Locality, Buffering, Caches Computer Operating Systems, Fall 2025

Instructors: Joel Ramirez

Head TAs: Maya Huizar Akash Kaukuntla

Vedansh Goenka Joy Liu

TAs:

Eric Zou Joseph Dattilo Aniket Ghorpade Shriya Sane

Zihao Zhou Eric Lee Shruti Agarwal Yemisi Jones

Connor Cummings Shreya Mukunthan Alexander Mehta Raymond Feng

Bo Sun Steven Chang Rania Souissi Rashi Agrawal

Sana Manesh



pollev.com/cis5480

If you were in charge of writing some software, what qualities would you prioritize?

Administrivia

PennOS

- Meet with your TAs by the end of this week!
- If you meet after Friday (unless there's no other availability) then you will use up your late tokens.
- Your TAs are going to be your life line if you have any questions regrading PennOS!

Check-In 07

- Make sure to get that in today!
- Due at 5PM today. So if you forgot, you still have a chance to turn it in!
- Recitation this week! Same time same place.
 - Led by Vedansh and Akash

Lecture Outline

- * Ideology
- Locality
- Buffering
- Caches



pollev.com/cis5480

If you were in charge of writing some software, what qualities would you prioritize?

Some Answers I'd Expect

- Efficiency
- Correctness
- Privacy
- Accessibility
- Security
- User-Friendliness
- Speech Rights
- Trustworthiness
- Maintainability
- Versatility

Efficiency is a big motivator

Both historically, and today

Can be seen in our own CS curriculum

- APCS or CS1: different sorting algorithms
- Space and time complexity -> O(n)
- Clock cycles and frequency in hardware
- Compiler optimization
- Programming approaches: recursion v. iteration

What do we mean by efficiency?

Well, what do you think we mean by it?

Efficiency to do ... what?

Why is efficiency so relevant?

Many ways to approach "making something efficient"

Constraint optimization problem

Some examples:

- Resource allocation
- Automation
- Production/development cycles

Is Efficiency the Most Important Thing™?

- How often do we discuss / take into consideration:
 - Correctness
 - Accessibility / Universal Design
 - Security
 - Maintainability (scalability, readability)

Software Bloat

- When successive features with minimal positive impact make the software run significantly slower or use more memory
- Affects performance, available space, and the security of your computer
- Wirth's Law (1995): software is getting slower more rapidly than hardware is becoming faster.

Examples of Software Bloat & Effects

- Apps growing larger in size over time
- Devices running slower as they age
 - Sometimes due to the battery or something wearing out
 - Sometimes cause the hardware is expected to do more and more
- Unwanted features (Grok on Twitter), Stories on Instagram (I still remember when this happened), Location Tracking on Instagram, DM feature on Spotify, etc etc etc.

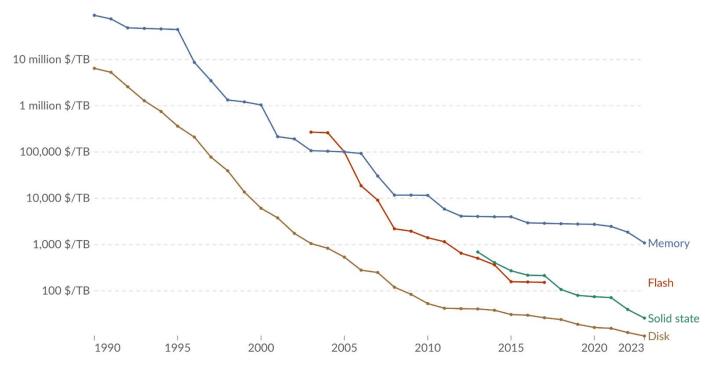
Why is Software Bloat a thing?

Historical price of computer memory and storage



This data is expressed in US dollars per terabyte (TB), adjusted for inflation. "Memory" refers to random access memory (RAM), "disk" to magnetic storage, "flash" to special memory used for rapid data access and rewriting, and "solid state" to solid-state drives (SSDs).

L18: Locality



Data source: John C. McCallum (2023); U.S. Bureau of Labor Statistics (2024)

OurWorldinData.org/technological-change | CC BY

Note: For each year, the time series shows the cheapest historical price recorded until that year. This data is expressed in constant 2020 US\$.

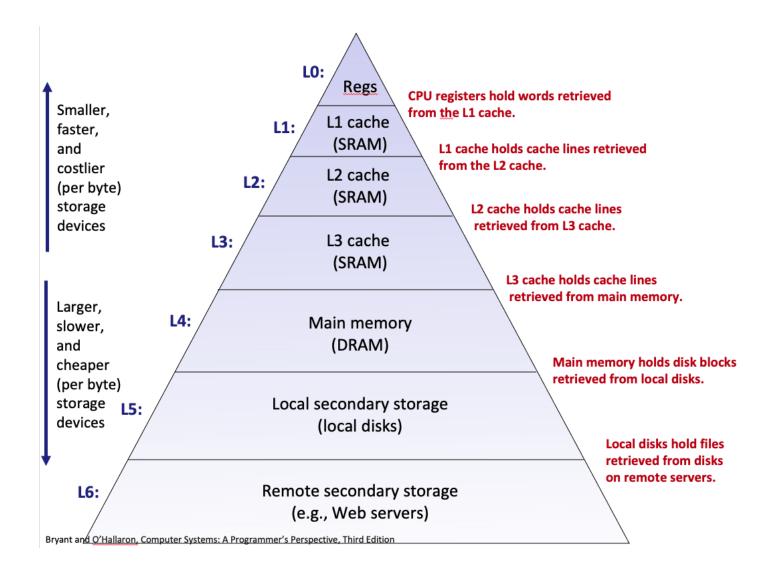
How is the OS responsible?

- Newer OS requiring more resources
- Can make it easier or harder to delete unwanted files
 - Deleting applications
 - This includes purging old data when upgrading the OS
- Installing features that you don't want
 - Apple intelligence is 4 GB
 - Microsoft Windows (Most of it)

Lecture Outline

- Ideology
- Locality
- Buffering
- Caches

Memory Hierarchy

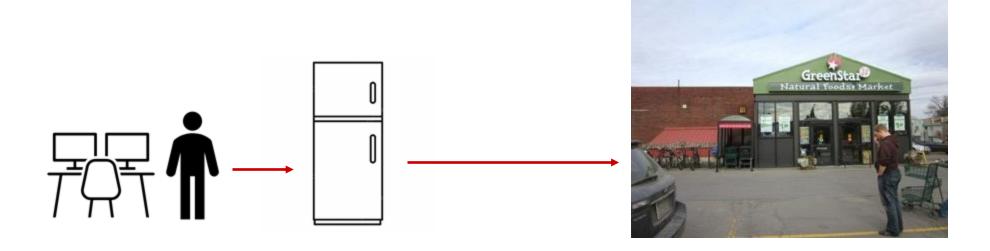


Principle of Locality

- The tendency for the Programs to access the same set of memory locations over a short period of time
- Two main types:
 - Temporal Locality: If we access a portion of memory, we will likely reference it again soon
 - **Spatial Locality**: If we access a portion of memory, we will likely reference memory close to it in the near future.
- Data that is accessed frequently can be stored in hardware that is quicker to access.

Locality Analogy

- If we are at home and we are hungry, where do we get food from?
 - We get it from our refrigerator!
 - If the refrigerator is empty, we go to the grocery store
 - When at the grocery store, we don't just get what we want right now, but also get other things we think we want in the near future (so that it will be in our fridge when we want it)



Lecture Outline

- Ideology
- Locality
- Buffering
- Caches



pollev.com/cis5480

If we compile this and run it, how many times is hello printed?

```
int main() {
   if (fork() == 0) {
      write(STDOUT_FILENO, "hello", 5);
   }
   if (fork() == 0) {
      write(STDOUT_FILENO, "hello", 5);
   }
   return EXIT_SUCCESS;
}
```



Raise Your Hands

If we compile this and run it, how many times is hello printed?

```
int main() {
   if (fork() == 0) {
      printf("hello");
   }
   if (fork() == 0) {
      printf("hello");
   }
   return EXIT_SUCCESS;
}
```

Poll Everywhere

If we compile this and run it, how many times is hello printed?

```
int main() {
   if (fork() == 0) {
      printf("hello\n");
   }
   if (fork() == 0) {
      printf("hello\n");
   }
   return EXIT_SUCCESS;
}
```

C stdio vs POSIX

- Why are we getting these different outputs?
- Let's start with the first two examples. Both use different ways of writing to standard out.
 - C stdio: user level portable library for standard input/output. Should work on any environment that has the C standard library
 - E.g. printf, fprintf, fputs, getline, etc.
 - POSIX C API: Portable Operating System Interface. Functions that are supported by many operating systems to support many OS-level concepts (Input/Output, networking, processes, threads...)

Buffered writing

- By default, C stdio uses buffering on top of POSIX:
 - When one writes with fwrite(), the data being written is copied into a buffer allocated by stdio inside your process' address space
 - As some point, once enough data has been written, the buffer will be "flushed" to the operating system.
 - When the buffer fills (often 1024 or 4096 bytes)
 - This prevents invoking the write system call and going to the filesystem too often

Arrow shows what has just finished executing.

```
int main(int argc, char** argv) {
   char *h = "hi"
   char *i = h + 1;

   FILE* fout = fopen("hi.txt", "w");

   // write "hi" one char at a time
   fwrite(h, sizeof(char), 1, fout);
   fwrite(i, sizeof(char), 1, fout);

   fclose(fout);
   return EXIT_SUCCESS;
}
```

buf h i

hi.txt (disk/OS)

Arrow shows what has just finished executing.

```
int main(int argc, char** argv) {
   char *h = "hi"
   char *i = h + 1;

FILE* fout = fopen("hi.txt", "w");

// write "hi" one char at a time
   fwrite(h, sizeof(char), 1, fout);
   fwrite(i, sizeof(char), 1, fout);

fclose(fout);
   return EXIT_SUCCESS;
}
```

DUT	
h	i
C stdi	o buffer for "hi.txt"
hi.txt	(disk/OS)

26

Arrow shows what has just finished executing.

```
int main(int argc, char** argv) {
   char *h = "hi"
   char *i = h + 1;

FILE* fout = fopen("hi.txt", "w");

// write "hi" one char at a time
fwrite(h, sizeof(char), 1, fout);
fwrite(i, sizeof(char), 1, fout);

fclose(fout);
   return EXIT_SUCCESS;
}
```

Store 'h' into buffer, so that we do not go to filesystem **yet** buf C stdio buffer hi.txt (disk/OS)

Arrow shows what has just finished executing.

```
int main(int argc, char** argv) {
   char *h = "hi"
   char *i = h + 1;

FILE* fout = fopen("hi.txt", "w");

// write "hi" one char at a time
   fwrite(h, sizeof(char), 1, fout);

fwrite(i, sizeof(char), 1, fout);

fclose(fout);
   return EXIT_SUCCESS;
}
```

Store 'i' into buffer, so that we do not go to filesystem **yet** buf C stdio buffer h hi.txt (disk/OS)

Arrow shows what has just finished executing.

```
int main(int argc, char** argv) {
   char *h = "hi"
   char *i = h + 1;

FILE* fout = fopen("hi.txt", "w");

// write "hi" one char at a time
   fwrite(h, sizeof(char), 1, fout);
   fwrite(i, sizeof(char), 1, fout);

fclose(fout);

return EXIT_SUCCESS;
}
```

uf		
h	i	

C stdio buffer

h	i	

When we call fclose, we deallocate and flush the buffer to disk hi.txt (disk/OS)



Arrow shows what has just finished executing.

CIS 4480 Fall 2025

```
int main(int argc, char** argv) {
   char *h = "hi"
   char *i = h + 1;

FILE* fout = fopen("hi.txt", "w");

// write "hi" one char at a time
   fwrite(h, sizeof(char), 1, fout);
   fwrite(i, sizeof(char), 1, fout);

fclose(fout);
   return EXIT_SUCCESS;
}
```

buf h i

hi.txt (disk/OS)

h i

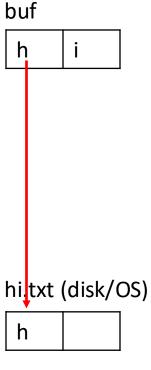
Arrow shows what has just finished executing.

```
int main(int argc, char** argv) {
   char buf[2] = {'h', 'i'};
   int fd = open("hi.txt", O_WRONLY | O_CREAT);

   // read "hi" one char at a time
   write(fd, &buf, sizeof(char));

   write(fd, &buf+1, sizeof(char));

   close(fd);
   return EXIT_SUCCESS;
}
```



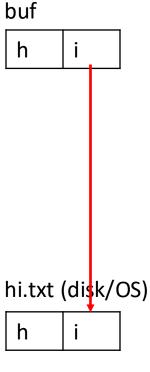
Arrow shows what has just finished executing.

```
int main(int argc, char** argv) {
  char buf[2] = {'h', 'i'};
  int fd = open("hi.txt", O_WRONLY | O_CREAT);

  // read "hi" one char at a time
  write(fd, &buf, sizeof(char));

  write(fd, &buf+1, sizeof(char));

  close(fd);
  return EXIT_SUCCESS;
}
```



Arrow shows what has just finished executing.

```
int main(int argc, char** argv) {
   char buf[2] = {'h', 'i'};
   int fd = open("hi.txt", O_WRONLY | O_CREAT);

   // read "hi" one char at a time
   write(fd, &buf, sizeof(char));

   write(fd, &buf+1, sizeof(char));

   close(fd);
   return EXIT_SUCCESS;
}
```

buf h i

Two OS/File system accesses instead of one

hi.txt (disk/OS)
h i

Buffered Reading

- By default, C stdio uses buffering on top of POSIX:
 - When one reads with fread(), a lot of data is copied into a buffer allocated by stdio inside your process' address space
 - Next time you read data, it is retrieved from the buffer
 - This avoids having to invoke a system call again
 - As some point, the buffer will be "refreshed":
 - When you process everything in the buffer (often 1024 or 4096 bytes)
 - Similar thing happens when you write to a file

Quick Aside: strace

- Basic usage: strace ./program [args]
- strace exec's a program and "traces" the system calls that it utilizes to execute.
- Will show you the system call, its arguments, return values and more!
 - "Arguments are printed in symbolic form with passion."
- Really useful tool. See example usage below.

```
$ strace ./cat /dev/null
...
open("/dev/null", O_RDONLY) = 3
...
```

We can use this to see if the past couple of slides are true or not.

Buffered Reading Example

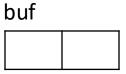
Arrow signifies what will be executed next

```
int main(int argc, char** argv) {
   char buf[2];
   FILE* fin = fopen("hi.txt", "rb");

   // read "hi" one char at a time
   fread(&buf, sizeof(char), 1, fin);

   fread(&buf+1, sizeof(char), 1, fin);

   fclose(fin);
   return EXIT_SUCCESS;
}
```



hi.txt (disk/OS)
h
i

Arrow signifies what will be executed next

```
int main(int argc, char** argv) {
   char buf[2];
   FILE* fin = fopen("hi.txt", "rb");

   // read "hi" one char at a time
   fread(&buf, sizeof(char), 1, fin);

   fread(&buf+1, sizeof(char), 1, fin);

   fclose(fin);
   return EXIT_SUCCESS;
}
```

	Copy	out w	<i>ı</i> hat
buf	wasr	eque	sted
C stdi	o buffe	er	_
h	i		
/	nd as r can fr		
hi.txt	(disk/	OS)	
h	i		

Arrow signifies what will be executed next

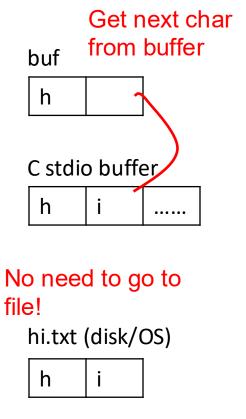
CIS 4480 Fall 2025

```
int main(int argc, char** argv) {
   char buf[2];
   FILE* fin = fopen("hi.txt", "rb");

   // read "hi" one char at a time
   fread(&buf, sizeof(char), 1, fin);

   fread(&buf+1, sizeof(char), 1, fin);

   fclose(fin);
   return EXIT_SUCCESS;
}
```



Arrow signifies what will be executed next

```
int main(int argc, char** argv) {
   char buf[2];
   FILE* fin = fopen("hi.txt", "rb");

   // read "hi" one char at a time
   fread(&buf, sizeof(char), 1, fin);

   fread(&buf+1, sizeof(char), 1, fin);

   fclose(fin);
   return EXIT_SUCCESS;
}
```

buf h i

C stdio buffer

h i

hi.txt (disk/OS)

h i

Arrow signifies what will be executed next

```
int main(int argc, char** argv) {
   char buf[2];
   FILE* fin = fopen("hi.txt", "rb");

   // read "hi" one char at a time
   fread(&buf, sizeof(char), 1, fin);

   fread(&buf+1, sizeof(char), 1, fin);

   fclose(fin);
   return EXIT_SUCCESS;
}
```

buf h i

hi.txt (disk/OS)

Why NOT Buffer?

- Reliability the buffer needs to be flushed
 - Loss of computer power = loss of data
 - "Completion" of a write (i.e. return from fwrite()) does not mean the data has actually been written

L18: Locality

- Performance buffering takes time
 - Copying data into the stdio buffer consumes CPU cycles and memory bandwidth
 - Can potentially slow down high-performance applications, like a web server or database ("zero-copy")
- When is buffering faster? Slower?

Many small writes

Large writes Singular write

Quick Aside: exit vs _exit

int exit(int status);

- Causes the process to terminate, but before it does it calls all "installed" exit functions.
 - Saw some of those in Penn-Shell...not sure why.
- Additionally, it flushes and closes all stdio buffers before the process terminates.
 - By default, "return" calls exit() somewhere down the line. If you'd like to see this, examine the obj file created by a binary.

```
int _exit(int status);
```

- Causes the process to terminate immediately. There is no flushing.
 - Check it out for yourself. Everything that uses fwrite will not actually write to the file.

Arrow signifies what will be executed next.

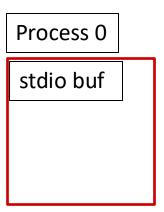
I execute processes in parallel and "in sync" for demonstration purposes

```
int main() {
    if (fork() == 0) {
        printf("hello");
    }
    if (fork() == 0) {
        printf("hello");
    }
    return EXIT_SUCCESS;
}
```



Arrow signifies what will be executed next. I execute processes in parallel and "in sync" for demonstration purposes

```
int main() {
   if (fork() == 0) {
      printf("hello");
   }
   if (fork() == 0) {
      printf("hello");
   }
   return EXIT_SUCCESS;
}
```





Arrow signifies what will be executed next.

I execute processes in parallel and "in sync" for demonstration purposes

```
Process 0
                                                    Process 1
int main() {
  if (fork() == 0) {
                                 stdio buf
                                                    stdio buf
    printf("hello");
                                                    hello
  if (fork() == 0) {
    printf("hello");
                                 Process 2
                                                    Process 3
  return EXIT_SUCCESS;
                                 stdio buf
                                                    stdio buf
                                                    hello
```

Arrow signifies what will be executed next.

I execute processes in parallel and "in sync" for demonstration purposes

```
Process 0
                                                       Process 1
   int main() {
     if (fork() == 0) {
                                     stdio buf
                                                       stdio buf
       printf("hello");
                                                       hello
     if (fork() == 0) {
       printf("hello");
                                    Process 2
                                                       Process 3
     return EXIT_SUCCESS;
                                     stdio buf
                                                        stdio buf
                                     hello
                                                       hello
                                                        hello
Hello is printed 4 times!
```

Fork Problem Explained (pt.2)

- Why did we get different outputs when printf printed a newline character after hello?
 - Only difference was:

```
printf("hello"); vs printf("hello\n");
```

- ❖ All we needed to do to get the expected output was add a \n. why?
- printf prints to stdout and by default stdout is line buffered because it refers to a TTY. Meaning it flushes the buffer on a newline character
 - If we ran ./prog > out.txt (redirect the output of the program to a file), we would get different output since buffering policy changes.

How to flush/modify the cstdio buffer

For C stdio:

```
•(int fflush(FILE* stream);
```

Flushes the stream to the OS/filesystem

- Has a family of related functions like setbuf(), setbuffer(), setlinebuf();
- Can set the stream to be unbuffered or a specified buffer

How to flush POSIX?

- When we write to a file with POSIX it is sent to the filesystem, is it immediately sent to disc? No
 - Well, we do have the block cache... so it may not be written to disc
 - Since all File I/O requests go to the file system, if another process accesses the same file, then it should see the data even if it is the block cache and not in disc.
 - If we lose power though...

How to flush POSIX to disk

- Two functions
 - •(int fsync(int fd);
 - Flushes all in-core data and metadata to the storage medium
 - int fdatasync(int fd);
 - Sends the file data to disk
 - Does not flush modified metadata unless necessary for data.

Blank slide

Transition Slide

Poll Everywhere

discuss

- Data Structures Review: I want to randomly generate a sequence of sorted numbers. To do this, we generate a random number and insert the number so that it remains sorted. Would a LinkedList or an ArrayList work better?
 - Assume we only use a Linear search!

e.g. if I have sequence [5, 9, 23] and I randomly generate 12, I will insert 12 between 9 and 23

- ❖ Part 2: Let's say we take the list from part 1, randomly generate an index and remove that index from the sequence until it is empty. Would this be faster on a LinkedList or an ArrayList?
 - Assume we only use a Linear search!

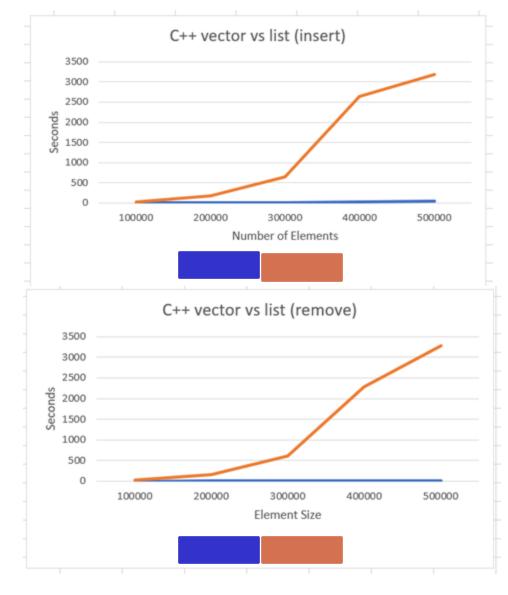
Lecture Outline

- Ideology
- Locality
- Buffering
- Caches

Answer:

- Run using c++:
- Terminology
 - ❖ Vector == ArrayList
 - List == LinkedList

 On Element size from 100,000 -> 500,000



Data Access Time

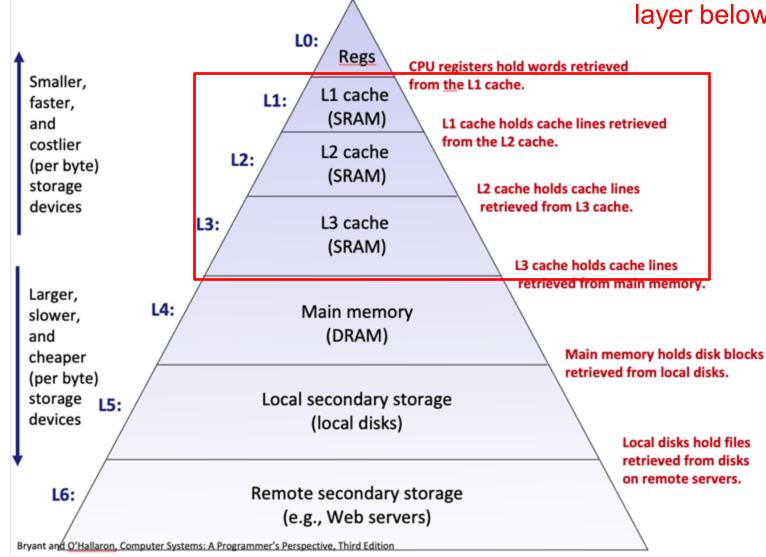
- Data is stored on a physical piece of hardware
- The distance data must travel on hardware affects how long it takes for that data to be processed
- Example: data stored closer to the CPU is quicker to access
 - We see this already with registers. Data in registers is stored on the chip and is faster to access than data in RAM

Memory Hierarchy so far

- So far, we know of three places where we store data
 - CPU Registers
 - Small storage size
 - Quick access time
 - Physical Memory
 - In-between registers and disk
 - Disk
 - Massive storage size
 - Long access time
- (Generally) as we go further from the CPU, storage space goes up, but access times increase

Memory Hierarchy

Each layer can be thought of as a "cache" of the layer below

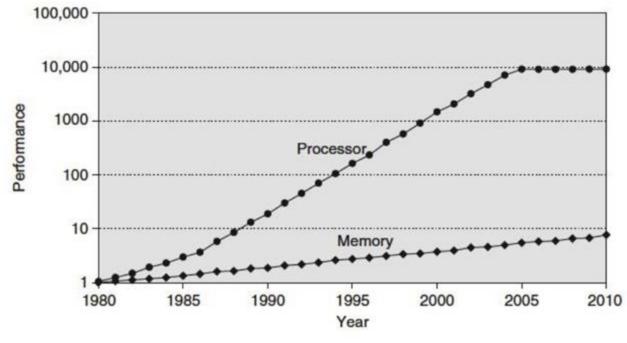


Cache

Pronounced "cash"

- English: A hidden storage space for equipment, weapons, valuables, supplies, etc.
- Computer: Memory with shorter access time used for the storage of data for increased performance. Data is usually either something frequently and/or recently used.
 - Physical memory is a "Cache" of page frames which may be stored on disk. (Instead of going to disk, we can go to physical memory which is quicker to access)

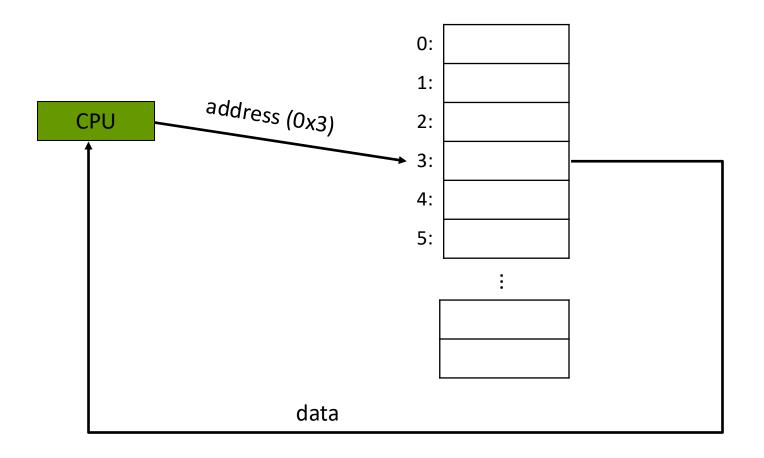
Processor Memory Gap



- Processor speed kept growing ~55% per year
- Time to access memory didn't grow as fast ~7% per year
- Memory access would create a bottleneck on performance
 - It is important that data is quick to access to get better CPU utilization

Memory (as we know it now)

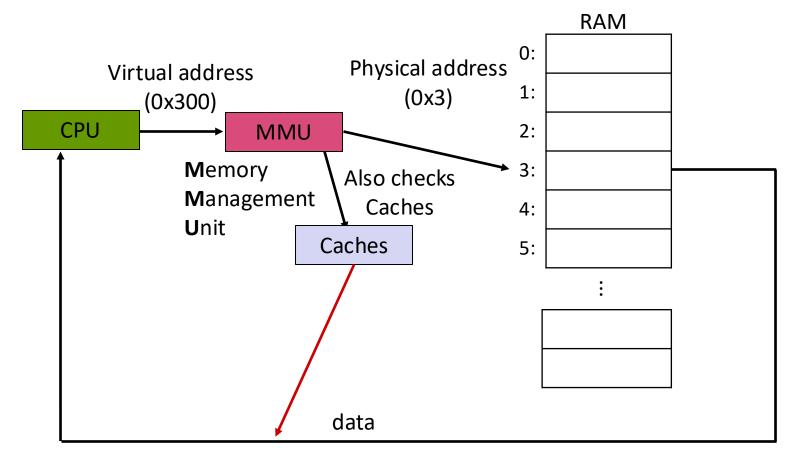
The CPU directly uses an address to access a location in memory



data

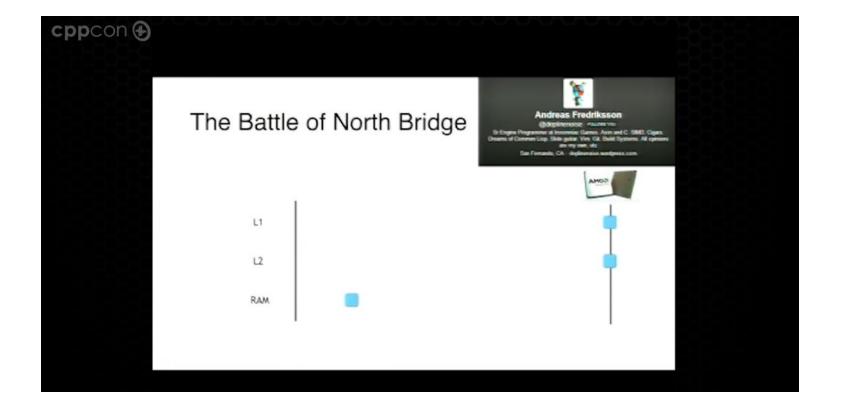
Memory (closer to reality)

Programs don't know about many of things going on under the hood with memory. They send an address to the MMU, and the MMU will help get the



Cache vs Memory Relative Speed

- Animation from Mike Acton's Cppcon 2014 talk on "data oriented design".
 - https://youtu.be/rX0ltVEVjHc?si=MRTeW3taRmRU1fpB&t=1830
 - Animation starts at 30:30, ends 31:07 ish



Cache Performance

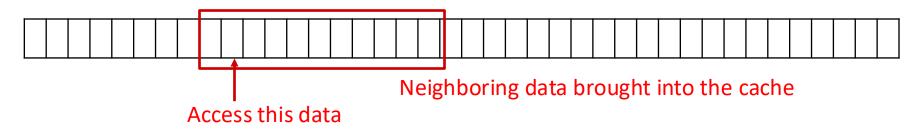
Accessing data in the cache allows for much better utilization of the CPU

Accessing data <u>not</u> in the cache can cause a bottleneck: CPU would have to wait for data to come from memory.

How is data loaded into a Cache?

Cache Lines

Imagine memory as a big array of data:



- We can split memory into 64-byte "lines" or "blocks" (64 bytes on most architectures)
- When we access data at an address, we bring the whole cache line (cache block) into the L1 Cache
 - Data next to address access is thus also brought into the cache!

Principle of Locality

- The tendency for the CPU to access the same set of memory locations over a short period of time
- Two main types:
 - Temporal Locality: If we access a portion of memory, we will likely reference it again soon
 - **Spatial Locality**: If we access a portion of memory, we will likely reference memory close to it in the near future.

Caches take advantage of these tendencies to help with cache management

Cache Replacement Policy

- Caches are small and can only hold so many cache lines inside it.
- When we access data not in the cache, and the cache is full, we must evict an existing entry.
- ❖ When we access a line, we can do a quick calculation on the address to determine which entry in the cache we can store it in. (Depending on architecture, 1 to 12 possible slots in the cache)
 - Cache's typically follow an LRU (Least Recently Used) on the entries a line can be stored in

LRU (Least Recently Used)

- If a cache line is used recently, it is likely to be used again in the near future
- Use past knowledge to predict the future
- Replace the cache line that has had the longest time since it was last used

Back to the Poll Questions

Data Structures Review: I want to randomly generate a sequence of sorted numbers. To do this, we generate a random number and insert the number so that it remains sorted. Would a LinkedList or an ArrayList work better?

Part 2: Let's say we take the list from part 1, randomly generate an index and remove that index from the sequence until it is empty. Would this be faster on a LinkedList or an ArrayList?

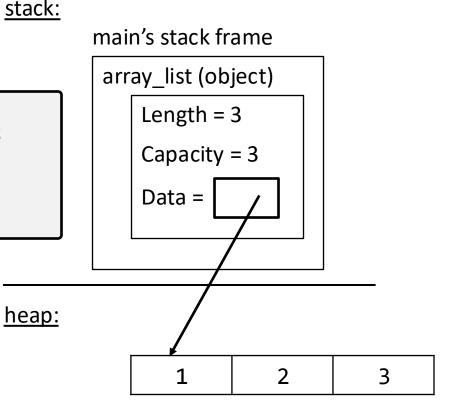
Data Structure Memory Layout

 Important to understanding the poll questions, we understand the memory layout of these data structures

ArrayList In C++:

```
int main() {
  vector<int> array_list {1, 2, 3};
  // ...
}
```

Elements are next to each other in memory ©



Data Structure Memory Layout

 Important to understanding the poll questions, we understand the memory layout of these data structures

stack: main's stack frame LinkedList In C++: linked_list (object) int main() { Length = 4 list<int> linked_list {1, 2, 3, 4}; tail = head = heap: Elements are not next to each other in memory 🖰

Poll Question: Explanation

- Vector wins in-part for a few reasons:
 - Less memory allocations
 - Integers are next to each other in memory, so they benefit from spatial complexity (and temporal complexity from being iterated through in order)
- Does this mean you should always use vectors?
 - No, there are still cases where you should use lists, but your default in C++, Rust, etc should be a vector
 - If you are doing something where performance matters, your best bet is to experiment try all options and analyze which is better.

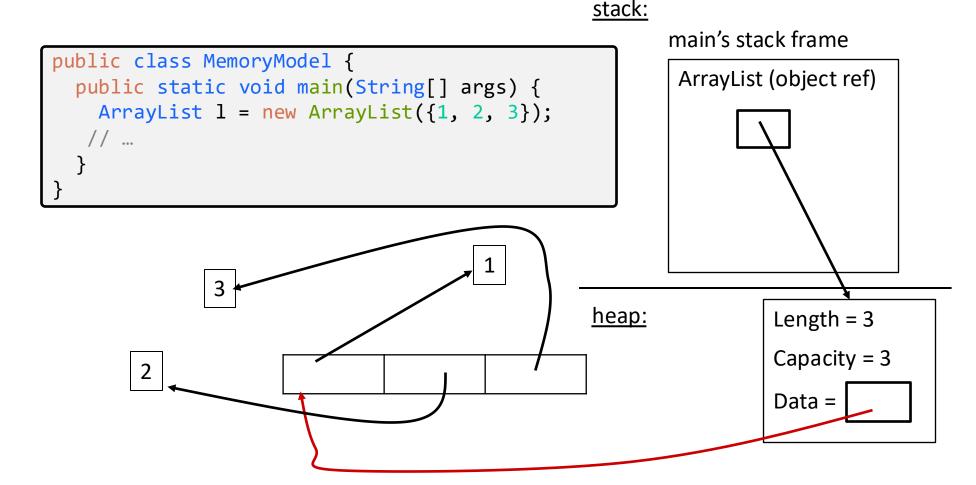
What about other languages?

- In C++ (and C, Rust, Zig ...) when you declare an object, you have an instance of that object. If you declare it as a local variable, it exists on the stack
- In most other languages (including Java, Python, etc.), the memory model is slightly different. Instead, all object variables are object references, that refer to an object on the heap

CIS 4480 Fall 2025

ArrayList in Java Memory Model

In Java, the memory model is slightly different. all object variables are object references, that refer to an object on the heap





pollev.com/cis5480

- Let's say I had a matrix (rectangular two-dimensional array) of integers, and I want the sum of all integers in it
- Would it be faster to traverse the matrix row-wise or column-wise?
 - row-wise (access all elements of the first row, then second)
 - column:-wise (access all elements of the first column, ...)

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



Poll Everywhere

- Let's say I had a matrix (rectangular two-dimensional array) of integers, and I want the sum of all integers in it
- Would it be faster to traverse the matrix row-wise or column-wise?
 - row-wise (access all elements of the first row, then second)
 - column:-wise (access all elements of the first column, ...)

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

Hint: Memory Representation in C & C++

Instruction Cache

- The CPU not only has to fetch data, but it also fetches instructions. There is a separate cache for this
 - which is why you may see something like L1I cache and L1D cache, for Instructions and Data respectively
- Consider the following three fake objects linked in inheritance

```
public class A {
  public void compute() {
     // ...
  }
}
```

```
public class B extends A {
  public void compute() {
     // ...
  }
}

public class C extends A {
  public void compute() {
     // ...
  }
}
```

Instruction Cache

Consider this code

```
public class ICacheExample {
  public static void main(String[] args) {
    ArrayList<A> l = new ArrayList<A>();
    // ...
    for (A item : l) {
        item.compute();
      }
  }
}
```

- When we call item.compute that could invoke A's compute,
 B's compute or C's compute
- Constantly calling different functions, may not utilizes instruction cache well

```
public class A {
  public void compute() {
     // ...
  }
}
```

```
public class B extends A {
  public void compute() {
      // ...
  }
}

public class C extends A {
  public void compute() {
      // ...
  }
}
```

Instruction Cache

Consider this code new code: makes it so we always do
 A.compute() -> B.compute() -> C.compute()

Instruction Cache is happier with this

```
public class ICacheExample {
  public static void main(String[] args) {
    ArrayList<A> la = new ArrayList<A>();
    ArrayList<B> lb = new ArrayList<B>();
    ArrayList<C> lc = new ArrayList<C>();
   for (A item : la) {
       item.compute();
    for (B item : lb) {
       item.compute();
    for (C item : lc) {
       item.compute();
```

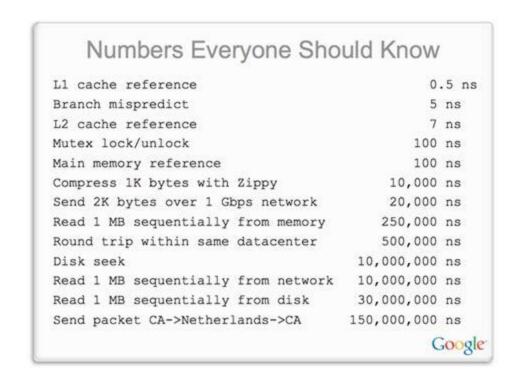
There is a set of numbers that called "numbers everyone you should know"

L18: Locality

- From Jeff Dean in 2009
- Numbers are out of date but the relative orders of magnitude are about the same

- More up to date numbers:
 - https://colin-

scott.github.io/personal_website/research/interactive_latency.html



Some more stuff ©

- Scott Meyers: CPU Caches and Why You Care
 - https://www.youtube.com/watch?v=WDIkqP4JbkE
- What Every Programmer Should Know About Caches:
 - https://people.freebsd.org/~lstewart/articles/cpumemory.pdf