Lecture 03: Concurrency Sources and Challenges

CS343 – Operating Systems Branden Ghena – Spring 2022

Some slides borrowed from: Stephen Tarzia (Northwestern), and UC Berkeley CS61C and CS162

Northwestern

Administrivia

- PCLab release will be delayed a couple days
 - I plan to release it sometime late on Saturday (instead of late today)
 - Tuesday's lecture will cover what you need to get started
 - Also I want to make some improvements to it

- Partner survey on Campuswire
 - Fill it out only if you do NOT have a partner, but want one

Today's Goals

• Describe where and why concurrency and parallelism are involved in computing.

• Be disappointed by performance limits on concurrency.

• Understand purpose and challenges of interrupts and signals.

• Introduce concept of data races as a concurrency problem.

Outline

Threads

- Need for Parallelism
- Processor Concurrency
- Concurrency Challenges
 - Amdahl's Law
 - Data Races

Processes and threads

- A process could have multiple threads
 - Each with its own registers and stack



Thread use case: web server

- Example: Web server
 - Receives multiple simultaneous requests
 - Reads web pages from disk to satisfy each request



Web server option 3: multi-threaded web server

• One thread per request. Thread handles only that request.



- Easy to program (maybe), and fast!
 - State is stored in the stacks of each thread and the thread scheduler
 - Simple to program if they are independent...

More Practical Motivation

Back to Jeff Dean's "Numbers Everyone Should Know"

Handle I/O in separate thread, avoid blocking other progress

L1 cache reference	0.5 1	ns
Branch mispredict	5 ns	
L2 cache reference	7 ns	
Mutex lock/unlock	25 ns	
Main memory reference	100 ns	
Compress 1K bytes with Zippy	3,000 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	٦
Read 1 MB sequentially from disk	20,000,000 ns	
Send packet CA->Netherlands->CA	150,000,000 ns	

Models for thread libraries: User Threads

- Thread scheduling is implemented within the process
 - OS only knows about the process, not the threads

- Upsides
 - Works on any hardware or OS
 - Performance is better when creating and switching
- Downsides
 - A system call in any thread blocks all threads



Models for thread libraries: Kernel Threads

- Thread scheduling is implemented by the operating system
 - OS manages the threads within each process



Threads versus Processes

Threads

• pthread_create()

- Creates a thread
- *Shares* all memory with all threads of the process.
- Scheduled independently of parent
- pthread_join()
 - Waits for a particular thread to finish
- Can communicate by reading/writing (shared) global variables.

Processes

• fork()

- Creates a single-threaded process
- Copies all memory from parent
 - Can be quick using copy-on-write
- Scheduled independently of parent
- •waitpid()
 - Waits for a particular child process to finish
- Can communicate by setting up shared memory, pipes, reading/writing files, or using sockets (network).

POSIX Threads Library: pthreads

- <u>https://man7.org/linux/man-pages/man7/pthreads.7.html</u>
- int pthread_create(pthread_t *thread, const pthread_attr_t *attr,
 void *(*start_routine)(void*), void *arg);
 - thread is created executing *start_routine* with *arg* as its sole argument.
 - return is implicit call to pthread_exit
- void pthread_exit(void *value_ptr);
 - terminates the thread and makes *value_ptr* available to any successful join

int pthread_join(pthread_t thread, void **value_ptr);

- suspends execution of the calling thread until the target *thread* terminates.
- On return with a non-NULL *value_ptr* the value passed to <u>*pthread_exit()*</u> by the terminating thread is made available in the location referenced by *value_ptr*.

Pthread system call example

• What happens when pthread_create() is called in a process?

```
Library:

int pthread_create(...) {

Do some work like a normal function

Put syscall number into register <----- clone (56) syscall on Linux

Put args into registers

Special trap instruction
```

Kernel:

Get args from regs Do the work to spawn the new thread Store return value in %eax

Get return values from regs
Do some more work like a normal function
};

Threads Example

```
include <stdio.h>
#include <stdlib.h>
#include <pthread.h>
#include <string.h>
int common = 162;
void *threadfun(void *threadid)
  long tid = (long)threadid;
  printf("Thread #%lx stack: %lx common: %lx (%d)\n", tid,
         (unsigned long) &tid, (unsigned long) &common, common++);
 pthread_exit(NULL);
int main (int argc, char *argv[])
{
  long t;
 int nthreads = 2;
 if (argc > 1) {
    nthreads = atoi(argv[1]);
  3
  pthread_t *threads = malloc(nthreads*sizeof(pthread_t));
  printf("Main stack: %lx, common: %lx (%d)\n",
         (unsigned long) &t, (unsigned long) &common, common);
  for(t=0; t<nthreads; t++){</pre>
   int rc = pthread_create(&threads[t], NULL, threadfun, (void *)t);
   if (rc){
      printf("ERROR; return code from pthread_create() is %d\n", rc);
      exit(-1);
    }
  }
  for(t=0; t<nthreads; t++){</pre>
    pthread_join(threads[t], NULL);
  3
  pthread_exit(NULL);
                                /* last thing in the main thread */
```

14

Threads Example

- Reads N from process
 arguments
- Creates N threads
- Each one prints a number, then increments it, then exits
- Main process waits for all of the threads to finish

```
include <stdio.h>
#include <stdlib.h>
#include <pthread.h>
#include <string.h>
int common = 162;
void *threadfun(void *threadid)
  long tid = (long)threadid;
  printf("Thread #%lx stack: %lx common: %lx (%d)\n", tid,
         (unsigned long) &tid, (unsigned long) &common, common++);
 pthread_exit(NULL);
int main (int argc, char *argv[])
  long t:
  int nthreads = 2;
  if (argc > 1) {
    nthreads = atoi(argv[1]);
  pthread_t *threads = malloc(nthreads*sizeof(pthread_t));
  printf("Main stack: %lx, common: %lx (%d)\n",
         (unsigned long) &t.(unsigned long) &common. common):
  for(t=0; t<nthreads; t++){</pre>
   int rc = pthread_create(&threads[t], NULL, threadfun, (void *)t);
    if (rc){
      printf("ERROR; return code from pthread_create() is %d\n", rc);
      exit(-1);
  for(t=0; t<nthreads; t++){</pre>
    pthread_join(threads[t], NULL);
 pthread_exit(NULL);
                                 /* last thing in the main thread */
                                                                         15
```

Threads Example

[(base) CullerMac19:code04 culler\$./pthread 4
Main stack: 7ffee2c6b6b8, common: 10cf95048 (162)
Thread #1 stack: 70000d83bef8 common: 10cf95048 (162)
Thread #3 stack: 70000d941ef8 common: 10cf95048 (164)
Thread #2 stack: 70000d8beef8 common: 10cf95048 (165)
Thread #0 stack: 70000d7b8ef8 common: 10cf95048 (163)

```
include <stdio.h>
#include <stdlib.h>
#include <pthread.h>
#include <string.h>
int common = 162;
void *threadfun(void *threadid)
  long tid = (long)threadid;
  printf("Thread #%lx stack: %lx common: %lx (%d)\n", tid,
         (unsigned long) &tid, (unsigned long) &common, common++);
  pthread_exit(NULL);
int main (int argc, char *argv[])
{
  long t;
 int nthreads = 2;
 if (argc > 1) {
    nthreads = atoi(argv[1]);
  3
  pthread_t *threads = malloc(nthreads*sizeof(pthread_t));
  printf("Main stack: %lx, common: %lx (%d)\n",
         (unsigned long) &t, (unsigned long) &common, common);
  for(t=0; t<nthreads; t++){</pre>
   int rc = pthread_create(&threads[t], NULL, threadfun, (void *)t);
    if (rc){
      printf("ERROR; return code from pthread_create() is %d\n", rc);
      exit(-1);
  }
  for(t=0; t<nthreads; t++){</pre>
    pthread_join(threads[t], NULL);
  3
  pthread_exit(NULL);
                                /* last thing in the main thread */
```

16

Ch	eck your understanding	<pre>include <stdio.h> #include <stdlib.h> #include <pthread.h> #include <string.h></string.h></pthread.h></stdlib.h></stdio.h></pre>
(base) C Main sta Thread # Thread # Thread # Thread #	ullerMac19:code04 culler\$./pthread 4 ck: 7ffee2c6b6b8, common: 10cf95048 (162) 1 stack: 70000d83bef8 common: 10cf95048 (162) 3 stack: 70000d941ef8 common: 10cf95048 (164) 2 stack: 70000d8beef8 common: 10cf95048 (165) 0 stack: 70000d7b8ef8 common: 10cf95048 (163)	<pre>int common = 162; void *threadfun(void *threadid) { long tid = (long)threadid; printf("Thread #%lx stack: %lx common: %lx (%d)\n", tid,</pre>
1.	How many threads are in this program?	<pre>int main (int argc, char *argv[]) { long t; int nthreads = 2;</pre>
2.	Does the main thread join with the threads in the same order that they were created?	<pre>if (argc > 1) { nthreads = atoi(argv[1]); } pthread_t *threads = malloc(nthreads*sizeof(pthread_t)); printf("Main stack: %lx, common: %lx (%d)\n",</pre>
3.	Do the threads exit in the same order they were created?	<pre>(unsigned long) &t,(unsigned long) &common, common); for(t=0; t<nthreads; t++){<br="">int rc = pthread_create(&threads[t], NULL, threadfun, (void *)t); if (rc){ printf("ERROR; return code from pthread_create() is %d\n", rc); if (rc)</nthreads;></pre>
4.	If we run the program again, would the result change?	exit(-1); } }
		<pre>for(t=0; t<nthreads; *="" in="" last="" main="" null);="" pre="" pthread_exit(null);="" pthread_join(threads[t],="" t++){="" the="" thing="" thread="" }="" }<=""></nthreads;></pre>

include <stdio.h> #include <stdlib.h> Check your understanding #include <pthread.h> #include <string.h> int common = 162; (base) CullerMac19:code04 culler\$./pthread 4 void *threadfun(void *threadid) Main stack: 7ffee2c6b6b8, common: 10cf95048 (162) long tid = (long)threadid; Thread #1 stack: 70000d83bef8 common: 10cf95048 (162) printf("Thread #%lx stack: %lx common: %lx (%d)\n", tid, Thread #3 stack: 70000d941ef8 common: 10cf95048 (164) (unsigned long) &tid, (unsigned long) &common, common++); Thread #2 stack: 70000d8beef8 common: 10cf95048 (165) pthread_exit(NULL); Thread #0 stack: 70000d7b8ef8 common: 10cf95048 (163) int main (int argc, char *argv[]) How many threads are in this 1. { program? **Five** long t; int nthreads = 2; if (argc > 1) { Does the main thread join with the threads in the same order 2. nthreads = atoi(argv[1]); 3 pthread_t *threads = malloc(nthreads*sizeof(pthread_t)); that they were created? **Yes** printf("Main stack: %lx, common: %lx (%d)\n", (unsigned long) &t, (unsigned long) &common, common); Do the threads exit in the 3. for(t=0; t<nthreads; t++){</pre> same order they were int rc = pthread_create(&threads[t], NULL, threadfun, (void *)t); if (rc){ created? Maybé?? printf("ERROR; return code from pthread_create() is %d\n", rc); exit(-1); If we run the program again, would the result change? 4. } **Possibly!** for(t=0; t<nthreads; t++){</pre> pthread_join(threads[t], NULL); pthread_exit(NULL); /* last thing in the main thread */

18

Outline

- Threads
- Need for Parallelism
- Processor Concurrency
- Concurrency Challenges
 - Amdahl's Law
 - Data Races

It's the mid 1990s and you work at Microsoft.

You need to double the speed of Excel in two years.

What do you do?

It's the mid 1990s and you work at Microsoft.

You need to double the speed of Excel in two years.

What do you do? **Take a vacation**

Moore's Law – CPU transistors counts



Data source: Wikipedia (wikipedia.org/wiki/Transistor_count) Year in which the microchip was first introduced OurWorldinData.org – Research and data to make progress against the world's largest problems. Licensed under CC-BY by the authors Hannah Ritchie and Max Roser.

Processors kept getting faster too



Power is a major limiting factor on speed

- We could make processors go very fast
 - But doing so uses more and more power
- More power means more heat generated
 - And chips typically work up to around 100°C
 - Hotter than that and things stop working
- We add heat sinks and fans and water coolers to keep chips cool
 But it's hard to remove heat quickly enough from chips
- So, power consumption ends up limiting processor speed

Denard Scaling

- Moore's Law corollary: Denar Scaling
 - As transistors get smaller, the power density stays the same
 - Which is to say that the power-per-transistor decreases!
- Making the processor clock speed faster uses more power
 - But the two balance out for roughly net even power
 - So not only do we get *more* transistors, but chip speed can be *faster* too

- From our Excel example:
 - In two years, new hardware would run the existing software twice as fast

Then they stopped getting faster



~2006: Leakage current becomes significant

Now smaller transistors doesn't mean lower power

So... now what?

In summary:

- Making transistors smaller doesn't make them lower power,
- so if we were to make them faster, they would take more power,
- which will eventually lead to our processors melting...
- and because of that, we can't reliably make performance better by waiting for clock speeds to increase.

How do we continue to get better computation performance?

Exploit parallelism!



Parallelism Analogy

- I want to peel 100 potatoes as fast as possible:
 - I can learn to peel potatoes faster

OR

- I can get 99 friends to help me
- Whenever one result doesn't depend on another, doing the task in parallel can be a big win!

Parallelism versus Concurrency

Two processes A and B





Parallelism versus Concurrency

- Parallelism
 - Two things happen strictly simultaneously
- Concurrency
 - More general term
 - Two things happen in the same time window
 - Could be simultaneous, could be interleaved
 - Concurrent execution occurs whenever two processes are both active



Outline

- Threads
- Need for Parallelism
- Processor Concurrency
- Concurrency Challenges
 - Amdahl's Law
 - Data Races

What are the hardware sources of concurrency?

Model of a processor



CPU



But instructions don't always have to be executed in order



Doesn't have to go after the movq instructions because it uses different registers

We can apply the multiprogramming approach of executing this addq while the movq is waiting on memory.

Out-of-order machines


Out-of-order processors obey normal execution results

- Initial thoughts on out-of-order execution
 - •
 - The processor could be executing my program in order it feels like?!!
 - How do I possibly reason about anything?
- Answer: the processor promises to have the same results as if things were done in the normal order.



Multiple threads might rely on memory ordering

- The processor can't account for multiple threads though
- If memory results are shared by two threads, the processor might mess something up for you.



• What will Thread 1 print?

Multiple threads might rely on memory ordering

- The processor can't account for multiple threads though
- If memory results are shared by two threads, the processor might mess something up for you.



This can be addressed with memory barriers

• What will Thread 1 print? Could be 42. Could be 0.

How else do processors employ concurrency?

Goal: Make computer faster by performing multiple tasks

Solutions:

- 1. Use multiple cores to run multiple tasks in parallel
- 2. Run multiple tasks on a single core concurrently

How else do processors employ concurrency?

Goal: Make computer faster by performing multiple tasks

Solutions:

1. Use multiple cores to run multiple tasks in parallel

2. Run multiple tasks on a single core concurrently

Multiprocessor Systems (in pictures)



Multiprocessor Systems (in words)

- A computer system with at least 2 processors or *cores*
 - Each core has its own registers
 - Each core executes independent instruction streams
 - Processors share the same system memory
 - But use different parts of it
 - Communication possible through memory accesses
- Deliver high throughput for independent jobs via task-level parallelism

Multiprocessor Example

Run Chrome and Spotify simultaneously

- Each are separate programs
- Each has a different memory space
- Each can run on a separate core

Don't even need to communicate...

Note: OS can fake this by interleaving processes, but hardware can make it actually simultaneous How else do processors employ concurrency?

Goal: Make computer faster by performing multiple tasks

Solutions:

1. Use multiple cores to run multiple tasks in parallel

2. Run multiple tasks on a single core concurrently

Basic idea: Processor resources are expensive and should not be left idle

Long memory latency to memory on cache miss?

- Hardware switches threads to bring in other useful work while waiting for cache miss
- Cost of thread context switch must be much less than cache miss latency

 Switching threads is less expensive than processes because they share memory

Hardware support for multithreading



- Two copies of PC and Registers inside processor hardware
- Looks like two processors to software (hardware thread 0, hardware thread 1)
- Control logic decides which thread to execute an instruction from next



Multithreading versus Multicore

- Multithreading => Better utilization
 - \approx 5% more hardware for \approx 1.3x better performance?
 - Gets to share ALUs, caches, memory controller
- Multicore => Duplicate processors
 - \approx 50% more hardware for \approx 2x better performance?
 - Share some caches (L2 cache, L3 cache), memory controller
- Modern machines do both
 - Multiple cores with multiple threads per core

My desktop computer



Raspberry Pi 4

Quad core processor

- One thread per core
- 3-way superscalar pipeline
- L1 Cache
 - 32 KiB 2-way set associative data cache
 - 48 KiB 3-way set associative instruction cache
 - Per core
- L2 Cache
 - 512 KiB to 4 MiB (shared)
- RAM 1-4 GB

\$35

Literally all computers are doing parallelism these days

Back up to the OS perspective

- Modern operating systems must manage concurrency
 - Both parallel operation and interleaving operations
- Concurrency is valuable
 - Performance gains are the reason

Break + Open Question

- How many cores/threads does your processor support?
 - Windows: Task Manager -> Performance -> CPU
 - MacOS: About this Mac -> System Report -> Hardware
 - M1 processor only does 1 thread per core
 - Linux: In terminal: lscpu
 - Android/iOS: You'll need to google it

Outline

- Threads
- Need for Parallelism
- Processor Concurrency
- Concurrency Challenges
 - Amdahl's Law
 - Data Races

Challenges to concurrency

Concurrency is great! We can do so many things!!

But what's the downside...?

- 1. How much speedup can we get from it?
- 2. How hard is it to write parallel programs?

Challenges to concurrency

Concurrency is great! We can do so many things!!

But what's the downside...?

1. How much speedup can we get from it?

2. How hard is it to write parallel programs?

Speedup Example



Imagine a program that takes 100 seconds to run

- 95 seconds in the blue part
- 5 seconds in the green part

Speedup from improvements

95 c	5 s Speedup w Improveme	Speedup with	improvement
33 5		Improvement	Execution time with

 $5 s \rightarrow 1 s$: Speedup = 100/96 = 1.042

 $5 s \rightarrow 0.001s$: Speedup = 100/95.001 = 1.053

The impact of a performance improvement is relative to the importance of the part being improved!

Execution time without



F = Fraction of execution time speed up

S = Scale of improvement

Example: 2x improvement to 25% of the program

$$\frac{1}{0.75 + \frac{0.25}{2}} = \frac{1}{0.75 + 0.125} = 1.14$$



Speedup with mprovement =
$$\frac{1}{(0.75) + (0.25/20)}$$

= 1.311

Amdahl's (heartbreaking) Law (in pictures)

• The amount of speedup that can be achieved through parallelism is limited by the non-parallel portion of your program!



Amdahl's (heartbreaking) Law (in words)

- Amdahl's Law tells us that to achieve linear speedup with more processors:
 - *none* of the original computation can be serial (non-parallelizable)
- To get a speedup of 90 from 100 processors, the percentage of the original program that could be scalar would have to be 0.1% or less

Speedup = 1/(.001 + .999/100) = 90.99

50%	50%
-----	-----

Speedup with improvement = $\frac{1}{(1-F) + (F/S)}$

- Suppose a program spends 50% of its time in a square root routine.
- How much must you speed up square root to make the program run 2x faster?

(A)	10
(B)	20
(C)	100
(D)	None of the above

= $\overline{(1-F)+(F/S)}$

ppose a program spends ! w much must you speed i

50%

Break + Question

50%

•	Suppose a	program spe	nds 50% of it	s time in a s	square	root routine.
					- <u>-</u>	

• How much must you speed up square root to make the program run 2x faster?



Speedup = 1 / [(1 - F) + (F/S)]2 = 1 / [(1 - 0.5) + (0.5/S)]S = $0.5 / ((1/2) - 0.5) = \infty$

Speedup with

improvement

The square root would need to decrease to nothing before you got 2x speedup

Outline

- Threads
- Need for Parallelism
- Processor Concurrency

Concurrency Challenges

- Amdahl's Law
- Data Races

Challenges to concurrency

Concurrency is great! We can do so many things!!

But what's the downside...?

- 1. How much speedup can we get from it?
- 2. How hard is it to write parallel programs?

Concurrency problem: data races

Consider two threads with a shared global variable: int count = 0



count could end up with a final value of 1 or 2. How?

Concurrency problem: data races

Consider two threads with a shared global variable: int count = 0

Thread 1:	Thread 2:	Assuming "count" is in memory location 0x8049a1c
<pre>void thread_fn(){ mov \$0x8049a1c, %edi mov (%edi), %eax add \$0x1, %eax mov %eax, (%edi) }</pre>	<pre>void thread_fn(){ mov \$0x8049a1c, %edi mov (%edi), %eax add \$0x1, %eax mov %eax, (%edi) }</pre>	

count could end up with a final value of 1 or 2. How? These instructions could be interleaved in any way.

Data race example

Assuming "count" is in memory location pointed to by **%edi**

	Thread 1	Thread 2	Thread 1	Thread 2
Time	mov (%edi), %eax		mov (%edi), %eax	
	add \$0x1, %eax			mov (%edi), %eax
	mov %eax, (\$edi)			add \$0x1, %eax
		mov (%edi), %eax		mov %eax, (%edi)
		add \$0x1, %eax	add \$0x1, %eax	
+		mov %eax, (%edi)	mov %eax, (%edi)	

Final value of count: 2

Final value of count: 1

Data race explanation

- Thread scheduling is **non-deterministic**
 - There is no guarantee that any thread will go first or last or not be interrupted at any point
- If different threads write to the same variable
 - The final value of the variable is also non-deterministic
 - This is a *data race*

Check your understanding: data races with multiple threads

Consider three threads with a shared global variable: int count = 0

Thread 1:	Thread 2:	Thread 3:
<pre>void main(){ count += 2; }</pre>	<pre>void main(){ count -= 2; }</pre>	<pre>void main(){ count += 3; }</pre>

What are the possible values of count?

Check your understanding: data races with multiple threads

Consider three threads with a shared global variable: int count = 0

Thread 1:	Thread 2:	Thread 3:
<pre>void main(){ count += 2; }</pre>	<pre>void main(){ count -= 2; }</pre>	<pre>void main(){ count += 3; }</pre>

What are the possible values of count?

-2, 0, 1, 2, 3, 5

How are you supposed to reason about this?! Need mechanisms for sharing memory.

Outline

- Threads
- Need for Parallelism
- Processor Concurrency
- Concurrency Challenges
 - Amdahl's Law
 - Data Races