The Hardware/Software Interface

Memory & Caches I

Instructors:

Justin Hsia, Amber Hu

Teaching Assistants:

Anthony Mangus

Grace Zhou

Jiuyang Lyu

Kurt Gu

Mendel Carroll

Naama Amiel

Rose Maresh

Violet Monserate

Divya Ramu

Jessie Sun

Kanishka Singh

Liander Rainbolt

Ming Yan

Pollux Chen

Soham Bhosale



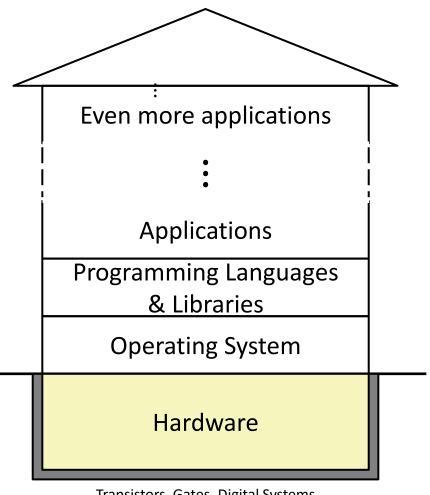
Relevant Course Information

- Mid-quarter Survey due tonight
- HW14 due tonight, HW15 due Monday, HW16 due Wednesday
- Lab 3 due next Friday (11/7)
 - Lots of resources: HW15, Section 06, GDB Stack Tutorial lesson
- Midterm grading is ongoing, will be released when available
 - Midterm clobber policy allows you to overwrite your midterm score!
 - Solutions file will be posted, Gradescope rubric items will have more scoring details

House of Computing Check-In

- * Topic Group 3: Scale & Coherence
 - Caches, Memory Allocation, Processes, Virtual Memory

- How do we maintain logical consistency in the face of more data and more processes?
 - How do we support data access, including dynamic requests, across multiple processes?
 - How do we support control flow both within many processes and things external to the computer?



Transistors, Gates, Digital Systems

Physics

Lecture Outline (1/4)

- * IEC Prefixes
- Caches and Cache Mechanics
- Cache Performance Metrics
- The Memory Hierarchy

Prefixes (Review)

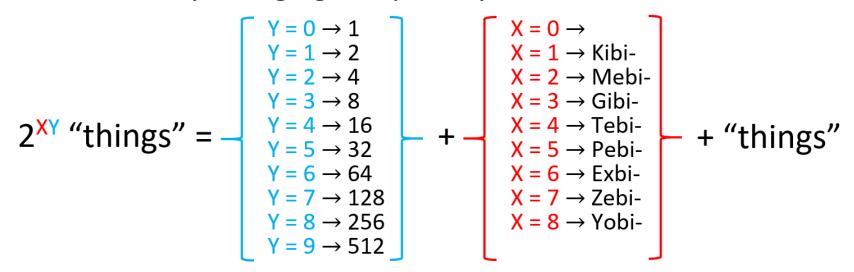
- Here focusing on large numbers (exponents > 0)
- SI prefixes are ambiguous if base 10 or 2
 - Note that $10^3 \approx 2^{10}$
- IEC prefixes are unambiguously base 2

SIZE PREFIXES (10^x for Disk, Communication; 2^x for Memory)

SI Size	Prefix	Symbol	IEC Size	Prefix	Symbol
10^{3}	Kilo-	K	2 ¹⁰	Kibi-	Ki
10 ⁶	Mega-	M	2 ²⁰	Mebi-	Mi
10 ⁹	Giga-	G	2 ³⁰	Gibi-	Gi
10^{12}	Tera-	T	2 ⁴⁰	Tebi-	Ti
10^{15}	Peta-	P	2 ⁵⁰	Pebi-	Pi
10^{18}	Exa-	Е	2 ⁶⁰	Exbi-	Ei
10^{21}	Zetta-	Z	2 ⁷⁰	Zebi-	Zi
10^{24}	Yotta-	Y	2 ⁸⁰	Yobi-	Yi

Large Powers of 2 and Units (Review)

- ❖ Because IEC prefixes are powers of 2¹⁰, we can convert any large power of 2 as follows:
 - Note that we are only changing the quantity and the units remain the same



- Examples:
 - 2³² bits into IEC:
 - How many address bits to use 13.2 TiB of memory?

How to Remember?

- Will be given to you on Final reference sheet
- Mnemonics
 - There unfortunately isn't one well-accepted mnemonic
 - But that shouldn't stop you from trying to come with one!
 - Killer Mechanical Giraffe Teaches Pet, Extinct Zebra to Yodel
 - Kirby Missed Ganondorf Terribly, Potentially Exterminating Zelda and Yoshi
 - xkcd: Karl Marx Gave The Proletariat Eleven Zeppelins, Yo
 - https://xkcd.com/992/
 - Post your best on Ed Discussion!

Polling Questions (1/2)

- Convert the following to or from IEC:
 - 512 Ki-books
 - 2²⁷ caches

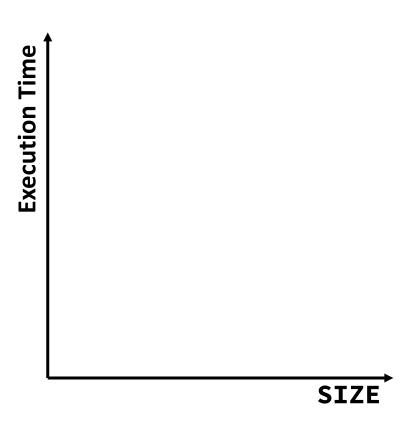
Lecture Outline (2/4)

- * IEC Prefixes
- Caches and Cache Mechanics
- Cache Performance Metrics
- The Memory Hierarchy

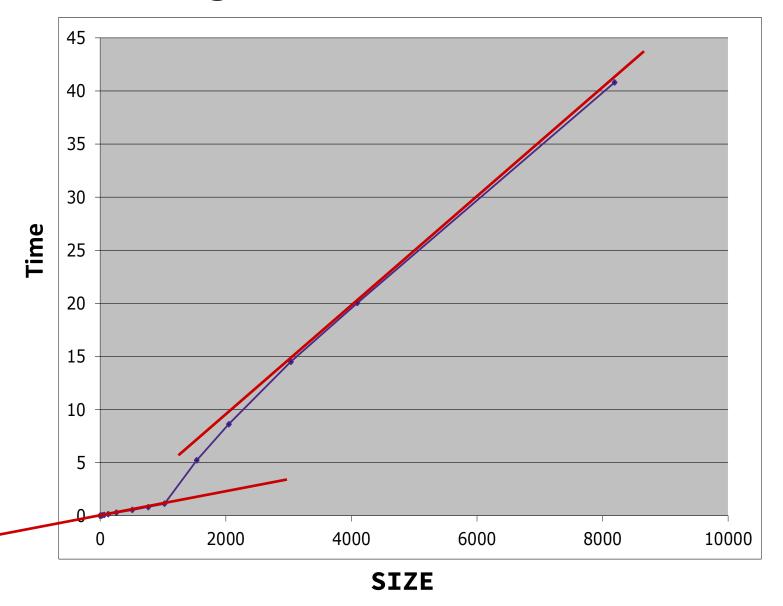
How does execution time grow with SIZE?

```
int array[SIZE];
int sum = 0;

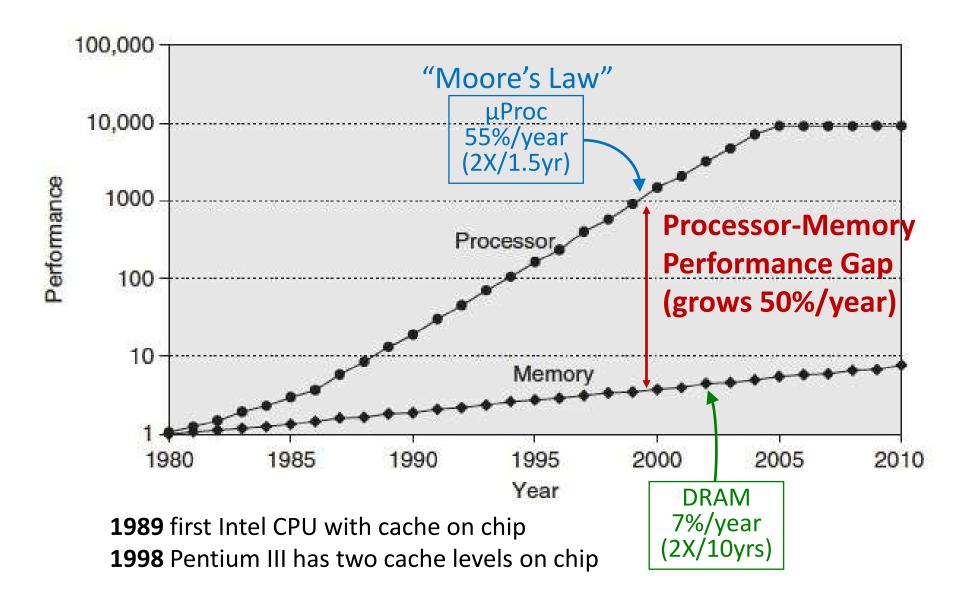
for (int i = 0; i < 200000; i++) {
   for (int j = 0; j < SIZE; j++) {
      sum += array[j];
   }
}</pre>
```



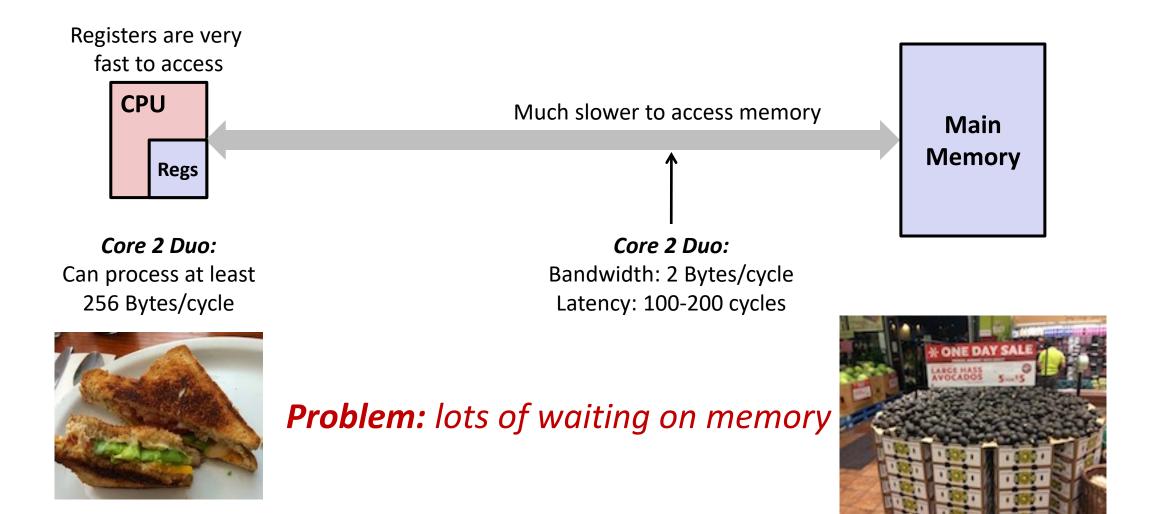
Actual Cache Timing Data



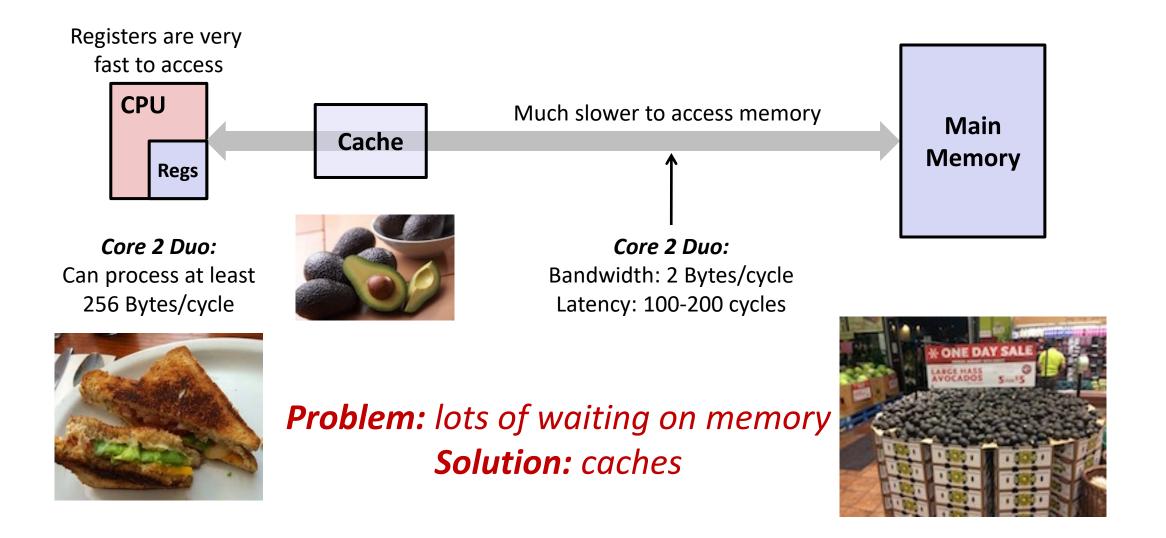
Processor-Memory Gap



Processor-Memory Bottleneck

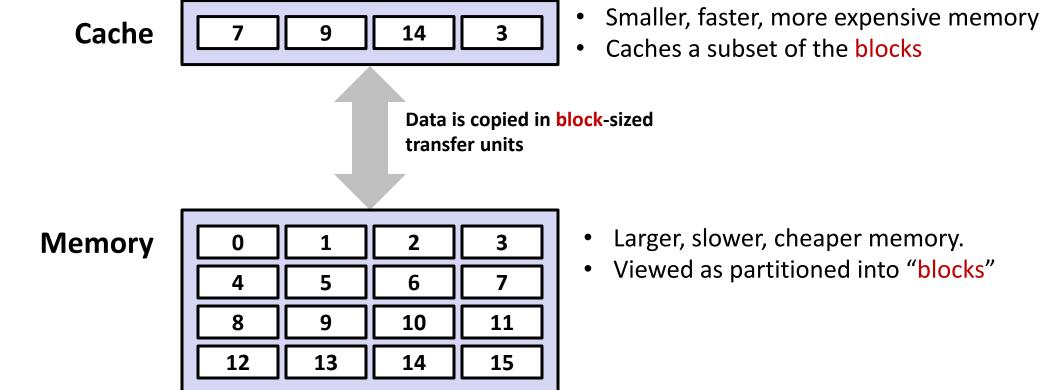


Processor-Memory Bottleneck Fix



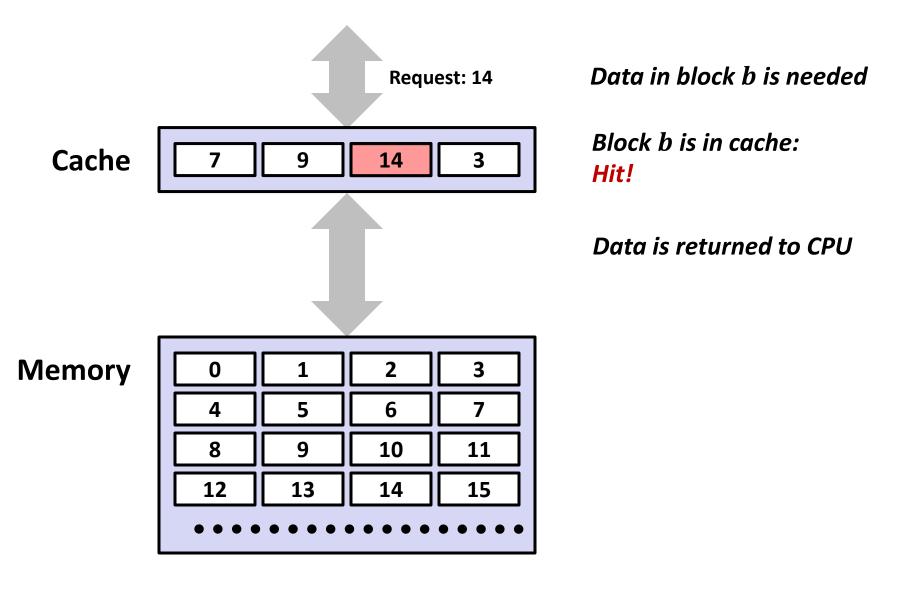
Cache 6

- Pronunciation: "cash"
 - We abbreviate this as "\$"
- English: A hidden storage space for provisions, weapons, and/or treasures
- Computer: Memory with short access time used for the storage of frequently or recently used instructions (i-cache/I\$) or data (d-cache/D\$)
 - More generally: Used to optimize data transfers between any system elements with different characteristics (network interface cache, I/O cache, etc.)

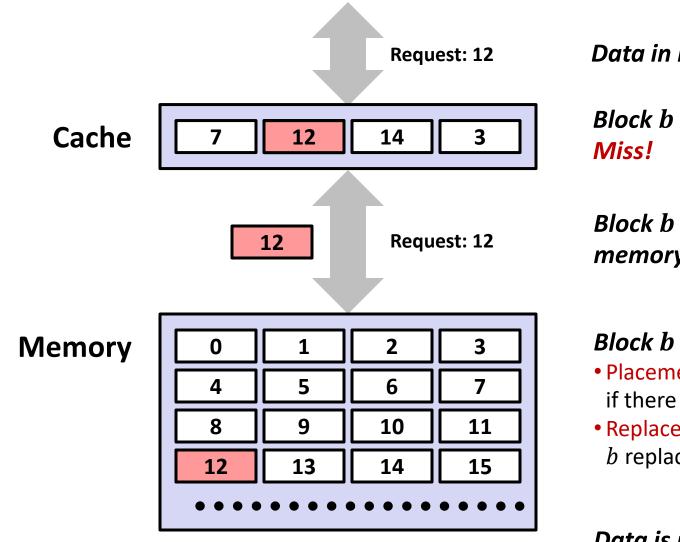


CSE351, Autumn 2025

General Cache Concepts: Hit (Review)



General Cache Concepts: Miss (Review)



Data in block b is needed

Block b is not in cache: Miss!

Block b is fetched from memory

Block b is stored in cache

- Placement policy: Where should b go if there is room in the cache
- Replacement policy: Which block should
 b replace if the ache is full

Data is returned to CPU

Why Caches Work: Locality (Review)

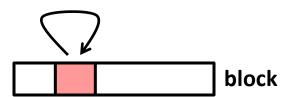
Locality: Programs tend to use data and instructions with addresses near or equal to those they have used recently

Why Caches Work: Temporal Locality (Review)

 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



Why Caches Work: Spatial Locality (Review)

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

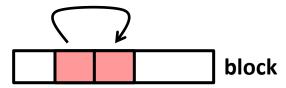
Temporal locality:

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





How do caches take advantage of this?

Locality in Practice

```
sum = 0;
for (i = 0; i < n; i++) {
   sum += a[i];
}
return sum;</pre>
```

Data:

Temporal: sum referenced in each iteration

Spatial: consecutive elements of array a [] accessed

Instructions:

Temporal: cycle through loop repeatedly

Spatial: reference instructions in sequence

Locality Example #1 Code

```
int sum_array_rows(int a[M][N]) {
  int i, j, sum = 0;

for (i = 0; i < M; i++)
  for (j = 0; j < N; j++)
    sum += a[i][j];

return sum;
}</pre>
```

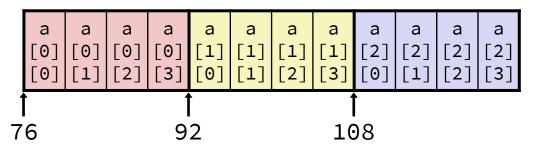
Locality Example #1 Access Pattern

```
int sum_array_rows(int a[M][N]) {
  int i, j, sum = 0;

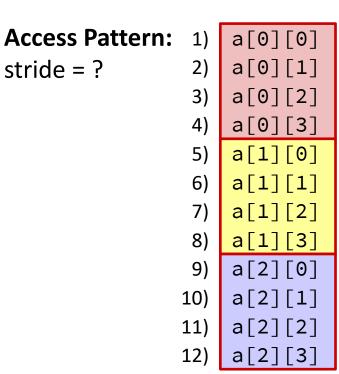
  for (i = 0; i < M; i++)
    for (j = 0; j < N; j++)
      sum += a[i][j];

  return sum;
}</pre>
```

Layout in Memory (arbitrary starting address)



M = 3, N=4 a[0][0] a[0][1] a[0][2] a[0][3] a[1][0] a[1][1] a[1][2] a[1][3] a[2][0] a[2][1] a[2][2] a[2][3] Access Pattern: 1) a[0][0]

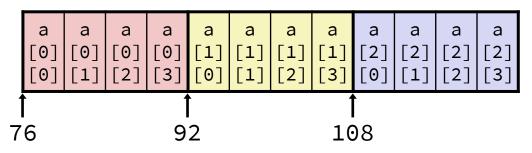


```
int sum_array_cols(int a[M][N]) {
  int i, j, sum = 0;

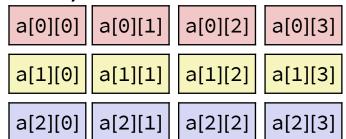
for (j = 0; j < N; j++)
  for (i = 0; i < M; i++)
    sum += a[i][j];

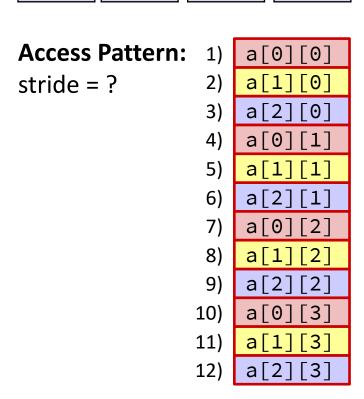
return sum;
}</pre>
```

Layout in Memory (arbitrary starting address)



M = 3, N=4





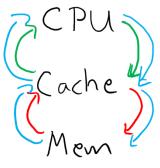
CSE351, Autumn 2025

Lecture Outline (3/4)

- * IEC Prefixes
- Caches and Cache Mechanics
- Cache Performance Metrics
- The Memory Hierarchy

Cache Performance Parameters (Review)

- Huge difference between a cache hit and a cache miss
 - Could be 100x speed difference between accessing cache and main memory (measured in clock cycles)
- Hit Time (HT)
 - Time to deliver a block in the cache to the processor
 - Includes time to determine whether the block is in the cache
- Hit takes HT Miss takes HT+MP



- Miss Penalty (MP)
 - Additional time required because of a miss
- Miss Rate (MR)
 - Fraction of memory references not found in cache (misses / accesses) = 1 Hit Rate

Cache Performance Measurement (Review)

 Average Memory Access Time (AMAT): average time to access memory considering both hits and misses

AMAT = Hit time + Miss rate × Miss penalty = HT + MR × MP

HT and MP generally fixed, so minimize MR

- Example: Assume HT of 1 clock cycle and MP of 100 clock cycles
 - 97%: AMAT =
 - 99%: AMAT =
 - 99% hit rate twice as good as 97% hit rate!!?!

Polling Questions (2/2)

Processor specs: 200 ps clock, MP of 50 clock cycles, MR of 0.02 misses/instruction, and HT of 1 clock cycle

AMAT =

- Which improvement would be best?
 - A. 190 ps clock
 - B. Miss penalty of 40 clock cycles
 - C. MR of 0.015 misses/instruction

Lecture Outline (4/4)

- * IEC Prefixes
- Caches and Cache Mechanics
- Cache Performance Metrics
- The Memory Hierarchy

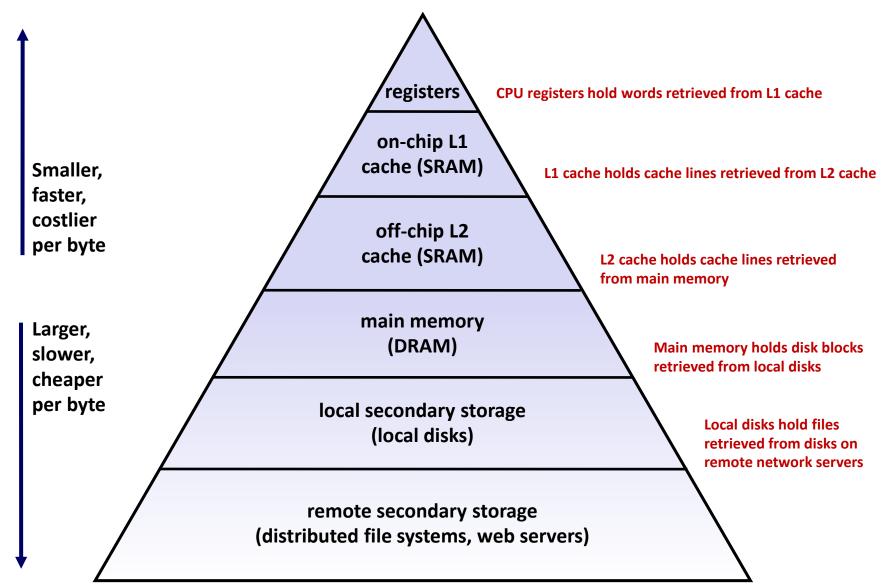
More than one cache? (Review)

- Why would we want to do that?
 - Avoid going to memory!
- Extra considerations for multilevel caching:
 - Block size may differ between different levels of cache
 - A memory access can miss in one level and then hit in the next, causing the block to be copied into the higher level
 - AMAT, however, is computed for your overall system caching, taking all levels of cache into account

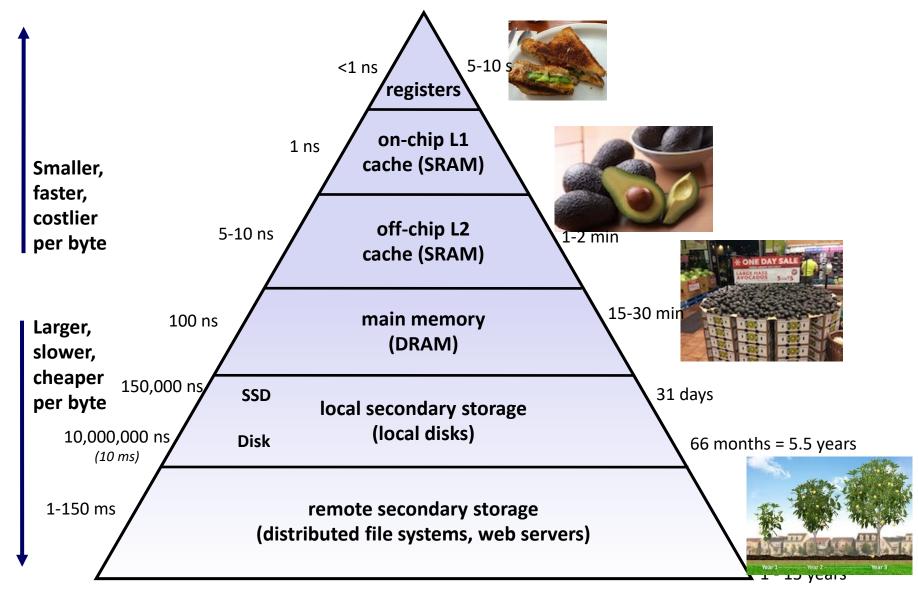
Memory Hierarchies (Review)

- Some fundamental properties of hardware and software systems:
 - The gaps between memory technology speeds are widening
 - True for: registers
 ⇔ cache, cache
 ⇔ main memory, main memory
 ⇔ disk, etc.
 - Faster storage technologies almost always cost more per byte and have lower capacity
 - Well-written programs tend to exhibit good locality
- These properties complement each other beautifully and suggest an approach for organizing memory and storage systems known as a memory hierarchy
 - For each level k, the faster, smaller device at level k serves as a cache for the larger, slower device at level k+1

An Example Memory Hierarchy (Review)



An Example Memory Hierarchy: Analogy Revisited



Summary (1/3)

IEC prefixes are unambiguously powers of 2:

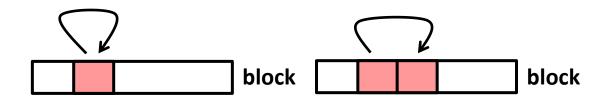
SIZE PREFIXES (10^x for Disk, Communication; 2^x for Memory)

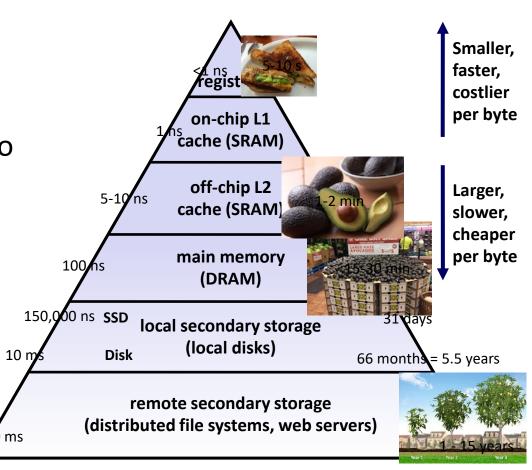
SI Size	Prefix	Symbol	IEC Size	Prefix	Symbol
10 ³	Kilo-	K	2 ¹⁰	Kibi-	Ki
10 ⁶	Mega-	M	2 ²⁰	Mebi-	Mi
10 ⁹	Giga-	G	2 ³⁰	Gibi-	Gi
10 ¹²	Tera-	T	2 ⁴⁰	Tebi-	Ti
10 ¹⁵	Peta-	P	2 ⁵⁰	Pebi-	Pi
10 ¹⁸	Exa-	Е	2 ⁶⁰	Exbi-	Ei
10 ²¹	Zetta-	Z	2 ⁷⁰	Zebi-	Zi
10 ²⁴	Yotta-	Y	2 ⁸⁰	Yobi-	Yi

$$2^{XY} \text{ "things"} = - \begin{bmatrix} Y = 0 \rightarrow 1 \\ Y = 1 \rightarrow 2 \\ Y = 2 \rightarrow 4 \\ Y = 3 \rightarrow 8 \\ Y = 4 \rightarrow 16 \\ Y = 5 \rightarrow 32 \\ Y = 6 \rightarrow 64 \\ Y = 7 \rightarrow 128 \\ Y = 8 \rightarrow 256 \\ Y = 9 \rightarrow 512 \end{bmatrix} + - \begin{cases} X = 0 \rightarrow \\ X = 1 \rightarrow \text{Kibi-} \\ X = 2 \rightarrow \text{Mebi-} \\ X = 3 \rightarrow \text{Gibi-} \\ X = 4 \rightarrow \text{Tebi-} \\ X = 5 \rightarrow \text{Pebi-} \\ X = 6 \rightarrow \text{Exbi-} \\ X = 7 \rightarrow \text{Zebi-} \\ X = 8 \rightarrow \text{Yobi-} \end{cases}$$

Summary (2/3)

- Memory Hierarchy
 - Successively higher levels contain "most used" data from lower levels
 - Caches are intermediate storage levels used to optimize data transfers between any system elements with different characteristics
 - Exploits temporal and spatial locality:





Summary (3/3)

- Cache Performance
 - Ideal case: found in cache (cache hit), return requested data immediately
 - Bad case: not found in cache (cache miss), search in next level
 - Bring entire *cache block* containing requested data into this cache once found
 - Average Memory Access Time (AMAT) = HT + MR × MP
 - Hurt by Miss Rate and Miss Penalty

