

# An Introduction to Parallel Programming:

**A hands-on Introduction** (using OpenMP)

**Tim Mattson**

**Human Learning Group**

(retired from Intel Aug'2023)

With a lot of help from Helen He, Alice Koniges, David Eder, and many more

# Introduction

I'm just a simple kayak instructor

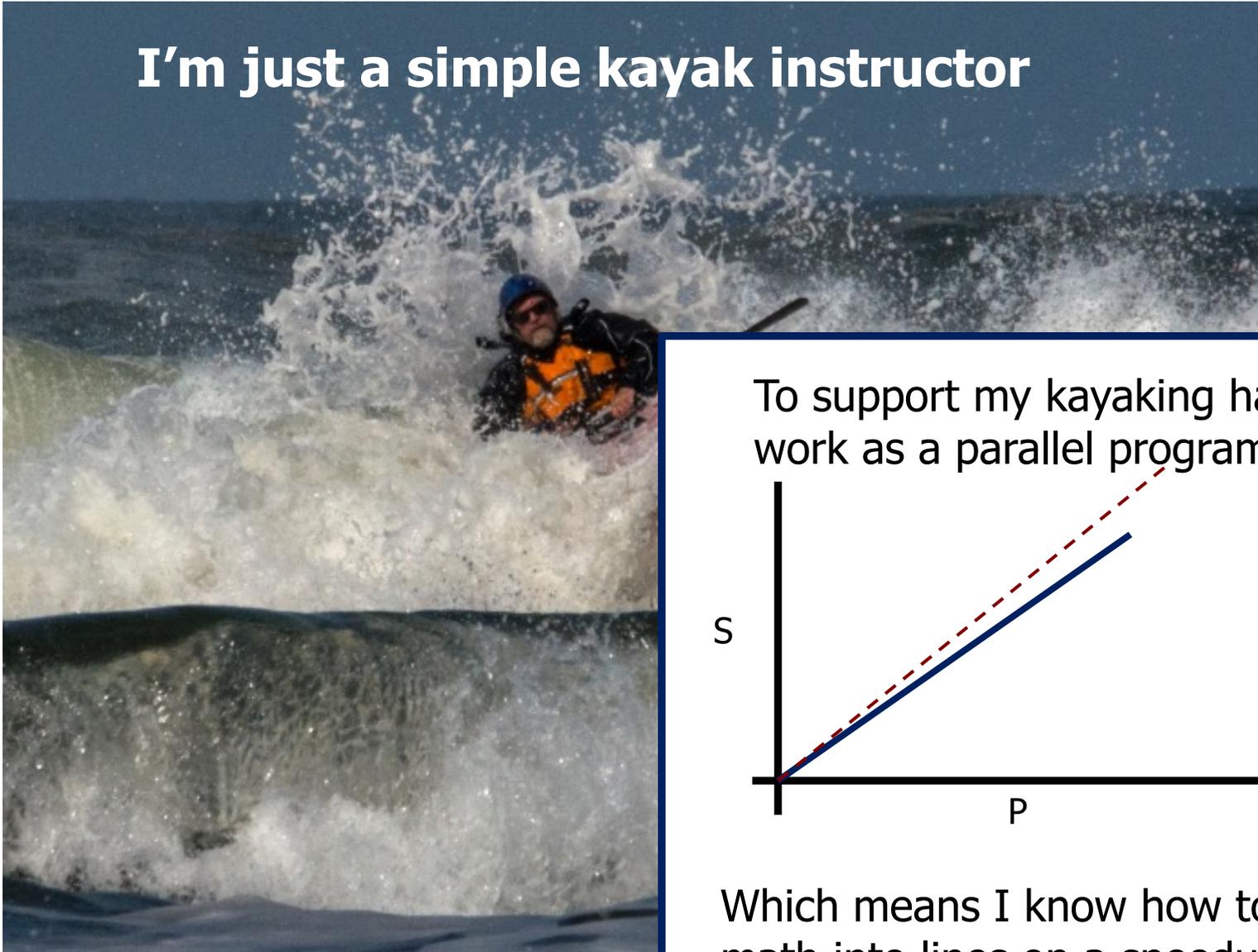
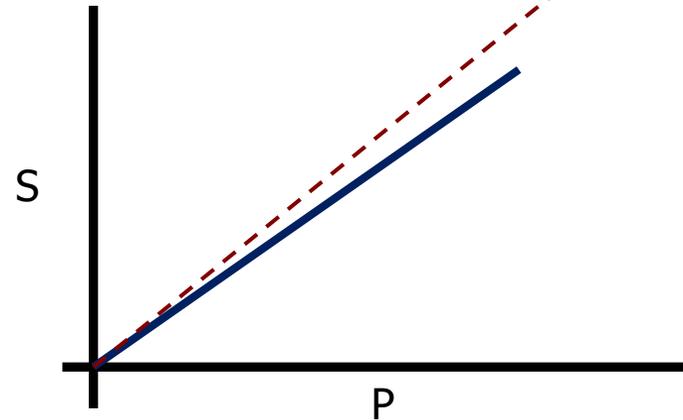


Photo © by Greg Clopton, 2014

To support my kayaking habit, I work as a parallel programmer



Which means I know how to turn math into lines on a speedup plot

# Preliminaries: Part 1

- Our plan for the day .. Active learning!
  - We will mix short lectures with short exercises.
- Please follow these simple rules
  - Do the exercises that we assign and then change things around and experiment.
    - Embrace active learning!
  - Don't cheat: Do Not look at the solutions before you complete an exercise ... even if you get really frustrated.

# Use homebrew to install gnu compilers on your Apple laptop

I tested this on a new  
(July 2023) MacBook  
Air with an Apple M2  
CPU

Warning: Xcode uses the name gcc for Apple's clang compiler.  
Use Homebrew to load a real, gcc compiler.

- Download Xcode. Be sure to choose the command line tools option.
- Go to the homebrew web site (brew.sh). Cut and paste the command near the top of the page to install homebrew (in /opt/homebrew):

```
/bin/bash -c "$(curl -fsSL https://raw.githubusercontent.com/Homebrew/install/HEAD/install.sh)"
```

- Add /opt/homebrew/bin to your path. I did this by adding the following line to .zshrc

```
% export PATH=/opt/homebrew/bin:$PATH
```

- Install the latest gcc compiler

```
% brew install gcc
```

- This will install the compiler in /opt/homebrew/bin. Check /opt/homebrew/bin to see which gcc compiler was installed. In my case, it installed gcc-13
- Test the compiler (and the openmp option) with a simple hello world program

```
% gcc-13 -fopenmp hello.c
```

# OpenMP Compilers on Apple Laptops: MacPorts

I have not tested this in a long time.  
I greatly prefer homebrew.

But if you prefer MacPorts, this procedure  
should work.

- To use OpenMP on your Apple laptop:
- Download Xcode. Be sure to choose the command line tools option.
- Download and use MacPorts to install the latest gnu compilers.

```
sudo port selfupdate
```

Update to latest version of  
MacPorts

```
sudo port install gcc13
```

Grab version 13 gnu  
compilers (5-10 mins)

```
port select --list gcc
```

List versions of gcc on your  
system

```
sudo port select --set gcc mp-gcc13
```

Select the mp enabled version of  
the most recent gcc release

```
gcc -fopenmp hello.c
```

Test the installation with a simple  
program

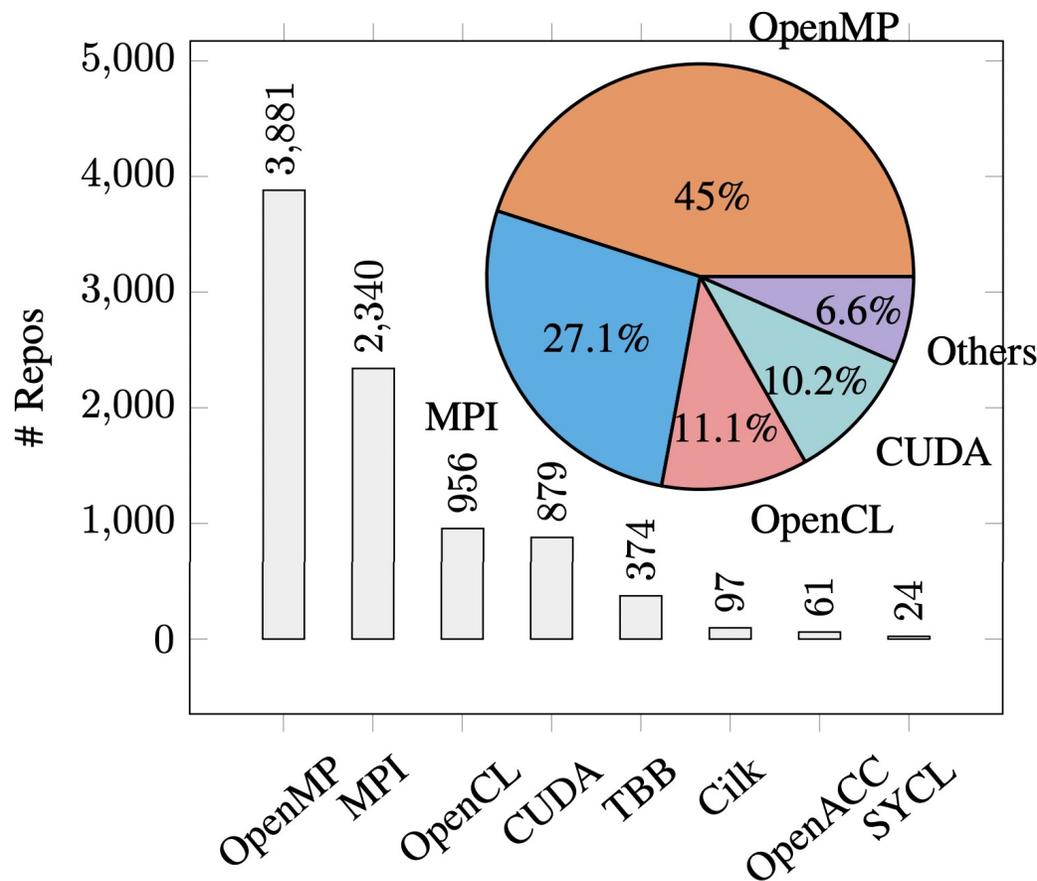
**The best way to master parallel computing ...**

**start with a simple approach to parallelism and build an intellectual foundation by writing parallel code.**

**... and the simplest API for parallelism is?**

# Outline

- ➔ Introduction to OpenMP
- Creating Threads
- Synchronization
- Parallel Loops
- Data Environment
- Memory Model
- OpenMP stuff we didn't cover
- Recap



In a dataset (HPCorpus) of all C/C++/Fortran github repositories from 2013-2023, OpenMP was found to be the most popular parallel programming model

Aggregate numbers over all repositories from 2013 to 2023

Quantifying OpenMP: Statistical insights into usage and adoption, Tal Kadosh, et al., HPEC'2023, <https://arxiv.org/abs/2308.08002>

# OpenMP\* Overview

```
C$OMP FLUSH
```

```
#pragma omp critical
```

```
#pragma omp single
```

```
C$OMP THREADPRIVATE (/ABC/)
```

```
C$OMP ATOMIC
```

```
CALL OMP_SET_NUM_THREADS(10)
```

## *OpenMP: An API for Writing Parallel Applications*

- A set of compiler directives and library routines for parallel application programmers
- Originally ... Greatly simplifies writing multithreaded programs in Fortran, C and C++
- Later versions ... supports non-uniform memories, vectorization and GPU programming

```
#pragma omp parallel for private(A, B)
```

```
C$OMP PARALLEL REDUCTION (+: A, B)
```

```
C$OMP PARALLEL COPYIN (/blk/)
```

```
C$OMP DO lastprivate(XX)
```

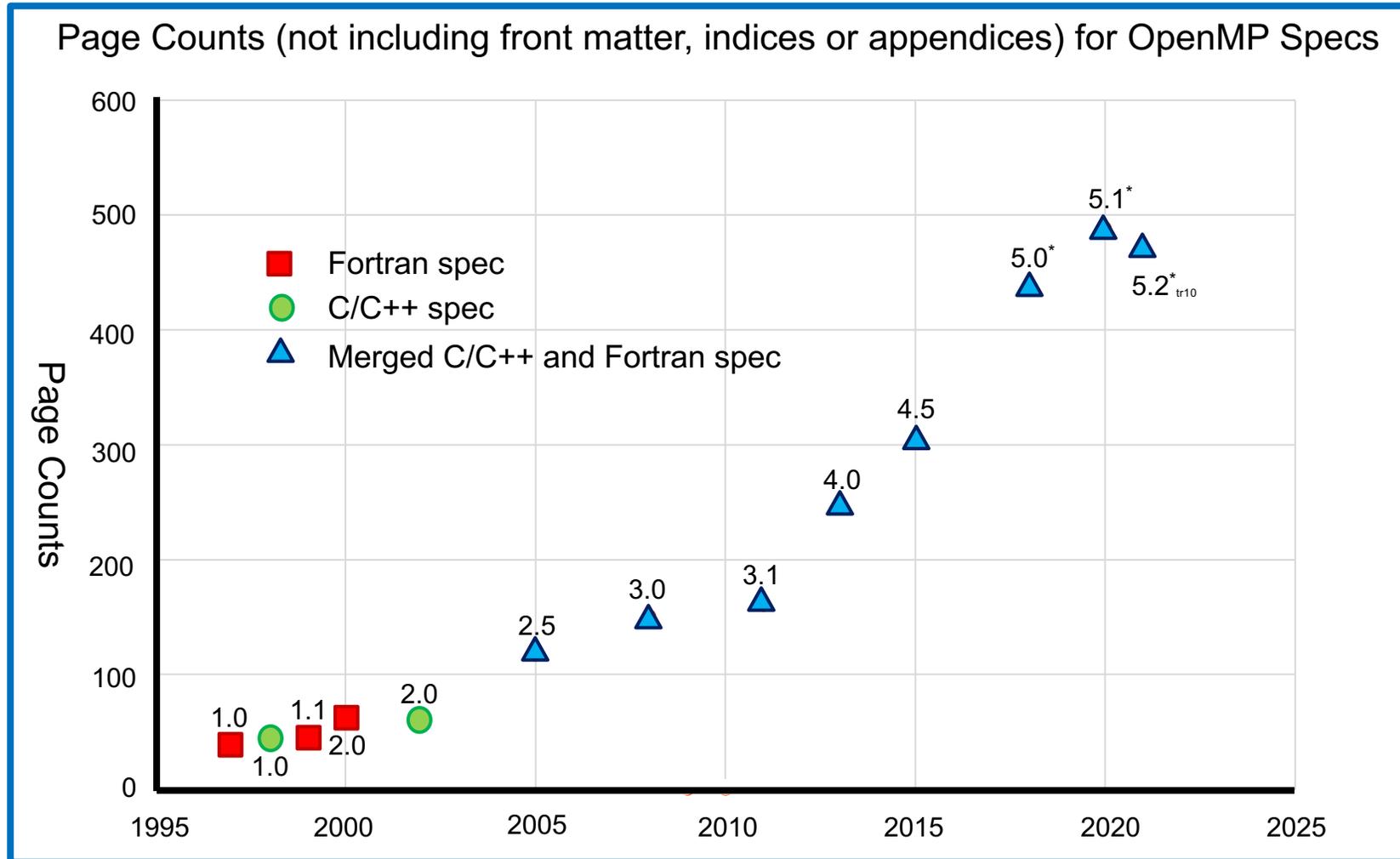
```
#pragma omp atomic seq_cst
```

```
Nthrds = OMP_GET_NUM_PROCS()
```

```
omp_set_lock(lck)
```

# The Growth of Complexity in OpenMP

Our goal in 1997 ... A simple interface for application programmers

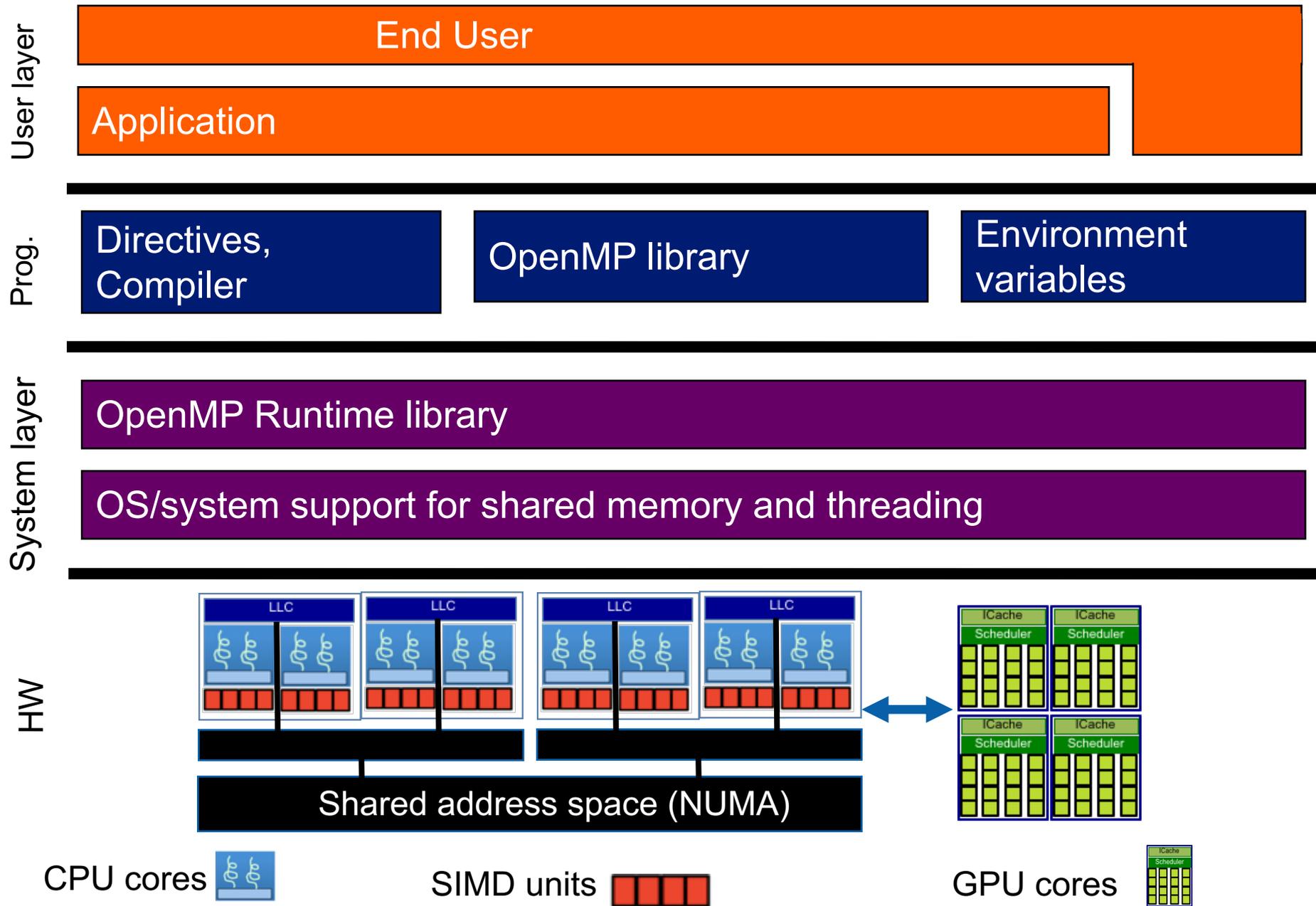


The full spec is overwhelming. We focus on the Common Core: the 21 items most people restrict themselves to

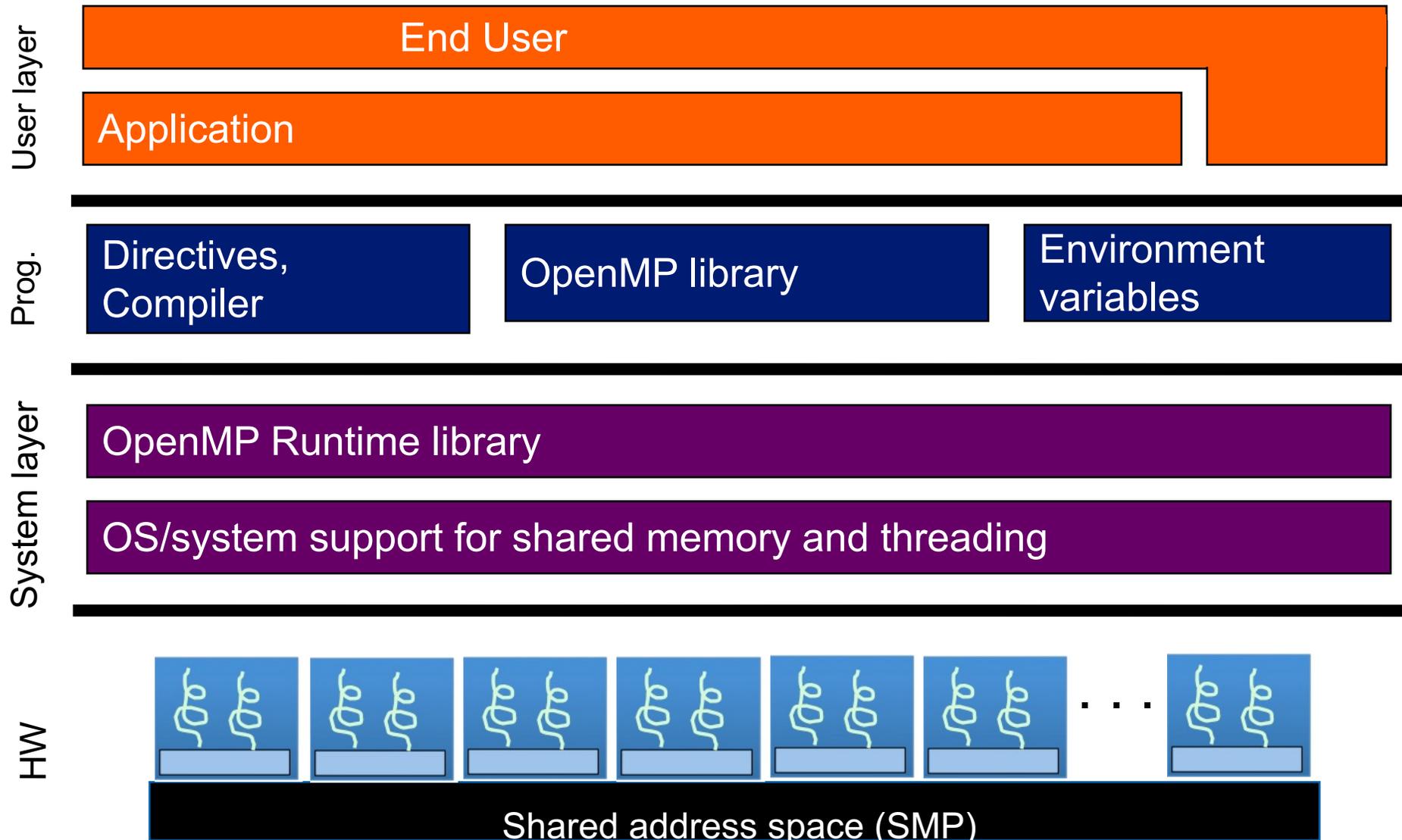
# The OpenMP Common Core: Most OpenMP programs only use these 21 items

OpenMP pragma, function, or clause	Concepts
#pragma omp parallel	Parallel region, teams of threads, structured block, interleaved execution across threads.
void omp_set_thread_num() int omp_get_thread_num() int omp_get_num_threads()	Default number of threads and internal control variables. SPMD pattern: Create threads with a parallel region and split up the work using the number of threads and the thread ID.
double omp_get_wtime()	Speedup and Amdahl's law. False sharing and other performance issues.
setenv OMP_NUM_THREADS N	Setting the internal control variable for the default number of threads with an environment variable
#pragma omp barrier #pragma omp critical	Synchronization and race conditions. Revisit interleaved execution.
#pragma omp for #pragma omp parallel for	Worksharing, parallel loops, loop carried dependencies.
reduction(op:list)	Reductions of values across a team of threads.
schedule (static [,chunk]) schedule(dynamic [,chunk])	Loop schedules, loop overheads, and load balance.
shared(list), private(list), firstprivate(list)	Data environment.
default(none)	Force explicit definition of each variable's storage attribute
nowait	Disabling implied barriers on workshare constructs, the high cost of barriers, and the flush concept (but not the flush directive).
#pragma omp single	Workshare with a single thread.
#pragma omp task #pragma omp taskwait	Tasks including the data environment for tasks.

# OpenMP Basic Definitions: Basic Solution Stack



# OpenMP Basic Definitions: Basic Solution Stack



For the OpenMP Common Core, we focus on Symmetric Multiprocessor Case ....  
i.e., lots of threads with “equal cost access” to memory

# OpenMP Basic Syntax

- Most of OpenMP happens through compiler directives.

C and C++	Fortran
Compiler directives	
<i><b>#pragma omp construct [clause [clause]...]</b></i>	<i><b>!\$OMP construct [clause [clause] ...]</b></i>
Example	
<i><b>#pragma omp parallel private(x)</b></i> {  }	<i><b>!\$OMP PARALLEL PRIVATE(X)</b></i>  <i><b>!\$OMP END PARALLEL</b></i>
Function prototypes and types:	
<i><b>#include &lt;omp.h&gt;</b></i>	<i><b>use OMP_LIB</b></i>

- Most OpenMP constructs apply to a “structured block”.
  - **Structured block**: a block of one or more statements with one point of entry at the top and one point of exit at the bottom.
  - It’s OK to have an exit() within the structured block.

# Exercise, Part A: Hello World

## Verify that your environment works

- Write a program that prints “hello world”.

```
#include <stdio.h>
int main()
{

    printf(" hello ");
    printf(" world \n");

}
```

# Exercise, Part B: Hello World

## Verify that your OpenMP environment works

- Write a multithreaded program that prints “hello world”.

```
#include <omp.h>
#include <stdio.h>
int main()
{
    #pragma omp parallel
    {

        printf(" hello ");
        printf(" world \n");

    }
}
```

### Switches for compiling and linking

**gcc -fopenmp**      **Gnu (Linux, OSX)**

**cc -fopenmp**

**icc -fopenmp**      **Intel (Linux, OSX)**

# Solution

## A Multi-Threaded “Hello World” Program

- Write a multithreaded program where each thread prints “hello world”.

```
#include <omp.h>
#include <stdio.h>
int main()
{
#pragma omp parallel
{
    printf(" hello ");
    printf(" world \n");
}
}
```

OpenMP include file

Parallel region with default number of threads

End of the Parallel region

### Sample Output:

```
hello hello world
world
hello hello world
world
```

The statements are interleaved based on how the operating schedules the threads

# **A brief digression on the terminology of parallel computing**

# Let's agree on a few definitions:

- **Computer:**

- A machine that transforms *input values* into *output values*.
- Typically, a computer consists of Control, Arithmetic/Logic, and Memory units.
- The transformation is defined by a stored **program** (von Neumann architecture).

- **Task:**

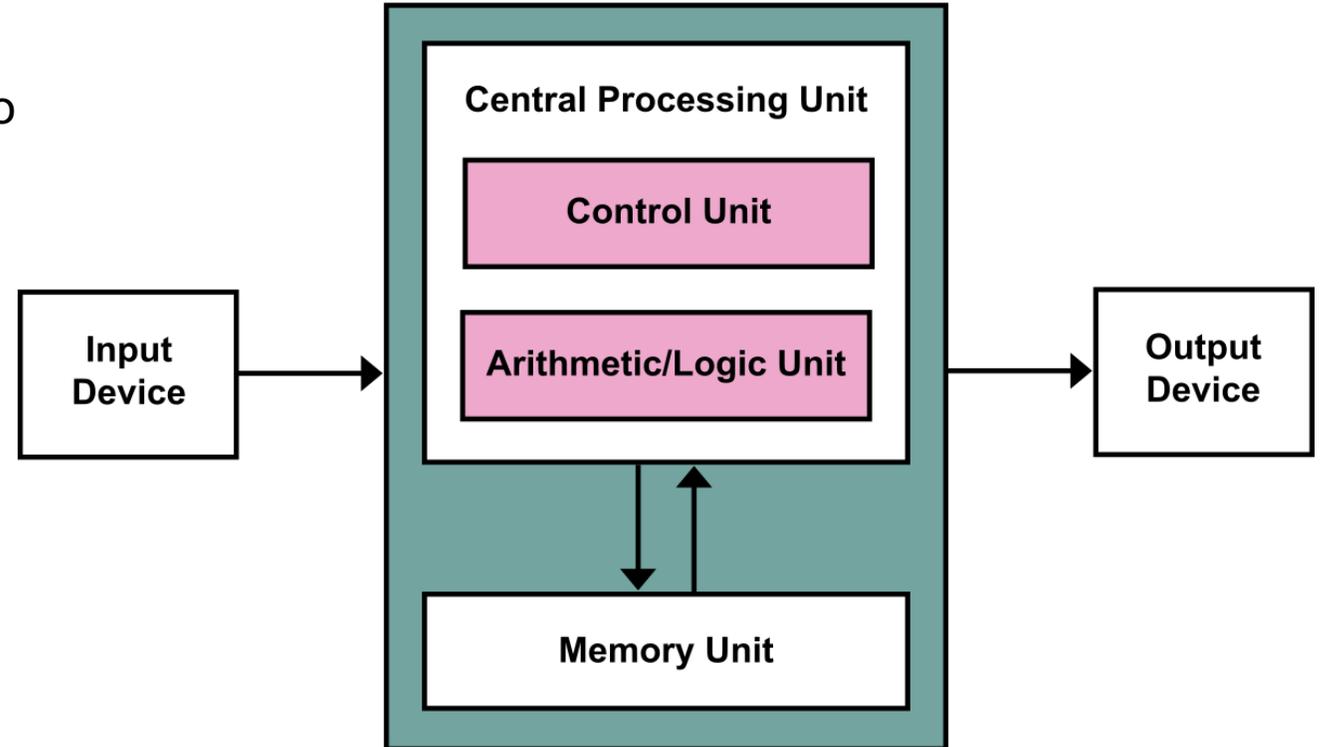
- A sequence of instructions plus a data environment. A program is composed of one or more tasks.

- **Active task:**

- A task that is available to be scheduled for execution. When the task is moving through its sequence of instructions, we say it is making **forward progress**

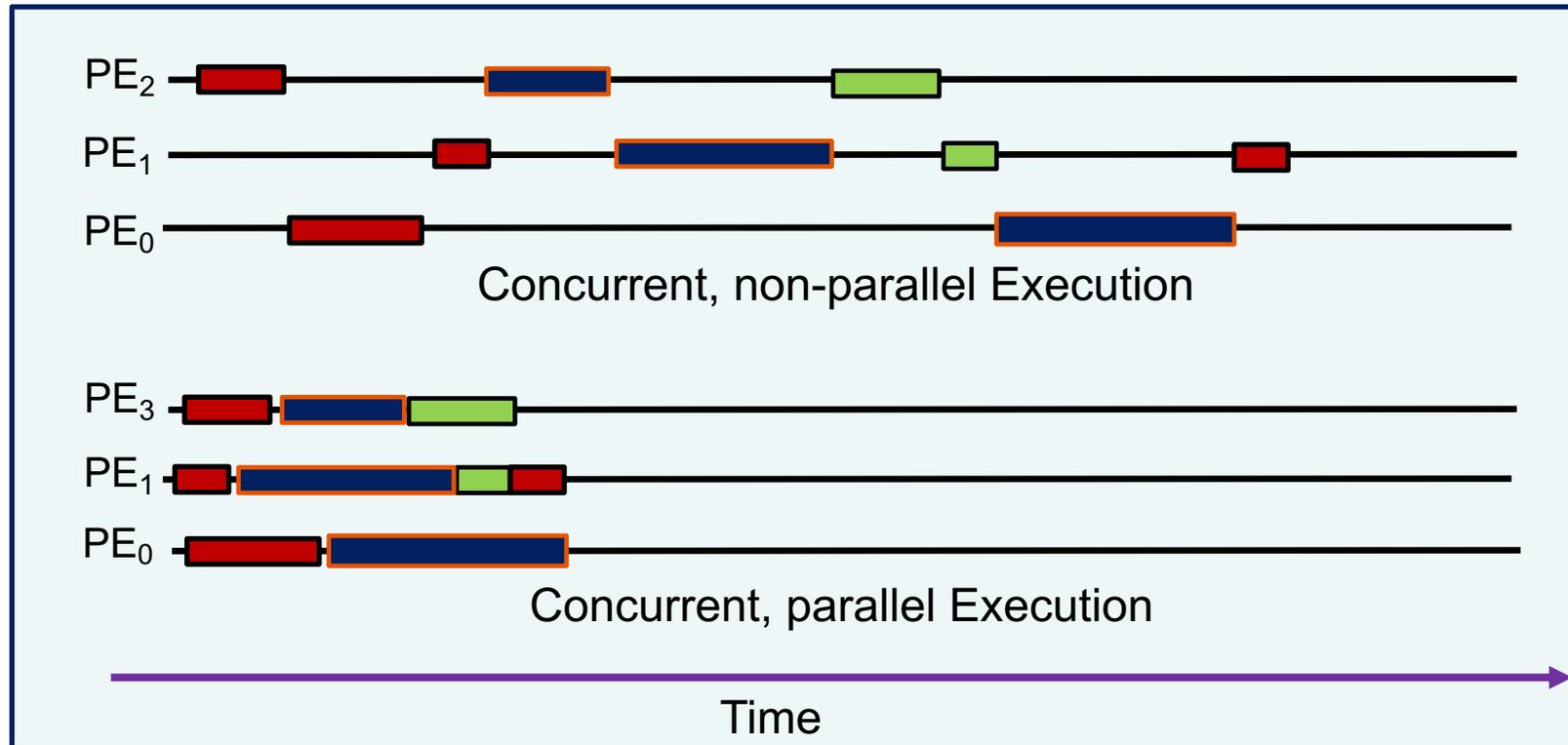
- **Fair scheduling:**

- When a scheduler gives each active task an equal *opportunity* for execution.



# Concurrency vs. Parallelism

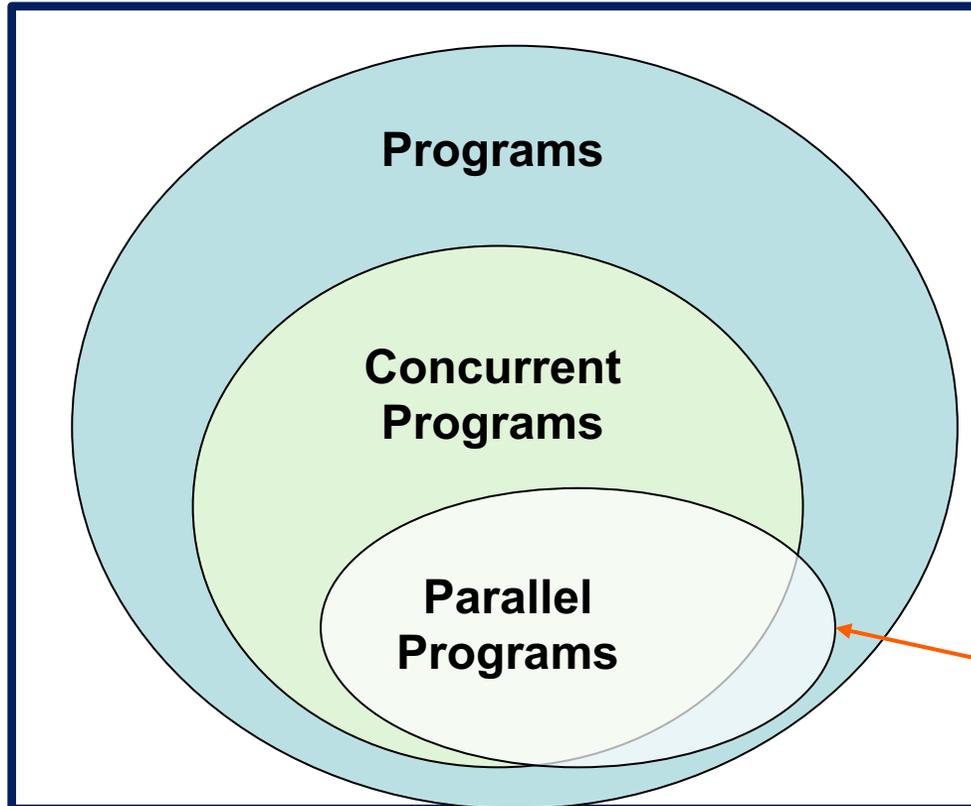
- Two important definitions:
  - Concurrency: A condition of a system in which multiple tasks are active and unordered. If **scheduled fairly**, they can be described as **logically** making **forward progress** at the same time.
  - Parallelism: A condition of a system in which multiple tasks are **actually** making **forward progress** at the same time.



PE = Processing Element

# Concurrency vs. Parallelism

- Two important definitions:
  - Concurrency: A condition of a system in which multiple tasks are active and unordered. If **scheduled fairly**, they can be described as **logically** making **forward progress** at the same time.
  - Parallelism: A condition of a system in which multiple tasks are **actually** making **forward progress** at the same time.



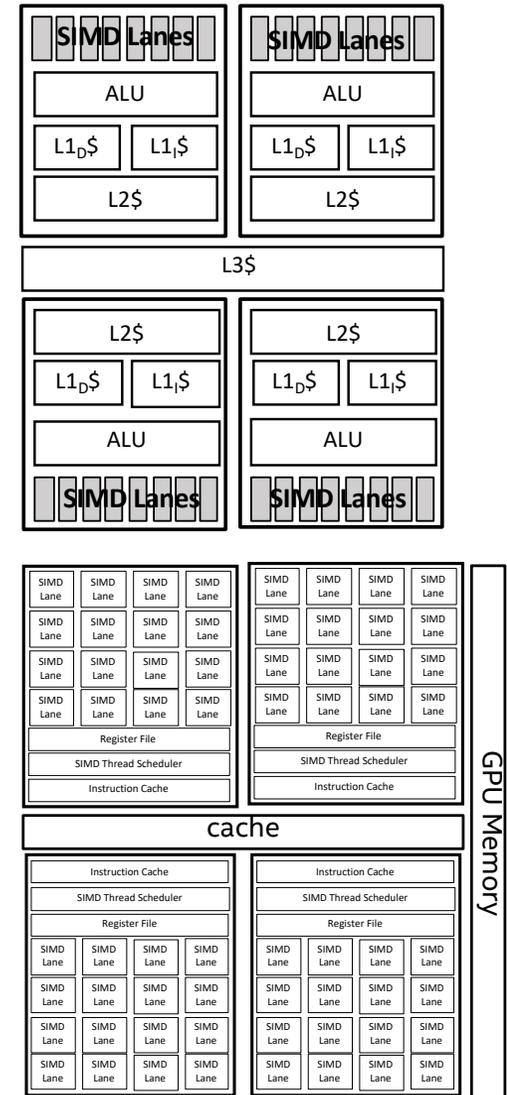
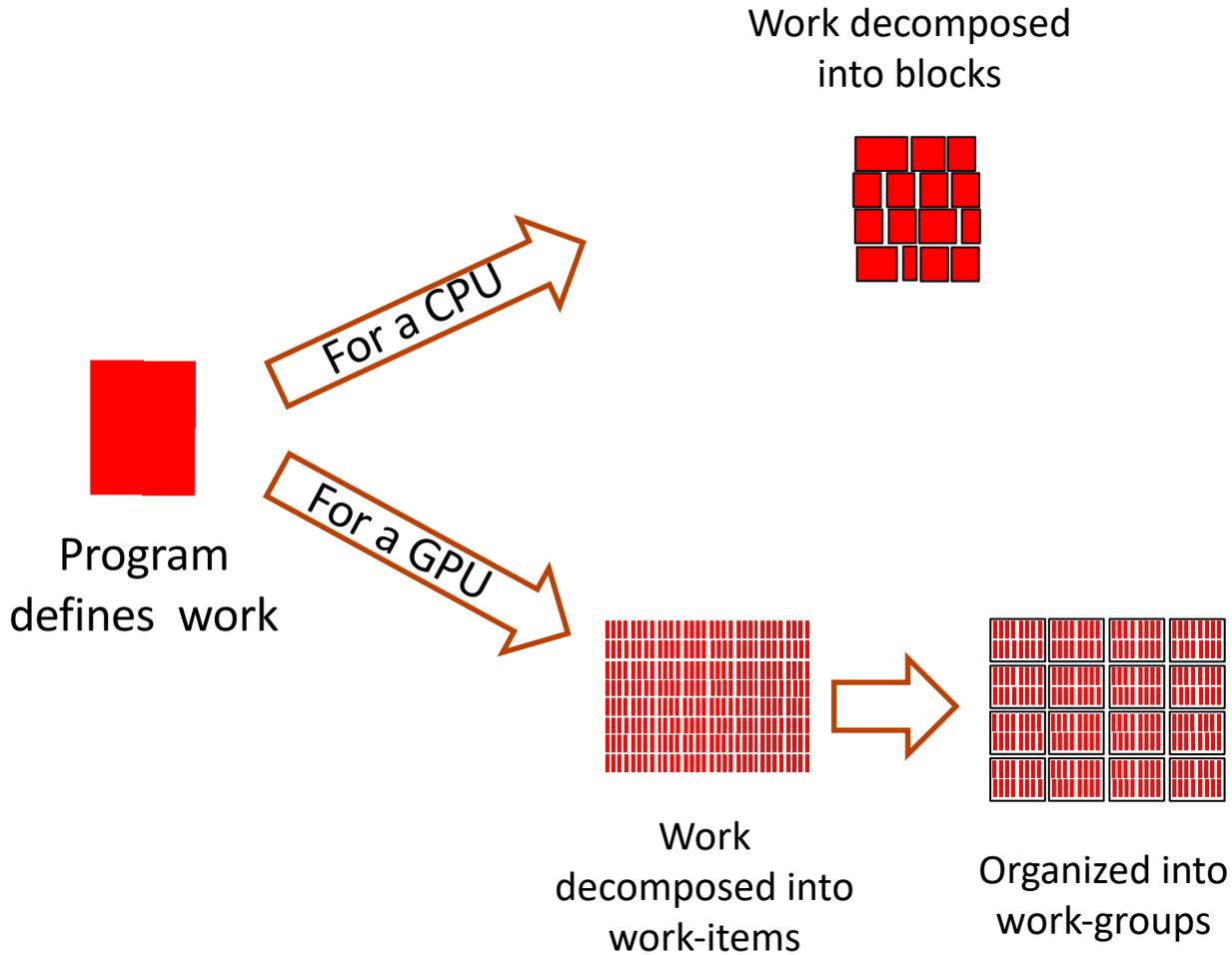
In most cases, parallel programs exploit concurrency in a problem to run tasks on multiple processing elements

We use Parallelism to:

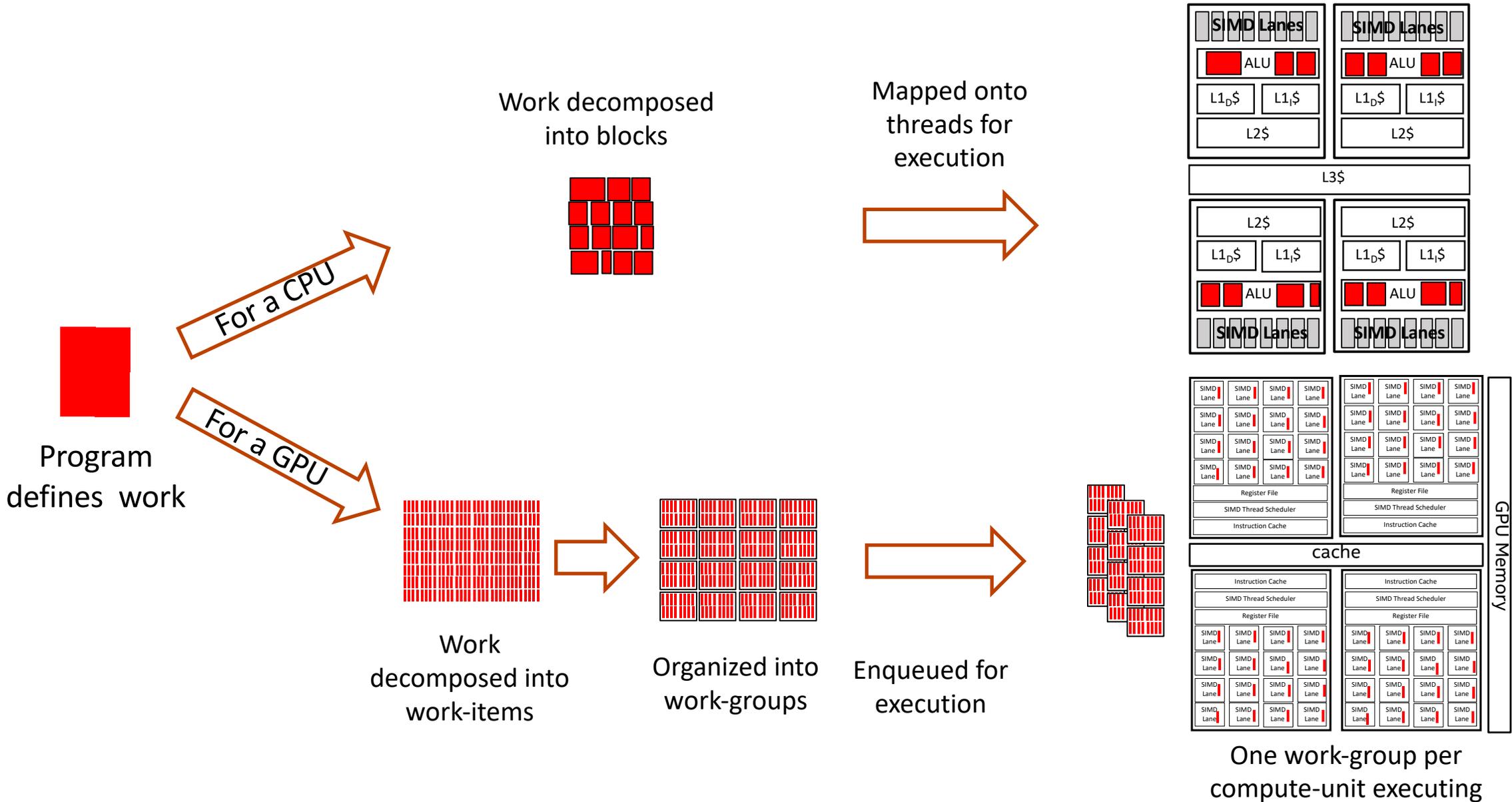
- Do more work in less time
- Work with larger problems

If tasks execute in “lock step” they are not concurrent, but they are still parallel.  
Example ... a SIMD unit.

# Executing a program on CPUs and GPUs

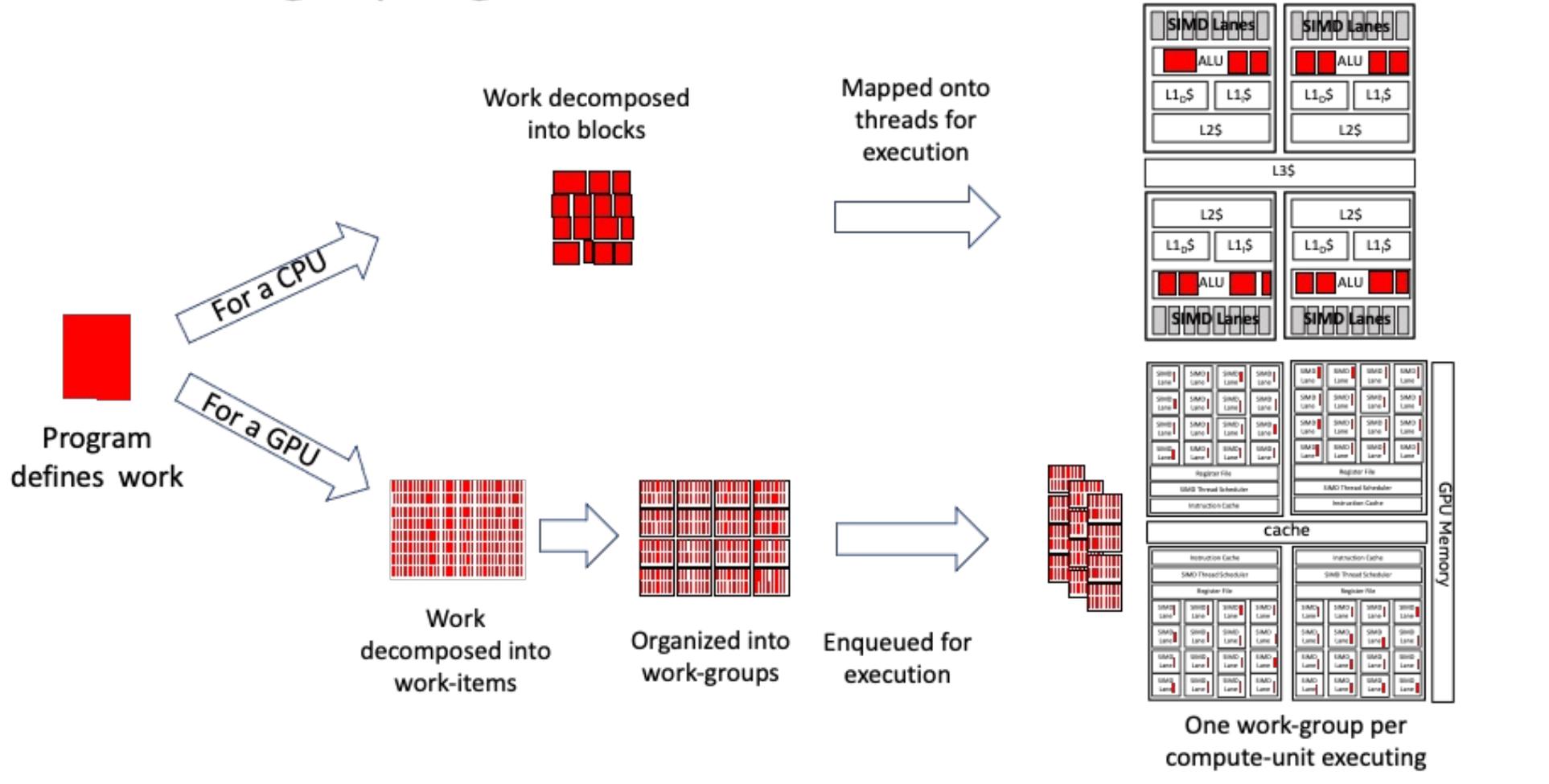


# Executing a program on CPUs and GPUs



# CPU/GPU execution modes

## Executing a program on CPUs and GPUs



For a CPU, the threads are all active and able to make forward progress.

For a GPU, any given work-group might be in the queue waiting to execute.

# Outline

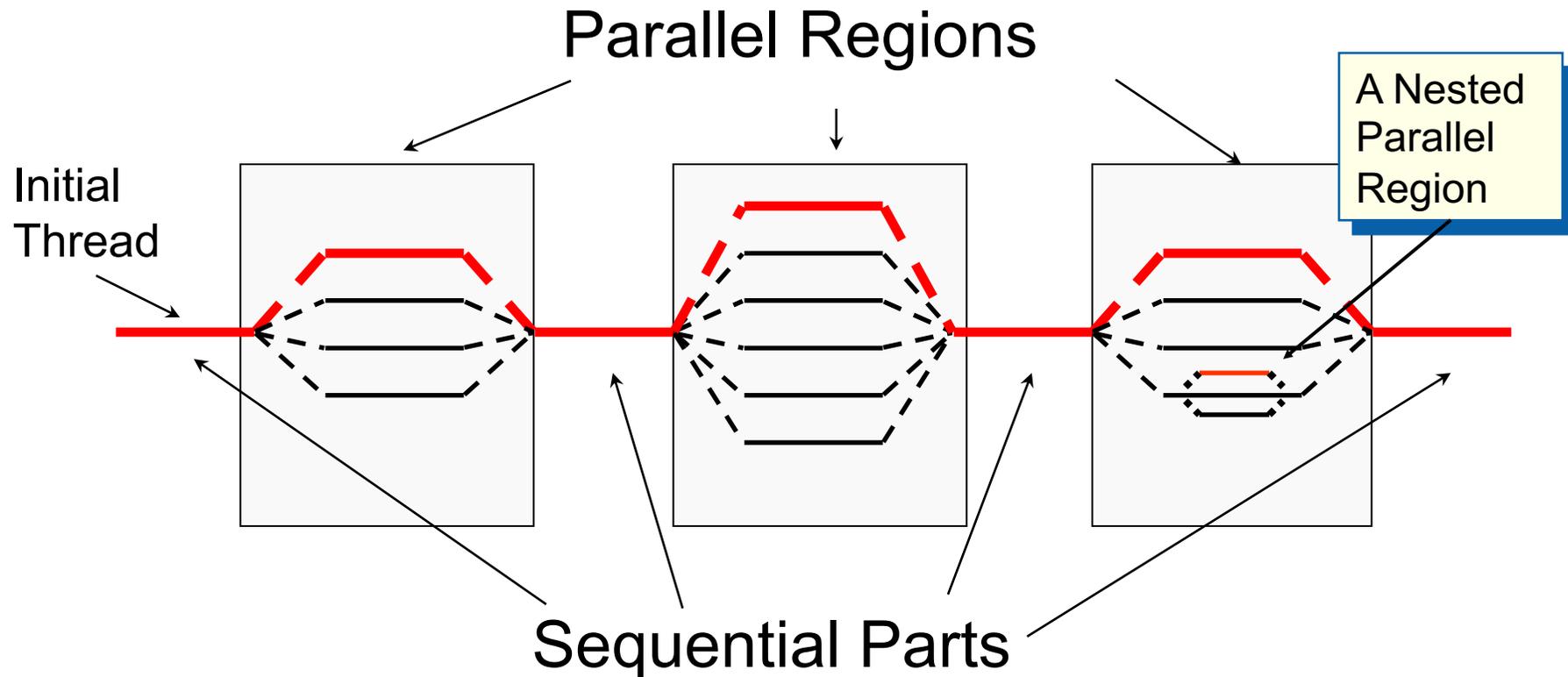
- Introduction to OpenMP
- • Creating Threads
- Synchronization
- Parallel Loops
- Data Environment
- Memory Model
- OpenMP stuff we didn't cover
- Recap

# OpenMP Execution model:

The fork-join model is also used with POSIX threads and `std::thread` in C++

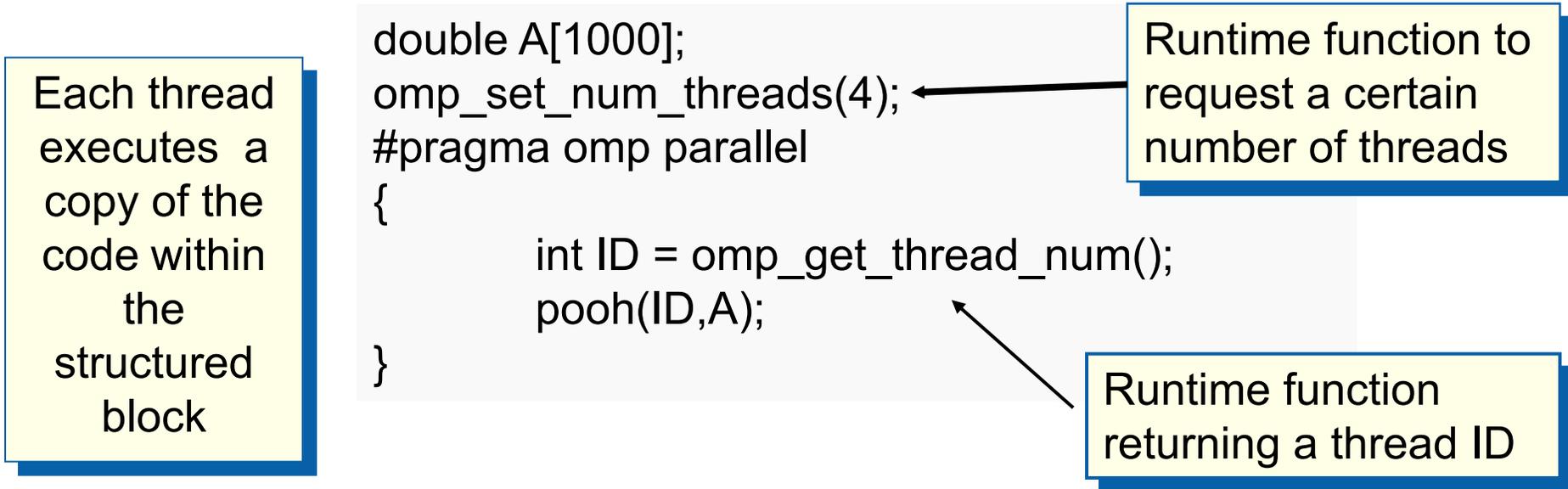
## Fork-Join Parallelism:

- ◆ Initial thread spawns a team of threads as needed.
- ◆ Parallelism added incrementally until performance goals are met, i.e., the sequential program evolves into a parallel program.



# Thread Creation: Parallel Regions

- You create threads in OpenMP with the parallel construct.
- For example, to create a 4 thread Parallel region:

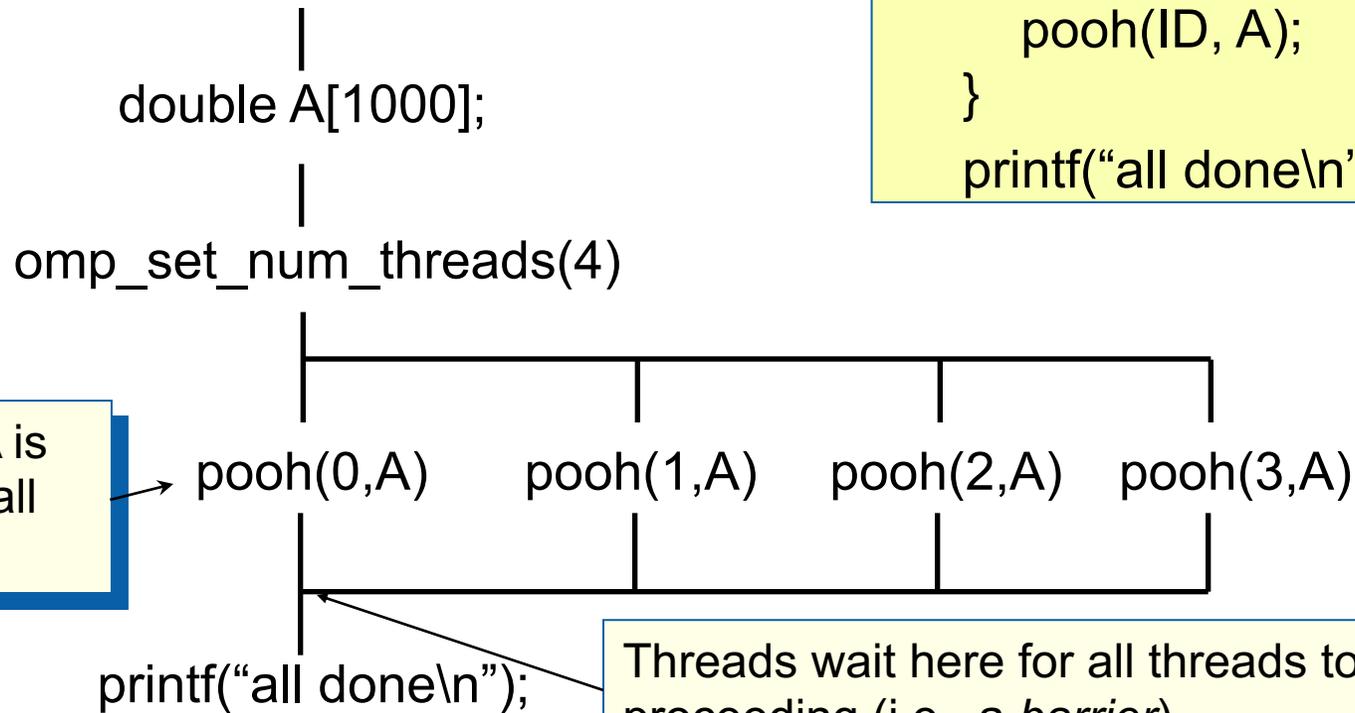


- Each thread calls `pooh(ID,A)` for `ID = 0 to 3`

# Thread Creation: Parallel Regions Example

- Each thread executes the same code redundantly.

```
double A[1000];
omp_set_num_threads(4);
#pragma omp parallel
{
    int ID = omp_get_thread_num();
    pooh(ID, A);
}
printf("all done\n");
```



A single copy of A is shared between all threads.

Threads wait here for all threads to finish before proceeding (i.e., a *barrier*)

# Thread creation: How many threads did you actually get?

- Request a number of threads with `omp_set_num_threads()`
- The number requested may not be the number you actually get.
  - An implementation may silently give you fewer threads than you requested.
  - Once a team of threads has launched, it will not be reduced.

Each thread executes a copy of the code within the structured block

```
double A[1000];
omp_set_num_threads(4);
#pragma omp parallel
{
    int ID = omp_get_thread_num();
    int nthrds = omp_get_num_threads();
    pooh(ID,A);
}
```

Runtime function to request a certain number of threads

Runtime function to return actual number of threads in the team

- Each thread calls `pooh(ID,A)` for `ID = 0` to `nthrds-1`



# Serial PI Program

```
static long num_steps = 100000;
double step;
int main ()
{
    double x, pi, sum = 0.0;

    step = 1.0/(double) num_steps;

    for (int i=0;i< num_steps; i++){
        x = (i+0.5)*step;
        sum = sum + 4.0/(1.0+x*x);
    }
    pi = step * sum;
}
```

# Serial PI Program

```
#include <omp.h>
static long num_steps = 100000;
double step;
int main ()
{
    double x, pi, sum = 0.0;

    step = 1.0/(double) num_steps;
double tdata = omp_get_wtime();
    for (int i=0;i< num_steps; i++){
        x = (i+0.5)*step;
        sum = sum + 4.0/(1.0+x*x);
    }
    pi = step * sum;
tdata = omp_get_wtime() - tdata;
    printf(" pi = %f in %f secs\n",pi, tdata);
}
```

The library routine `get_omp_wtime()` is used to find the elapsed "wall time" for blocks of code



# Hints: the Parallel Pi Program

- Use a parallel construct:

```
#pragma omp parallel
```

- The challenge is to:

- divide loop iterations between threads (use the thread ID and the number of threads).
- Create an accumulator for each thread to hold partial sums that you can later combine to generate the global sum.

- In addition to a parallel construct, you will need the runtime library routines

- `omp_set_num_threads();`
- `int omp_get_num_threads();`
- `int omp_get_thread_num();`
- `double omp_get_wtime();`









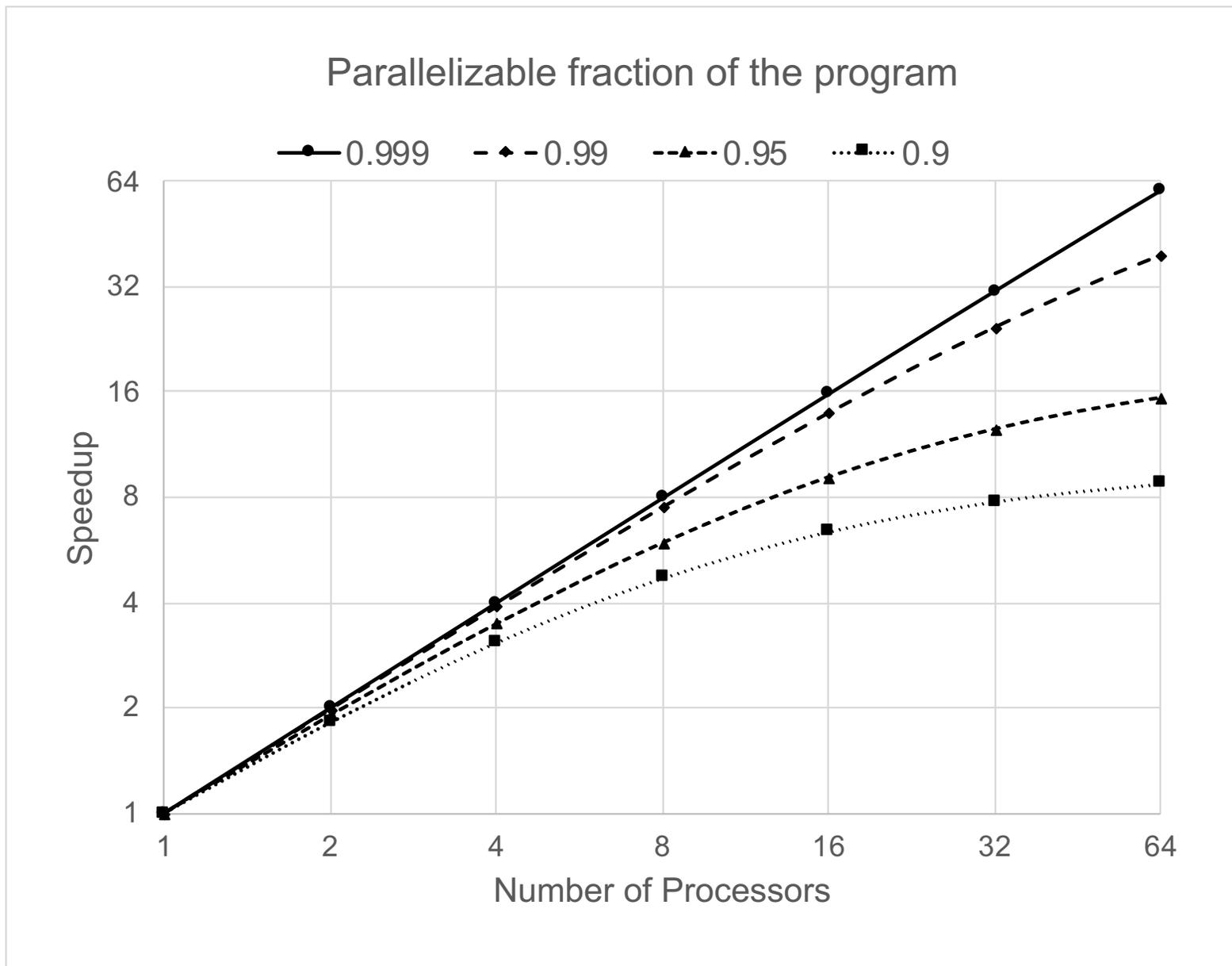
**A brief digression to talk about  
performance issues in parallel  
programs**







# Amdahl's Law



# So now you should understand my silly introduction slide.

## Introduction

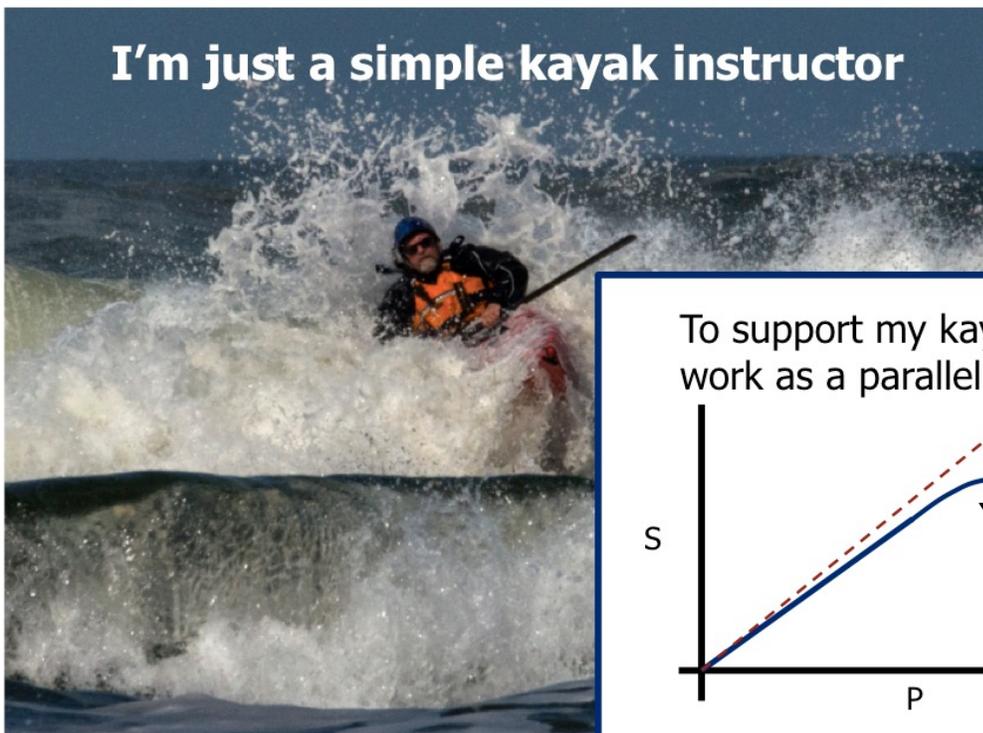
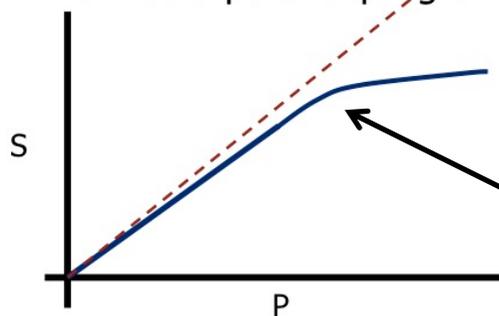


Photo © by Greg Clopton, 2014

We measure our success as parallel programmers by how close we come to ideal linear speedup.

To support my kayaking habit I work as a parallel programmer



Which means I know how to turn math into lines on a speedup plot

A good parallel programmer always figures out when you fall off the linear speedup curve and why that has occurred.



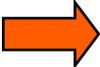








# Outline

- Introduction to OpenMP
- Creating Threads
-  • Synchronization
- Parallel Loops
- Data Environment
- Memory Model
- OpenMP stuff we didn't cover
- Recap

# Synchronization

Synchronization is used to impose order constraints and to protect access to shared data

- High level synchronization included in the common core:
  - critical
  - barrier
- Other, more advanced, synchronization operations:
  - atomic
  - ordered
  - flush
  - locks (both simple and nested)













# Example: Using a critical section to remove impact of false sharing

```
#include <omp.h>
static long num_steps = 100000;    double step;
#define NUM_THREADS 2
void main ()
{ int nthreads; double pi=0.0;      step = 1.0/(double) num_steps;
  omp_set_num_threads(NUM_THREADS);
  #pragma omp parallel
  {
    int i, id, nthrds;  double x, sum;
    id = omp_get_thread_num();
    nthrds = omp_get_num_threads();
    if (id == 0)  nthreads = nthrds;
    for (i=id, sum=0.0; i< num_steps; i=i+nthrds) {
      x = (i+0.5)*step;
      #pragma omp critical
      sum += 4.0/(1.0+x*x);
    }
  }
}
```

What would happen if you put the critical section inside the loop?













# Reduction

- How do we handle this case?

```
double ave=0.0, A[MAX];
for (int i=0;i< MAX; i++) {
    ave + = A[i];
}
ave = ave/MAX;
```

- We are combining values into a single accumulation variable (ave) ... there is a true dependence between loop iterations that can't be trivially removed.
- This is a very common situation ... it is called a "reduction".
- Support for reduction operations is included in most parallel programming environments.



# Reduction ... with arrays

- OpenMP reduction clause:  
reduction (op : list)
- Inside a parallel or a work-sharing construct:
  - A local copy of each list variable is made and initialized depending on the “op” (e.g. 0 for “+”).
  - Updates occur on the local copy.
  - Local copies are reduced into a single value and combined with the original global value.
- The variables in “list” must be shared in the enclosing parallel region.

```
double arr[N] = {0.0};  
#pragma omp parallel for reduction(+:arr[0:N])  
for (int j=0; j< M; j++) {  
    double val = A_Function(j);  
    for (int i=0; i<N; i++){  
        arr[i] += SomeFunction(i,val);  
    }  
}
```

Indicates an array section with N elements starting at element 0









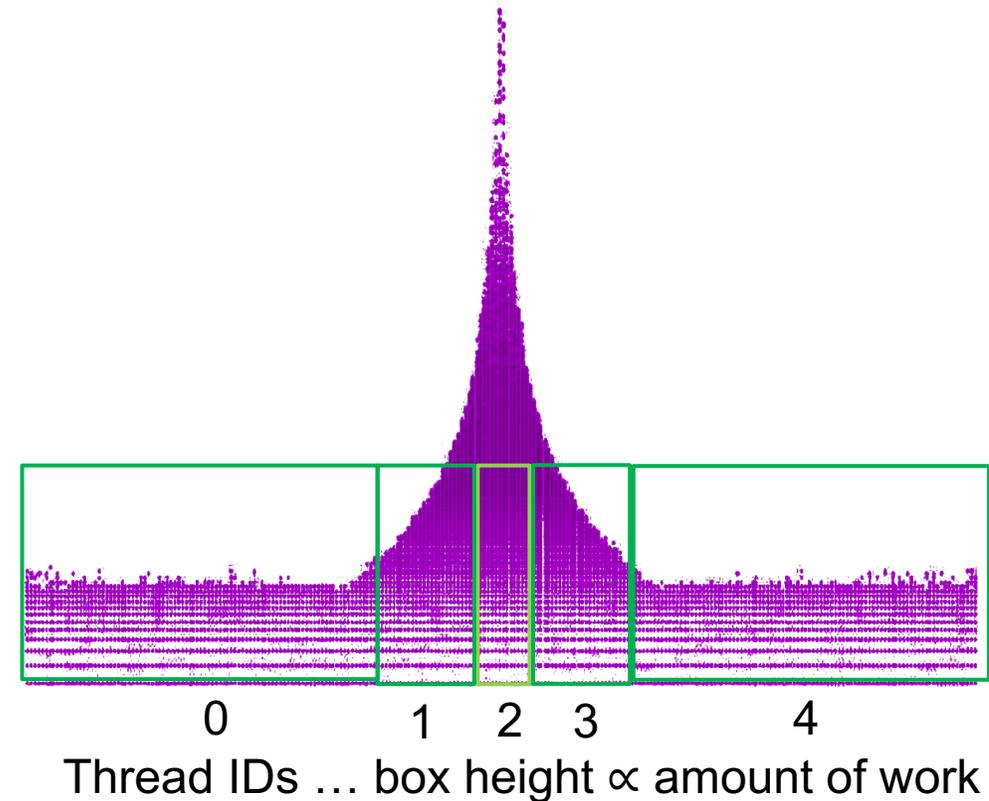






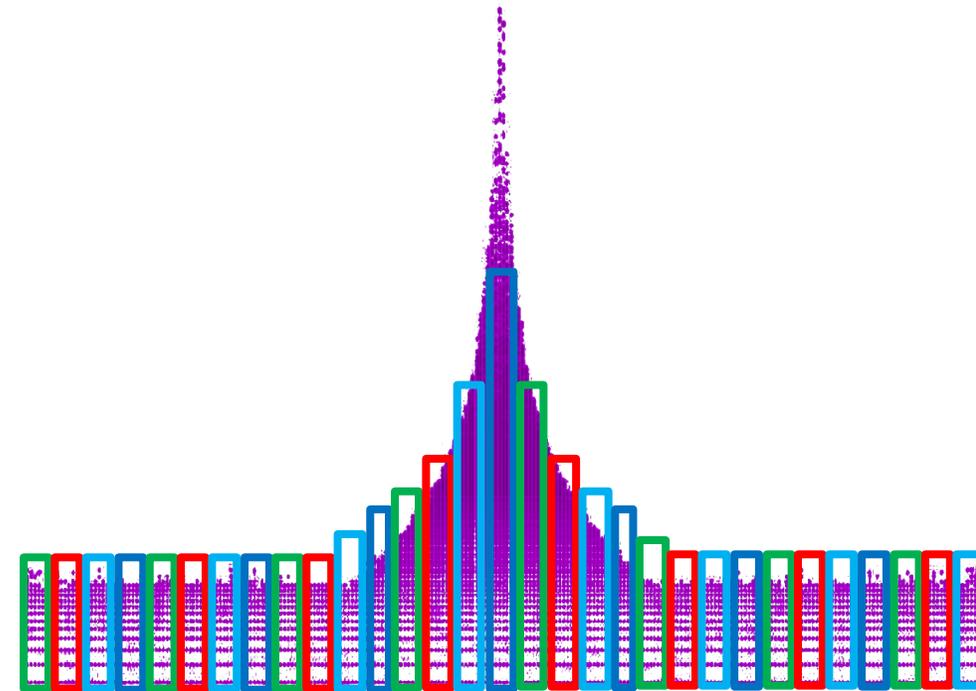
# Load Balancing

- A parallel job isn't done until the last thread is finished
- The work in our problem is unevenly distributed spatially.
- A key part of parallel programming is to design how you partition the work between threads so every thread has about the same amount of work.
- This topic is referred to as Load Balancing.
- In this case we adjusted the size of each chunk to equalize the work assigned to each thread.
  - Getting the right sized chunks for a variable partitioning (as done here) can be really difficult.



# Load Balancing

- A parallel job isn't done until the last thread is finished
- An easier path to Load Balancing.
  - Over-decompose the problem into small, fine-grained chunks
  - Spread the chunks out among the threads (in this case using a cyclic distribution)
  - The work is spread out and statistically, you are likely to get a good distribution of work



Colors mapped to 4 different Threads

- 0 
- 1 
- 2 
- 3 

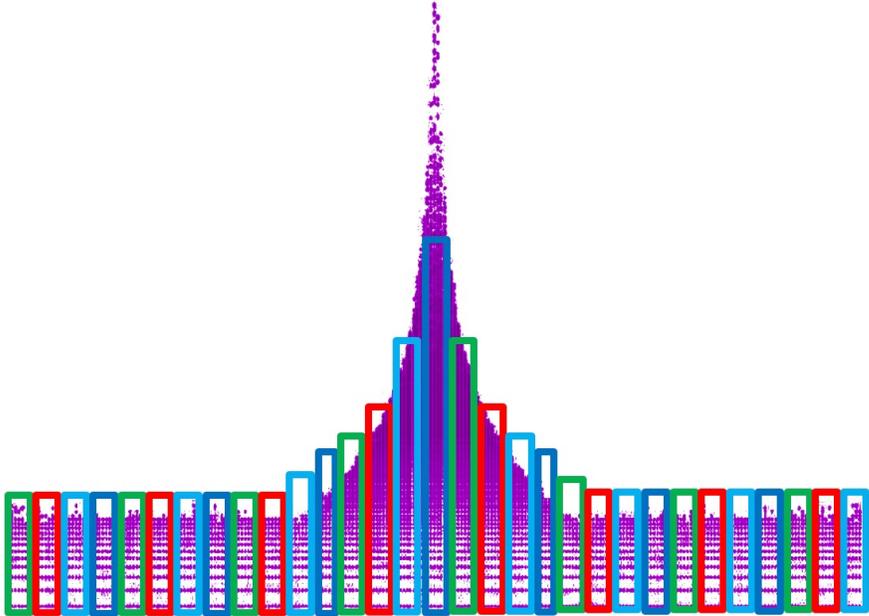




# Loop Worksharing Constructs: The schedule clause

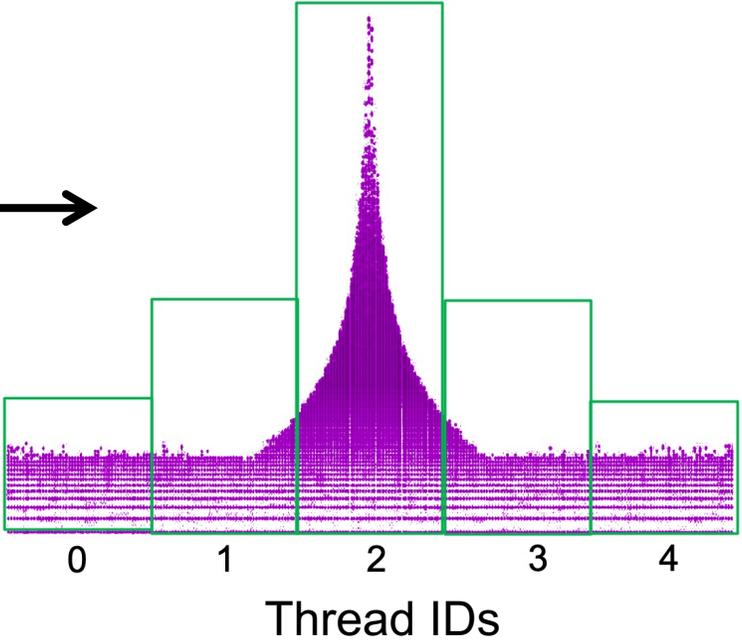
- The schedule clause ... most common cases:

```
#pragma omp parallel for schedule (static)
```



Colors mapped to 4 different Threads

- 0
- 1
- 2
- 3



```
Int small = 8; // loop iterations, i.e., width of boxes in the figure  
#pragma omp parallel for schedule (static, small)
```





# Outline

- Introduction to OpenMP
- Creating Threads
- Synchronization
- Parallel Loops
-  • Data Environment
- Memory Model
- OpenMP stuff we didn't cover
- Recap







# Data Sharing: Private clause

- `private(var)` creates a new local copy of `var` for each thread.

```
int N = 1000;
extern void init_arrays(int N, double *A, double *B, double *C);
```

```
void example () {
    int i, j;
    double A[N][N], B[N][N], C[N][N];
    init_arrays(N, *A, *B, *C);
```

```
    #pragma omp parallel for private(j)
    for (i = 0; i < 1000; i++)
        for(j = 0; j<1000; j++)
            C[i][j] = A[i][j] + B[i][j];
}
```

OpenMP makes the loop control index on the parallel loop (i) private by default ... but not for the second loop (j)

# Data Sharing: Private clause

- `private(var)` creates a new local copy of `var` for each thread.
  - The value of the private copies is uninitialized
  - The value of the original variable is unchanged after the region

```
void wrong() {  
    int tmp = 0;  
    #pragma omp parallel for private(tmp)  
    for (int j = 0; j < 1000; ++j)  
        tmp += j;  
    printf(“%d\n”, tmp);  
}
```

When you need to refer to the variable `tmp` that exists prior to the construct, we call it the **original variable**.

`tmp` is 0 here

`tmp` was not initialized



# Exercise: Mandelbrot set area

- The supplied program (mandel.c) computes the area of a Mandelbrot set.
- The program has been parallelized with OpenMP, but we were lazy and didn't do it right.
- Find and fix the errors.
- Once you have a working version, try to optimize the program.

```
#pragma omp parallel           #pragma omp parallel private (list)
#pragma omp for               #pragma omp parallel shared (list)
#pragma omp parallel for      #pragma omp parallel firstprivate (list)
#pragma omp critical          #pragma omp parallel default(none)
int omp_get_num_threads();    #pragma omp for reduction(op:list)
int omp_get_thread_num();
double omp_get_wtime();
```

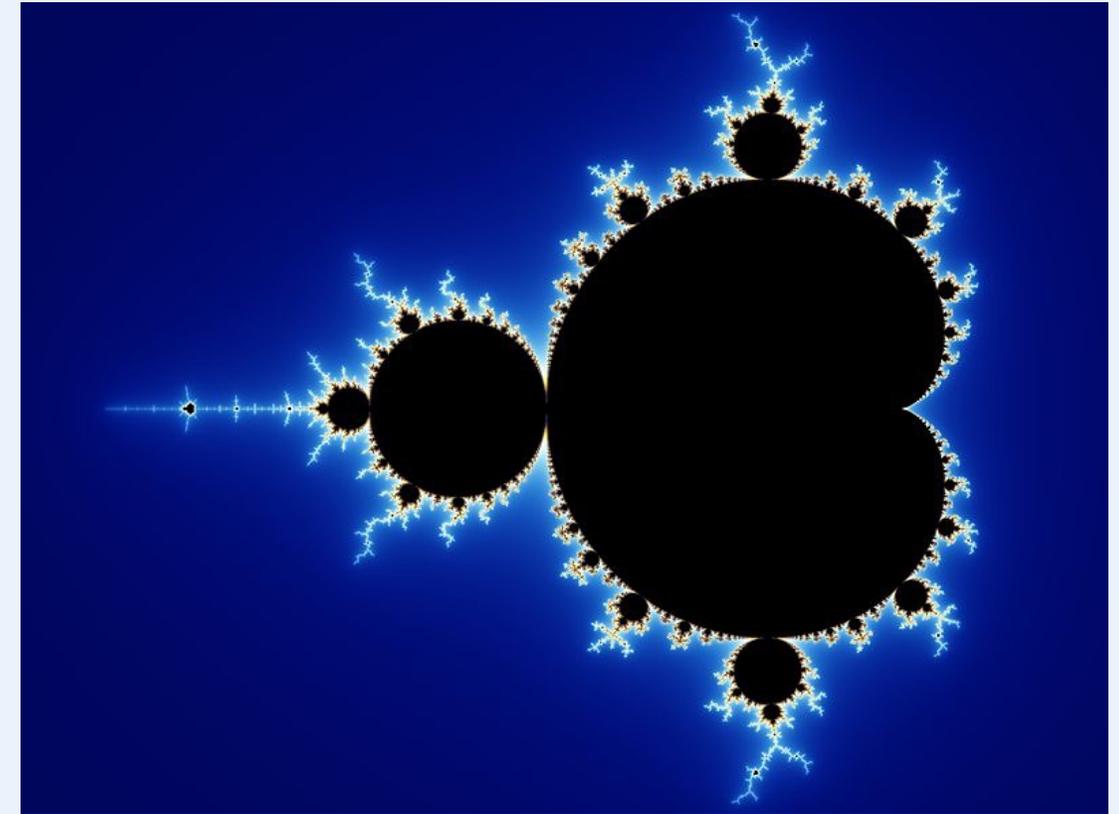


Image Source: Created by Wolfgang Beyer with the program Ultra Fractal 3. - Own work, CC BY-SA 3.0, <https://commons.wikimedia.org/w/index.php?curid=321973>

The Mandelbrot set ... The points,  $c$ , for which the following iterative map converges

$$z_{n+1} = z_n^2 + c$$

With  $z_n$  and  $c$  as complex numbers and  $z_0 = 0$ .





# Outline

- Introduction to OpenMP
- Creating Threads
- Synchronization
- Parallel Loops
- Data Environment
-  • Memory Model
- OpenMP stuff we didn't cover
- Recap







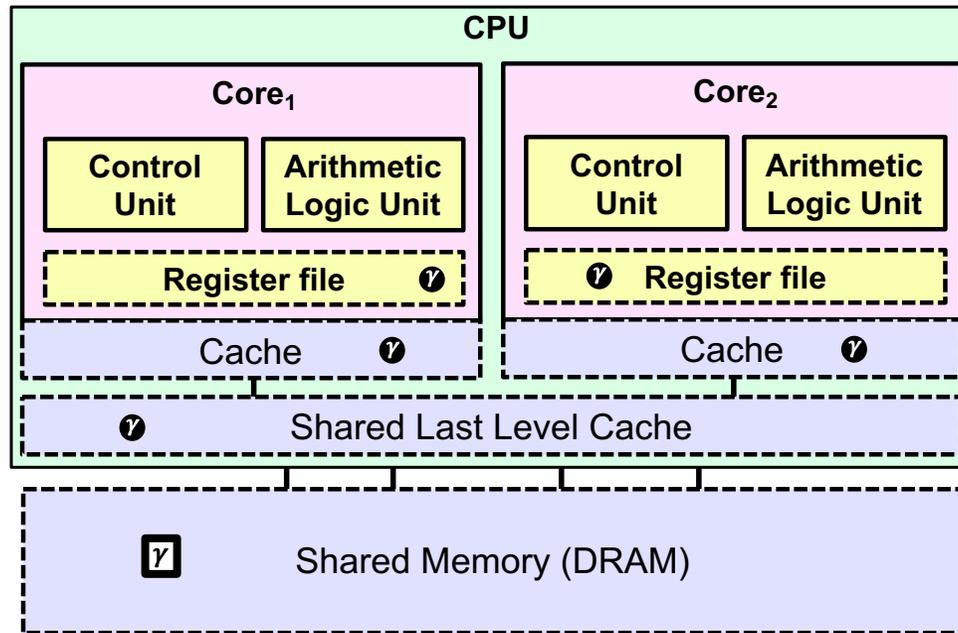


# Memory Models ...

- The fundamental issue is how do the values of variables across the memory hierarchy interact with the statements executed by two or more threads?
- Two options:

## 1. Sequential Consistency

- Threads execute and the associated loads/stores appear in some order defined by the semantically allowed interleaving of program statements.
- **All threads see the same interleaved order of loads and stores**



## 2. Relaxed Consistency

- Threads execute and the associated loads/stores appear in some order defined by the semantically allowed interleaving of program statements.
- **Threads may see different orders of loads and stores**

Most (if not all) multithreading programming models assume **relaxed consistency**. Maintaining sequential consistency across the full program-execution adds too much synchronization overhead.

# Why did this program fail?

Two issues:

(1) Can thread 1 fail to see updates to **flag**?

(2) Can **flag** be 1 while **answer** is still 0?

```
#include <stdio.h>
#include <omp.h>
#define COUNT 1000000
int main()
{
    int answer = 0, flag= 0,err=0;
    for (int i=0; i<COUNT; i++) {
        flag = 0; answer=0;
        #pragma omp parallel shared(flag,answer) num_threads(2)
        {
            int id = omp_get_thread_num();
            if (id == 0) {
                answer = 42;
                flag = 1;
            }
            else if (id == 1){
                while (flag == 0) { }
                if(answer!=42) err++;
            }
        }
    }
    return 0;
}
```

The compiler can reorder statements, so **flag** is set to 1 before **answer** is set to **42**

Thread 1 can load **flag** from the register file. It may not even go to cache (let alone memory) to see an updated value.

Regardless of how the compiler orders stores to **answer** and **flag**, thread 1 may see a different order than thread 0

# Why did this program fail?

Two issues:

(1) Can thread 1 fail to see updates to **flag**?

(2) Can **flag** be 1 while **answer** is still 0?

```
#include <stdio.h>
#include <omp.h>
#define COUNT 1000000
int main()
{
    int answer = 0, flag= 0, err=0;
    for (int i=0; i<COUNT; i++) {
        flag = 0; answer=0;
        #pragma omp parallel shared(flag,answer) num_threads(2)
        {
            int id = omp_get_thread_num();
            if (id == 0) {
                answer = 42;
                flag = 1;
            }
            else if (id == 1){
                while (flag == 0) { }
                if(answer!=42) err++;
            }
        }
    }
    return 0;
}
```

The compiler can reorder statements, so **flag** is set to 1 before **answer** is set to **42**

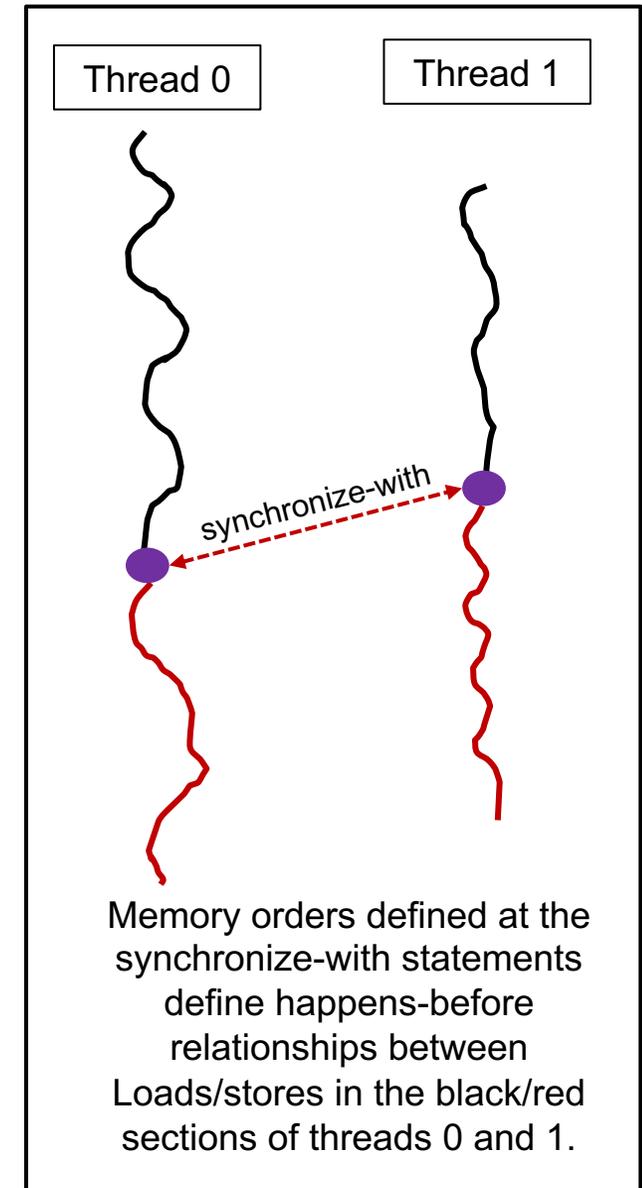
We need to enforce ordering constraints between the concurrent threads ... we need to consider the memory model and put the right synchronization constructs in place.

Thread 1 can load **flag** from the register file. It may not even go to cache (let alone memory) to see an updated value.

Regardless of how the compiler orders stores to **answer** and **flag**, thread 1 may see a different order than thread 0

# Memory Models: *Happens-before* and *synchronized-with* relations

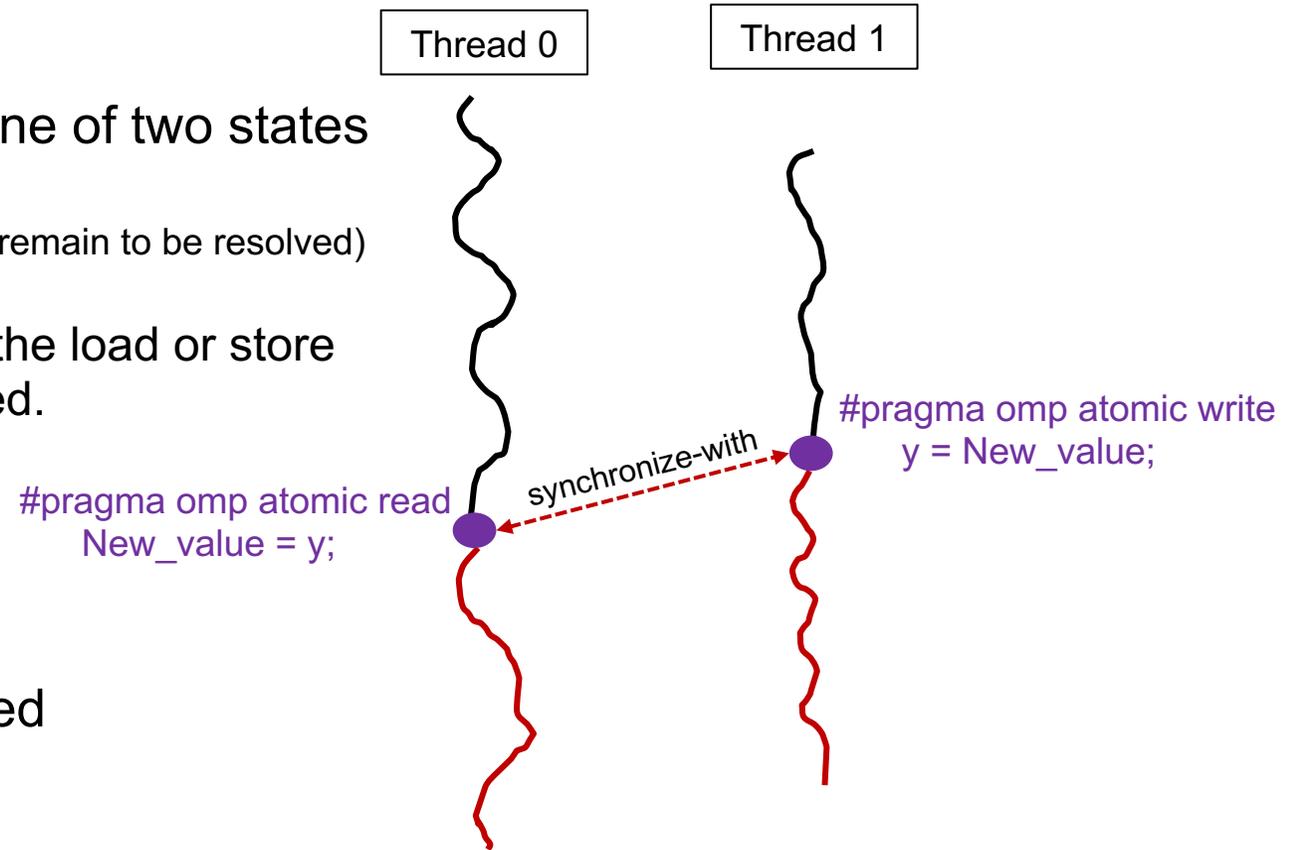
- Single thread execution:
  - Program order ... Loads and stores appear to occur in the order defined by the program's semantics. If you can't observe it, however, compilers can reorder instructions to maximize performance.
- Multithreaded execution ... concurrency in action
  - The compiler doesn't understand instruction-ordering across threads ... loads/stores to shared memory across threads can expose ambiguous orders of loads and stores
  - Instructions between threads are unordered except when specific ordering constraints are imposed, i.e., **synchronization**.
  - Synchronization lets us force that some instructions **happens-before** other instructions
- Two parts to synchronization:
  - A **synchronize-with** relationship exists at statements in 2 or more threads at which memory order constraints can be established.
  - **Memory order**: defines the view of loads/stores on either side of a synchronized-with operations.



# Atomic Operations and Synchronized-with

- An atomic operation can only be observed in one of two states
  - The operation has not happened yet
  - The operation has happened and is complete (no side-effects remain to be resolved)
- For example, on an atomic load or store operation, the load or store has happened and is complete, or it has not occurred.

- A **synchronized-with** relationship is established between a pair of atomic operations.
- The variables involved are visible to the programmer (such as with atomic constructs) or the variables are internal to a high level synchronization construct (barrier, critical, locks, etc).

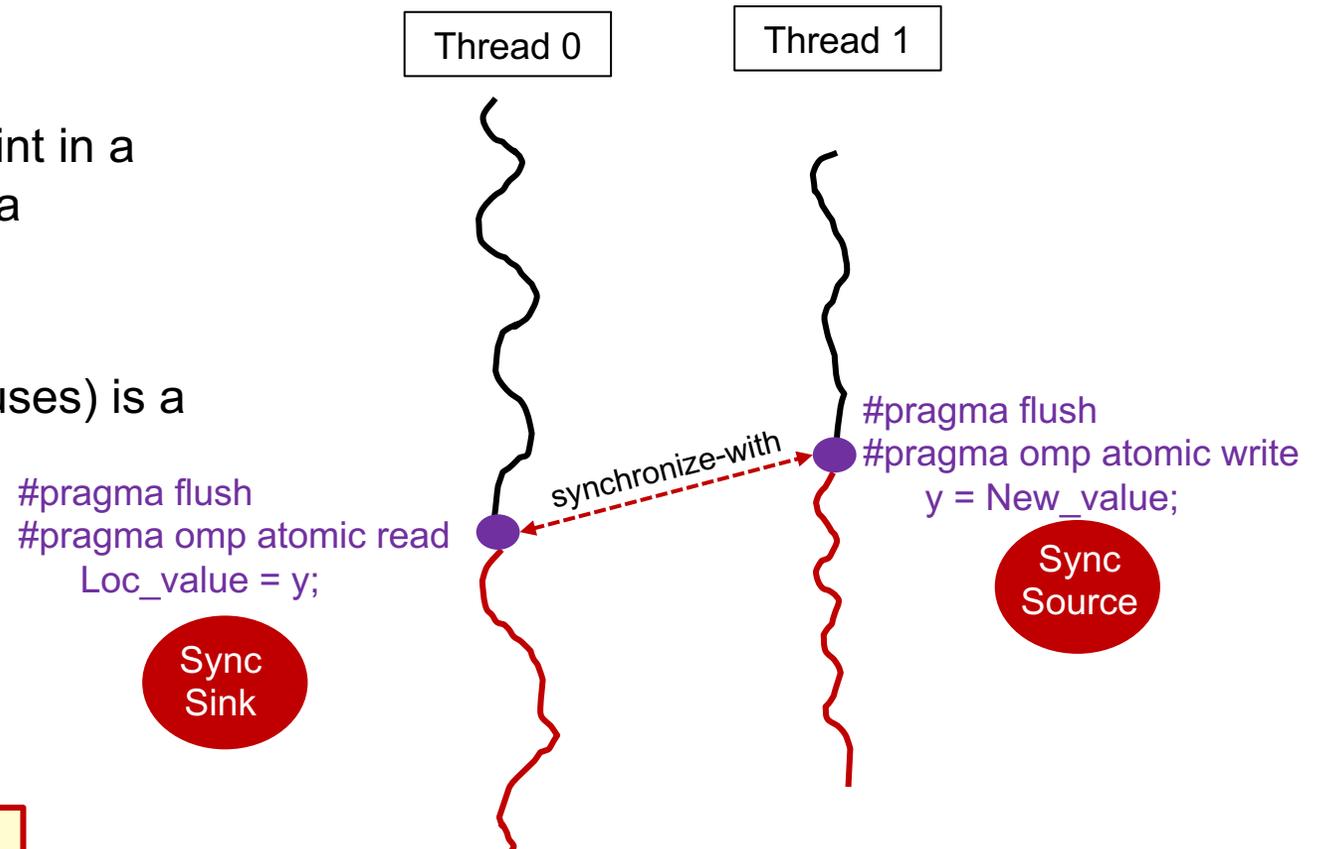


# Memory orders

- Memory orders establish which loads and stores can be moved around synchronized-with relations.
- The key construct is **flush**. ... flush defines a point in a program at which a thread is guaranteed to see a consistent view of memory.
- The default case for flush (i.e., no additional clauses) is a **strong flush**:
  - Previous read/writes by this thread have completed and are visible to other threads
  - No subsequent read/writes by this thread have occurred

A strong flush on its own does NOT define a synchronization point. The flush only addresses memory orders.

To synchronize threads, you need a synchronized-with relation which in this case, comes from an atomic write paired with an atomic read



Memory orders defined at the synchronize-with statements define happens-before relationships between Loads/stores in the black/red sections of threads 0 and 1.

Black operations on Thread 1 happen-before Red operations on thread 0.





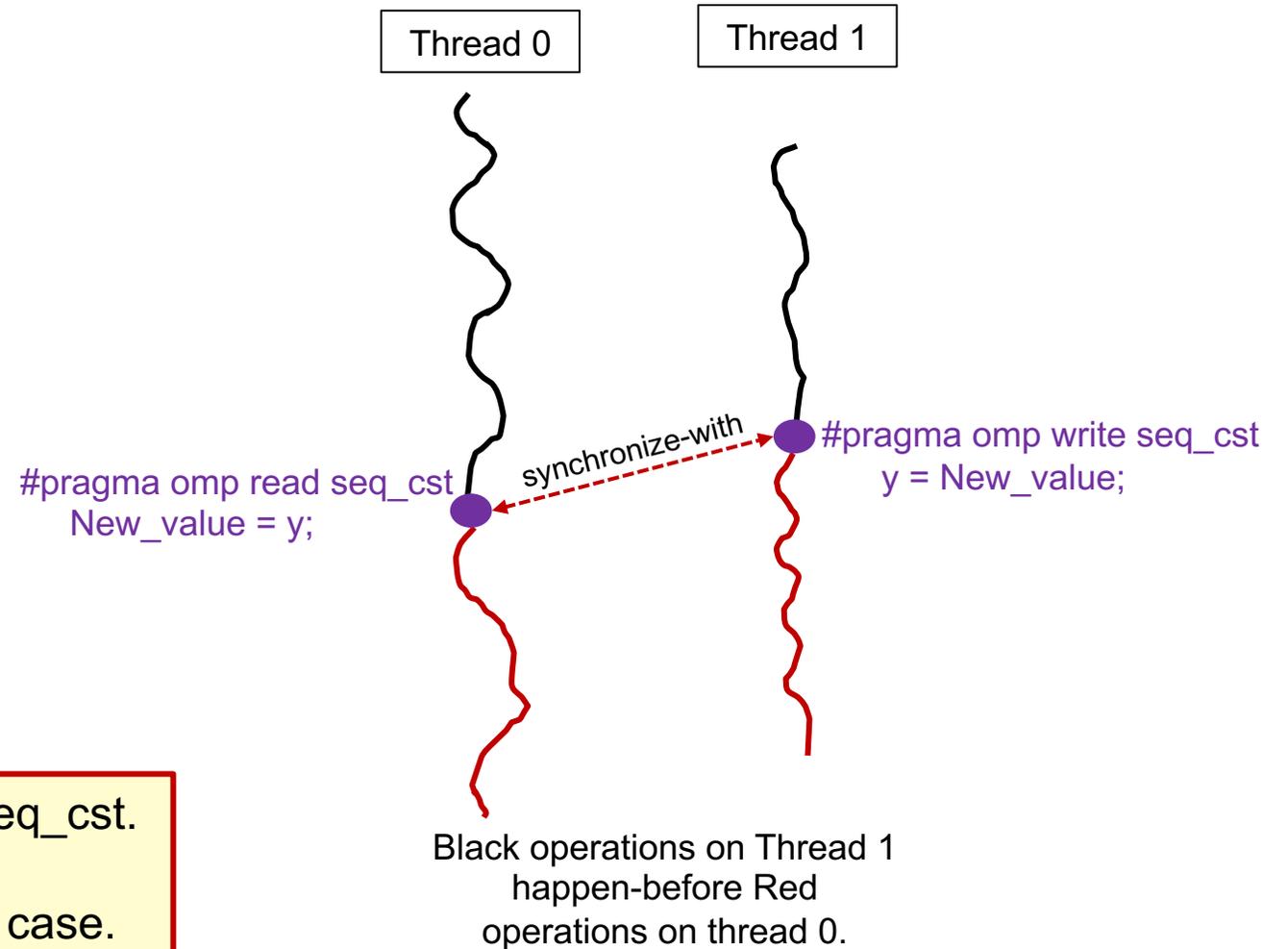
# producer/consumer program correctly synchronized

```
#include <stdio.h>
#include <omp.h>
#define COUNT 1000000
int main()
{
    int answer = 0, flag= 0,err=0;
    #pragma omp parallel shared(flag,answer) num_threads(2)
    {
        int id = omp_get_thread_num();
        if (id == 0) {
            answer = 42;
            #pragma omp write release
            flag = 1;
        }
        else if (id == 1){
            int fetch = 0;
            while (fetch == 0) {
                #pragma omp atomic read acquire
                fetch = flag;
            }
            if(answer!=42) err++;
        }
    }
    return 0;
}
```

# Other Memory orders in OpenMP

- Other OpenMP memory orders

- **acq\_rel**: Applies acquire and release memory order constraints at a single point in a program's execution.
- **seq\_cst**: sequential consistency. All data accessible to a thread are written to memory, subsequent writes are set to load from memory (akin to the strong flush)



The most important memory order to use is seq\_cst.

It can be more expensive, but it is the safest case.

# Keep it simple ... let OpenMP take care of Flushes for you

- A flush operation is implied by OpenMP constructs ...
  - at entry/exit of parallel regions
  - at implicit and explicit barriers
  - at entry/exit of critical regions

This has not been a detailed discussion of the full OpenMP memory model. The goal was to explain how memory models work and to understand the subset of features people commonly use.

- OpenMP programs that:
  - Do not use non-sequentially consistent atomic constructs;
  - Do not rely on the accuracy of a false result from `omp_test_lock` and `omp_test_nest_lock`; and
  - Correctly avoid data races

... behave as though operations on shared variables were simply interleaved in an order consistent with the order in which they are performed by each thread. The relaxed consistency model is invisible for such programs, and any explicit flushes in such programs are redundant.

## WARNING:

If you find yourself wanting to write code with explicit flushes, stop and get help. It is very difficult to manage flushes on your own. Even experts often get them wrong.

This is why we defined OpenMP constructs to automatically apply flushes most places where you really need them.





# Outline

- Introduction to OpenMP
- Creating Threads
- Synchronization
- Parallel Loops
- Data Environment
- Memory Model
-  • OpenMP stuff we didn't cover
- Recap



## **... But more importantly, we covered the most of the concepts you need to understand parallel computing**

- Parallelism vs concurrency ... and how it maps onto GPUs and CPUs
- How we talk about performance: Speedup and Amdahl's law
- SPMD and Loop level parallelism design patterns
- Load Balancing
- False sharing, race conditions, and other challenges of shared address spaces
- Memory models and the challenges of working with them
- Synchronization: mutual exclusion, barriers, and atomic operations with flush



# GPU programming: Vector addition with CUDA

Avoid CUDA ... it is proprietary and will lock you in to NVIDIA

```
// Compute sum of length-N vectors: C = A + B
void __global__
vecAdd (float* a, float* b, float* c, int N) {
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    if (i < N) c[i] = a[i] + b[i];
}

int main () {
    int N = ... ;
    float *a, *b, *c;
    cudaMalloc (&a, sizeof(float) * N);
    // ... allocate other arrays (b and c), fill with data

    // Use thread blocks with 256 threads each
    vecAdd <<< (N+255)/256, 256 >>> (a, b, c, N);
}
```

CUDA kernel as function

Unified shared memory ... allocate on host, visible on device too

Enqueue the kernel to execute on the Grid

# GPU programming: Vector addition with SYCL

```
// Compute sum of length-N vectors: C = A + B
#include <CL/sycl.hpp>

int main () {
    int N = ... ;
    float *a, *b, *c;
    sycl::queue q;
    *a = (float *)sycl::malloc_shared(N * sizeof(float), q);
    // ... allocate other arrays (b and c), fill with data

    q.parallel_for(sycl::range<1>{N},
         [=](sycl::id<1> i) {
            c[i] = a[i] + b[i];
        });
    q.wait();
}
```

Create a queue  
for SYCL  
commands

Unified shared  
memory ... allocate  
on host, visible on  
device too

Kernel as a C++  
Lambda function

[=] means capture external  
variables by value.

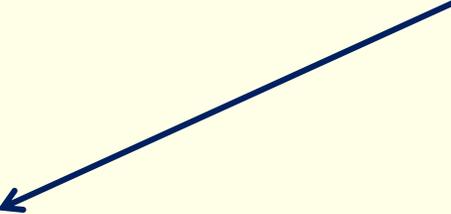
# Vector addition with OpenMP

- Let's add two vectors together ...  $C = A + B$

```
int main(){
    float *a, *b, *c;  int n = 10000;
    // allocate and fill a and b

    #pragma omp target map(to:a(0:n), b(0:n)) map(tofrom:c(0:n))
    #pragma omp loop
    for (i=0; i<n; i++)
        c[i] = a[i] + b[i];
}
```

This code is portable  
... it works with or  
without Unified  
Shared Memory



# Outline

- Introduction to OpenMP
- Creating Threads
- Synchronization
- Parallel Loops
- Data Environment
- Memory Model
- OpenMP stuff we didn't cover
- • Recap

# OpenMP Organizations

- OpenMP Architecture Review Board (ARB) URL, the “owner” of the OpenMP specification:

[www.openmp.org](http://www.openmp.org)

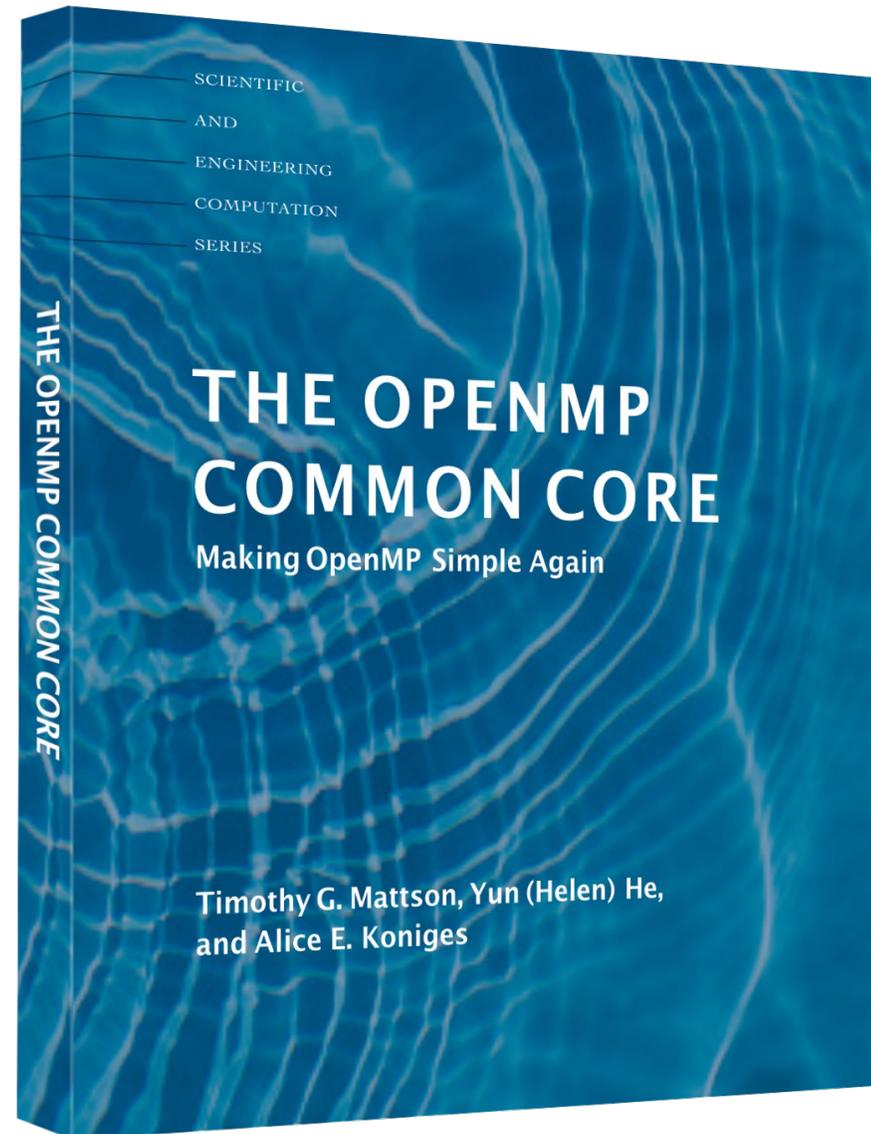
Get involved, get your organization to join the ARB.

Help define the future of OpenMP



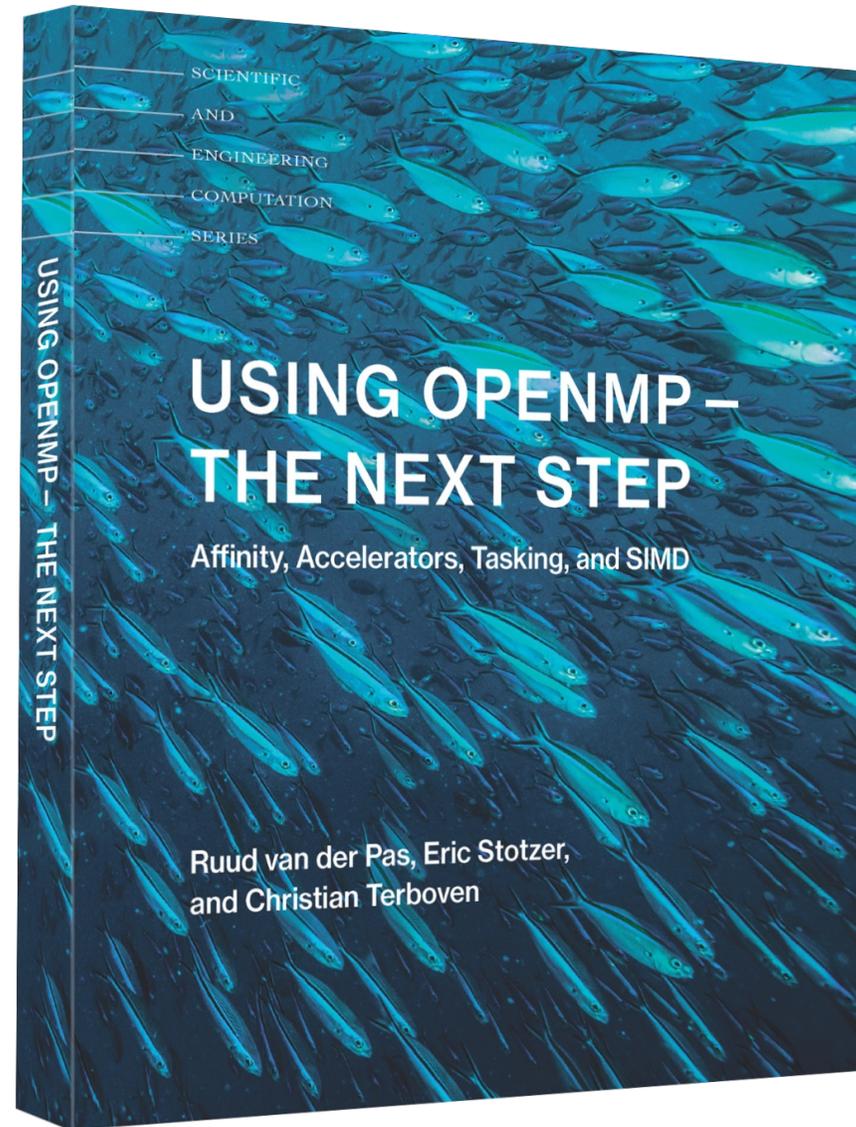
# To learn OpenMP:

- An exciting new book that Covers the Common Core of OpenMP plus a few key features beyond the common core that people frequently use
- It's geared towards people learning OpenMP, but as one commentator put it ... **everyone at any skill level should read the memory model chapters.**
- Available from MIT Press



# Books about OpenMP

A great book that covers  
OpenMP features beyond  
OpenMP 2.5

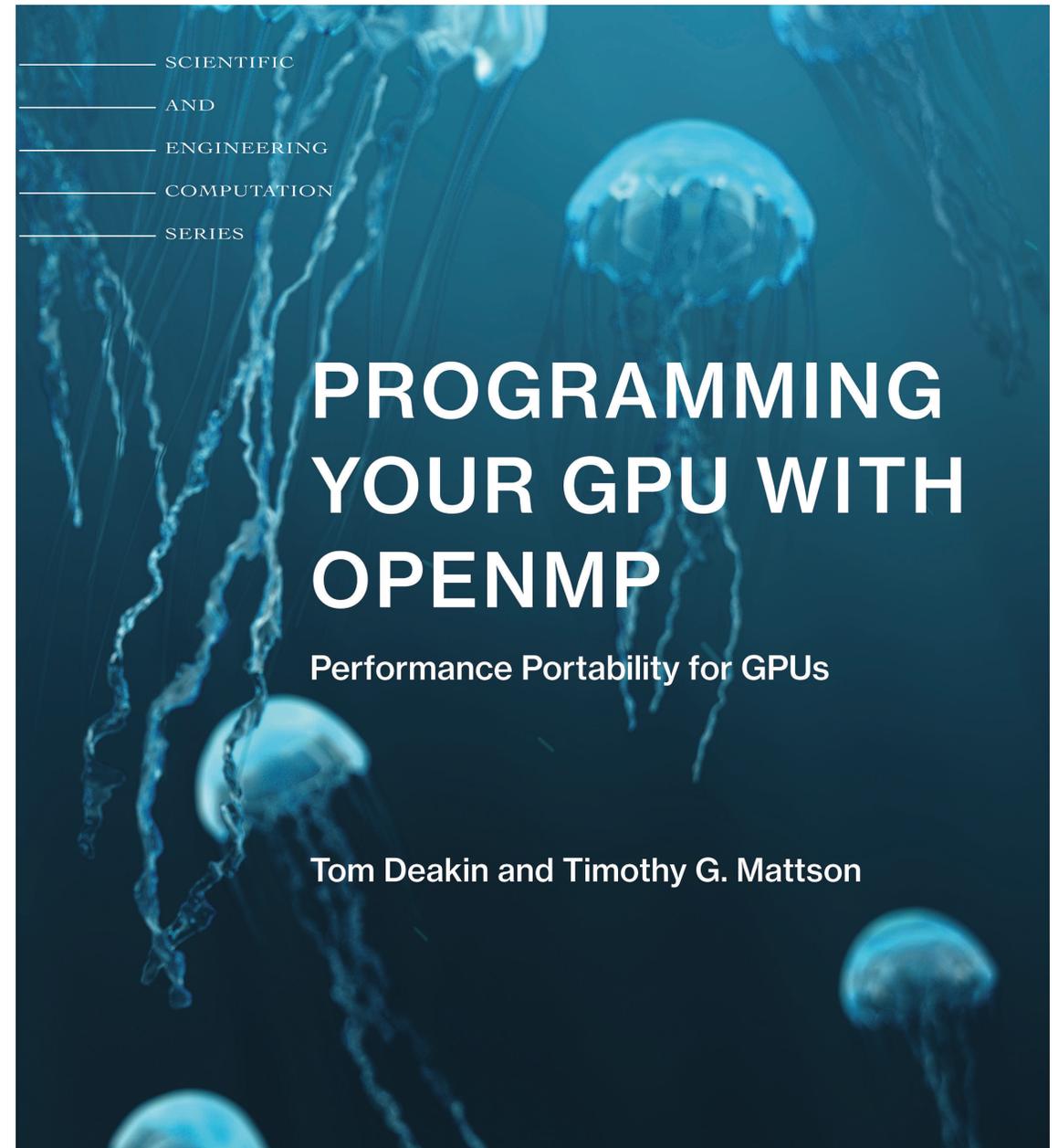


# Books about OpenMP

The latest book on OpenMP ...

Came out in early November 2023.

A book about how to use OpenMP to program a GPU.





# Backup content

- ➔ • Extra exercises
  - Irregular parallelism and OpenMP tasks
  - Worksharing Revisited
  - Synchronization Revisited: Additional options for Mutual exclusion

# Additional exercises for consolidation

- We provide three exercises you can use to solidify your understanding of OpenMP.
  1. Optimizing the Mandelbrot-area program.
  2. Building a histogram in parallel: This exercise explores different ways to synchronize code and how they impact performance. Given the importance of synchronization in multithreaded programming, this is an important skill to refine.
  3. Traversing linked lists: OpenMP added tasks to solve what we call irregular parallelism (i.e., they do not map onto basic for-loops). Parallelizing the traversal of linked lists with using tasks requires a lot of creativity. There are several ways to do this each with different impacts on performance.
- In some cases, the exercises content we presented, but didn't include in prior exercises. Hence, expect to need to review the lecture slides to figure out some key details.







# Backup content

- Extra exercises

- • Irregular parallelism and OpenMP tasks

- Worksharing Revisited

- Synchronization Revisited: Additional options for Mutual exclusion

# Irregular Parallelism

- Let's call a problem “irregular” when one or both of the following hold:
  - Data Structures are sparse or involve indirect memory references
  - Control structures are not basic for-loops
- Example: Traversing Linked lists:

```
p = listhead ;  
while (p) {  
    process(p) ;  
    p=p->next ;  
}
```

- Using what we've learned so far, traversing a linked list in parallel using OpenMP is difficult.

# Exercise: Traversing linked lists

- Consider the program linked.c
  - Traverses a linked list computing a sequence of Fibonacci numbers at each node.
- Parallelize this program selecting from the following list of constructs:

```
#pragma omp parallel
#pragma omp for
#pragma omp parallel for
#pragma omp for reduction(op:list)
#pragma omp critical
int omp_get_num_threads();
int omp_get_thread_num();
double omp_get_wtime();
schedule(static[,chunk]) or schedule(dynamic[,chunk])
private(), firstprivate(), default(none)
```

- Hint: Just worry about the while loop that is timed inside main(). You don't need to make any changes to the "list functions"

# Linked Lists with OpenMP (without tasks)

- See the file solutions/linked\_notasks.c

```
while (p != NULL) {
    p = p->next;
    count++;
}
struct node *parr = (struct node*) malloc(count*sizeof(struct node));
p = head;
for(i=0; i<count; i++) {
    parr[i] = p;
    p = p->next;
}
#pragma omp parallel
{
    #pragma omp for schedule(static,1)
    for(i=0; i<count; i++)
        processwork(parr[i]);
}
```

Count number of items in the linked list

Copy pointer to each node into an array

Process nodes in parallel with a for loop

Number of threads	Schedule	
	Default	Static,1
1	48 seconds	45 seconds
2	39 seconds	28 seconds

# Linked Lists with OpenMP (without tasks)

- See the file solutions/linked\_notasks.c

```
while (p != NULL) {
    p = p->next;
    count++;
}
struct node *parr = (struct node*) malloc(count*sizeof(struct node));
p = head;
for(i=0; i<count; i++) {
    parr[i] = p;
    p = p->next;
}
#pragma omp parallel
{
    #pragma omp for schedule(static,1)
    for(i=0; i<count; i++)
        processwork(parr[i]);
}
```

Count number of items in the linked list

Copy pointer to each node into an array

Process nodes in parallel with a for loop

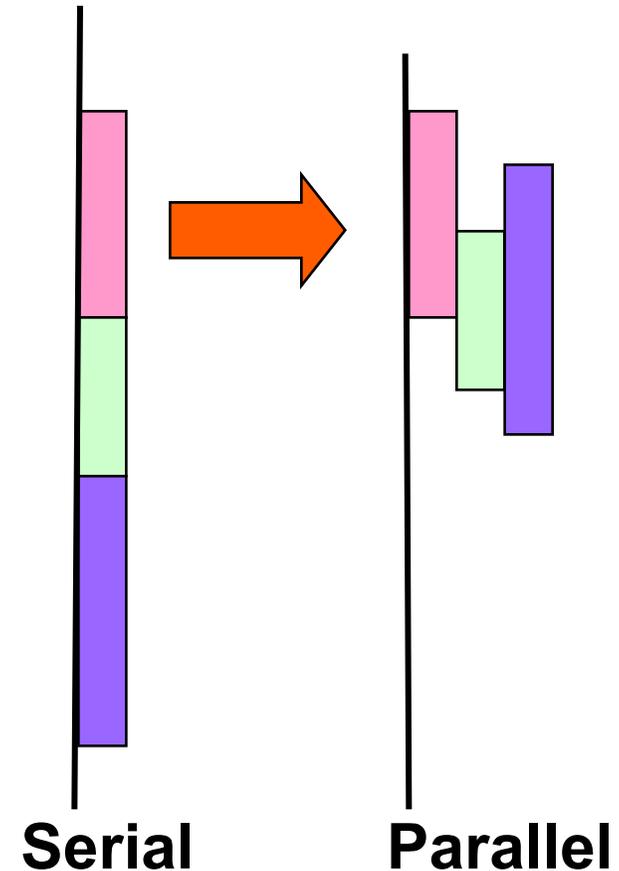
Number of threads	Schedule	
	Default	Static,1
1	48 seconds	45 seconds
2	39 seconds	28 seconds

With so much code to add and three passes through the data, this is really ugly.

There has got to be a better way to do this

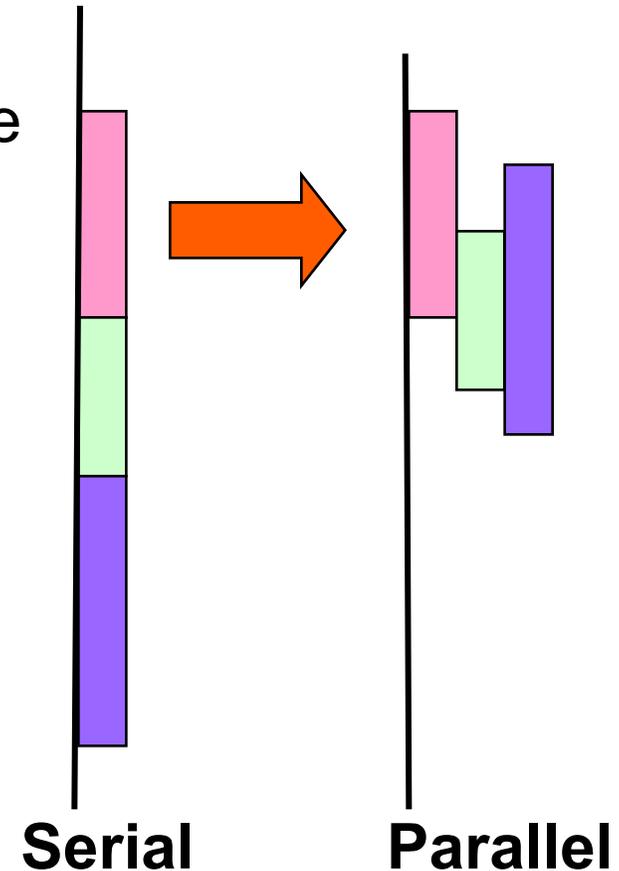
# What are Tasks?

- Tasks are independent units of work
- Tasks are composed of:
  - code to execute
  - data to compute with
- Threads are assigned to perform the work of each task.
  - The thread that encounters the task construct may execute the task immediately.
  - The threads may defer execution until later



# What are Tasks?

- The task construct includes a structured block of code
- Inside a parallel region, a thread encountering a task construct will package up the code block and its data for execution
- Tasks can be nested: i.e., a task may itself generate tasks.



A common Pattern is to have one thread create the tasks while the other threads wait at a barrier and execute the tasks

# Single Worksharing Construct

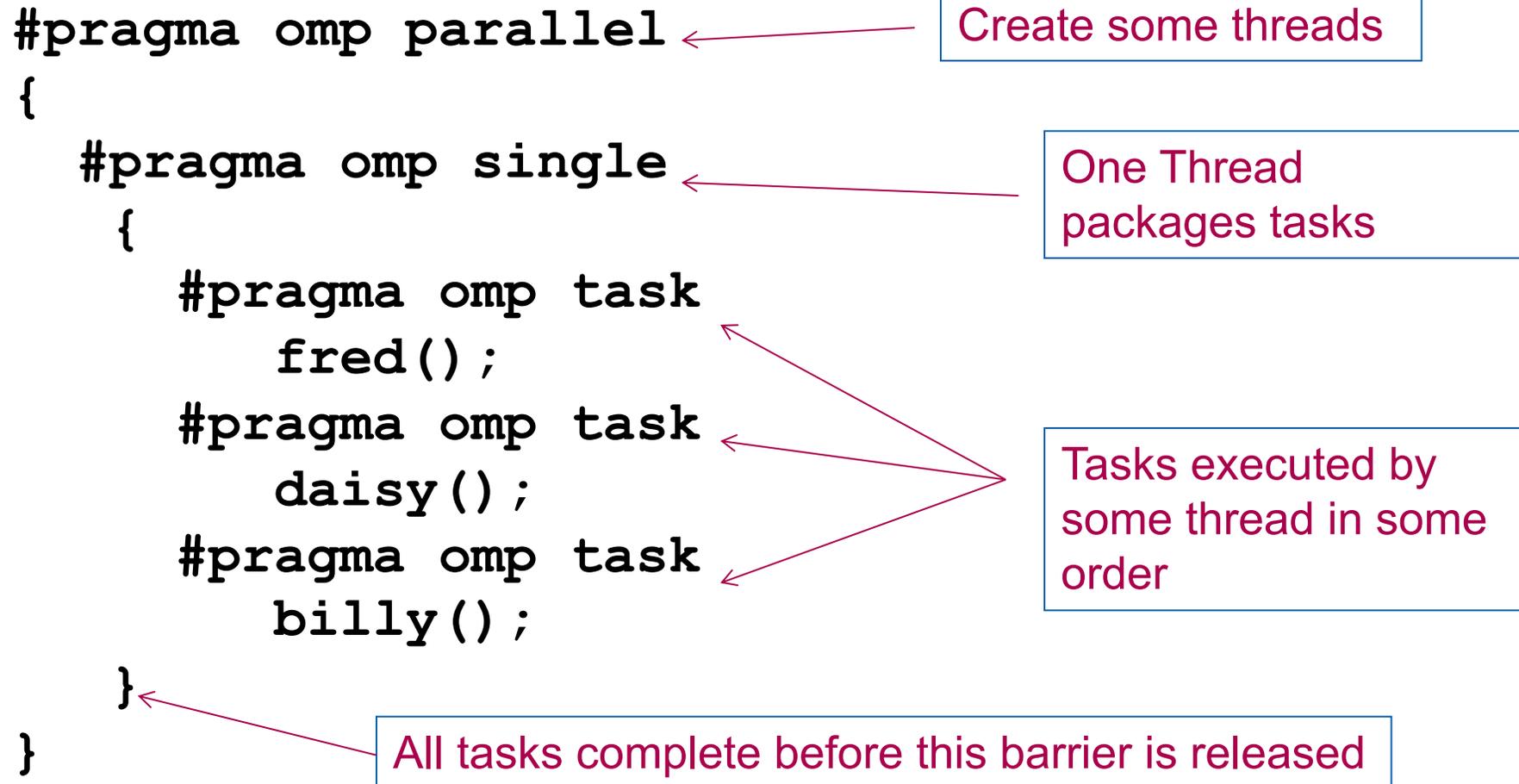
- The **single** construct denotes a block of code that is executed by only one thread (not necessarily the primary\* thread).
- A barrier is implied at the end of the single block (can remove the barrier with a *nowait* clause).

```
#pragma omp parallel
{
    do_many_things();
    #pragma omp single
    { exchange_boundaries(); }
    do_many_other_things();
}
```

\*This used to be called the “master thread”. The term “master” has been deprecated in OpenMP 5.1 and replaced with the term “primary”.

# Task Directive

```
#pragma omp task [clauses]  
    structured-block
```



# Exercise: Simple tasks

- Write a program using tasks that will “randomly” generate one of two strings:
  - “I think “ “race” “car” “s are fun”
  - “I think “ “car” “race” “s are fun”
- Hint: use tasks to print the indeterminate part of the output (i.e. the “race” or “car” parts).
- This is called a “Race Condition”. It occurs when the result of a program depends on how the OS schedules the threads.
- NOTE: A “data race” is when threads “race to update a shared variable”. They produce race conditions. Programs containing data races are undefined (in OpenMP but also ANSI standards C++’11 and beyond).

```
#pragma omp parallel
```

```
#pragma omp task
```

```
#pragma omp single
```

# Racey Cars: Solution

```
#include <stdio.h>
#include <omp.h>
int main()
{ printf("I think");
  #pragma omp parallel
  {
    #pragma omp single
    {
      #pragma omp task
      printf(" car");
      #pragma omp task
      printf(" race");
    }
  }
  printf("s");
  printf(" are fun!\n");
}
```

# Data Scoping with Tasks

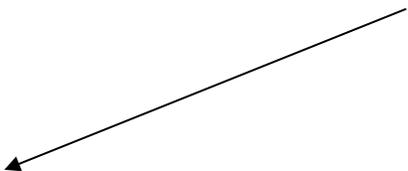
- Variables can be shared, private or firstprivate with respect to task
- These concepts are a little bit different compared with threads:
  - If a variable is **shared** on a task construct, the references to it inside the construct are to the storage with that name at the point where the task was encountered
  - If a variable is **private** on a task construct, the references to it inside the construct are to new uninitialized storage that is created when the task is executed
  - If a variable is **firstprivate** on a construct, the references to it inside the construct are to new storage that is created and initialized with the value of the existing storage of that name when the task is encountered

# Data Scoping Defaults

- The behavior you want for tasks is usually firstprivate, because the task may not be executed until later (and variables may have gone out of scope)
  - Variables that are private when the task construct is encountered are firstprivate by default
- Variables that are shared in all constructs starting from the innermost enclosing parallel construct are shared by default

```
#pragma omp parallel shared(A) private(B)
{
    ...
    #pragma omp task
    {
        int C;
        compute(A, B, C);
    }
}
```

A is shared  
B is firstprivate  
C is private



# Exercise: Traversing linked lists

- Consider the program `linked.c`
  - Traverses a linked list computing a sequence of Fibonacci numbers at each node.
- Parallelize this program selecting from the following list of constructs:

```
#pragma omp parallel
#pragma omp single
#pragma omp task
int omp_get_num_threads();
int omp_get_thread_num();
double omp_get_wtime();
private(), firstprivate()
```

- Hint: Just worry about the contents of `main()`. You don't need to make any changes to the "list functions"

# Parallel Linked List Traversal

```
#pragma omp parallel
{
    #pragma omp single
    {
        p = listhead ;
        while (p) {
            #pragma omp task firstprivate(p)
            {
                process (p) ;
            }
            p=next (p) ;
        }
    }
}
```

Only one thread  
packages tasks

makes a copy of **p**  
when the task is  
packaged

# When/Where are Tasks Complete?

- At thread barriers (explicit or implicit)
  - all tasks generated inside a region must complete at the next barrier encountered by the threads in that region. Common examples:
    - **Tasks generated inside a single construct:** all tasks complete before exiting the barrier on the single.
    - **Tasks generated inside a parallel region:** all tasks complete before exiting the barrier at the end of the parallel region.
- At taskwait directive
  - i.e. Wait until all tasks defined in the current task have completed.  
`#pragma omp taskwait`
  - Note: applies only to tasks generated in the current task, not to “descendants” .

# Example

```
#pragma omp parallel
{
    #pragma omp single
    {
        #pragma omp task
        fred();
        #pragma omp task
        daisy();
        #pragma omp taskwait
        #pragma omp task
        billy();
    }
}
```

**fred()** and **daisy()** must complete before **billy()** starts, but this does not include tasks created inside **fred()** and **daisy()**

**All tasks** including those created inside **fred()** and **daisy()** must complete before exiting this barrier

# Example

```
#pragma omp parallel
{
    #pragma omp single nowait
    {
        #pragma omp task
            fred();
        #pragma omp task
            daisy();
        #pragma omp taskwait
        #pragma omp task
            billy();
    }
}
```

The barrier at the end of the single is expensive and not needed since you get the barrier at the end of the parallel region. So use `nowait` to turn it off.

**All tasks** including those created inside `fred()` and `daisy()` must complete before exiting this barrier

# Example: Fibonacci numbers

```
int fib (int n)
{
    int x,y;
    if (n < 2) return n;

    x = fib(n-1);
    y = fib (n-2);
    return (x+y);
}
```

```
int main()
{
    int NW = 5000;
    fib(NW);
}
```

- $F_n = F_{n-1} + F_{n-2}$
- Inefficient  $O(2^n)$  recursive implementation!

# Parallel Fibonacci

```
int fib (int n)
{  int x,y;
   if (n < 2) return n;
```

```
#pragma omp task shared(x)
```

```
  x = fib(n-1);
```

```
#pragma omp task shared(y)
```

```
  y = fib (n-2);
```

```
#pragma omp taskwait
```

```
  return (x+y);
```

```
}
```

```
Int main()
```

```
{  int NW = 5000;
```

```
  #pragma omp parallel
```

```
  {
```

```
    #pragma omp single
```

```
      fib(NW);
```

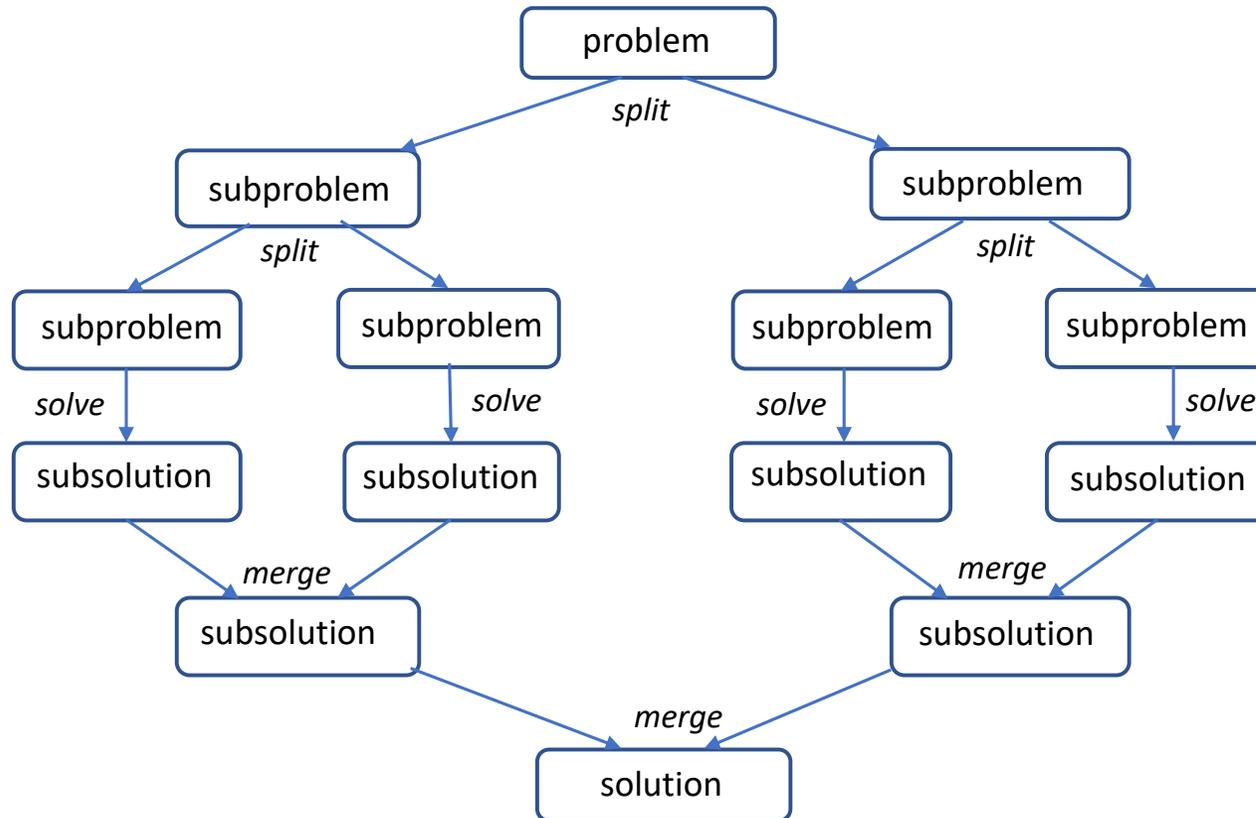
```
  }
```

```
}
```

- Binary tree of tasks
- Traversed using a recursive function
- A task cannot complete until all tasks below it in the tree are complete (enforced with taskwait)
- **x, y** are local, and so by default they are private to current task
  - must be shared on child tasks so they don't create their own firstprivate copies at this level!

# Divide and Conquer

- Split the problem into smaller sub-problems; continue until the sub-problems can be solved directly



- 3 Options for parallelism:
  - Do work as you split into sub-problems
  - Do work only at the leaves
  - Do work as you recombine



# Program: OpenMP tasks

```
include <omp.h>
static long num_steps = 100000000;
#define MIN_BLK 10000000
double pi_comp(int Nstart,int Nfinish,double step)
{ int i,iblk;
  double x, sum = 0.0,sum1, sum2;
  if (Nfinish-Nstart < MIN_BLK){
    for (i=Nstart;i< Nfinish; i++){
      x = (i+0.5)*step;
      sum = sum + 4.0/(1.0+x*x);
    }
  }
  else{
    iblk = Nfinish-Nstart;
    #pragma omp task shared(sum1)
      sum1 = pi_comp(Nstart,      Nfinish-iblk/2,step);
    #pragma omp task shared(sum2)
      sum2 = pi_comp(Nfinish-iblk/2, Nfinish,      step);
    #pragma omp taskwait
      sum = sum1 + sum2;
  }return sum;
}
```

```
int main ()
{
  int i;
  double step, pi, sum;
  step = 1.0/(double) num_steps;
  #pragma omp parallel
  {
    #pragma omp single
      sum =
        pi_comp(0,num_steps,step);
  }
  pi = step * sum;
}
```

# Results\*: Pi with tasks

threads	1 <sup>st</sup> SPMD	SPMD critical	PI Loop	Pi tasks
1	1.86	1.87	1.91	1.87
2	1.03	1.00	1.02	1.00
3	1.08	0.68	0.80	0.76
4	0.97	0.53	0.68	0.52

\*Intel compiler (icpc) with no optimization on Apple OS X 10.7.3 with a dual core (four HW thread) Intel® Core™ i5 processor at 1.7 Ghz and 4 Gbyte DDR3 memory at 1.333 Ghz.

# Using Tasks

- Don't use tasks for things already well supported by OpenMP
  - e.g. standard do/for loops
  - the overhead of using tasks is greater
  
- Don't expect miracles from the runtime
  - best results usually obtained where the user controls the number and granularity of tasks

# Backup content

- Extra exercises
- Irregular parallelism and OpenMP tasks
-  • Worksharing Revisited
- Synchronization Revisited: Additional options for Mutual exclusion

# The Loop Worksharing Constructs

- The loop worksharing construct splits up loop iterations among the threads in a team

```
#pragma omp parallel  
{  
  #pragma omp for  
    for (I=0;I<N;I++){  
      NEAT_STUFF(I);  
    }  
}
```

The variable I is made “private” to each thread by default. You could do this explicitly with a “private(I)” clause

Loop construct name:

- C/C++: for
- Fortran: do

# Loop Worksharing Constructs: The schedule clause

- The schedule clause affects how loop iterations are mapped onto threads
  - **schedule(static [,chunk])**
    - Deal-out blocks of iterations of size “chunk” to each thread.
  - **schedule(dynamic[,chunk])**
    - Each thread grabs “chunk” iterations off a queue until all iterations have been handled.
  - **schedule(guided[,chunk])**
    - Threads dynamically grab blocks of iterations. The size of the block starts large and shrinks down to size “chunk” as the calculation proceeds.
  - **schedule(runtime)**
    - Schedule and chunk size taken from the OMP\_SCHEDULE environment variable (or the runtime library) ... vary schedule without a recompile!
  - **Schedule(auto)**
    - Schedule is left up to the runtime to choose (does not have to be any of the above).

OpenMP 4.5 added modifiers monotonic, nonmonotonic and simd.

# Loop Worksharing Constructs: The schedule clause

Schedule Clause	When To Use
<b>STATIC</b>	Pre-determined and predictable by the programmer
<b>DYNAMIC</b>	Unpredictable, highly variable work per iteration
<b>GUIDED</b>	Special case of dynamic to reduce scheduling overhead
<b>AUTO</b>	When the runtime can “learn” from previous executions of the same loop

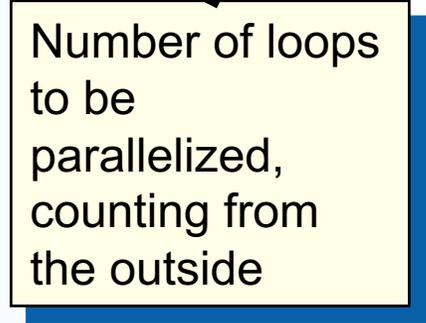
Least work at runtime : scheduling done at compile-time

Most work at runtime : complex scheduling logic used at run-time

# Nested Loops

- For perfectly nested rectangular loops we can parallelize multiple loops in the nest with the collapse clause:

```
#pragma omp parallel for collapse(2)  
for (int i=0; i<N; i++) {  
    for (int j=0; j<M; j++) {  
        .....  
    }  
}
```



Number of loops to be parallelized, counting from the outside

- Will form a single loop of length  $N \times M$  and then parallelize that.
- Useful if  $N$  is  $O(\text{no. of threads})$  so parallelizing the outer loop makes balancing the load difficult.

# Sections Worksharing Construct

- The *Sections* worksharing construct gives a different structured block to each thread.

```
#pragma omp parallel
{
    #pragma omp sections
    {
        #pragma omp section
            x_calculation();
        #pragma omp section
            y_calculation();
        #pragma omp section
            z_calculation();
    }
}
```

By default, there is a barrier at the end of the “omp sections”. Use the “nowait” clause to turn off the barrier.

# Array Sections with Reduce

```
#include <stdio.h>
#define N 100
void init(int n, float (*b)[N]);
int main(){
int i,j; float a[N], b[N][N]; init(N,b);
for(i=0; i<N; i++) a[i]=0.0e0;
```

Works the same as any other reduce ... a private array is formed for each thread, element wise combination across threads and then with original array at the end

```
#pragma omp parallel for reduction(+:a[0:N]) private(j)
for(i=0; i<N; i++){
    for(j=0; j<N; j++){
        a[j] += b[i][j];
    }
}
printf(" a[0] a[N-1]: %f %f\n", a[0], a[N-1]);
return 0;
```

# Exercise

- Go back to your parallel mandel.c program.
- Using what we've learned in this block of slides can you improve the runtime?

# Optimizing mandel.c

```
wtime = omp_get_wtime();
#pragma omp parallel for collapse(2) schedule(runtime) firstprivate(eps) private(j,c)
for (i=0; i<NPOINTS; i++) {
  for (j=0; j<NPOINTS; j++) {
    c.r = -2.0+2.5*(double)(i)/(double)(NPOINTS)+eps;
    c.i = 1.125*(double)(j)/(double)(NPOINTS)+eps;
    testpoint(c);
  }
}
wtime = omp_get_wtime() - wtime;
```

```
$ export OMP_SCHEDULE="dynamic,100"
$ ./mandel_par
```

default schedule	0.48 secs
schedule(dynamic,100)	0.39 secs
collapse(2) schedule(dynamic,100)	0.34 secs

Four threads on a dual core Apple laptop (Macbook air ... 2.2 Ghz Intel Core i7 with 8 GB memory) and the gcc version 9.1. Times are the minimum time from three runs

# Backup content

- Extra exercises
- Irregular parallelism and OpenMP tasks
- Worksharing Revisited
-  • Synchronization Revisited: Additional options for Mutual exclusion

# Synchronization

Synchronization is used to impose order constraints between threads and to protect access to shared data

- High level synchronization included in the common core:

- critical
  - barrier

Covered earlier

- Other, more advanced, synchronization operations:

- atomic
  - ordered
  - flush

Covered in this section

- locks (both simple and nested)

# Synchronization: Atomic

- Atomic provides mutual exclusion but only applies to the update of a memory location (the update of X in the following example)

```
#pragma omp parallel
{
    double B;
    B = DOIT();

    #pragma omp atomic
        X += big_ugly(B);
}
```

# Synchronization: Atomic

- Atomic provides mutual exclusion but only applies to the update of a memory location (the update of X in the following example)

```
#pragma omp parallel
{
    double B, tmp;
    B = DOIT();
    tmp = big_ugly(B);
    #pragma omp atomic
        X += tmp;
}
```

Atomic only protects the read/update of X











# Histogram Program: Critical section

- A critical section means that only one thread at a time can update a histogram bin ... but this effectively serializes the loops and adds huge overhead as the runtime manages all the threads waiting for their turn for the update.

```
#pragma omp parallel for
for(i=0;i<NVALS;i++){
    ival = (int) x[i];
    #pragma omp critical
        hist[ival]++;
}
```

Easy to write and correct, but terrible performance

