

Caches & Scheduling

Computer Systems Programming, Spring 2024

Instructor: Travis McGaha

TAs:

Ash Fujiyama

Lang Qin

CV Kunjeti

Sean Chuang

Felix Sun

Serena Chen

Heyi Liu

Yuna Shao

Kevin Bernat

Administrivia

- ❖ HW1 was due this Friday
 - Already out
 - Everything you need has been covered
- ❖ HW2 to be released soon
 - Due after break
 - We expect you to look at it and try some of it (maybe implement one of the threading components) before the exam
 - We do not expect you to work on it over break
- ❖ Travis still recovering from a stomach virus :(

Administrivia

- ❖ Midterm Exam: Wednesday February 28th 7-9 pm in Towne 100
 - Please contact Travis if you cannot make it at that time

- ❖ Tentative Final Exam Time:
 - Friday May 5th 3pm-5pm in Towne 100
 - Not confirmed, but this is likely it



pollev.com/tqm

❖ Any questions?

Lecture Outline

- ❖ **Intro to Caches**
- ❖ Scheduling

 **Poll Everywhere**pollev.com/tqm

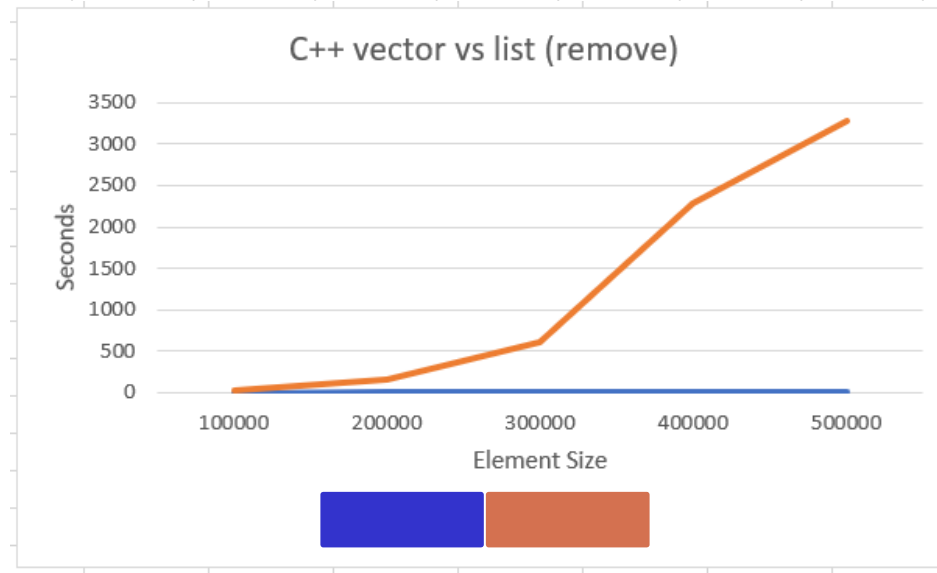
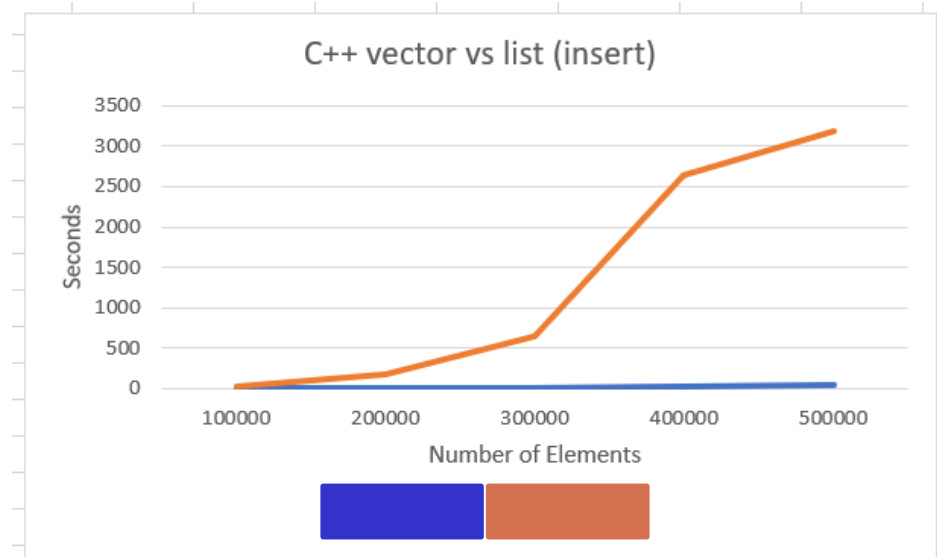
- ❖ Data Structures Review: I want to randomly generate a sequence of sorted numbers. To do this, we generate a random number and insert the number so that it remains sorted. Would a LinkedList or an ArrayList work better?

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?

Answer:

- ❖ I ran this in C++ on this laptop:
- ❖ Terminology
 - Vector == ArrayList
 - List == LinkedList
- ❖ On Element size from 100,000 -> 500,000

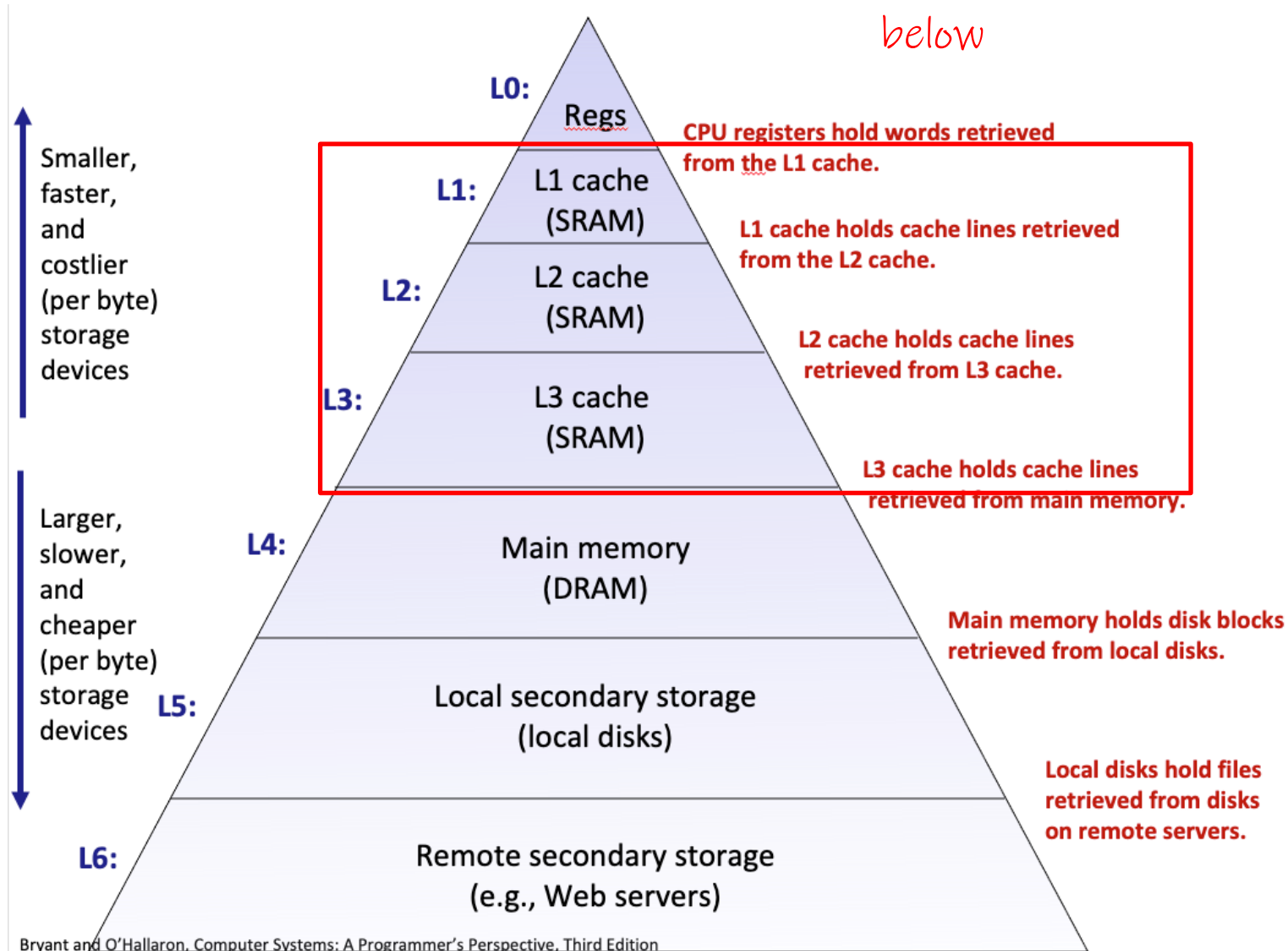


Data Access Time

- ❖ 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 registers

Memory Hierarchy

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

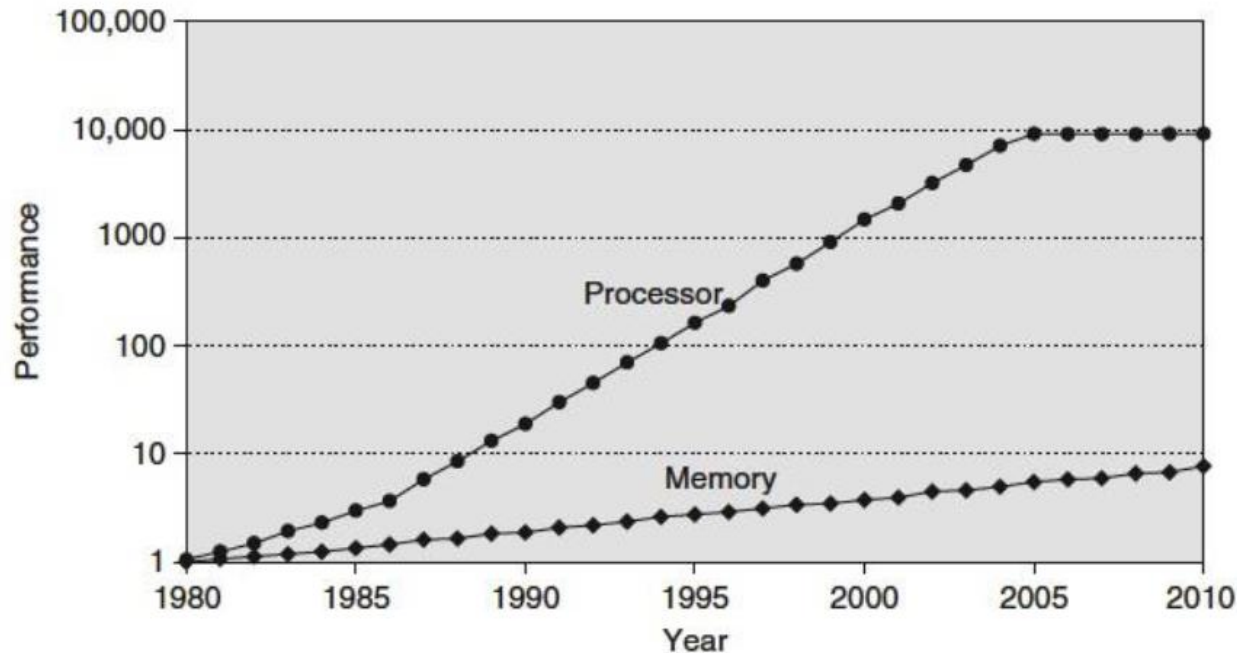


Memory Hierarchy so far

- ❖ 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

Processor Memory Gap



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

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)

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

cppcon

The Battle of North Bridge

Andreas Fredriksson
@afredriksson

Dr. Engine Programmer at Electronic Arts, Ages and C, SMD, Cigars, Director of Common Lisp, Shell games, Vex, G4, Build Systems, All aspects of my own, etc.
San Francisco, CA - afredriksson@electronicarts.com

L1

L2

RAM

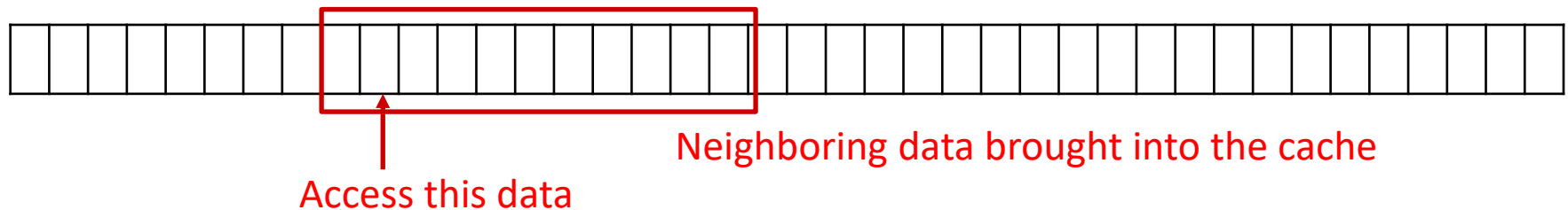
AMD

Cache Performance

- ❖ Accessing data in the cache allows for much better utilization of the CPU
- ❖ Accessing data **not** 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)
 - This means bottom 6 bits of an address are the offset into a line
 - The top 58 bits of the address specify the “line” number
- ❖ 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!

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

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

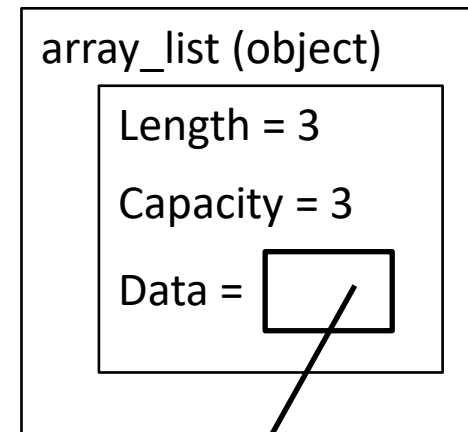
- ❖ Important to understanding the poll questions, we understand the memory layout of these data structures

- ❖ ArrayList In C++:

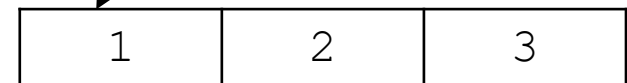
```
int main() {
    vector<int> array_list {1, 2, 3};
    // ...
}
```

stack:

main's stack frame



heap:



Data Structure Memory Layout

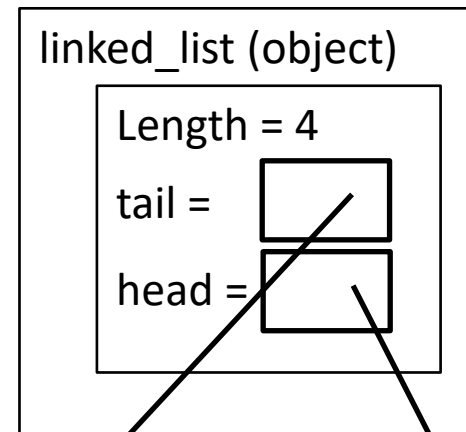
- ❖ Important to understanding the poll questions, we understand the memory layout of these data structures

stack:

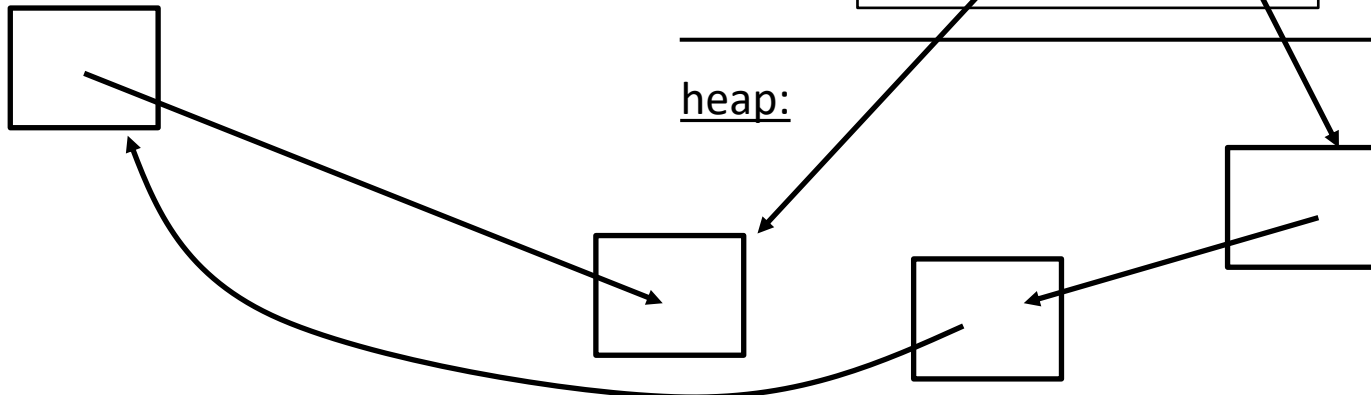
❖ LinkedList In C++:

```
int main() {
    list<int> linked_list {1, 2, 3, 4};
    // ...
}
```

main's stack frame



heap:



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 locality (and temporal locality 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?

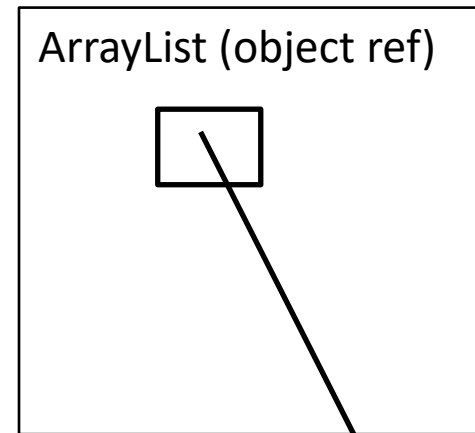
- ❖ In C++ (and C, Rust, Zig ...) when you declare an object, you have an instance of that object. If you declare it as a local variable, it **exists on the stack**
- ❖ In most other languages (including Java, Python, etc.), the memory model is slightly different. Instead, **all object variables are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

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

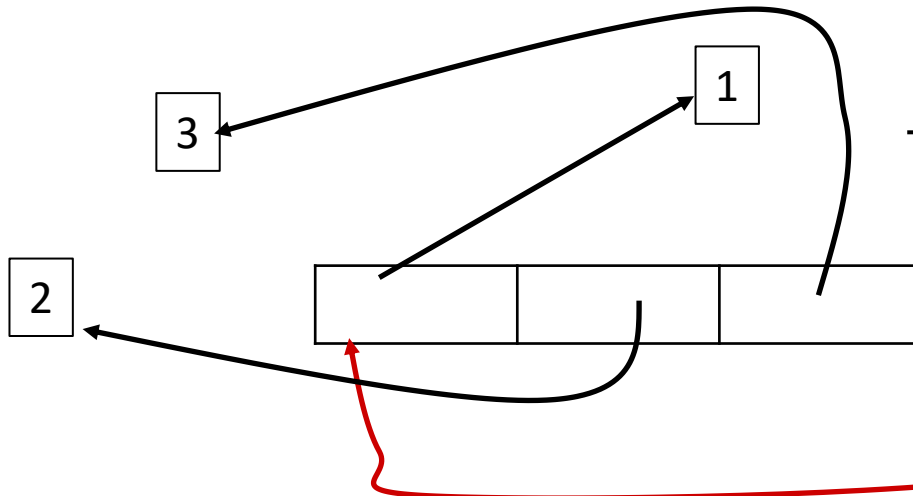
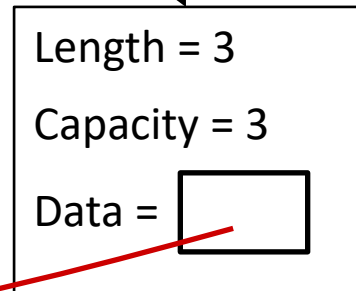
stack:

main's stack frame



```
public class MemoryModel {
    public static void main(String[] args) {
        ArrayList l = new ArrayList({1, 2, 3});
        // ...
    }
}
```

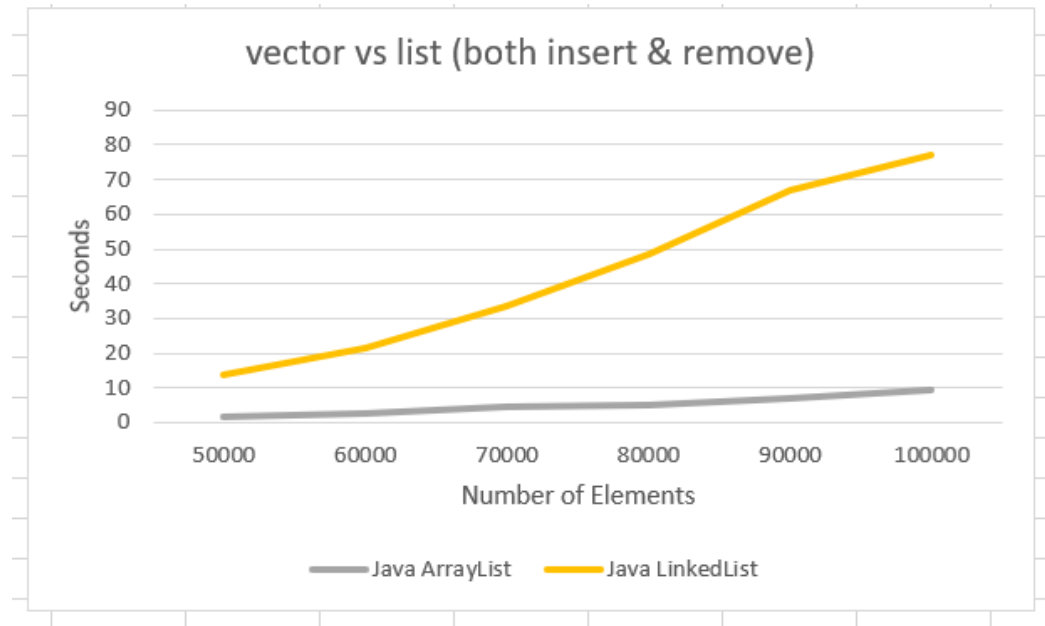
heap:



Does Caching apply to Java?

❖ I believe so, yes. Doing the same experiment in java got:

❖ Note: did this on smaller number of elements.
50,000 -> 100,000





Poll Everywhere

pollev.com/tqm

- ❖ 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**pollev.com/tqm

- ❖ 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, ...)

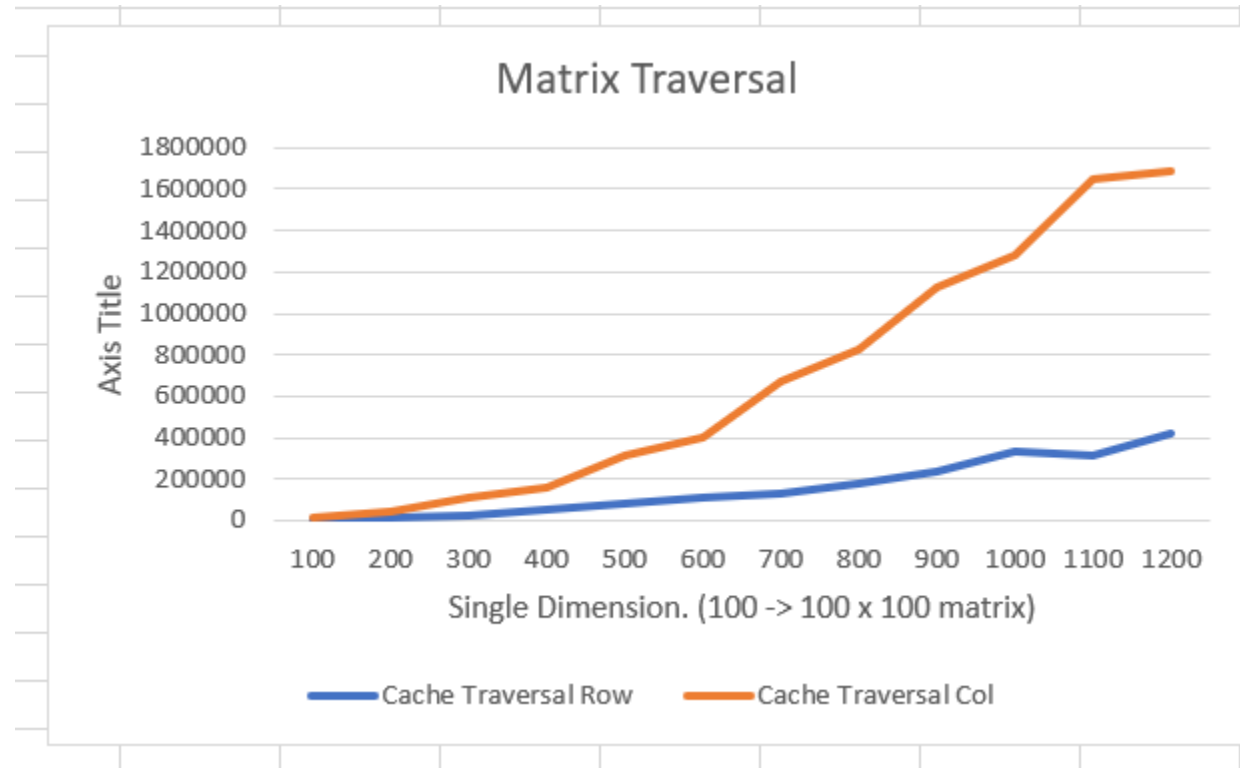
1	5	8	10
11	2	6	9
14	12	3	7
0	15	13	4

Hint: Memory Representation in C & C++

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

Experiment Results

❖ I ran this in C:



❖ Row traversal is better since it means you can take advantage of the cache

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();
        }
    }
}
```

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

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

```
public class C extends A {
    public void 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 utilize instruction cache well

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();
        }
    }
}
```

Lecture Outline

- ❖ Intro to Caches
- ❖ **Scheduling**
 - **FCFS**
 - **SJF**
 - **RR**
 - **RR Variants**

OS as the Scheduler

- ❖ The scheduler is code that is part of the kernel (OS)

- ❖ The scheduler runs when a thread:
 - starts (“arrives to be scheduled”),
 - Finishes
 - Blocks (e.g., waiting on something, usually some form of I/O)
 - Has run for a certain amount of time

- ❖ It is responsible for scheduling threads
 - Choosing which one to run
 - Deciding how long to run it

Scheduler Terminology

- ❖ The scheduler has a scheduling algorithm to decide what runs next.

- ❖ Algorithms are designed to consider many factors:
 - Fairness: Every program gets to run
 - Liveness: That “something” will eventually happen
 - Throughput: amount of work completed over an interval of time
 - Wait time: Average time a “task” is “alive” but not running
 - Turnaround time: time between task being ready and completing
 - Response time: time it takes between task being ready and when it can take user input
 - Etc...

Goals

- ❖ The scheduler will have various things to prioritize
- ❖ Some examples:
 - ❖ Minimizing wait time
 - Get threads started as soon as possible
 - ❖ Minimizing latency
 - Quick response times and task completions are preferred
 - ❖ Maximizing throughput
 - Do as much work as possible per unit of time
 - ❖ Maximizing fairness
 - Make sure every thread can execute fairly
- ❖ These goals depend on the system and can conflict

Scheduling: Other Considerations

- ❖ It takes time to context switch between threads
 - Could get more work done if thread switching is minimized
- ❖ Scheduling takes resources
 - It takes time to decide which thread to run next
 - It takes space to hold the required data structures
- ❖ Different tasks have different priorities
 - Higher priority tasks should finish first

Types of Scheduling Algorithms

- ❖ **Non-Preemptive:** if a thread is running, it continues to run until it completes or until it gives up the CPU
 - First come first serve (FCFS)
 - Shortest Job First (SJF)

- ❖ **Preemptive:** the thread may be interrupted after a given time and/or if another thread becomes ready
 - Round Robin
 - Priority Round Robin
 - ...

First Come First Serve (FCFS)

- ❖ Idea: Whenever a thread is ready, schedule it to run until it is finished (or blocks).
- ❖ Maintain a queue of ready threads
 - Threads go to the back of the queue when it arrives or becomes unblocked
 - The thread at the front of the queue is the next to run

Example of FCFS

1 CPU
 Job 2 arrives slightly after job 1.
 Job 3 arrives slightly after job 2

- ❖ Example workload with three “jobs”:
 Job 1: 24 time units; Job 2: 3 units; Job 3: 3 units

- ❖ FCFS schedule:



- ❖ Total waiting time: $0 + 24 + 27 = 51$
- ❖ Average waiting time: $51/3 = 17$
- ❖ Total turnaround time: $24 + 27 + 30 = 81$
- ❖ Average turnaround time: $81/3 = 27$

 **Poll Everywhere**pollev.com/tqm

- ❖ What are the advantages/disadvantages/concerns with **First Come First Serve**

- ❖ Things a scheduler should prioritize:
 - Minimizing wait time
 - Minimizing Latency
 - Maximizing fairness
 - Maximizing throughput
 - Task priority
 - Cost to schedule things
 - Cost to context Switch

- ❖ Imagine we have 1 core, and tasks of various lengths...

FCFS Analysis

❖ Advantages:

- Simple, low overhead
- Hard to screw up the implementation
- Each thread will DEFINITELY get to run eventually.

❖ Disadvantages

- Doesn't work well for interactive systems
- Throughput can be low due to long threads
- Large fluctuations in average turn around time
- Priority not taken into considerations

Shortest Job First (SJF)

- ❖ Idea: variation on FCFS, but have the tasks with the smallest CPU-time requirement run first
 - Arriving jobs are instead put into the queue depending on their run time, shorter jobs being towards the front
 - Scheduler selects the shortest job (1st in queue) and runs till completion

Example of SJF

1 CPU
 Job 2 arrives slightly after job 1.
 Job 3 arrives slightly after job 2

- ❖ Same example workload with three “jobs”:
 Job 1: 24 time units; Job 2: 3 units; Job 3: 3 units

- ❖ FCFS schedule:



- ❖ Total waiting time: $6 + 0 + 3 = 9$
- ❖ Average waiting time: 3
- ❖ Total turnaround time: $30 + 3 + 6 = 39$
- ❖ Average turnaround time: $39/3 = 13$

 **Poll Everywhere**pollev.com/tqm

- ❖ What are the advantages/disadvantages/concerns with **Shortest Job First**

- ❖ Things a scheduler should prioritize:
 - Minimizing wait time
 - Minimizing Latency
 - Maximizing fairness
 - Maximizing throughput
 - Task priority
 - Cost to schedule things
 - Cost to context Switch

- ❖ Imagine we have 1 core, and tasks of various lengths...

Types of Scheduling Algorithms

- ❖ **Non-Preemptive:** if a thread is running, it continues to run until it completes or until it gives up the CPU
 - First come first serve (FCFS)
 - Shortest Job First (SJF)

- ❖ **Preemptive:** the thread may be interrupted after a given time and/or if another thread becomes ready
 - Round Robin
 - Priority Round Robin
 - ...

Round Robin

- ❖ Sort of a preemptive version of FCFS
 - Whenever a thread is ready, add it to the end of the queue.
 - Run whatever job is at the front of the queue
- ❖ BUT only let it run for a fixed amount of time (quantum).
 - If it finishes before the time is up, schedule another thread to run
 - If time is up, then send the running thread back to the end of the queue.

Example of Round Robin

- ❖ Same example workload:

Job 1: 24 units, Job 2: 3 units, Job 3: 3 units

- ❖ RR schedule with time quantum=2:



- ❖ Total waiting time: $(0 + 4 + 2) + (2 + 4) + (4 + 3) = 19$
 - Counting time spent waiting between each “turn” a job has with the CPU
- ❖ Average waiting time: $19/3$ (~ 6.33)
- ❖ Total turnaround time: $30 + 9 + 10 = 49$
- ❖ Average turnaround time: $49/3$ (~ 16.33)



Poll Everywhere

pollev.com/tqm

- ❖ What are the advantages/disadvantages/concerns with Round Robin
- ❖ Things a scheduler should prioritize:
 - Minimizing wait time
 - Minimizing Latency
 - Maximizing fairness
 - Maximizing throughput
 - Task priority
 - Cost to schedule things
 - Cost to context Switch
- ❖ Imagine we have 1 core, and tasks of various lengths...

Round Robin Analysis

❖ Advantages:

- Still relatively simple
- Can work for interactive systems

❖ Disadvantages

- If quantum is too small, can spend a lot of time context switching
- If quantum is too large, approaches FCFS
- Still assumes all processes have the same priority.

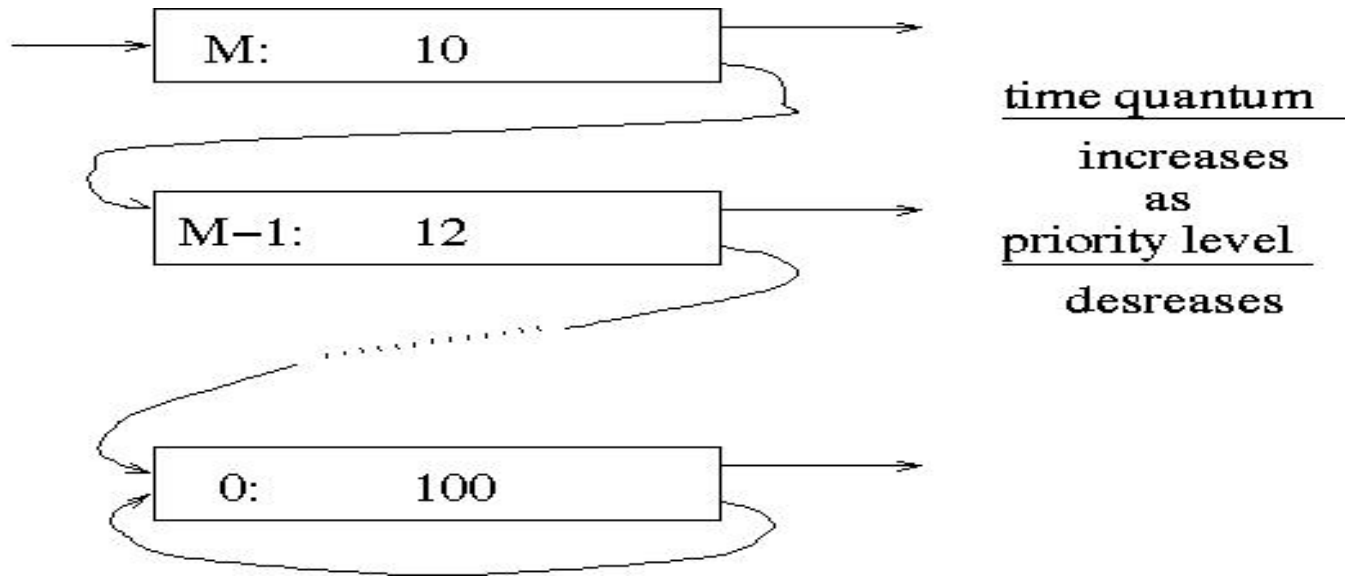
❖ Rule of thumb:

- Choose a unit of time so that most jobs (80-90%) finish in one usage of CPU time

RR Variant: Priority Round Robin

- ❖ Same idea as round robin, but with multiple queues for different priority levels.
- ❖ Scheduler chooses the first item in the highest priority queue to run
- ❖ Scheduler only schedules items in lower priorities if all queues with higher priority are empty.

RR Variant: Multi Level Feedback



- ❖ Each priority level has a ready queue, and a time quantum
- ❖ Thread enters highest priority queue initially, and lower queue with each timer interrupt
- ❖ If a thread voluntarily stops using CPU before time is up, it is moved to the end of the current queue
- ❖ Bottom queue is standard Round Robin
- ❖ Thread in a given queue not scheduled until all higher queues are empty

Multi Level Feedback Analysis

- ❖ Threads with high I/O bursts are preferred
 - Makes higher utilization of the I/O devices
 - Good for interactive programs (keyboard, terminal, mouse is I/O)
- ❖ Threads that need the CPU a lot will sink to lower priority, giving shorter threads a chance to run
- ❖ Still have to be careful in choosing time quantum
- ❖ Also have to be careful in choosing how many layers

Multi Level Feedback Variants: Priority

- ❖ Can assign tasks different priority levels upon initiation that decide which queue it starts in
 - E.g. the scheduler should have higher priority than HelloWorld.java
- ❖ Update the priority based on recent CPU usage rather than overall cpu usage of a task
 - Makes sure that priority is consistent with recent behavior
- ❖ Many others that vary from system to system

Why did we talk about this?

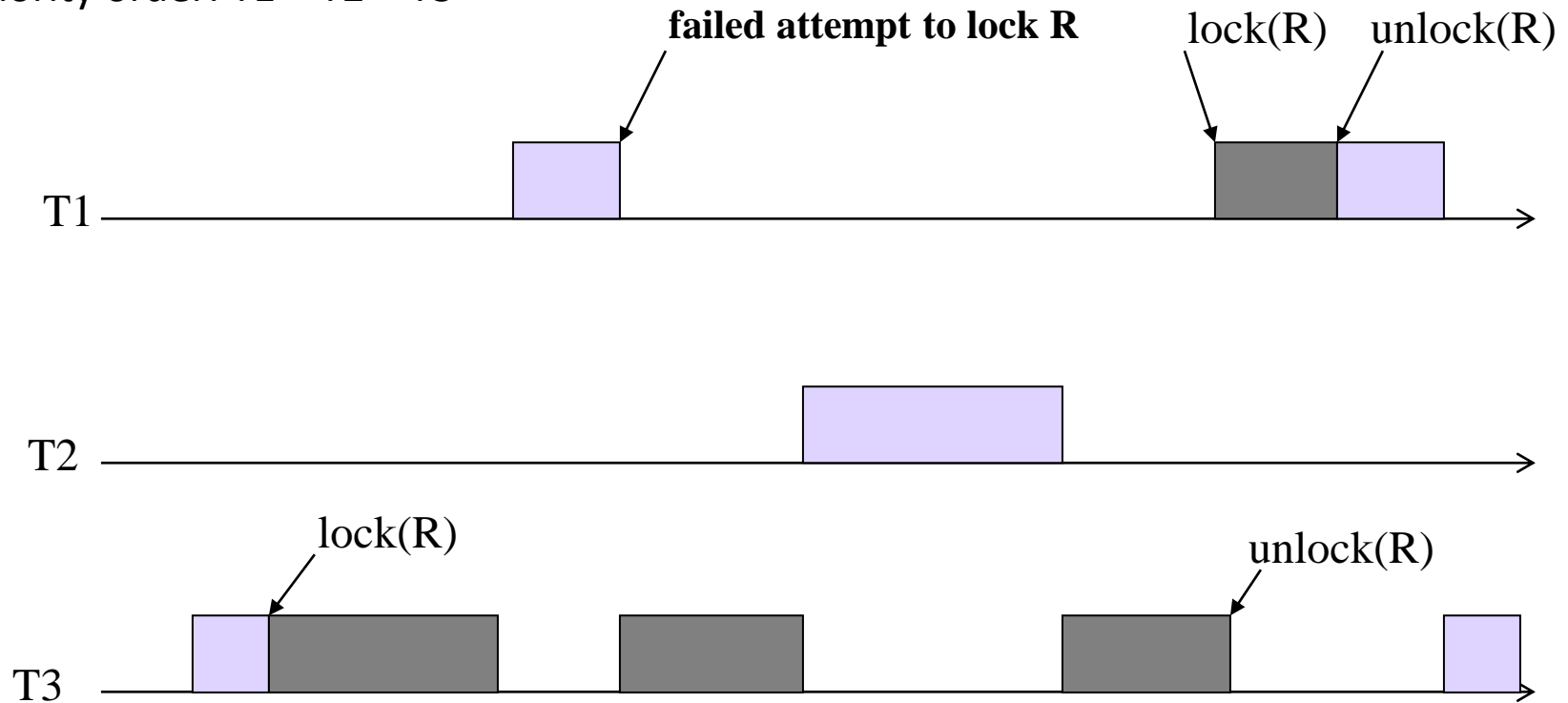
- ❖ Scheduling is fundamental towards how computer can multi-task
- ❖ This is a great example of how “systems” intersects with algorithms :)
- ❖ It shows up occasionally in the real world :)
 - Scheduling threads with priority with shared resources can cause a priority inversion, potentially causing serious errors.

What really happened on Mars Rover Pathfinder, Mike Jones.

<http://www.cs.cornell.edu/courses/cs614/1999sp/papers/pathfinder.html>

The Priority Inversion Problem

Priority order: $T1 > T2 > T3$



T2 is causing a higher priority task T1 wait !

More

- ❖ For those curious, there was a LOT left out

- ❖ RTOS (Real Time Operating Systems)
 - For real time applications
 - CRITICAL that data and events meet defined time constraints
 - Different focus in scheduling. Throughput is de-prioritized

- ❖ Fair-share scheduling
 - Equal distribution across different users instead of by processes

- ❖ Etc.

A little exam practice

Poll Everywhere

pollev.com/tqm

- ❖ The code below has three functions that could be executed in separate threads. Note that these are not thread entry points, just functions used by threads:

- Assume that "lock" has been initialized

- ❖ Thread-1 executes line 8 while Thread-2 executes line 21.

Choose one:

- Could lead to a race condition.
- There is no possible race condition.
- The situation cannot occur.

- ❖ Thread-1 executes line 15 while Thread-2 executes line 15.

Choose one:

- Could lead to a race condition.
- There is no possible race condition.
- The situation cannot occur.

```
1 // global variables
2 pthread_mutex_t lock;
3 int g = 0;
4 int k = 0;
5
6 void fun1() {
7     pthread_mutex_lock(&lock);
8     g += 3;
9     pthread_mutex_unlock(&lock);
10    k++;
11 }
12
13 void fun2(int a, int b) {
14     g += a;
15     a += b;
16     k = a;
17 }
18
19 void fun3() {
20     pthread_mutex_lock(&lock);
21     g = k + 2;
22     pthread_mutex_unlock(&lock);
23 }
```


Poll Everywhere

pollev.com/tqm

- ❖ The code below has three functions that could be executed in separate threads. Note that these are not thread entry points, just functions used by threads:

- Assume that "lock" has been initialized

- ❖ Thread-1 executes line 8 while Thread-2 executes line 21.

Choose one:

- Could lead to a race condition.
- There is no possible race condition.
- The situation cannot occur.

- ❖ Thread-1 executes line 15 while Thread-2 executes line 15.

Choose one:

- Could lead to a race condition.
- There is no possible race condition.
- The situation cannot occur.

```
1 // global variables
2 pthread_mutex_t lock;
3 int g = 0;
4 int k = 0;
5
6 void fun1() {
7     pthread_mutex_lock(&lock);
8     g += 3;
9     pthread_mutex_unlock(&lock);
10    k++;
11 }
12
13 void fun2(int a, int b) {
14     g += a;
15     a += b;
16     k = a;
17 }
18
19 void fun3() {
20     pthread_mutex_lock(&lock);
21     g = k + 2;
22     pthread_mutex_unlock(&lock);
23 }
```



Poll Everywhere

pollev.com/tqm

- ❖ The code below has three functions that could be executed in separate threads. Note that these are not thread entry points, just functions used by threads:

- Assume that "lock" has been initialized

- ❖ Thread-1 executes line 8 while Thread-2 executes line 14

Choose one:

- Could lead to a race condition.
- There is no possible race condition.
- The situation cannot occur.

- ❖ Thread-1 executes line 14 while Thread-2 executes line 16.

Choose one:

- Could lead to a race condition.
- There is no possible race condition.
- The situation cannot occur.

```
1 // global variables
2 pthread_mutex_t lock;
3 int g = 0;
4 int k = 0;
5
6 void fun1() {
7     pthread_mutex_lock(&lock);
8     g += 3;
9     pthread_mutex_unlock(&lock);
10    k++;
11 }
12
13 void fun2(int a, int b) {
14     g += a;
15     a += b;
16     k = a;
17 }
18
19 void fun3() {
20     pthread_mutex_lock(&lock);
21     g = k + 2;
22     pthread_mutex_unlock(&lock);
23 }
```

Poll Everywhere

pollev.com/tqm

- ❖ The code below has three functions that could be executed in separate threads. Note that these are not thread entry points, just functions used by threads:

- Assume that "lock" has been initialized

- ❖ Thread-1 executes line 8 while Thread-2 executes line 14

Choose one:

- Could lead to a race condition.
- There is no possible race condition.
- The situation cannot occur.

- ❖ Thread-1 executes line 14 while Thread-2 executes line 16.

Choose one:

- Could lead to a race condition.
- There is no possible race condition.
- The situation cannot occur.

```
1 // global variables
2 pthread_mutex_t lock;
3 int g = 0;
4 int k = 0;
5
6 void fun1() {
7     pthread_mutex_lock(&lock);
8     g += 3;
9     pthread_mutex_unlock(&lock);
10    k++;
11 }
12
13 void fun2(int a, int b) {
14     g += a;
15     a += b;
16     k = a;
17 }
18
19 void fun3() {
20     pthread_mutex_lock(&lock);
21     g = k + 2;
22     pthread_mutex_unlock(&lock);
23 }
```