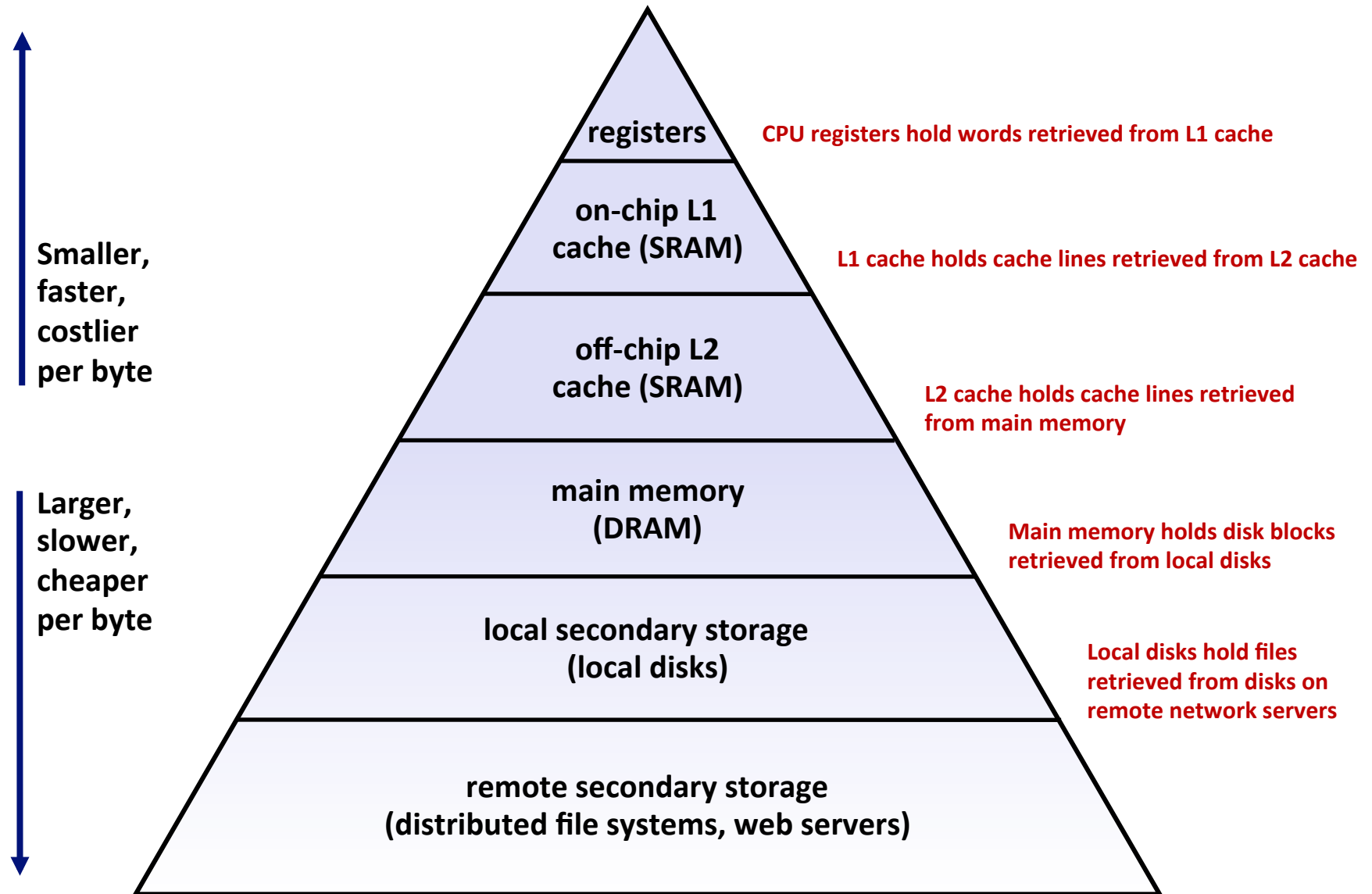  - Typically 50 - 200 cycles

# Memory Hierarchies

- **Some fundamental and enduring properties of hardware and software systems:**
  - Faster storage technologies almost always cost more per byte and have lower capacity
  - The gaps between memory technology speeds are widening
    - True for: registers ↔ cache, cache ↔ DRAM, DRAM ↔ disk, etc.
  - Well-written programs tend to exhibit good locality

- **These properties complement each other beautifully**

- **They suggest an approach for organizing memory and storage systems known as a <span style="color:red">memory hierarchy</span>**
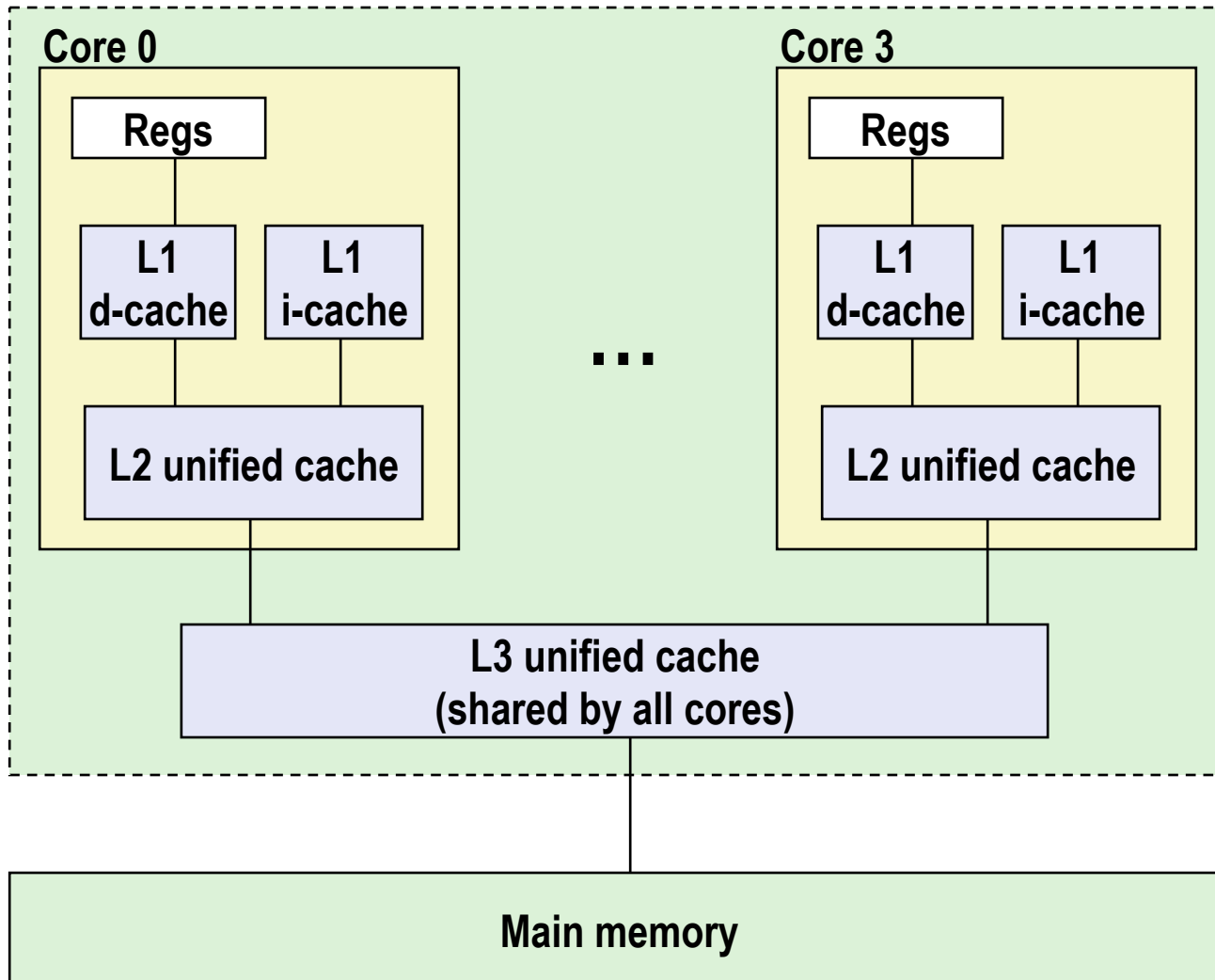
# Memory Hierarchies

- **Fundamental idea of a memory hierarchy:**
  - Each level k serves as a cache for the larger, slower, level k+1 below.

- **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.**

# An Example Memory Hierarchy



**Smaller, faster, costlier per byte**

**Larger, slower, cheaper per byte**

registers — CPU registers hold words retrieved from L1 cache

on-chip L1 cache (SRAM) — L1 cache holds cache lines retrieved from L2 cache

off-chip L2 cache (SRAM) — L2 cache holds cache lines retrieved from main memory

main memory (DRAM) — Main memory holds disk blocks retrieved from local disks

local secondary storage (local disks) — Local disks hold files retrieved from disks on remote network servers

remote secondary storage (distributed file systems, web servers)

Caches - Memory Hierarchy

# Intel Core i7 Cache Hierarchy

**Processor package**



**Core 0**

Regs

L1 d-cache | L1 i-cache

L2 unified cache

**Core 3**

Regs

L1 d-cache | L1 i-cache

L2 unified cache

...

L3 unified cache
(shared by all cores)

Main memory

**L1 i-cache and d-cache:**
32 KB, 8-way,
Access: 4 cycles

**L2 unified cache:**
256 KB, 8-way,
Access: 11 cycles

**L3 unified cache:**
8 MB, 16-way,
Access: 30-40 cycles

**Block size**: 64 bytes for all caches.