

# Threads

## Computer Systems Programming, Spring 2024

**Instructor:** Travis McGaha

**TAs:**

CV Kunjeti

Lang Qin

Felix Sun

Sean Chuang

Heyi Liu

Serena Chen

Kevin Bernat

Yuna Shao

# Administrivia

- ❖ HW1 is due this Friday
  - Already out
  - Everything you need has been covered
  - Auto-grader should be out sometime today
  
- ❖ HW2 to be released over the weekend
  
- ❖ Check-in was due before lecture today



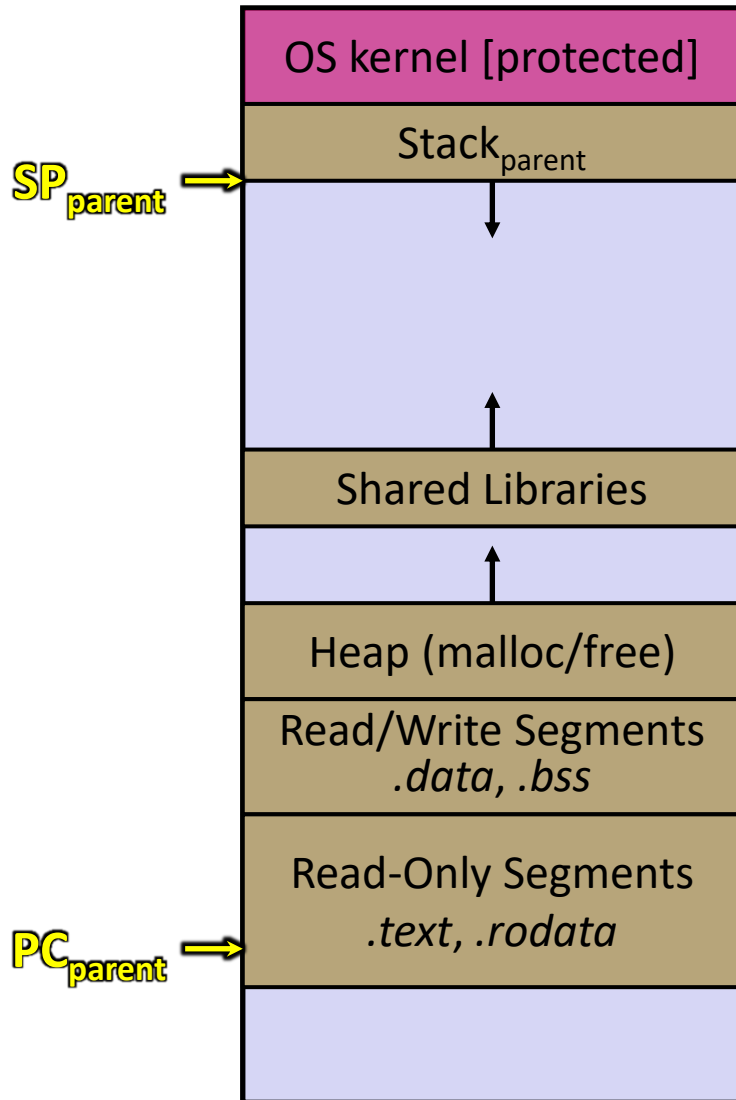
[pollev.com/tqm](https://pollev.com/tqm)

❖ Any questions?

# Lecture Outline

- ❖ **pthread review**
- ❖ Why threads?
  - Parallelism
  - Efficient use of System Resources
- ❖ Shared resources & data races
- ❖ Locks & mutexes

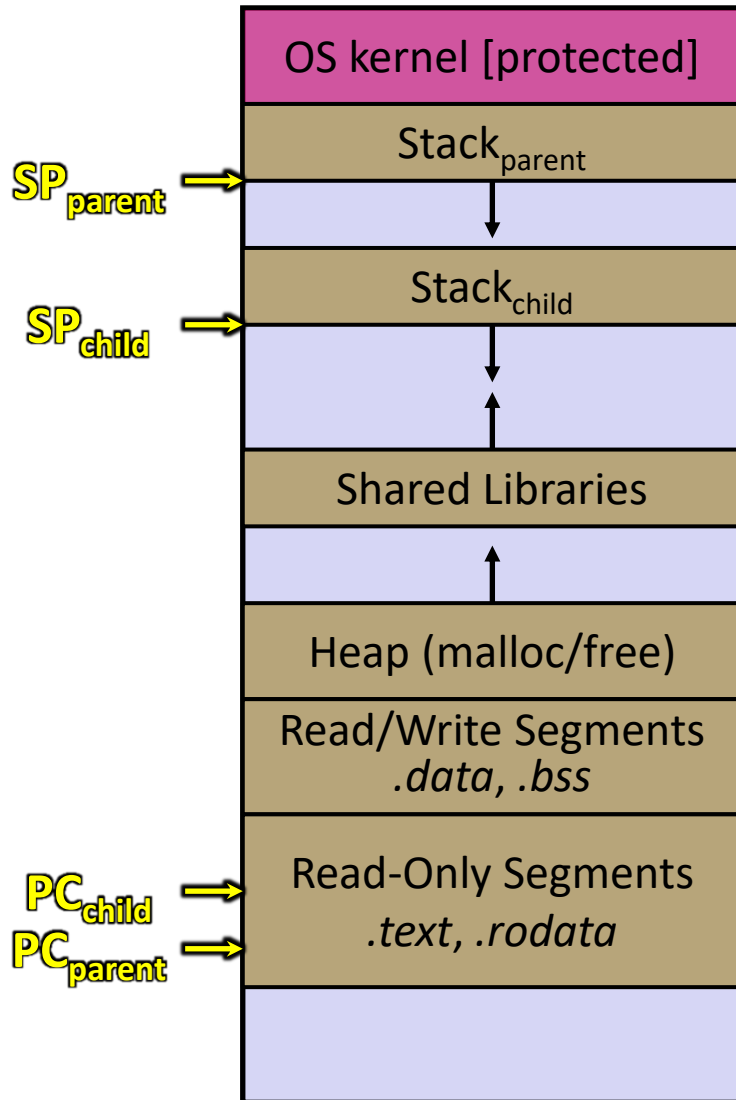
# Single-Threaded Address Spaces



## ❖ Before creating a thread

- One thread of execution running in the address space
  - One PC, stack, SP
- That main thread invokes a function to create a new thread
  - Typically `pthread_create()`

# Multi-threaded Address Spaces



## ❖ After creating a thread

- Two threads of execution running in the address space
  - Original thread (parent) and new thread (child)
  - New stack created for child thread
  - Child thread has its own *values* of the PC and SP
- Both threads share the other segments (code, heap, globals)
  - They can cooperatively modify shared data

# POSIX Threads (pthreads)

- ❖ The POSIX APIs for dealing with threads
  - Declared in `pthread.h`
    - Not part of the C/C++ language
  - To enable support for multithreading, must include `-pthread` flag when compiling and linking with `gcc` command
    - `g++ -g -Wall -std=c++23 -pthread -o main main.c`
  - Implemented in C
    - Must deal with C programming practices and style

# Creating and Terminating Threads

Output parameter.

Gives us a "thread\_descriptor"

```
❖ int pthread_create (
    pthread_t* thread,
    const pthread_attr_t* attr,
    void* (*start_routine) (void*)
    void* arg) ;
```

Function pointer!

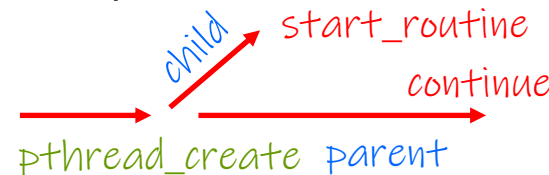
Takes & returns void\* to allow "generics" in C

Argument for the thread function

- Creates a new thread into `*thread`, with attributes `*attr` (`NULL` means default attributes)

- Returns `0` on success and an error number on error (can check against error constants)

- The new thread runs `start_routine` (`arg`)



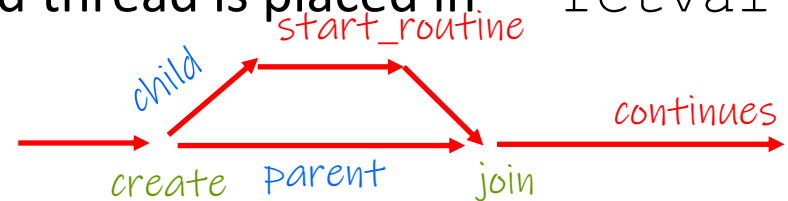


# What To Do After Forking Threads?

❖ `int pthread_join(pthread_t thread, void** retval);`

- Waits for the thread specified by `thread` to terminate
- The thread equivalent of `waitpid()`
- The exit status of the terminated thread is placed in `**retval`

Parent thread waits for child thread to exit, gets the child's return value, and child thread is cleaned up



# Thread Example

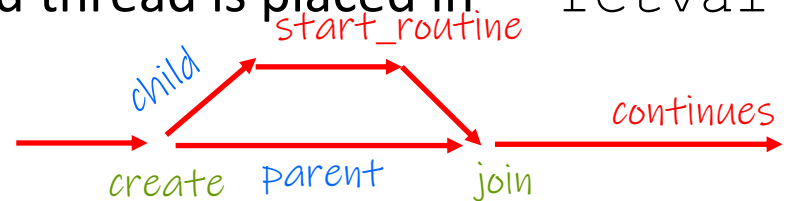
- ❖ See `cthreads.cpp`
  - How do you properly handle memory management?
    - Who allocates and deallocates memory?
    - How long do you want memory to stick around?

# What To Do After Forking Threads?

❖ `int pthread_join(pthread_t thread, void** retval);`

- Waits for the thread specified by `thread` to terminate
- The thread equivalent of `waitpid()`
- The exit status of the terminated thread is placed in `**retval`

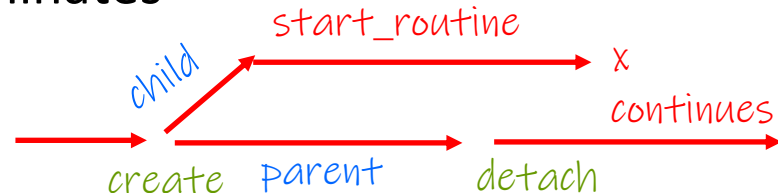
Parent thread waits for child thread to exit, gets the child's return value, and child thread is cleaned up



❖ `int pthread_detach(pthread_t thread);`

- Mark thread specified by `thread` as detached – it will clean up its resources as soon as it terminates

Detach a thread. Thread is cleaned up when it is finished



# Thread Examples

- ❖ See `cthreads.cpp`
  - How do you properly handle memory management?
    - Who allocates and deallocates memory?
    - How long do you want memory to stick around?
  
- ❖ See `exit_thread.cpp`
  - Do we need to join every thread we create?

## Discuss

## ❖ What are all possible outputs of this program?

```

void* thrd_fn(void* arg) {
    int* ptr = reinterpret_cast<int*>(arg);
    cout << *ptr << endl;
    return nullptr;
}

int main() {
    pthread_t thd1{};
    pthread_t thd2{};
    int x = 1;
    pthread_create(&thd1, nullptr, thrd_fn, &x);
    x = 2;
    pthread_create(&thd2, nullptr, thrd_fn, &x);

    pthread_join(thd1, nullptr);
    pthread_join(thd2, nullptr);
}

```

Are these output possible?

-----

1

2

-----

2

2

-----

1

1

-----

2

1

# Visualization

```
int main() {  
    int x = 1;  
    pthread_create(...);  
    x = 2;  
    pthread_create(...);  
  
    pthread_join(...);  
    pthread_join(...);  
}
```

```
thrd_fn() {  
    cout << *ptr ...;  
    return nullptr;  
}
```

```
thrd_fn() {  
    cout << *ptr ...;  
    return nullptr;  
}
```

# Visualization: Memory

- ❖ The variable `x` is shared across all threads.

```
main()
```

```
int x 1
```

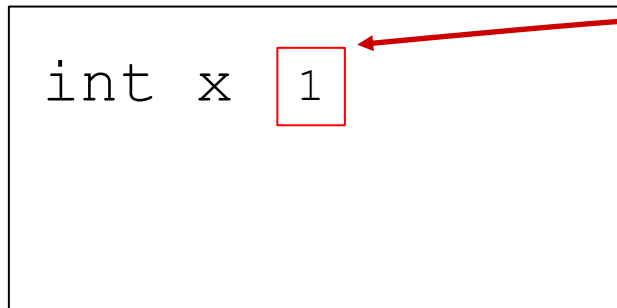
```
int main() {
    int x = 1;
    pthread_create(thd1);
    x = 2;
    pthread_create(thd2);

    pthread_join(thd1);
    pthread_join(thd2);
}
```

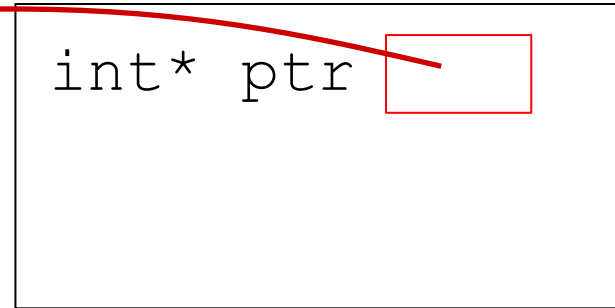
# Visualization: Memory

- ❖ The variable `x` is shared across all threads.

`main()`



`thd1`



```

int main() {
    int x = 1;
    pthread_create(thd1);
    x = 2;
    pthread_create(thd2);

    pthread_join(thd1);
    pthread_join(thd2);
}
    
```

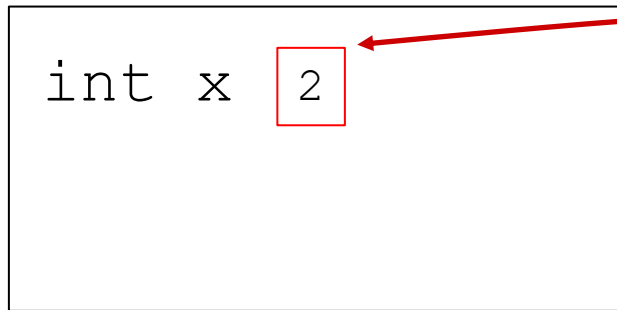
A red arrow points from the `pthread_create(thd1);` line to the `thd1` memory diagram above.



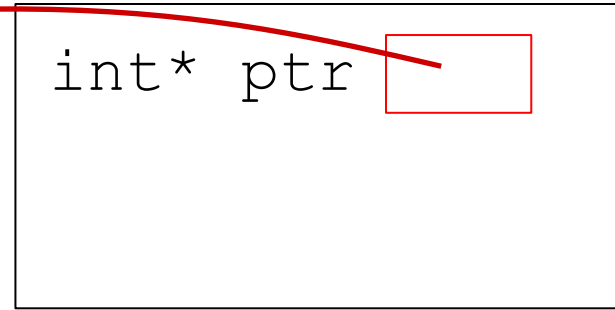
# Visualization: Memory

- ❖ The variable `x` is shared across all threads.

`main()`



`thd1`



```

int main() {
    int x = 1;
    pthread_create(thd1);
    x = 2;
    pthread_create(thd2);

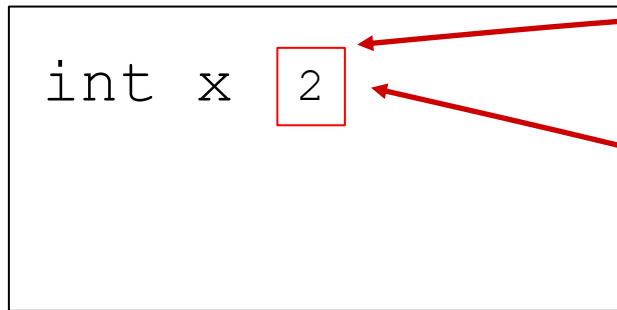
    pthread_join(thd1);
    pthread_join(thd2);
}
    
```

A red arrow points from the `x = 2;` line in the code block to the `2` in the `main()` memory diagram above.

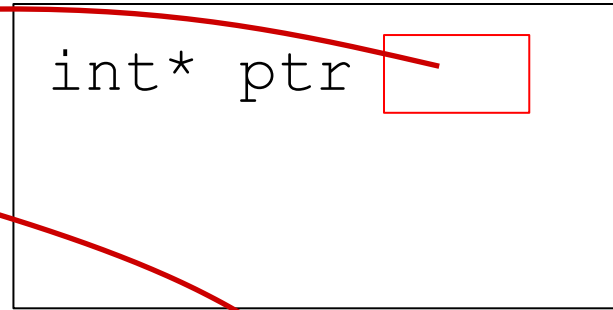
# Visualization: Memory

- ❖ The variable `x` is shared across all threads.

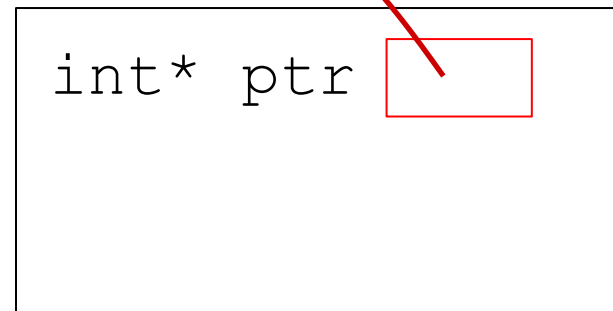
`main()`



`thd1`



`thd2`



```
int main() {
    int x = 1;
    pthread_create(thd1);
    x = 2;
    pthread_create(thd2);

    pthread_join(thd1);
    pthread_join(thd2);
}
```

# Visualization: Ordering

- ❖ Threads run concurrently; we can't be sure of the ordering of things across threads.

`main()`

```
int x = 1
```

`thd1`

`thd2`

# Visualization: Ordering

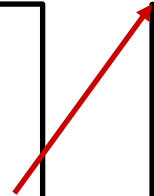
- ❖ Threads run concurrently; we can't be sure of the ordering of things across threads.

main()

```
int x = 1
create thd1
```

thd1

thd2



# Visualization: Ordering

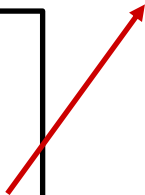
- ❖ Threads run concurrently; we can't be sure of the ordering of things across threads.

`main()`

```
int x = 1
create thd1
```

`thd1`

`thd2`



# Visualization: Ordering

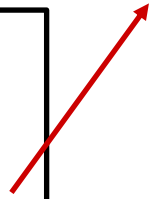
- ❖ Threads run concurrently; we can't be sure of the ordering of things across threads.

main()

```
int x = 1
create thd1
x = 2
```

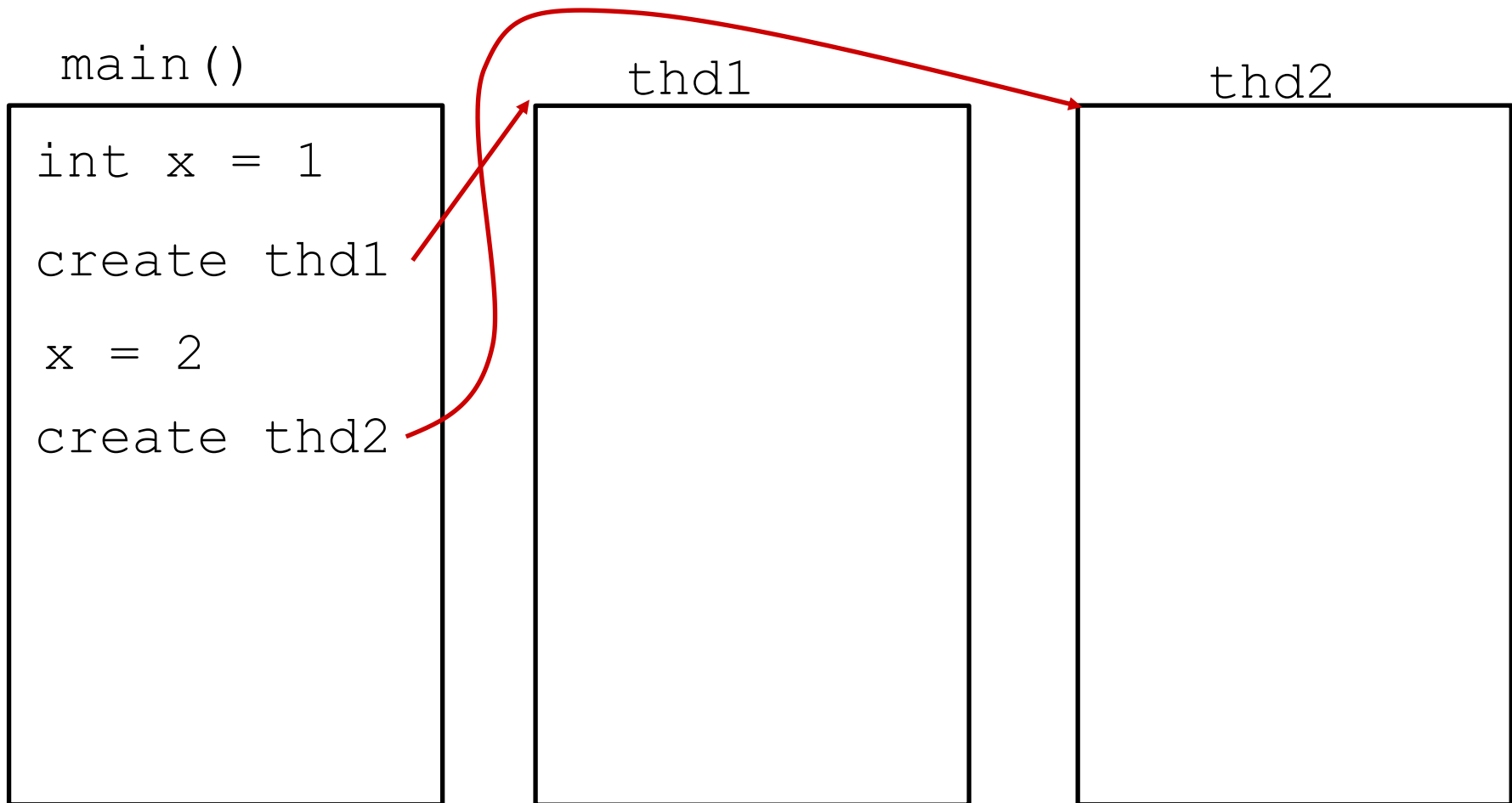
thd1

thd2



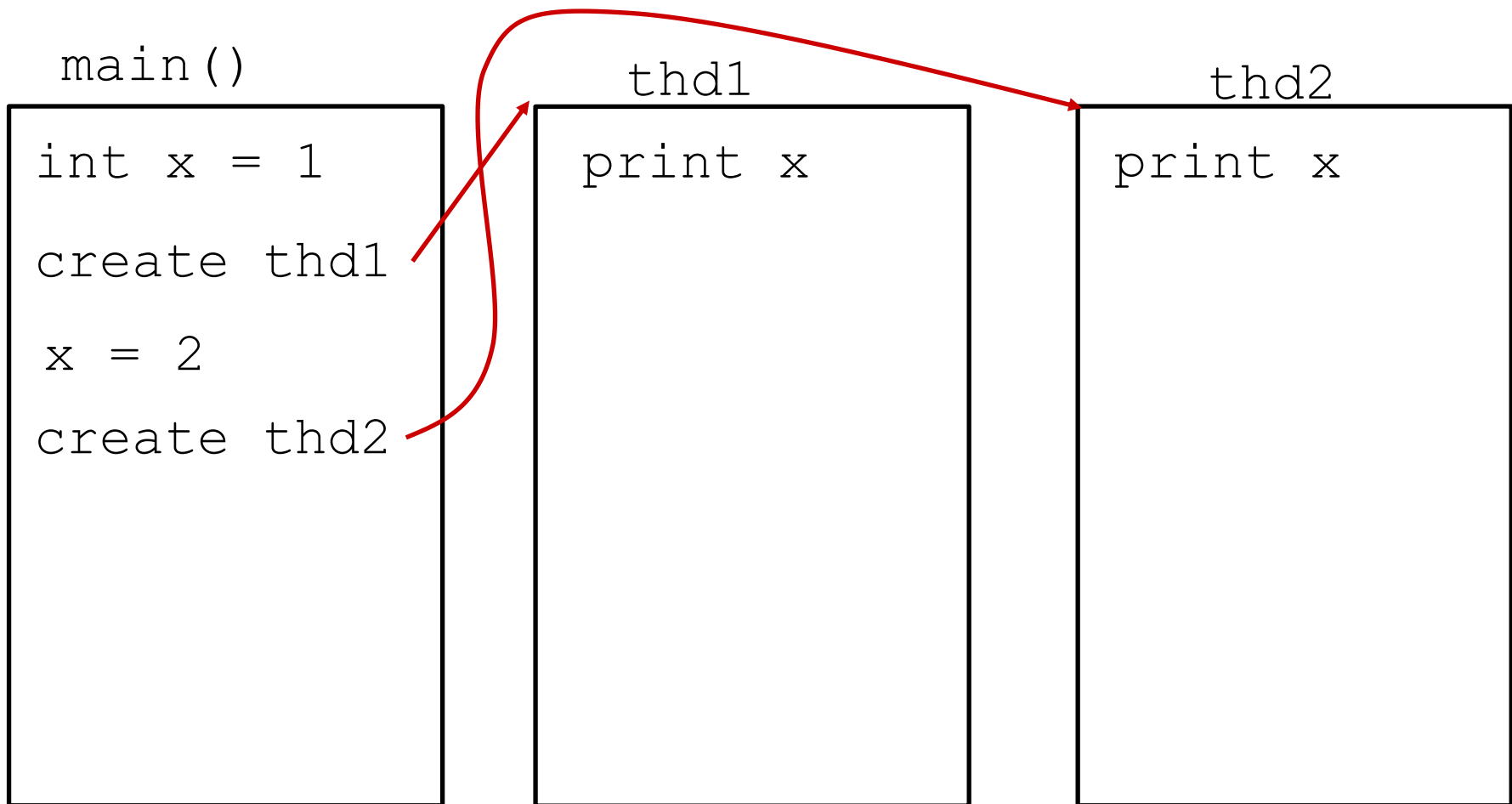
# Visualization: Ordering

- ❖ Threads run concurrently; we can't be sure of the ordering of things across threads.



# Visualization: Ordering

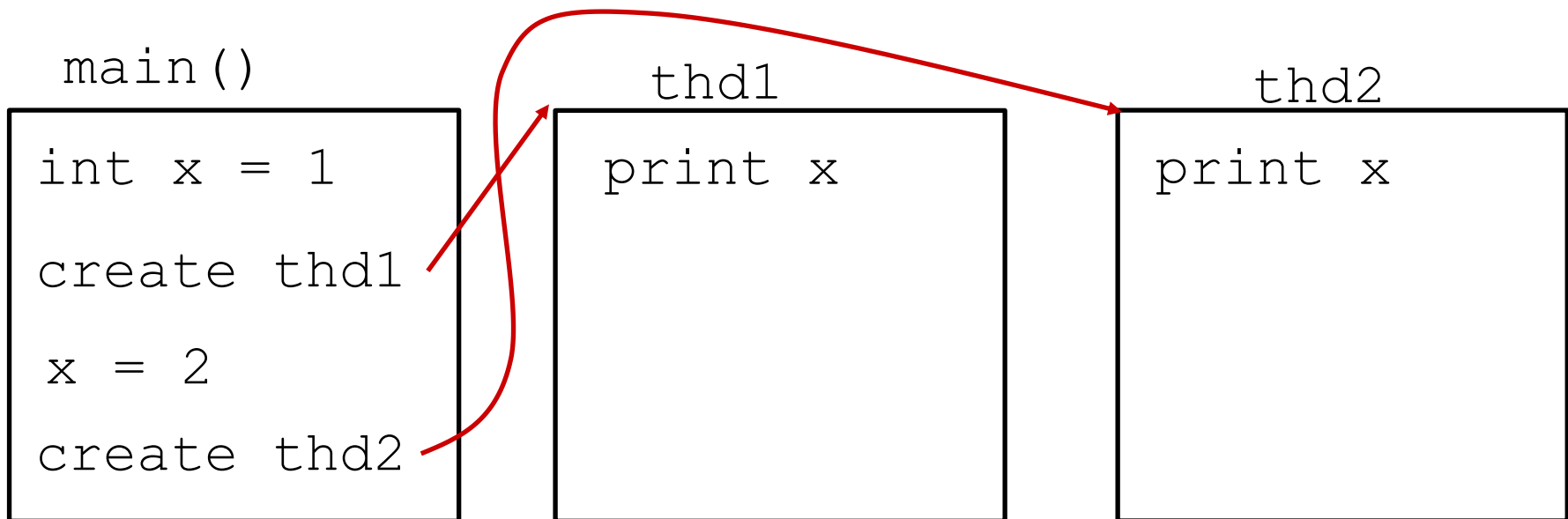
- ❖ Threads run concurrently; we can't be sure of the ordering of things across threads.





# Visualization: Ordering

- ❖ Threads run concurrently; we can't be sure of the ordering of things across threads.



We know that `x` is initialized to 1 before `thd1` is created

We know that `x` is set to 2 and `thd1` is created before `thd2` is created

Anything else that we know? **No**. Beyond those statements, we do not know the ordering of `main` and the threads running.

# Lecture Outline

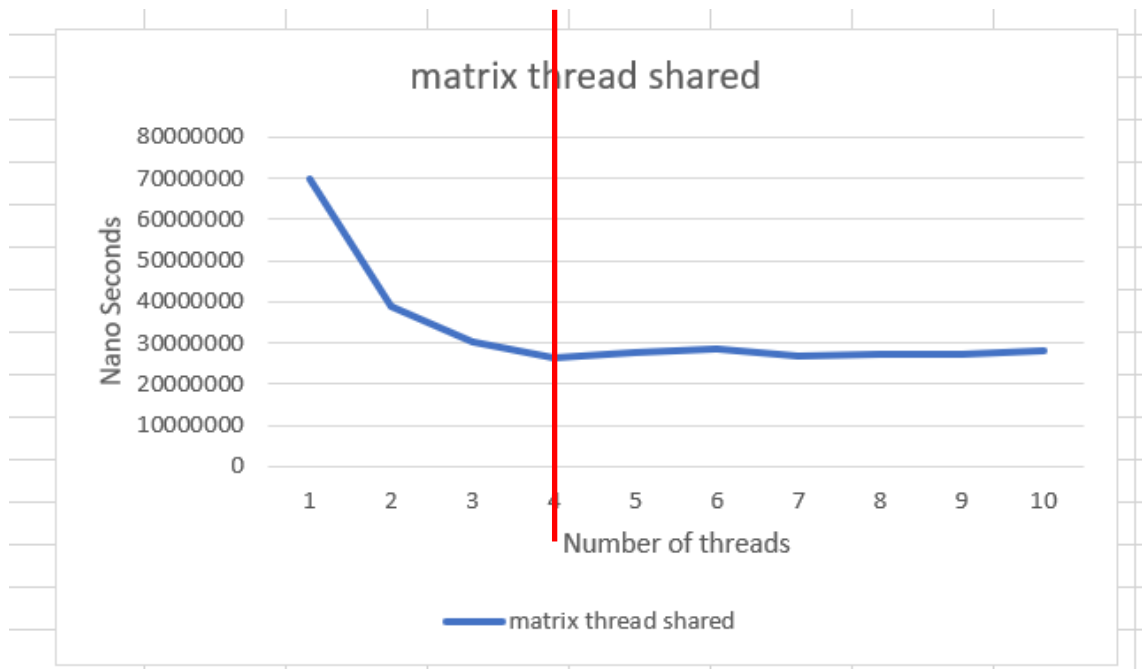
- ❖ pthreads review
- ❖ **Why threads?**
  - **Parallelism**
  - Efficient use of System Resources
- ❖ Shared resources & data races
- ❖ Locks & mutexes

# Parallelism

- ❖ You can gain performance by running things in parallel
  - Each thread can use another core and run code in parallel
- ❖ I have a 3800 x 3800 integer matrix, and I want to count the number of odd integers in the matrix

# Parallelism

- ❖ I have a 3800 x 3800 integer matrix, and I want to count the number of odd integers in the matrix
- ❖ I can speed this up by giving each thread a part of the matrix to check!
  - Works with threads since they share memory



*Diminishing returns*

*After 4 threads, no gain in speed*

*why? Machine run on only has 4 cores*

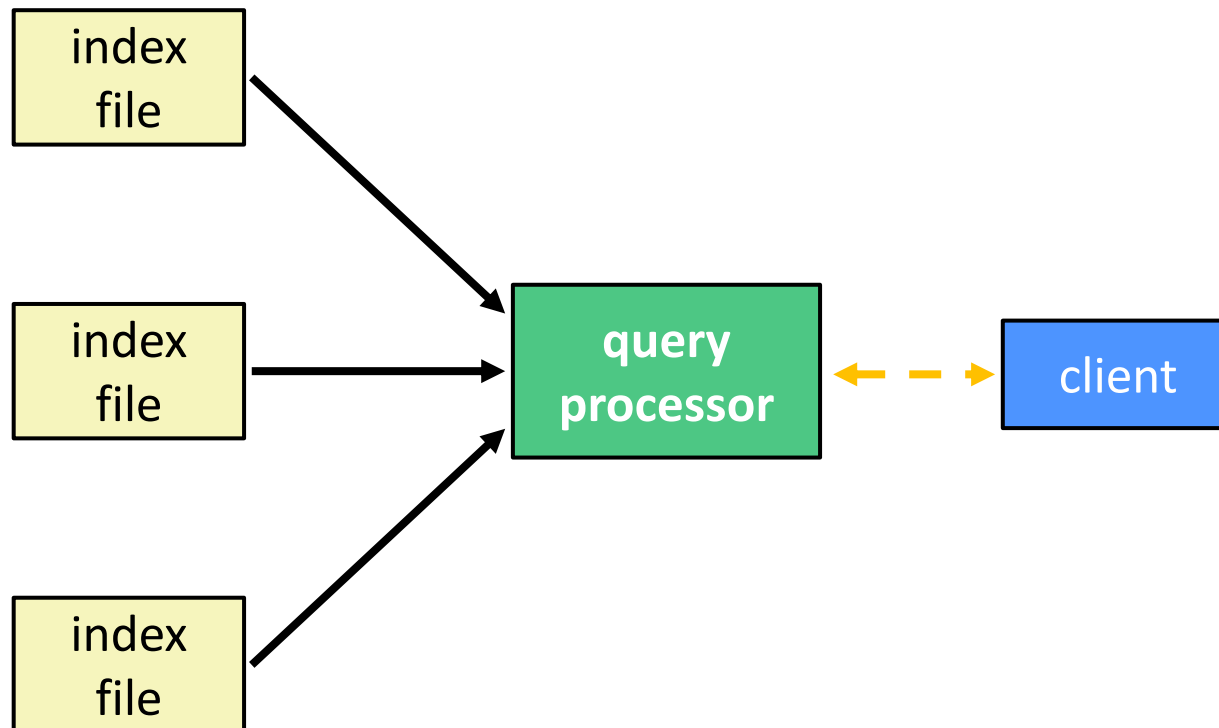
# Lecture Outline

- ❖ pthreads review
- ❖ **Why threads?**
  - Parallelism
  - **Efficient use of System Resources**
- ❖ Shared resources & data races
- ❖ Locks & mutexes

# Building a Web Search Engine

- ❖ We have:
  - A web index
    - A map from *<word>* to *<list of documents containing the word>*
    - This is probably *sharded* over multiple files
  - A query processor
    - Accepts a query composed of multiple words
    - Looks up each word in the index
    - Merges the result from each word into an overall result set

# Search Engine Architecture



# Search Engine (Pseudocode)

```

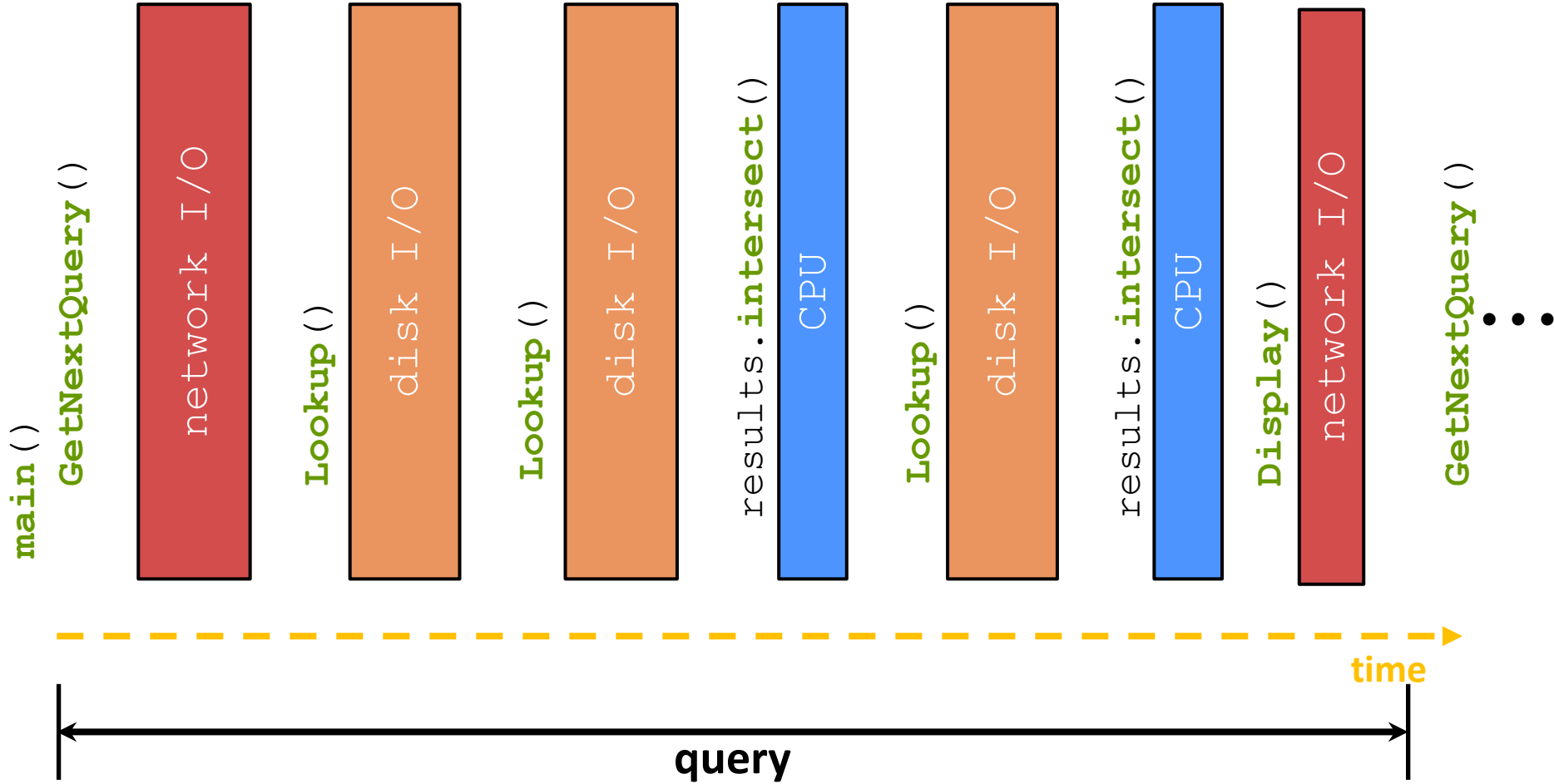
doclist Lookup(string word) {
    bucket = hash(word);
    hitlist = file.read(bucket); ← Disk I/O
    foreach hit in hitlist {
        doclist.append(file.read(hit)); ←
    }
    return doclist;
}

main() {
    SetupServerToReceiveConnections();
    while (1) {
        string query_words[] = GetNextQuery(); ← Network
        results = Lookup(query_words[0]); ← I/O
        foreach word in query[1..n] {
            results = results.intersect(Lookup(word));
        }
        Display(results); ← Network
    }
}
    I/O

```





# Execution Timeline: a Multi-Word Query



# What About I/O-caused Latency?

- ❖ Jeff Dean's "Numbers Everyone Should Know" (LADIS '09)

Numbers Everyone Should Know	
L1 cache reference	0.5 ns
Branch mispredict	5 ns
L2 cache reference	7 ns
Mutex lock/unlock	100 ns
Main memory reference	100 ns
Compress 1K bytes with Zip	10,000 ns
Send 2K bytes over 1 Gbps network	20,000 ns
Read 1 MB sequentially from memory	250,000 ns
Round trip within same datacenter	500,000 ns
Disk seek	10,000,000 ns
Read 1 MB sequentially from network	10,000,000 ns
Read 1 MB sequentially from disk	30,000,000 ns
Send packet CA->Netherlands->CA	150,000,000 ns

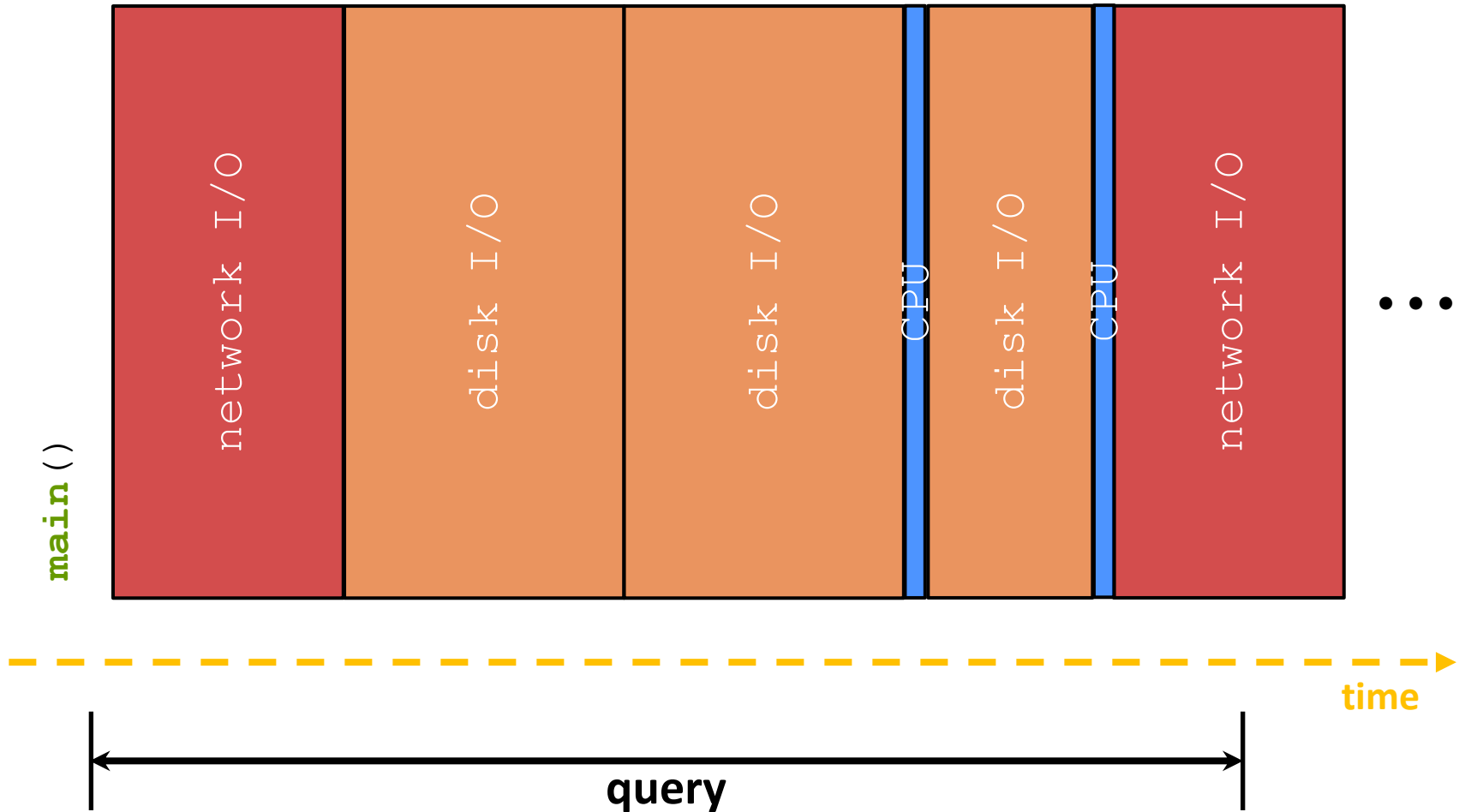


# Execution Timeline: To Scale

Model isn't perfect:

Technically also some cpu usage to setup I/O.

Network output also (probably) won't block program .....



# Multiple (Single-Word) Queries

# is the Query Number

#.a -> GetNextQuery ()

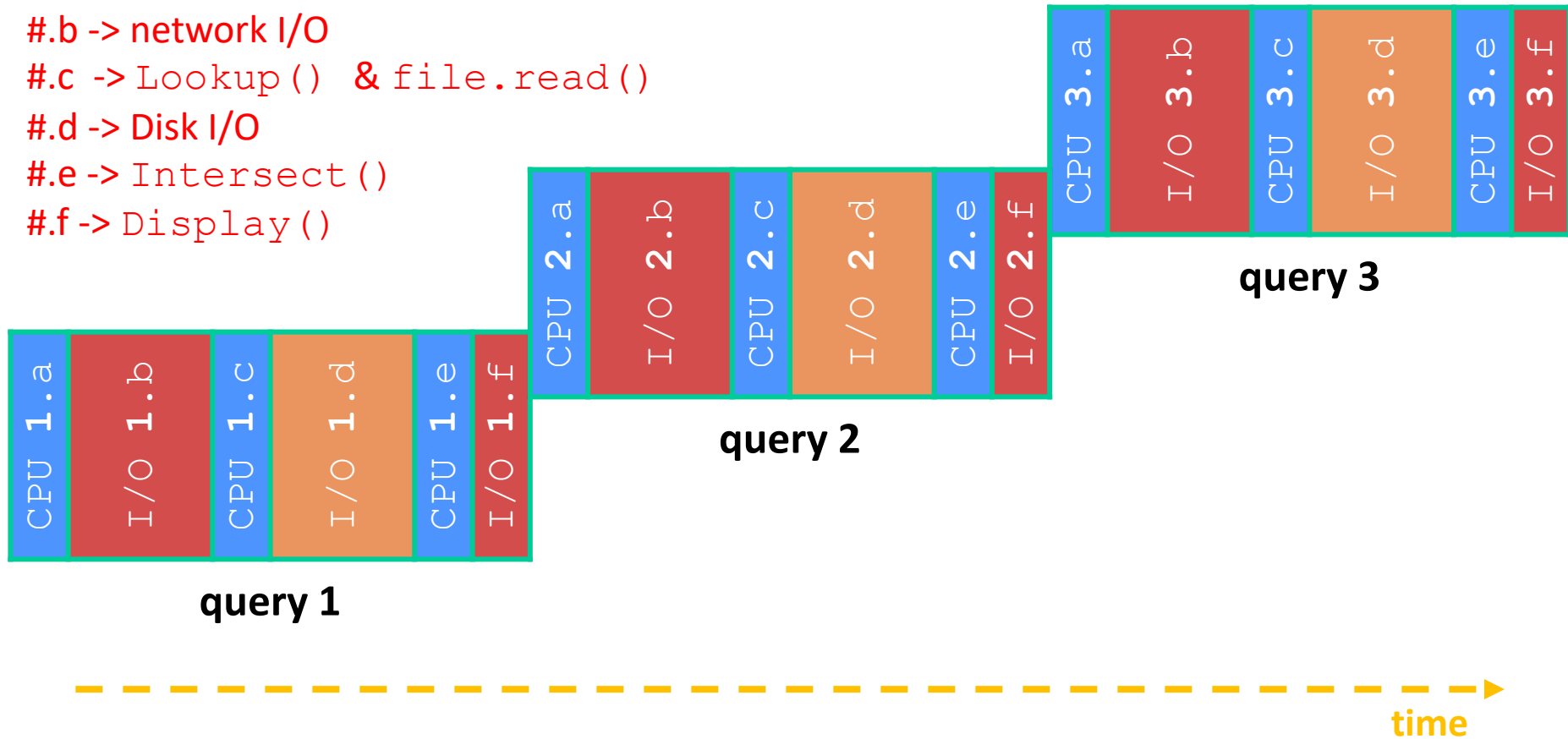
#.b -> network I/O

#.c -> Lookup () & file.read ()

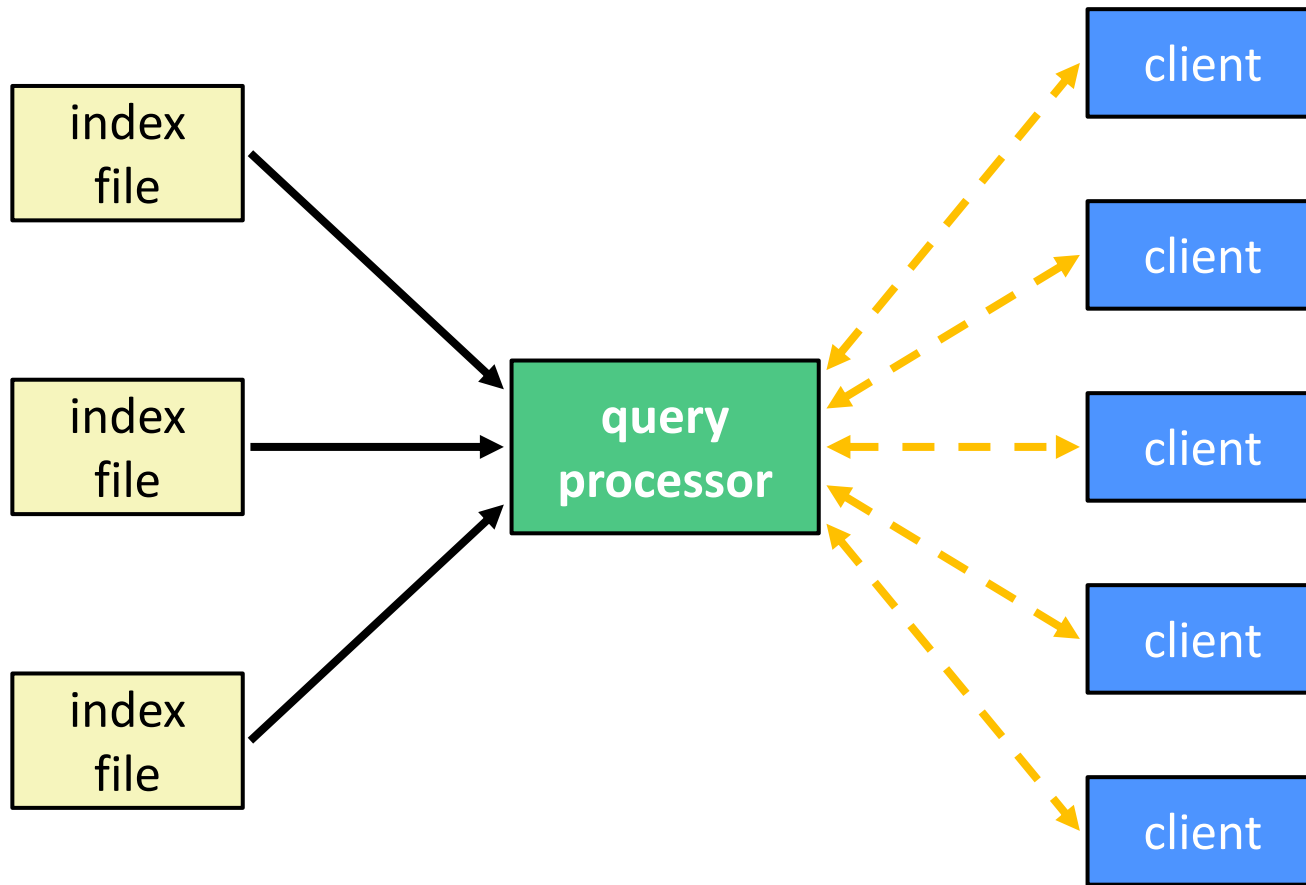
#.d -> Disk I/O

#.e -> Intersect ()

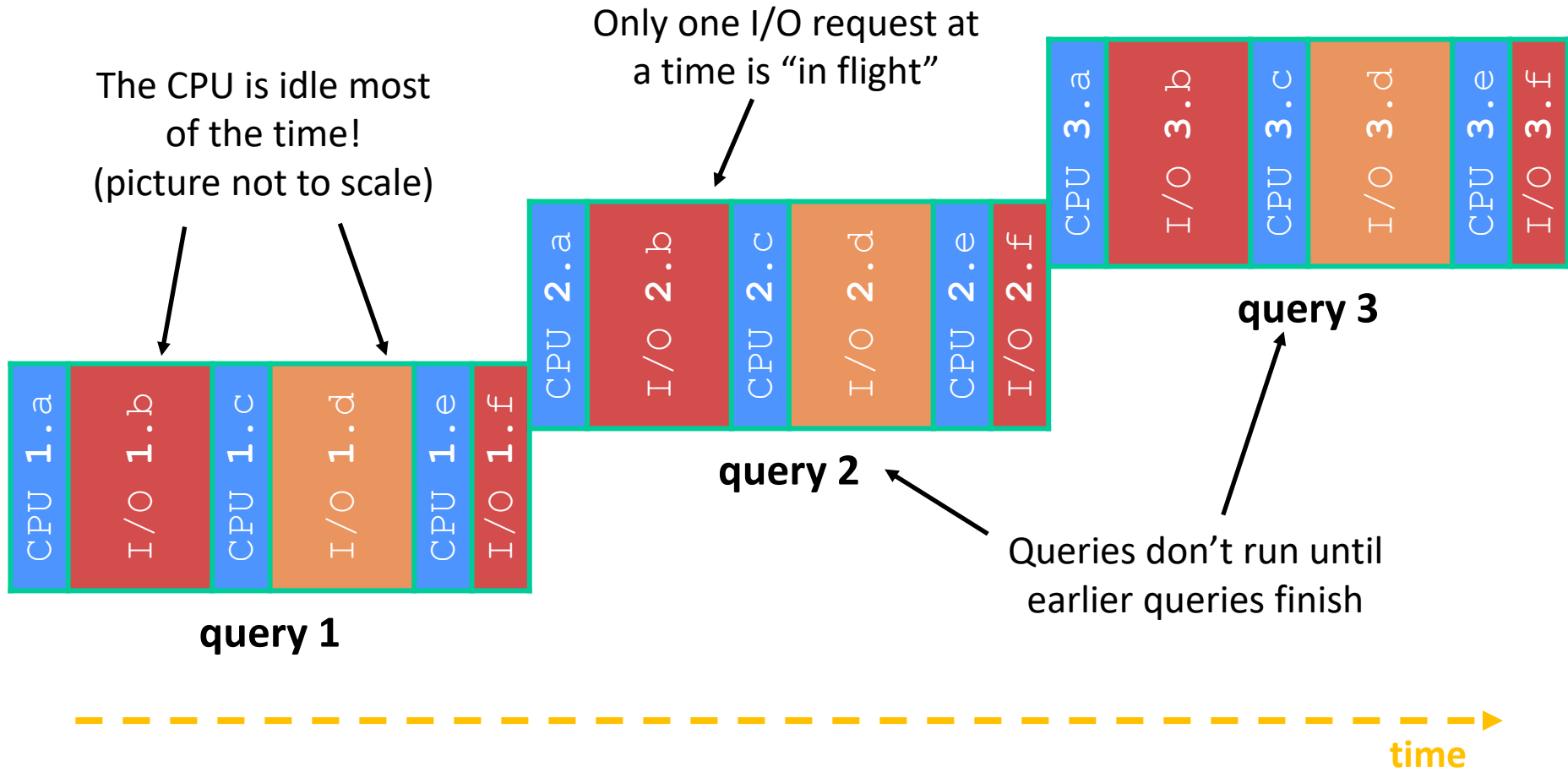
#.f -> Display ()



# Uh-Oh (1 of 2)



# Uh-Oh (2 of 2)



# Sequential Can Be Inefficient

- ❖ Only one query is being processed at a time
  - All other queries queue up behind the first one
  - And clients queue up behind the queries ...
- ❖ Even while processing one query, the CPU is idle the vast majority of the time
  - It is *blocked* waiting for I/O to complete
    - Disk I/O can be very, very slow (10 million times slower ...)
- ❖ At most one I/O operation is in flight at a time
  - Missed opportunities to speed I/O up
    - Separate devices in parallel, better scheduling of a single device, etc.

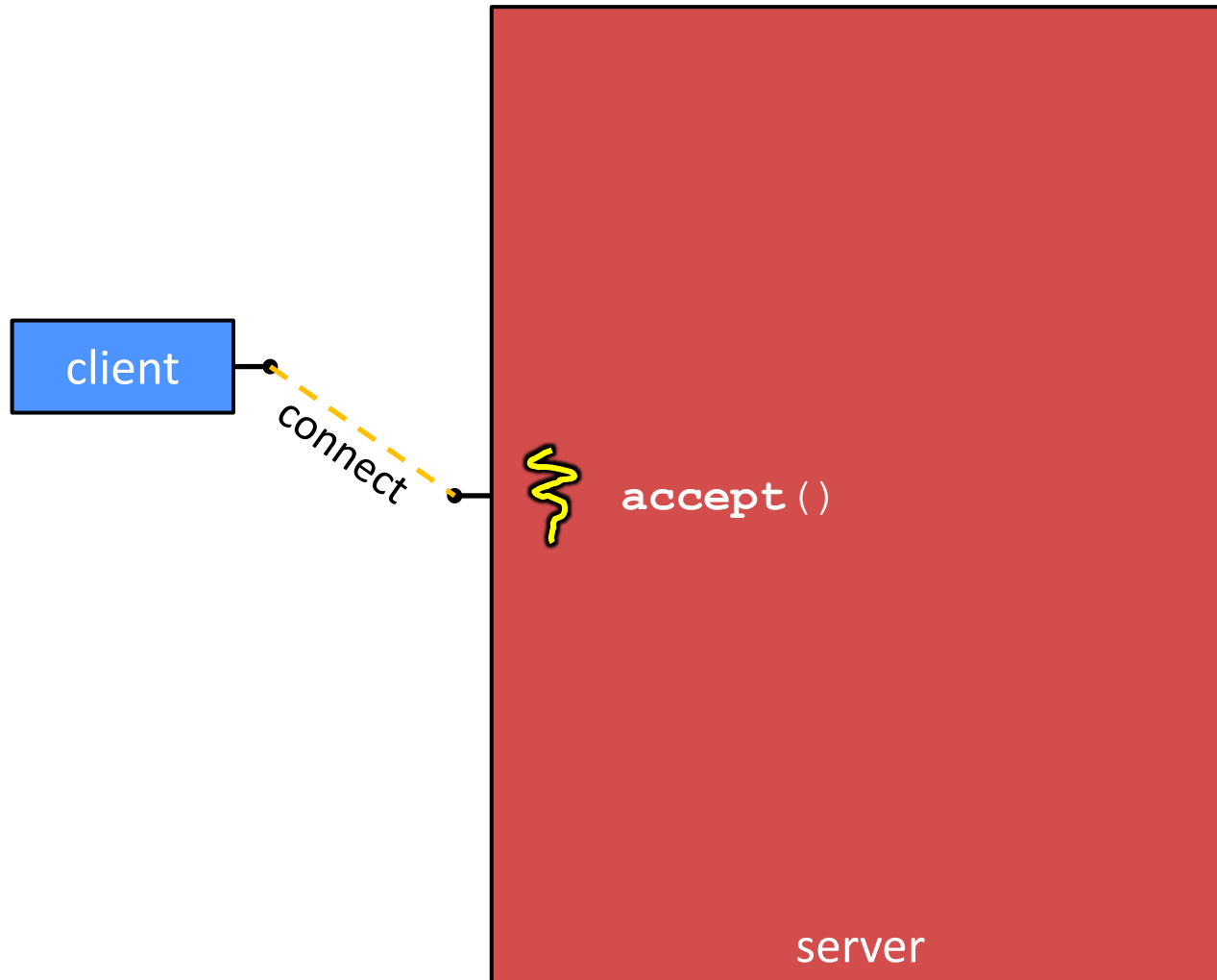
# A Concurrent Implementation

- ❖ Use multiple “workers”
  - As a query arrives, create a new “worker” to handle it
    - The “worker” reads the query from the network, issues read requests against files, assembles results and writes to the network
    - The “worker” uses blocking I/O; the “worker” alternates between consuming CPU cycles and blocking on I/O
  - The OS context switches between “workers”
    - While one is blocked on I/O, another can use the CPU
    - Multiple “workers” I/O requests can be issued at once
- ❖ So what should we use for our “workers”?

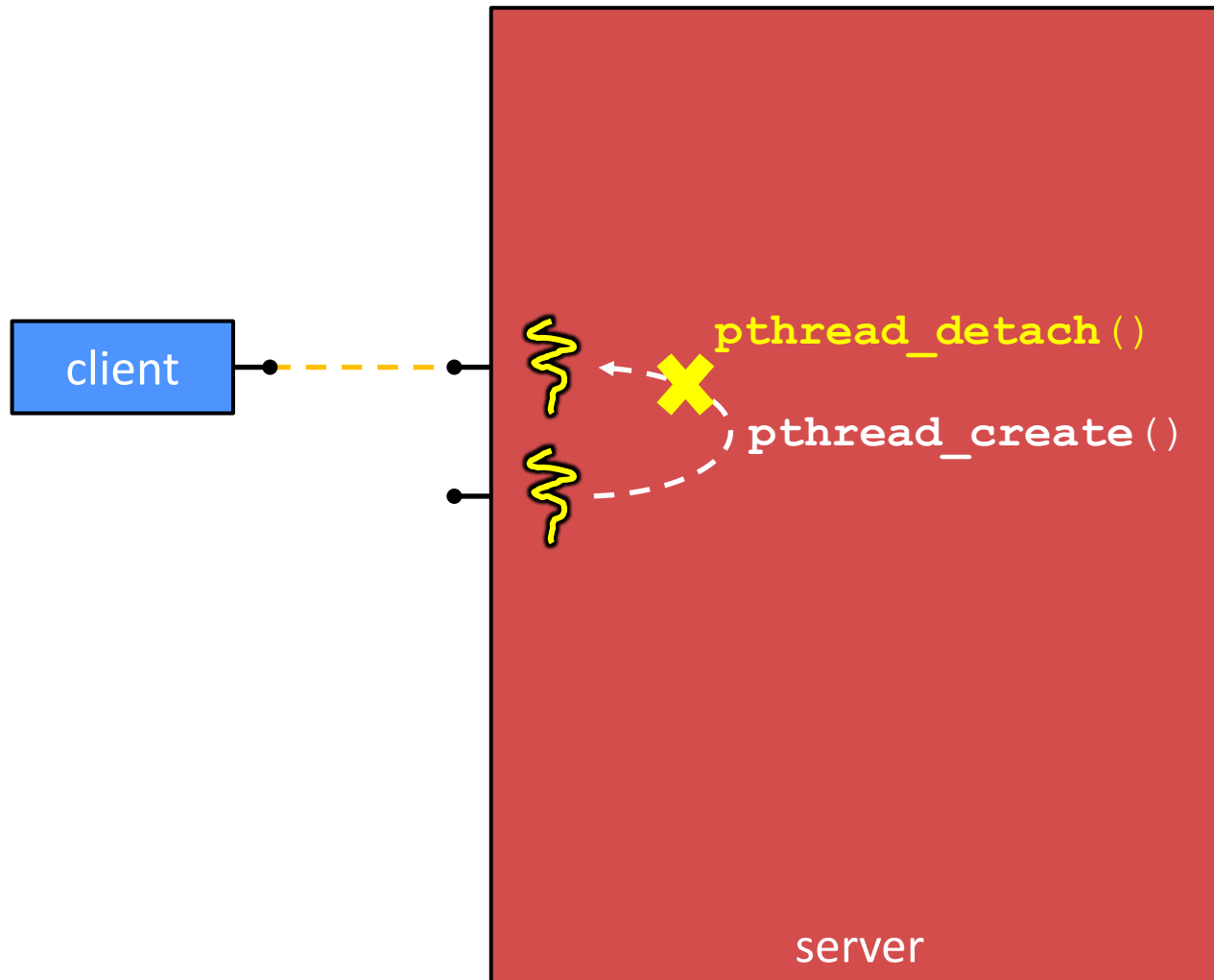
Threads!!!!



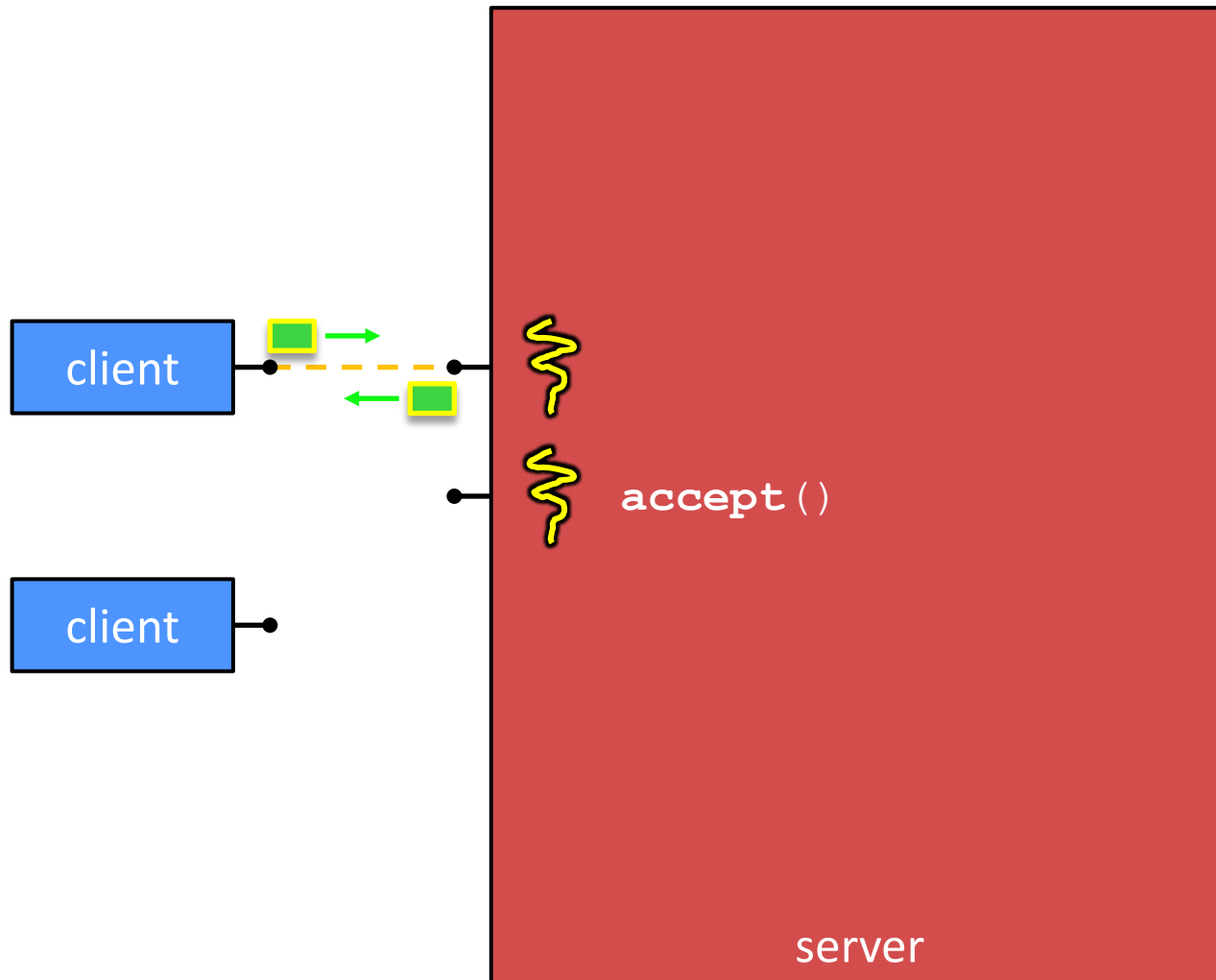
# Multithreaded Server



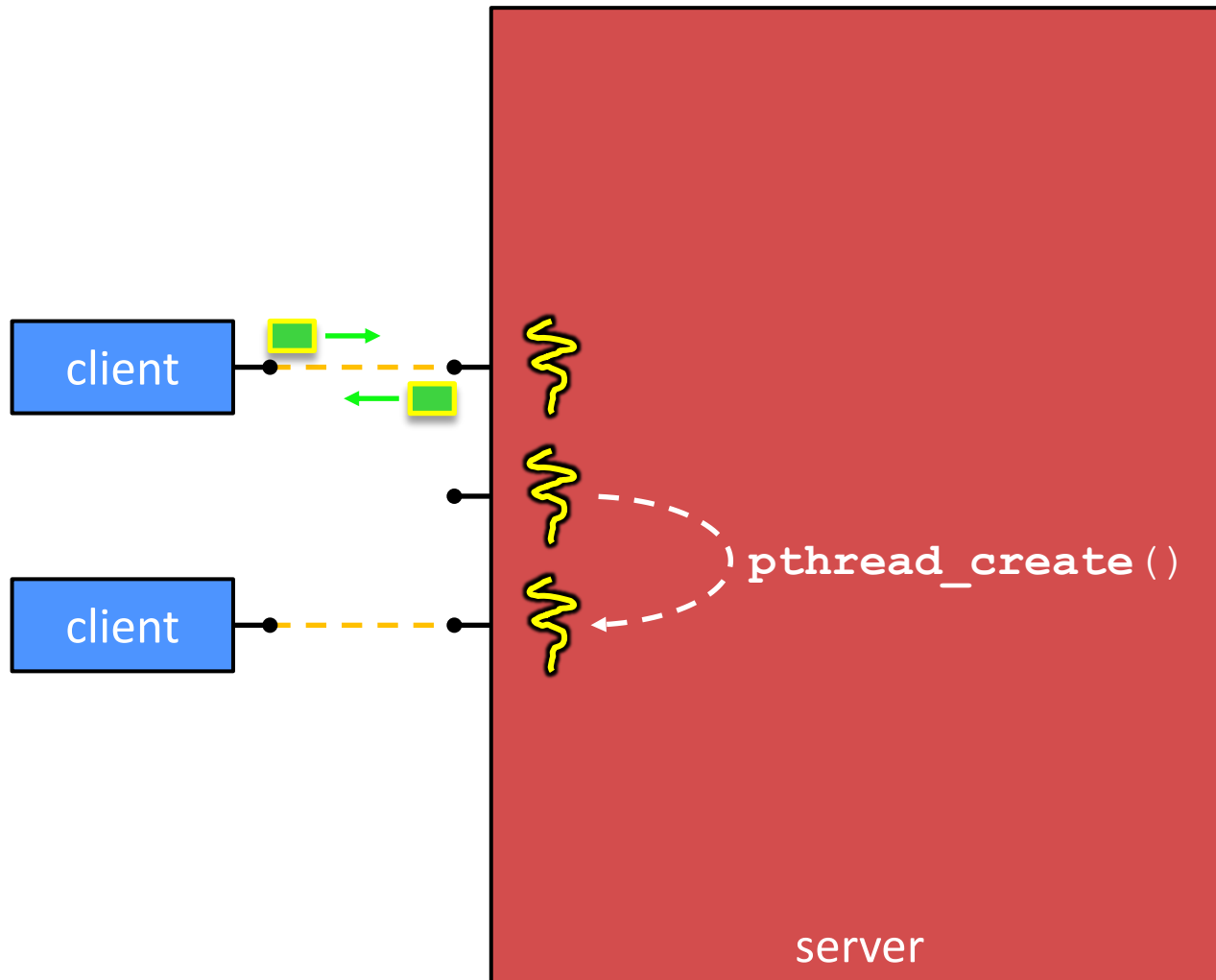
# Multithreaded Server



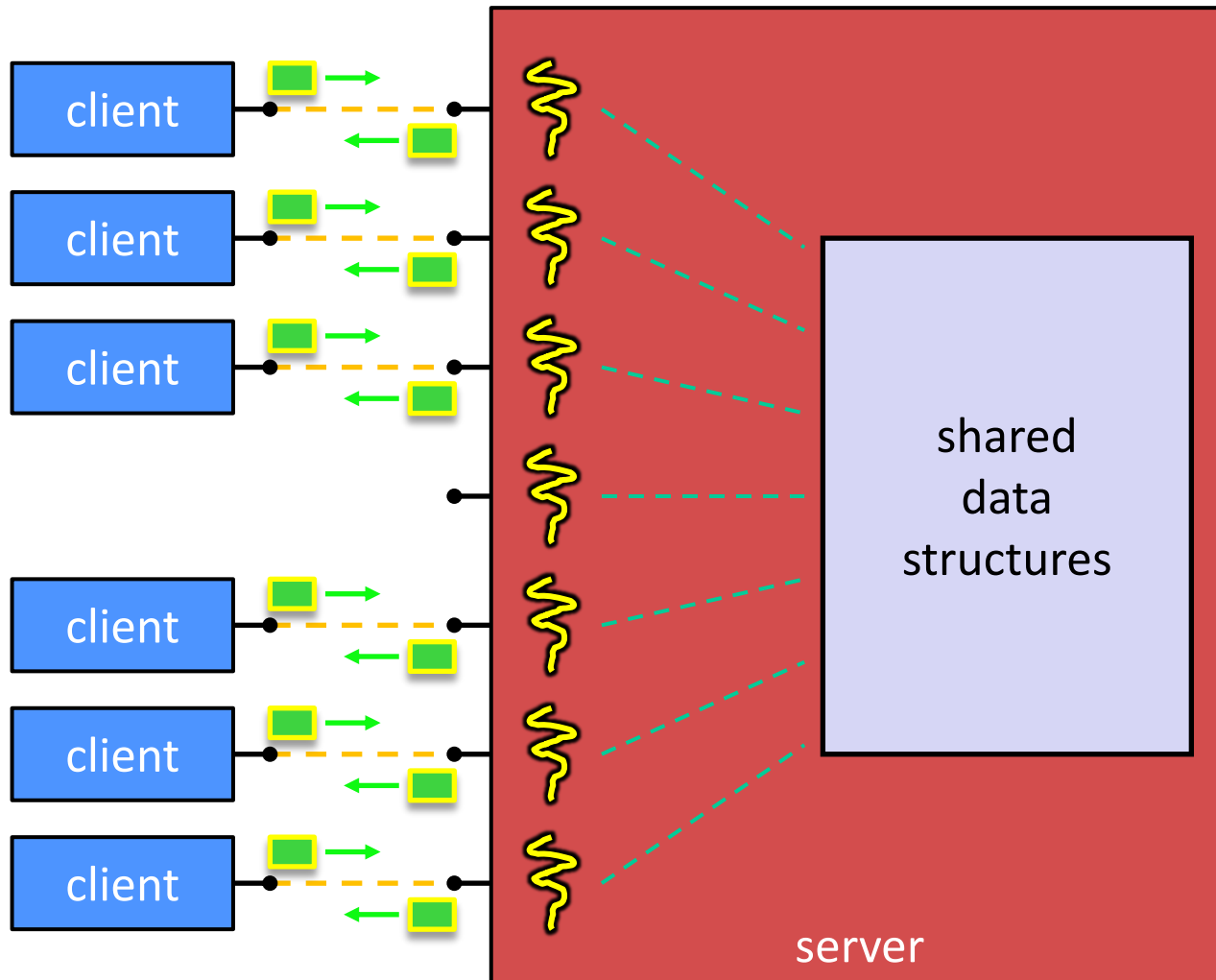
# Multithreaded Server



# Multithreaded Server

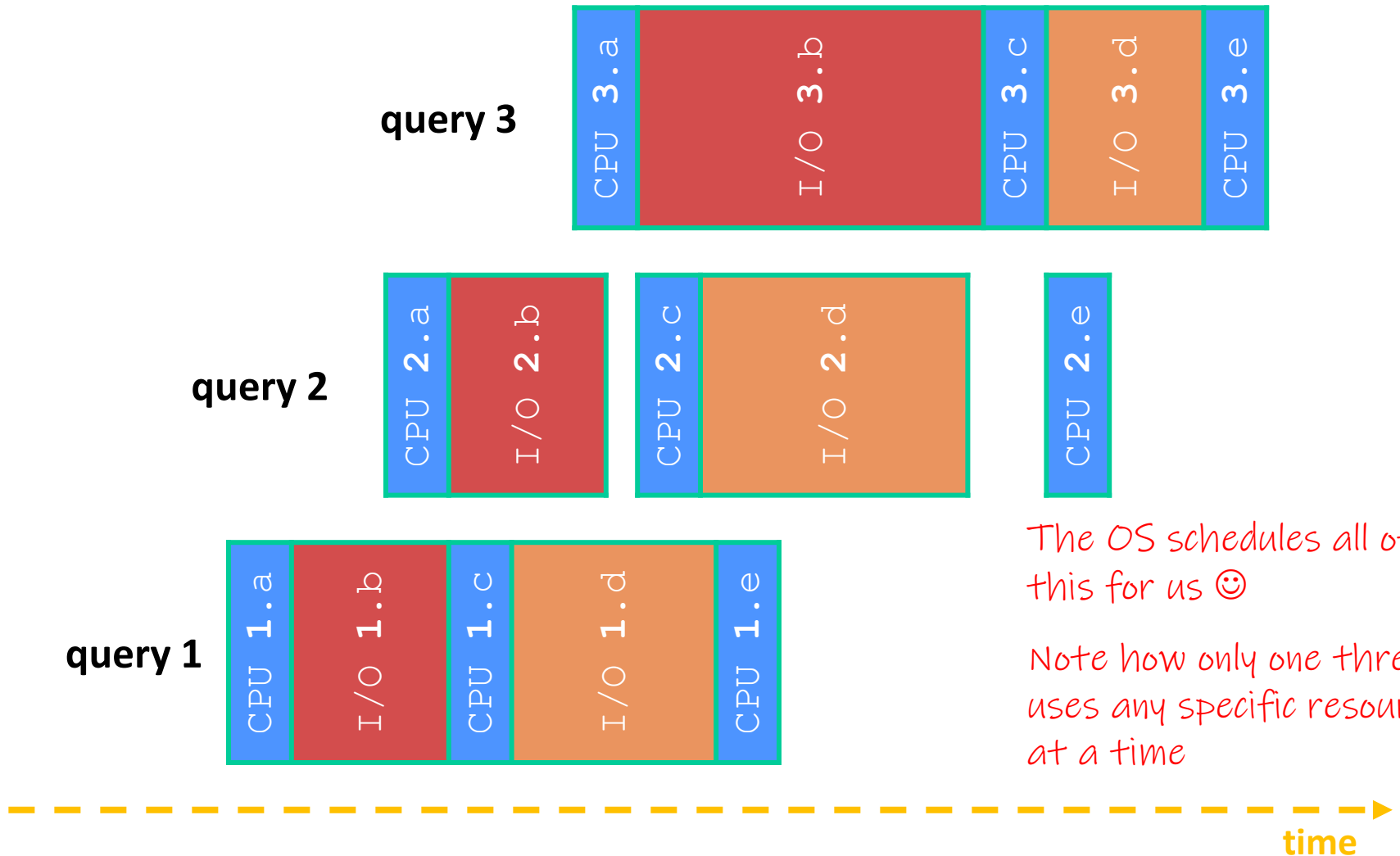


# Multithreaded Server



# Multi-threaded Search Engine (Execution)

*\*Running with 1 CPU*



*The OS schedules all of this for us 😊*


*Note how only one thread uses any specific resource at a time*

# Why Threads?

## ❖ Advantages:

- You (mostly) write sequential-looking code
- Threads can run in parallel if you have multiple CPUs/cores

## ❖ Disadvantages:

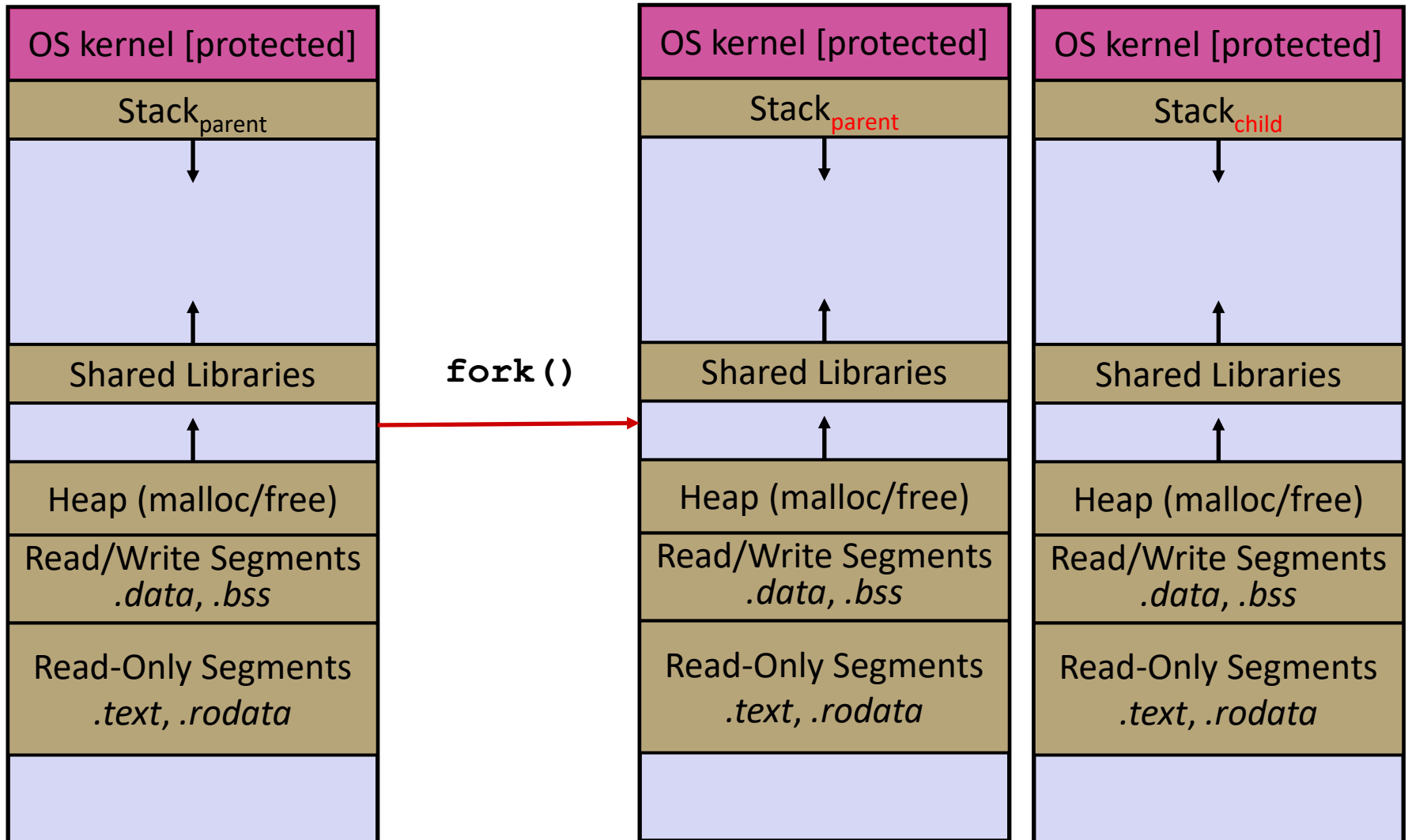
-  If threads share data, you need locks or other synchronization
  - Very bug-prone and difficult to debug
- Threads can introduce overhead
  - Lock contention, context switch overhead, and other issues
- Need language support for threads

# Threads vs. Processes

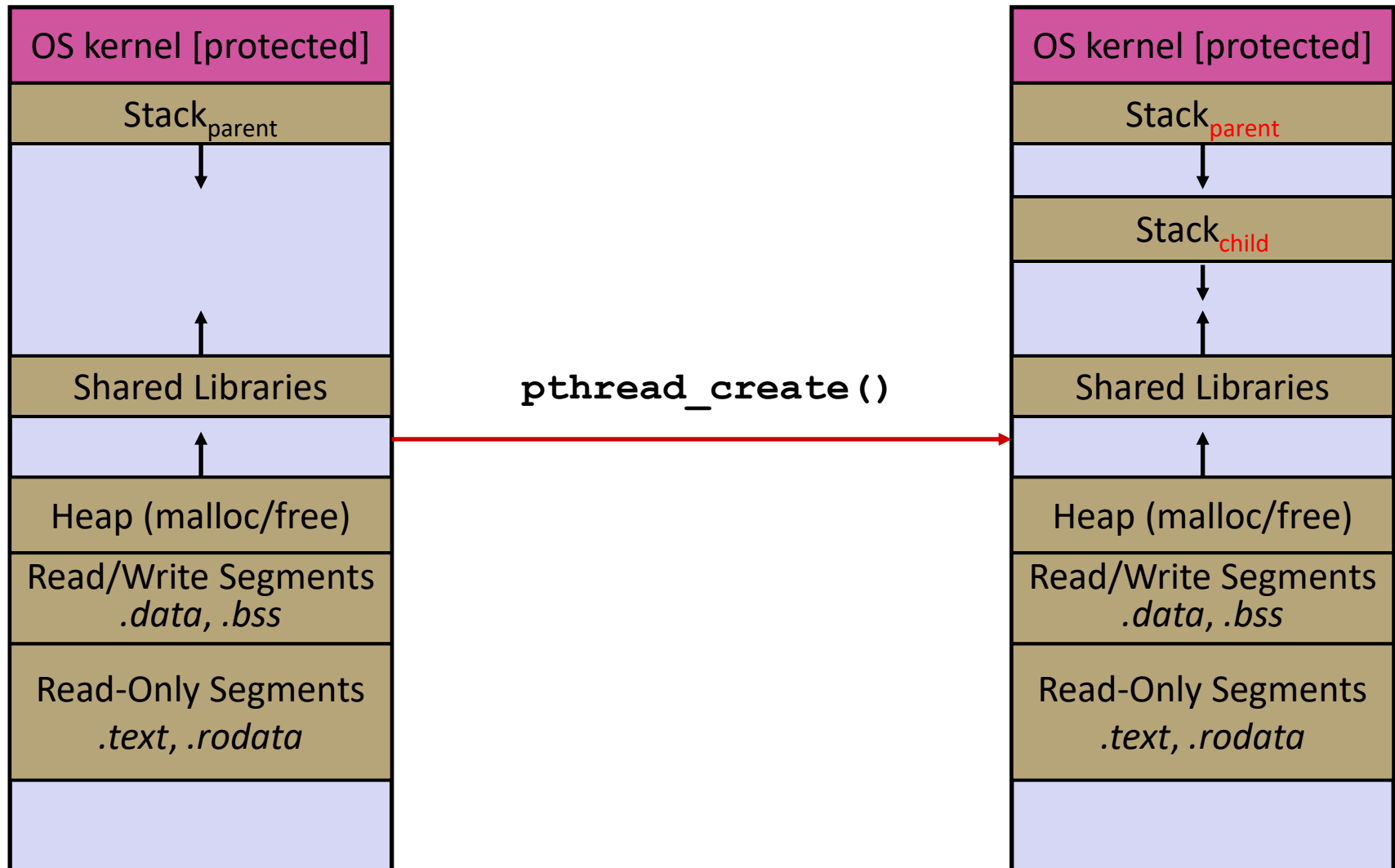
- ❖ In most modern OS's:
  - A Process has a unique: address space, OS resources, & security attributes
  - A Thread has a unique: stack, stack pointer, program counter, & registers
  - Threads are the *unit of scheduling* and processes are their *containers*; every process has at least one thread running in it



# Threads vs. Processes



# Threads vs. Processes



# Alternative: Processes

- ❖ What if we forked processes instead of threads?
- ❖ Advantages:
  - No shared memory between processes
  - No need for language support; OS provides “fork”
  - Processes are isolated. If one crashes, other processes keep going
- ❖ Disadvantages:
  - More overhead than threads during creation and context switching (Context switching == switching between threads/processes)
  - Cannot easily share memory between processes – typically communicate through the file system



# Poll Everywhere

[pollev.com/tqm](https://pollev.com/tqm)

- ❖ If I wanted to make a web browser, what concurrency model should I use?
  - Note that a web browser may need to request many resources over the network and combine them together to load a page
  
- A. **Do it sequentially**
- B. **Use threads**
- C. **Use processes**
- D. **We're lost...**

# Poll Everywhere

[pollev.com/tqm](https://pollev.com/tqm)

- ❖ If I wanted to make a web browser, what concurrency model should I use?
  - Note that a web browser may need to request many resources over the network and combine them together to load a page

**A. Do it sequentially**

*Concurrency will make more efficient use of time*

**B. Use threads**

*We will need to share the data we request across "workers"*

**C. Use processes**

**D. We're lost...**

*We want to be fast*

# Lecture Outline

- ❖ pthreads review
- ❖ Why threads?
  - Parallelism
  - Efficient use of System Resources
- ❖ **Shared resources & data races**
- ❖ Locks & mutexes

# Shared Resources

- ❖ Some resources are shared between threads and processes
  
- ❖ Thread Level:
  - Memory
  - Things shared by processes
  
- ❖ Process level
  - I/O devices
    - Files
    - terminal input/output
    - The network

*Issues arise when we try to shared things*

# Data Races

- ❖ Two memory accesses form a **data race** if different threads access the same location, and at least one is a write, and they occur one after another
  - Means that the result of a program can vary depending on chance (which thread ran first?)



# Data Race Example

- ❖ If your fridge has no milk, then go out and buy some more
  - What could go wrong?

```
if (!milk) {  
    buy milk  
}
```

- ❖ If you live alone:



- ❖ If you live with a roommate:





# Poll Everywhere

[pollev.com/tqm](https://pollev.com/tqm)

## ❖ Idea: leave a note!

- Does this fix the problem?

- A. Yes, problem fixed
- B. No, could end up with no milk
- C. No, could still buy multiple milk
- D. We're lost...

```
if (!note) {  
    if (!milk) {  
        leave note  
        buy milk  
        remove note  
    }  
}
```

# Poll Everywhere

[pollev.com/tqm](https://pollev.com/tqm)

## ❖ Idea: leave a note!

- Does this fix the problem?

*We can be interrupted  
between checking note and  
leaving note ☹️*

**A. Yes, problem fixed**

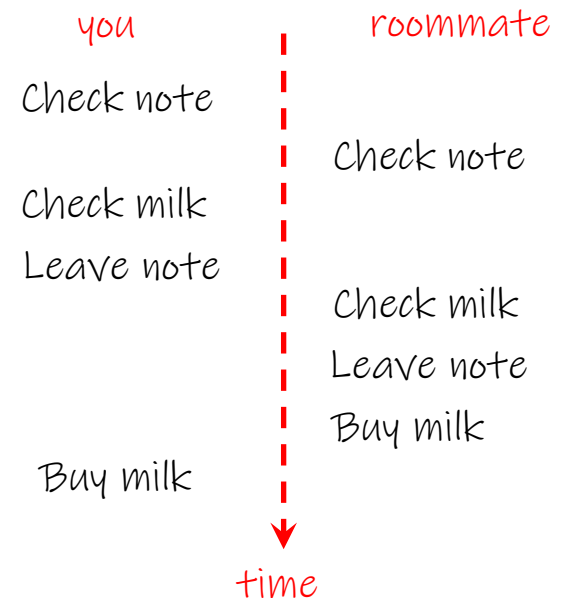
**B. No, could end up with no milk**

**C. No, could still buy multiple milk**

**D. We're lost...**

*\*There are other  
possible scenarios  
that result in  
multiple milks*

```
if (!note) {
    if (!milk) {
        leave note
        buy milk
        remove note
    }
}
```



# Threads and Data Races

- ❖ Data races might interfere in painful, non-obvious ways, depending on the specifics of the data structure
- ❖ Example: two threads try to read from and write to the same shared memory location
  - Could get “correct” answer
  - Could accidentally read old value
  - One thread’s work could get “lost”
- ❖ Example: two threads try to push an item onto the head of the linked list at the same time
  - Could get “correct” answer
  - Could get different ordering of items
  - Could break the data structure! ☠

# Increment Data Race

- ❖ What seems like a single operation `++sum total` is actually multiple operations in one. The increment looks something like this in assembly:

```
LOAD  sum_total into R0  
ADD   R0 R0 #1  
STORE R0 into sum_total
```

- ❖ What happens if we context switch to a different thread while executing these three instructions?
- ❖ **Reminder: Each thread has its own registers to work with. Each thread would have its own R0**

# Increment Data Race

- ❖ Consider that `sum_total` starts at 0 and two threads try to

execute `++sum total`      `sum_total = 0`

Thread 0      `R0 = 0`

**LOAD** `sum_total into R0`

Thread 1

# Increment Data Race

- ❖ Consider that `sum_total` starts at 0 and two threads try to

execute `++sum total`      `sum_total = 0`

Thread 0      `R0 = 0`

**LOAD** `sum_total into R0`

Thread 1      `R0 = 0`

**LOAD** `sum_total into R0`

# Increment Data Race

- ❖ Consider that `sum_total` starts at 0 and two threads try to

execute `++sum total`

`sum_total = 0`

Thread 0 `R0 = 0`

**LOAD** `sum_total into R0`

Thread 1 `R0 = 1`

**LOAD** `sum_total into R0`  
**ADD** `R0 R0 #1`



# Increment Data Race

- ❖ Consider that `sum_total` starts at 0 and two threads try to

execute `++sum total`

`sum_total = 1`

Thread 0 `R0 = 0`

**LOAD** `sum_total into R0`

Thread 1 `R0 = 1`

**LOAD** `sum_total into R0`

**ADD** `R0 R0 #1`

**STORE** `R0 into sum_total`

# Increment Data Race

- ❖ Consider that `sum_total` starts at 0 and two threads try to

execute `++sum total`

`sum_total = 1`

Thread 0 `R0 = 1`

```
LOAD sum_total into R0
```

```
ADD R0 R0 #1
```

Thread 1 `R0 = 1`

```
LOAD sum_total into R0
```

```
ADD R0 R0 #1
```

```
STORE R0 into sum_total
```

# Increment Data Race

- ❖ Consider that `sum_total` starts at 0 and two threads try to

execute `++sum total`

`sum_total = 1`

Thread 0 `R0 = 1`

**LOAD** `sum_total into R0`

**ADD** `R0 R0 #1`

**STORE** `R0 into sum_total`

Thread 1 `R0 = 1`

**LOAD** `sum_total into R0`

**ADD** `R0 R0 #1`

**STORE** `R0 into sum_total`

- ❖ With this example, we could get 1 as an output instead of 2, even though we executed `++sum_total` twice

# Lecture Outline

- ❖ pthreads review
- ❖ Why threads?
  - Parallelism
  - Efficient use of System Resources
- ❖ Shared resources & data races
- ❖ **Locks & mutexes**

# Synchronization

- ❖ **Synchronization** is the act of preventing two (or more) concurrently running threads from interfering with each other when operating on shared data
  - Need some mechanism to coordinate the threads
    - “Let me go first, then you can go”
  - Many different coordination mechanisms have been invented
  
- ❖ **Goals of synchronization:**
  - **Liveness** – ability to execute in a timely manner (informally, “something good eventually happens”)
  - **Safety** – avoid unintended interactions with shared data structures (informally, “nothing bad happens”)

# Lock Synchronization

- ❖ Use a “Lock” to grant access to a *critical section* so that only one thread can operate there at a time
  - Executed in an uninterruptible (*i.e.* *atomic*) manner

- ❖ Lock Acquire

- Wait until the lock is free, then take it

- ❖ Lock Release

- Release the lock
  - If other threads are waiting, wake exactly one up to pass lock to

- ❖ Pseudocode:

```

// non-critical code
lock.acquire();
// critical section
lock.release();
// non-critical code
    
```

# Milk Example – What is the Critical Section?

- ❖ What if we use a lock on the refrigerator?
  - Probably overkill – what if roommate wanted to get eggs?
  
- ❖ For performance reasons, only put what is necessary in the critical section
  - Only lock the milk
  - But lock *all* steps that must run uninterrupted (*i.e.* must run as an atomic unit)

```

fridge.lock()
if (!milk) {
    buy milk
}
fridge.unlock()
    
```



```

milk_lock.lock()
if (!milk) {
    buy milk
}
milk_lock.unlock()
    
```

# pthread and Locks

- ❖ Another term for a lock is a **mutex** (“mutual exclusion”)

- `pthread.h` defines datatype `pthread_mutex_t`

- ❖ 

```
int pthread_mutex_init(pthread_mutex_t* mutex,
                       const pthread_mutexattr_t* attr);
```

- Initializes a mutex with specified attributes

- ❖ 

```
int pthread_mutex_lock(pthread_mutex_t* mutex);
```

- Acquire the lock – blocks if already locked *Un-blocks when lock is acquired*

- ❖ 

```
int pthread_mutex_unlock(pthread_mutex_t* mutex);
```

- Releases the lock

- ❖ 

```
int pthread_mutex_destroy(pthread_mutex_t* mutex);
```

- “Uninitializes” a mutex – clean up when done



# pthread Mutex Examples

- ❖ See `total.cpp`
  - Data race between threads
- ❖ See `total_locking.cpp`
  - Adding a mutex fixes our data race
- ❖ How does `total_locking` compare to sequential code and to `total`?
  - Likely *slower* than both— only 1 thread can increment at a time, and must deal with checking the lock and switching between threads
  - One possible fix: each thread increments a local variable and then adds its value (once!) to the shared variable at the end
    - See `total_locking_better.cc`