



# An Introduction to Parallel Programming with OpenMP

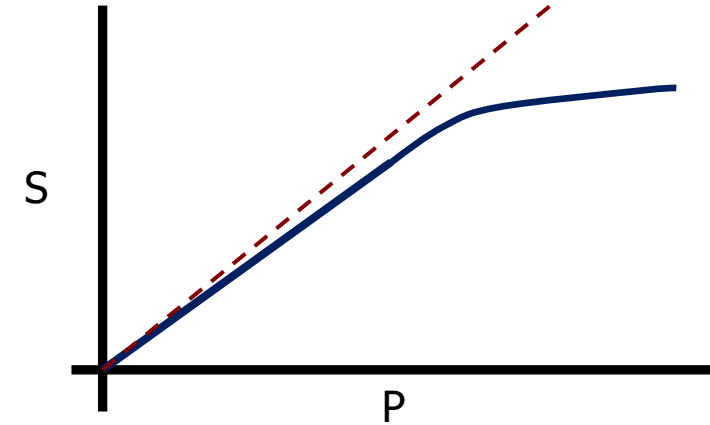
**Tim Mattson**  
Intel Corp.

All materials are on github. To download them:  
`git clone https://github.com/tgmattso/ParProgCourse.git`

# Introduction

I'm just a simple kayak instructor

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

# Disclaimer & Optimization Notice

INFORMATION IN THIS DOCUMENT IS PROVIDED “AS IS”. NO LICENSE, EXPRESS OR IMPLIED, BY ESTOPPEL OR OTHERWISE, TO ANY INTELLECTUAL PROPERTY RIGHTS IS GRANTED BY THIS DOCUMENT. INTEL ASSUMES NO LIABILITY WHATSOEVER AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY, RELATING TO THIS INFORMATION INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT OR OTHER INTELLECTUAL PROPERTY RIGHT.

Performance tests and ratings are measured using specific computer systems and/or components and reflect the approximate performance of Intel products as measured by those tests. Any difference in system hardware or software design or configuration may affect actual performance. Buyers should consult other sources of information to evaluate the performance of systems or components they are considering purchasing. For more information on performance tests and on the performance of Intel products, reference [www.intel.com/software/products](http://www.intel.com/software/products).

All rights reserved. Intel, the Intel logo, Xeon, Xeon Phi, VTune, and Cilk are trademarks of Intel Corporation in the U.S. and other countries.

\*Other names and brands may be claimed as the property of others.

## Optimization Notice

Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice.

Notice revision #20110804

# Preliminaries: Part 1

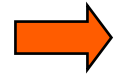
- Disclosures
  - The views expressed in this tutorial are those of the people delivering the tutorial.
    - We are not speaking for our employers.
    - We are not speaking for the OpenMP ARB
- We take these tutorials VERY seriously:
  - Help us improve ... tell us how you would make this tutorial better.

# Preliminaries: Part 2

- Our plan for the day .. Active learning!
  - We will mix short lectures with short exercises.
  - You will use your laptop to connect to a multiprocessor server.
- 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.

# Grab content from github

- Clone the parallel programming course git hub repository
  - `git clone https://github.com/tgmattso/ParProgCourse.git`
- We will use the exercises in the directory: Exercises/OpenMP
- These lectures assume familiarity with C. A simple C program the text of which explains all the C you need to know for this tutorial is included with the exercises
  - Exercises/learningC.c
- The slides for this lecture are in the file: Openmp\_Intro\_hands\_on.pdf
- There are also three other lectures we won't be covering ... for you to study on your own:
  - Par\_Comp\_Intro\_short.pdf,
  - MPI\_intro\_hands\_on.pdf,
  - Other\_par\_prog\_envs.pdf
- If we need to use slurm to submit jobs, this repository has information on how to use slurm on Adroit
  - [https://github.com/PrincetonUniversity/hpc\\_beginning\\_workshop/tree/2021fall/RC\\_example\\_jobs/cxx/multithreaded](https://github.com/PrincetonUniversity/hpc_beginning_workshop/tree/2021fall/RC_example_jobs/cxx/multithreaded)



- Introduction to OpenMP
- Creating Threads
- Synchronization
- Parallel Loops
- Data Environment
- Memory Model
- Irregular Parallelism and Tasks
- Recap
- Beyond the Common Core:
  - Worksharing Revisited
  - Synchronization Revisited: Options for Mutual exclusion
  - Thread Affinity and Data Locality
  - Thread Private Data
  - Memory Models and Point-to-Point Synchronization
  - Programming your GPU with OpenMP

# 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
- Greatly simplifies writing multi-threaded (MT) programs in Fortran, C and C++
- Also 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)

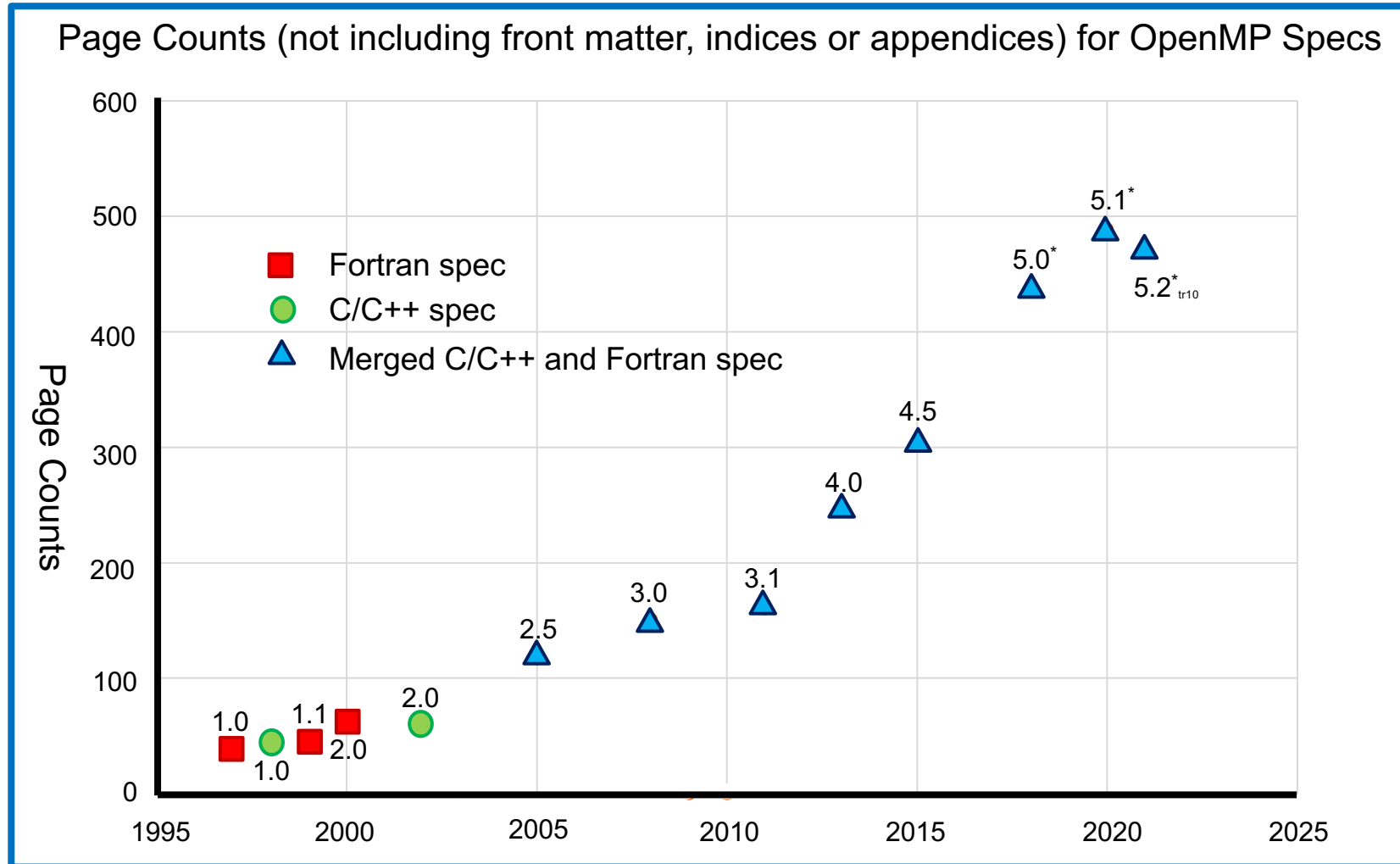
#pragma omp atomic seq\_cst

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

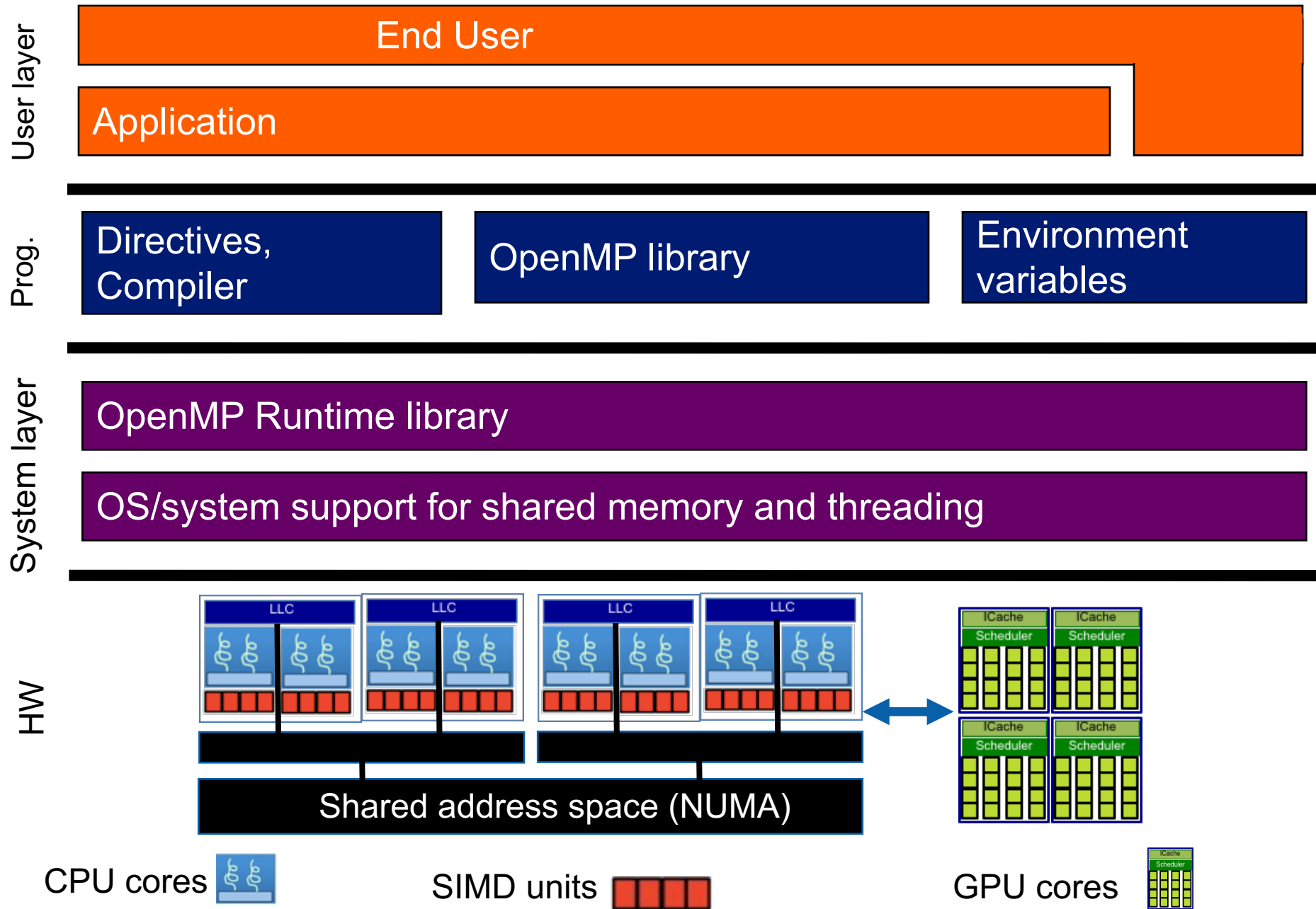


The full spec is overwhelming. We focus on the Common Core: the 21 items most people restrict themselves to

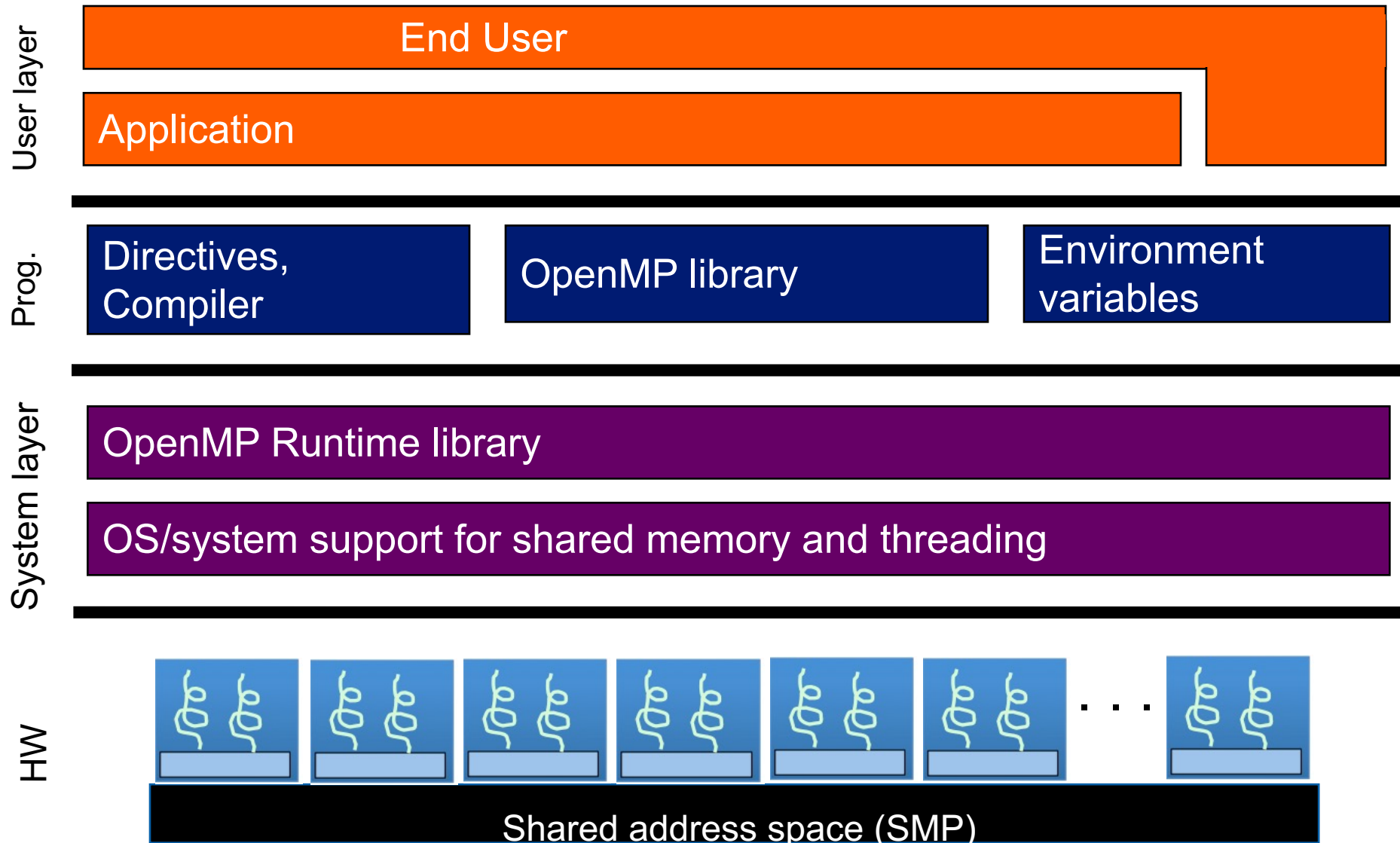
# The OpenMP Common Core: Most OpenMP programs only use these 21 items

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 the constructs in OpenMP are compiler directives.

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

- 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”.

`git clone https://github.com/tgmattso/OpenMPCommonCore.git`

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

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

}
```

- To download the slides:

<https://github.com/tgmattso/OmpCommonCore.git>

# Exercise, Part B: Hello World

## Verify that your OpenMP environment works

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

```
git clone https://github.com/tgmattso/OpenMP_Common_Core.git
```

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

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

    }
}
```

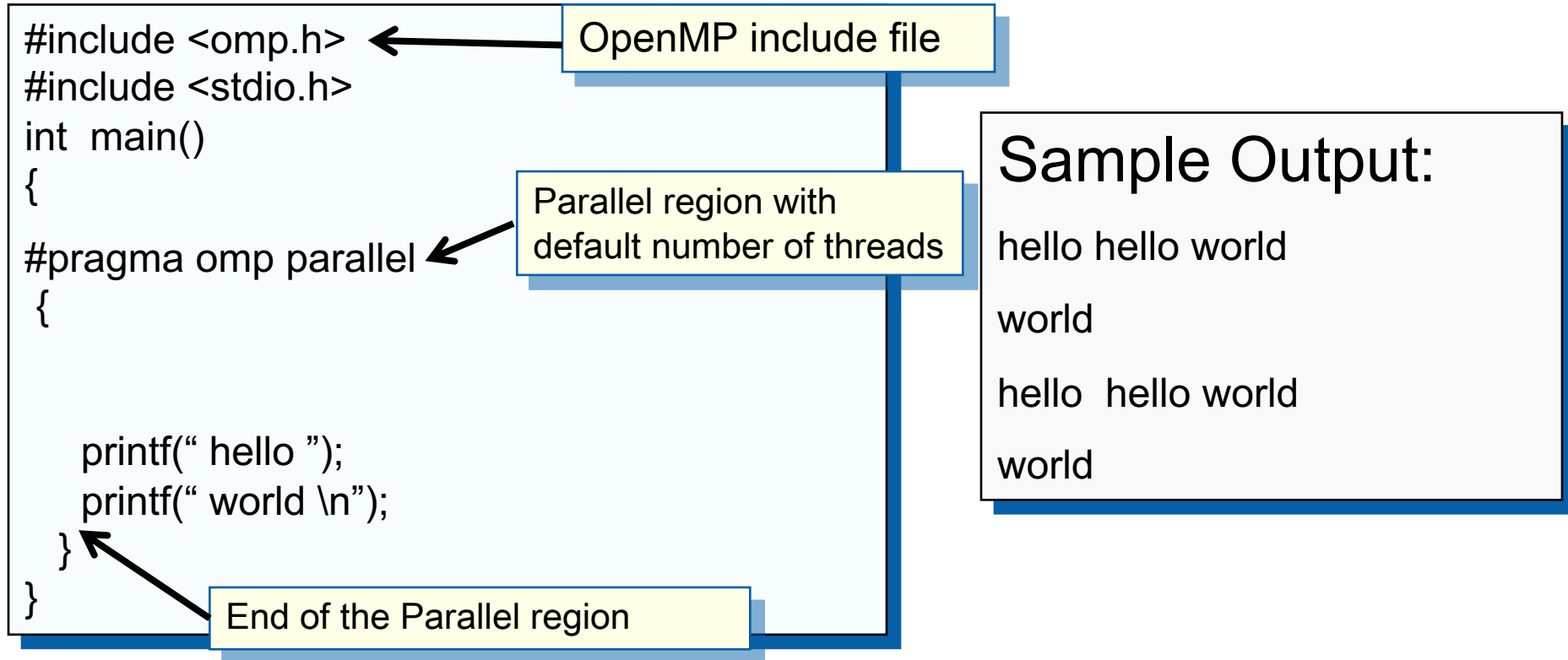
### Switches for compiling and linking

<b>gcc -fopenmp</b>	<b>Gnu (Linux, OSX)</b>
<b>cc -qopenmp</b>	<b>Intel (Linux@NERSC)</b>
<b>icl /Qopenmp</b>	<b>Intel (windows)</b>
<b>icc -fopenmp</b>	<b>Intel (Linux, OSX)</b>


# Solution

## A Multi-Threaded “Hello World” Program

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



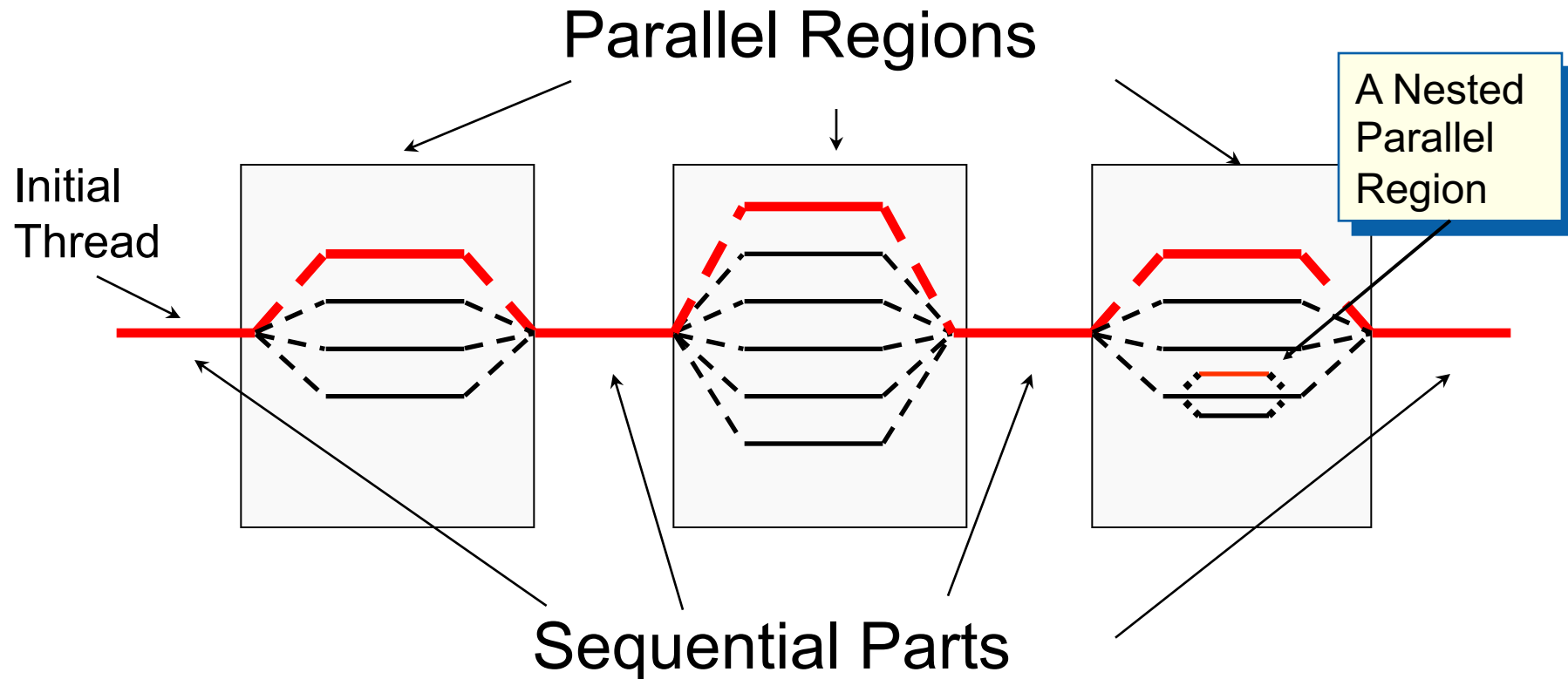
The statements are interleaved based on how the operating schedules the threads

- Introduction to OpenMP
-  • Creating Threads
- Synchronization
- Parallel Loops
- Data Environment
- Memory Model
- Irregular Parallelism and Tasks
- Recap
- Beyond the Common Core:
  - Worksharing Revisited
  - Synchronization Revisited: Options for Mutual exclusion
  - Thread Affinity and Data Locality
  - Thread Private Data
  - Memory Models and Point-to-Point Synchronization
  - Programming your GPU with OpenMP

# OpenMP Execution model:

## 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 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();  
    pooh(ID,A);  
}
```

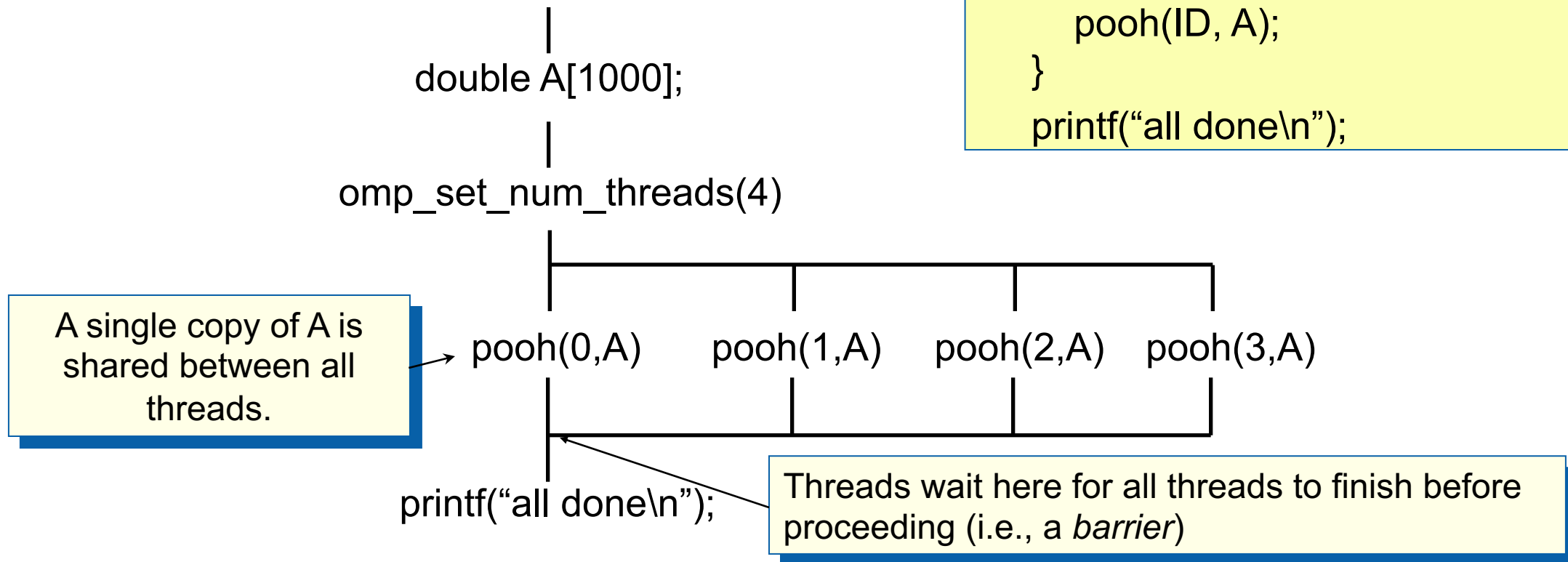
Runtime function to request a certain number of threads

Runtime function returning a thread ID

- 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");
```

# 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);
}
```

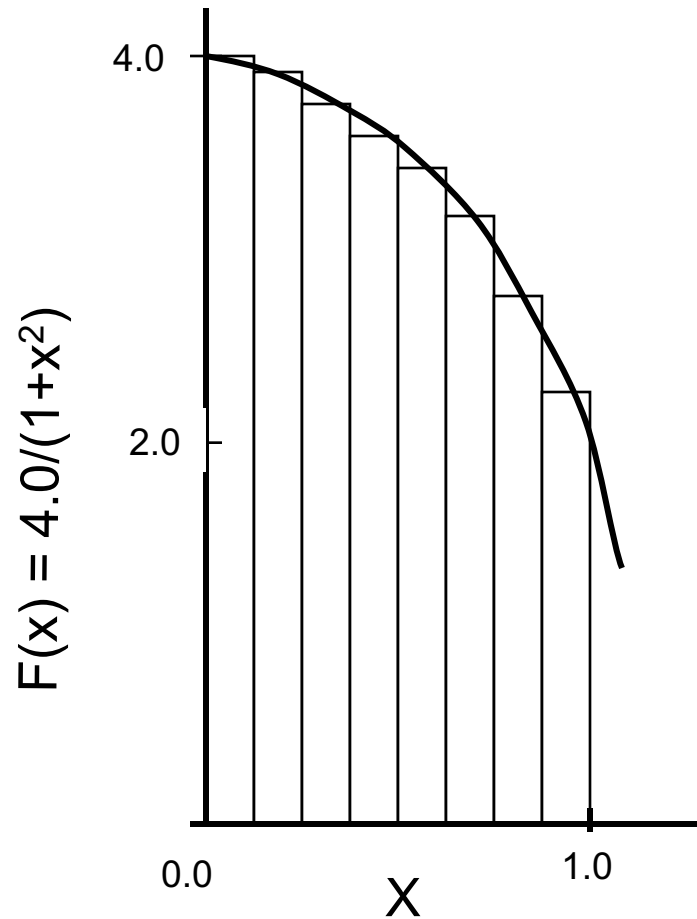
Runtime function to request a certain number of threads

Runtime function to return actual number of threads in the team

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

# An Interesting Problem to Play With

## Numerical Integration



Mathematically, we know that:

$$\int_0^1 \frac{4.0}{(1+x^2)} dx = \pi$$

We can approximate the integral as a sum of rectangles:

$$\sum_{i=0}^N F(x_i) \Delta x = \Delta x \sum_{i=0}^N F(x_i) \approx \pi$$

Where each rectangle has width  $\Delta x$  and height  $F(x_i)$  at the middle of interval  $i$ .

# Serial PI Program

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

    step = 1.0/(double) num_steps;

    for (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 ()
{
    int i;    double x, pi, sum = 0.0;

    step = 1.0/(double) num_steps;
double tdata = omp_get_wtime();
    for (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

# Exercise: the Parallel Pi Program

- Create a parallel version of the pi program using a parallel construct:

`#pragma omp parallel`

- Pay close attention to shared versus private variables.
- In addition to a parallel construct, you will need the runtime library routines

– `int omp_get_num_threads();`

Number of threads in the team

– `int omp_get_thread_num();`

Thread ID or rank

– `double omp_get_wtime();`

– `omp_set_num_threads();`

Time in seconds since a fixed point in the past

Request a number of threads in the team

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

# Example: A simple SPMD pi program

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

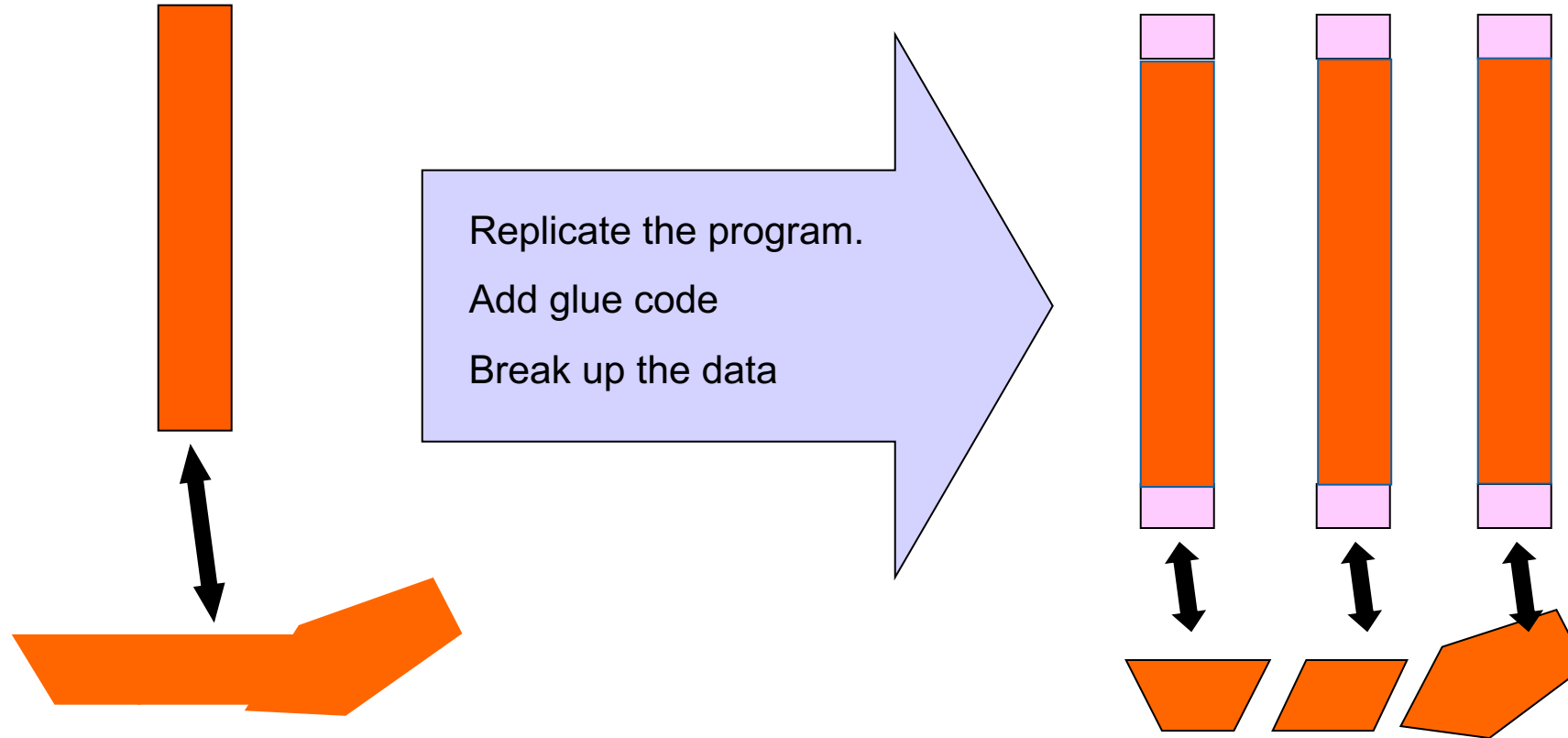
Promote scalar to an array dimensioned by number of threads to avoid race condition.

Only one thread should copy the number of threads to the global value to make sure multiple threads writing to the same address don't conflict.

This is a common trick in SPMD programs to create a cyclic distribution of loop iterations

# SPMD: Single Program Multiple Data

- Run the same program on  $P$  processing elements where  $P$  can be arbitrarily large.



- Use the rank ... an ID ranging from 0 to  $(P-1)$  ... to select between a set of tasks and to manage any shared data structures.

MPI programs almost always use this pattern ... it is probably the most commonly used pattern in the history of parallel programming.

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

# Consider performance of parallel programs

**Compute N independent tasks on one processor**

Load Data

Compute  $T_1$

...

Compute  $T_N$

Consume Results

$$\text{Time}_{\text{seq}}(1) = T_{\text{load}} + N * T_{\text{task}} + T_{\text{consume}}$$

**Compute N independent tasks with P processors**

Load Data

Compute  $T_1$

...

Compute  $T_N$

Consume Results

$$\text{Time}_{\text{par}}(P) = T_{\text{load}} + (N/P) * T_{\text{task}} + T_{\text{consume}}$$

Ideally Cut  
runtime by  $\sim 1/P$

*(Note: Parallelism  
only speeds-up the  
concurrent part)*

# Talking about performance

- Speedup: the increased performance from running on  $P$  processors.
- Perfect Linear Speedup: happens when no parallel overhead and algorithm is 100% parallel.
- Super-linear Speedup: typically due to cache effects ... i.e. as  $P$  grows, aggregate cache size grows so more of the problem fits in cache

$$S(P) = \frac{Time_{seq}(1)}{Time_{par}(P)}$$

$$S(P) = P$$

$$S(P) > P$$

# Amdahl's Law

- What is the maximum speedup you can expect from a parallel program?
- Approximate the runtime as a part that can be sped up with additional processors and a part that is fundamentally serial.

$$Time_{par}(P) = (serial\_fraction + \frac{parallel\_fraction}{P}) * Time_{seq}$$

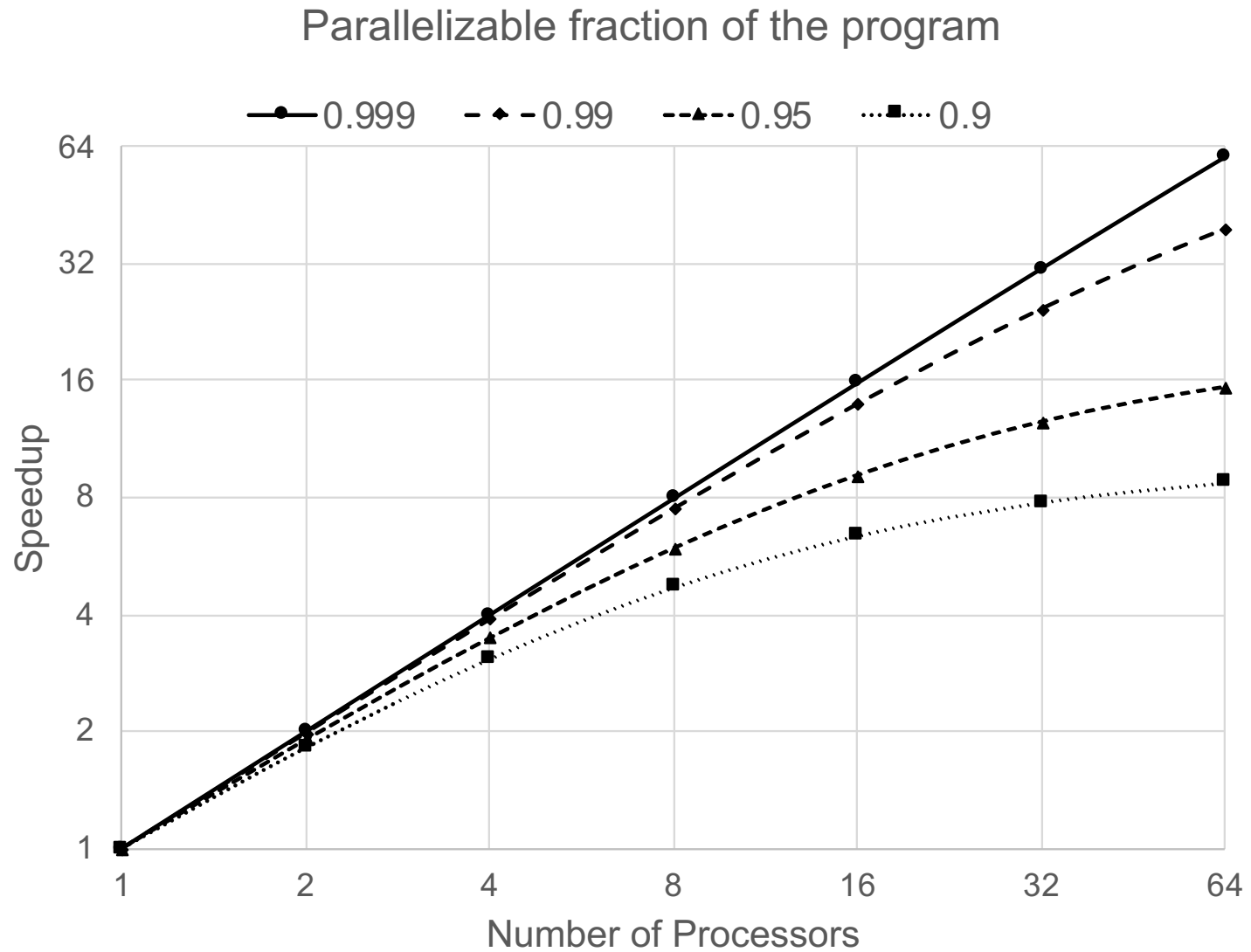
- If serial\_fraction is  $\alpha$  and parallel\_fraction is  $(1 - \alpha)$  then the speedup is:

$$S(P) = \frac{Time_{seq}}{Time_{par}(P)} = \frac{Time_{seq}}{(\alpha + \frac{1 - \alpha}{P}) * Time_{seq}} = \frac{1}{\alpha + \frac{1 - \alpha}{P}}$$

- If you had an unlimited number of processors:  $P \rightarrow \infty$

- The maximum possible speedup is:  $S = \frac{1}{\alpha}$  ← Amdahl's Law

# Amdahl's Law



So now you should understand my silly introduction slide.

## Introduction

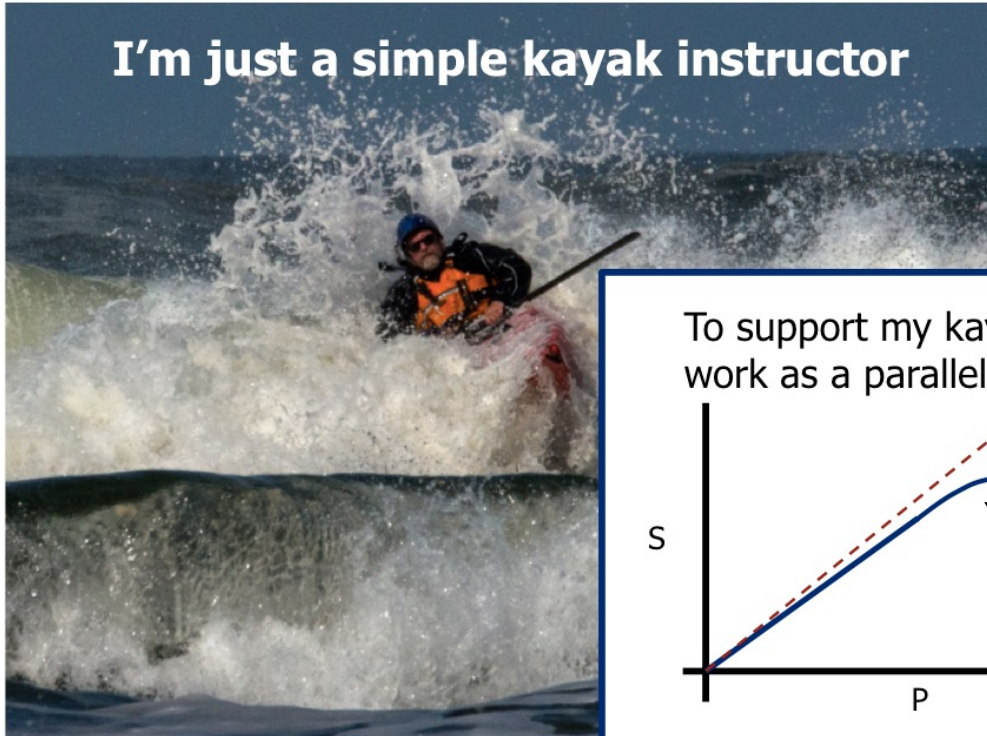
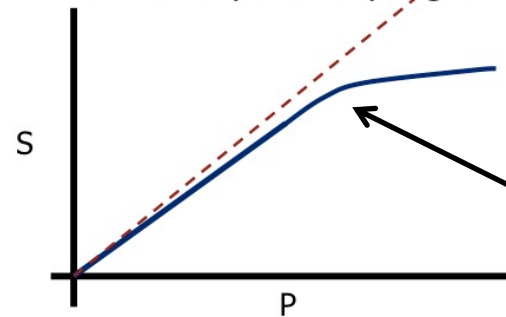


Photo © by Greg Clopton, 2014

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.

**Now that you understand how to  
think about parallel performance,  
lets get back to OpenMP**

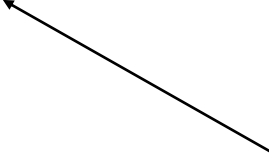
# Internal control variables and how to control the number of threads in a team

- We've used the following construct to control the number of threads. (e.g. to request 12 threads):
  - `omp_set_num_threads(12)`
- What does `omp_set_num_threads()` actually do?
  - It **resets** an **“internal control variable”** the system queries to define the default number of threads to request on subsequent parallel constructs.
- Is there an easier way to change this internal control variable ... perhaps one that doesn't require re-compilation? Yes.
  - When an OpenMP program starts up, it queries an environment variable `OMP_NUM_THREADS` and sets the appropriate internal control variable to the value of **`OMP_NUM_THREADS`**
  - For example, to set the initial, default number of threads to request in OpenMP from my apple laptop
    - > **`export OMP_NUM_THREADS=12`**

# Exercise

- Go back to your parallel pi program and explore how well it scales with the number of threads.
- Can you explain your performance with Amdahl's law? If not what else might be going on?

- `int omp_get_num_threads();`
- `int omp_get_thread_num();`
- `double omp_get_wtime();`
- `omp_set_num_threads();`
- `export OMP_NUM_THREADS = N`



An environment variable  
to set the default number  
of threads to request to N

# Results\*

- Original Serial pi program with 100000000 steps ran in 1.83 seconds.

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

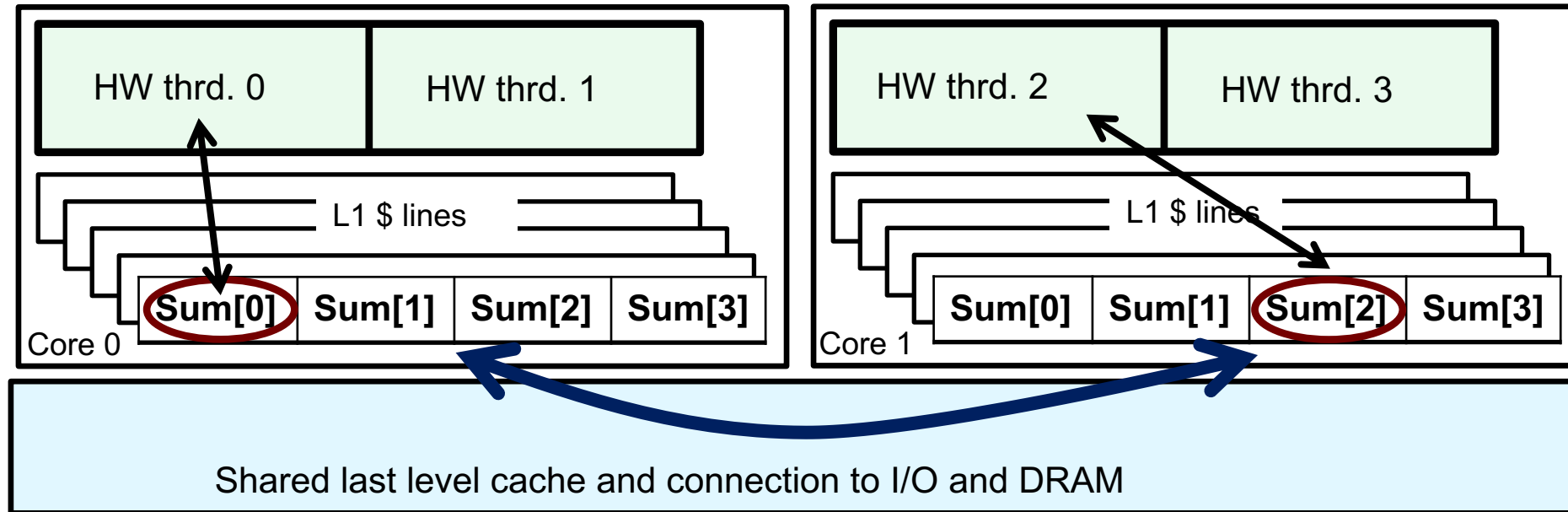
threads	1 <sup>st</sup> SPMD*
1	1.86
2	1.03
3	1.08
4	0.97

Intel compiler (icpc) with default optimization level (O2) 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.

\*SPMD: Single Program Multiple Data

# Why Such Poor Scaling? False Sharing

- If independent data elements happen to sit on the same cache line, each update will cause the cache lines to “slosh back and forth” between threads ... This is called **“false sharing”**.




- If you promote scalars to an array to support creation of an SPMD program, the array elements are contiguous in memory and hence share cache lines ... Results in poor scalability.
- Solution: Pad arrays so elements you use are on distinct cache lines.

# Example: Eliminate false sharing by padding the sum array

```
#include <omp.h>
static long num_steps = 100000;      double step;
#define NUM_THREADS 2
#define PAD 8      // assume 64 byte L1 cache line size
void main ()
{  int i, nthreads; double pi, sum[NUM_THREADS][PAD] ;
  step = 1.0/(double) num_steps;
  omp_set_num_threads(NUM_THREADS);
  #pragma omp parallel
  {
    int i, id, nthrds;
    double x;
    id = omp_get_thread_num();
    nthrds = omp_get_num_threads();
    if (id == 0) nthreads = nthrds;
    for (i=id, sum[id]=0.0; i< num_steps; i=i+nthrds) {
      x = (i+0.5)*step;
      sum[id][0] += 4.0/(1.0+x*x);
    }
  }
  for(i=0, pi=0.0; i<nthreads; i++) pi += sum[i][0] * step;
}
```

Pad the array so each  
sum value is in a  
different cache line



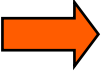
# Results\*: PI Program, Padded Accumulator

- Original Serial pi program with 100000000 steps ran in 1.83 seconds.

```
#include <omp.h>
static long num_steps = 100000;      double step;
#define NUM_THREADS 2
#define PAD 8      // assume 64 byte L1 cache line size
void main ()
{  int i, nthreads; double pi, sum[NUM_THREADS][PAD] ;
  step = 1.0/(double) num_steps;
  omp_set_num_threads(NUM_THREADS);
  #pragma omp parallel
  {
    int i, id,nthrds;
    double x;
    id = omp_get_thread_num();
    nthrds = omp_get_num_threads();
    if (id == 0)  nthreads = nthrds;
    for (i=id, sum[id]=0.0;i< num_steps; i=i+nthrds) {
      x = (i+0.5)*step;
      sum[id][0] += 4.0/(1.0+x*x);
    }
  }
  for(i=0, pi=0.0;i<nthreads;i++)pi += sum[i][0] * step;
}
```

threads	1 <sup>st</sup> SPMD	1 <sup>st</sup> SPMD padded
1	1.86	1.86
2	1.03	1.01
3	1.08	0.69
4	0.97	0.53

\*Intel compiler (icpc) with default optimization level (O2) 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.

- Introduction to OpenMP
- Creating Threads
-  • Synchronization
- Parallel Loops
- Data Environment
- Memory Model
- Irregular Parallelism and Tasks
- Recap
- Beyond the Common Core:
  - Worksharing Revisited
  - Synchronization Revisited: Options for Mutual exclusion
  - Thread Affinity and Data Locality
  - Thread Private Data
  - Memory Models and Point-to-Point Synchronization
  - Programming your GPU with OpenMP

# 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)

# Synchronization: critical

- Mutual exclusion: Only one thread at a time can enter a **critical** region.

Threads wait their turn  
– only one thread at a  
time calls consume()

```
float res;  
  
#pragma omp parallel  
{  float B;  int i, id, nthrds;  
    id = omp_get_thread_num();  
    nthrds = omp_get_num_threads();  
    B = big_SPMD_job(id, nthrds);  
    #pragma omp critical  
        res += consume (B);  
}
```

# Synchronization: barrier

- Barrier: a point in a program all threads must reach before any threads are allowed to proceed.
- It is a “stand alone” pragma meaning it is not associated with user code ... it is an executable statement.

```
double Arr[8], Brr[8];      int numthrds;

omp_set_num_threads(8)


#pragma omp parallel
{  int id, nthrds;

    id = omp_get_thread_num();
    nthrds = omp_get_num_threads();
    if (id==0) numthrds = nthrds;

    Arr[id] = big_ugly_calc(id, nthrds);

    #pragma omp barrier
    Brr[id] = really_big_and_ugly(id, nthrds, Arr);
}
```

Threads wait until all  
threads hit the barrier.  
Then they can go on.



# Exercise

- In your first Pi program, you probably used an array to create space for each thread to store its partial sum.
- If array elements happen to share a cache line, this leads to false sharing.
  - Non-shared data in the same cache line so each update invalidates the cache line ... in essence “sloshing independent data” back and forth between threads.
- Modify your “pi program” to avoid false sharing due to the partial sum array.

```
int omp_get_num_threads();
int omp_get_thread_num();
double omp_get_wtime();
omp_set_num_threads();
#pragma parallel
#pragma critical
```

# PI Program with False Sharing

- Original Serial pi program with 100000000 steps ran in 1.83 seconds.

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

Recall that promoting sum to an array made the coding easy, but led to false sharing and poor performance.

threads	1 <sup>st</sup> SPMD
1	1.86
2	1.03
3	1.08
4	0.97

\*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.

## 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;
      sum += 4.0/(1.0+x*x);
    }
    #pragma omp critical
    pi += sum * step;
  }
}
```

Create a scalar local to each thread to accumulate partial sums.

No array, so no false sharing.

Sum goes “out of scope” beyond the parallel region ... so you must sum it in here. Must protect summation into pi in a critical region so updates don’t conflict

# Results\*: pi program critical section

- Original Serial pi program with 100000000 steps ran in 1.83 seconds.

```
#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;
      sum += 4.0/(1.0+x*x);
    }
    #pragma omp critical
      pi += sum * step;
  }
}
```

threads	1 <sup>st</sup> SPMD	1 <sup>st</sup> SPMD padded	SPMD critical
1	1.86	1.86	1.87
2	1.03	1.01	1.00
3	1.08	0.69	0.68
4	0.97	0.53	0.53

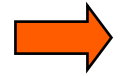
\*Intel compiler (icpc) with no optimization on Apple OS X 10.7.3 with a dual core (four HW thread) Intel® Core™ i5 processor at 1.7 Ghz and 4 Gbyte DDR3 memory at 1.333 Ghz.

## 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?

# Outline



- Introduction to OpenMP
- Creating Threads
- Synchronization
- Parallel Loops
- Data Environment
- Memory Model
- Irregular Parallelism and Tasks
- Recap
- Beyond the Common Core:
  - Worksharing Revisited
  - Synchronization Revisited: Options for Mutual exclusion
  - Thread Affinity and Data Locality
  - Thread Private Data
  - Memory Models and Point-to-Point Synchronization
  - Programming your GPU with OpenMP

# The Loop Worksharing Construct

- 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);
  }
}
```

Loop construct name:

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

The loop control index I is made  
“private” to each thread by default.

Threads wait here until all  
threads are finished with the  
parallel loop before any proceed  
past the end of the loop

# Loop Worksharing Construct

## A motivating example

Sequential code

```
for(i=0;i<N;i++) { a[i] = a[i] + b[i];}
```

OpenMP parallel region  
(SPMD Pattern)

```
#pragma omp parallel  
{  
    int id, i, Nthrds, istart, iend;  
    id = omp_get_thread_num();  
    Nthrds = omp_get_num_threads();  
    istart = id * N / Nthrds;  
    iend = (id+1) * (N / Nthrds)-1;  
    if (id == Nthrds-1)iend = N;  
    for(i=istart;i<iend;i++) { a[i] = a[i] + b[i];}  
}
```

OpenMP parallel region and  
a worksharing for construct

```
#pragma omp parallel  
#pragma omp for  
for(i=0;i<N;i++) { a[i] = a[i] + b[i];}
```

# 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.
- Example:
  - #pragma omp for schedule(dynamic, 10)

Schedule Clause	When To Use
<b>STATIC</b>	<b>Pre-determined and predictable by the programmer</b>
<b>DYNAMIC</b>	<b>Unpredictable, highly variable work per iteration</b>

Least work at runtime :  
scheduling done at  
compile-time

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

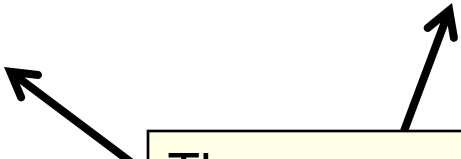
# Combined Parallel/Worksharing Construct

- OpenMP shortcut: Put the “parallel” and the worksharing directive on the same line

```
double res[MAX]; int i;  
#pragma omp parallel  
{  
    #pragma omp for  
    for (i=0;i< MAX; i++) {  
        res[i] = huge();  
    }  
}
```

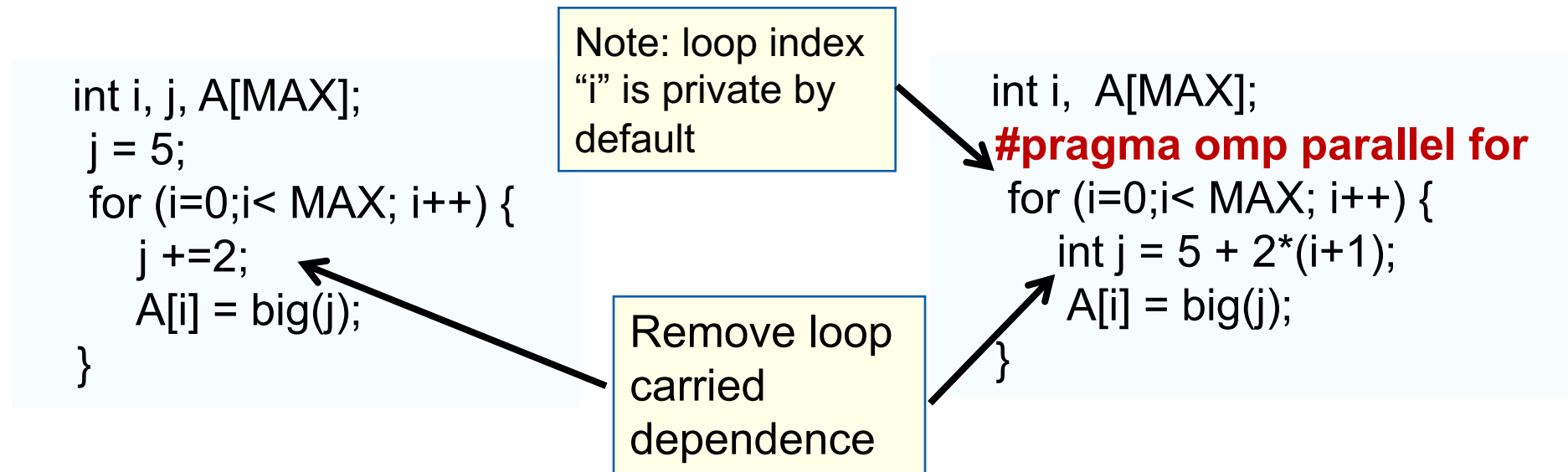
```
double res[MAX]; int i;  
#pragma omp parallel for  
    for (i=0;i< MAX; i++) {  
        res[i] = huge();  
    }
```

These are equivalent



# Working with loops

- Basic approach
  - Find compute intensive loops
  - Make the loop iterations independent ... So they can safely execute in any order without loop-carried dependencies
  - Place the appropriate OpenMP directive and test



# Reduction

- How do we handle this case?

```
double ave=0.0, A[MAX];  
int i;  
for (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

- 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 ave=0.0, A[MAX];  int i;  
#pragma omp parallel for reduction (+:ave)  
  for (i=0;i< MAX; i++) {  
    ave + = A[i];  
  }  
ave = ave/MAX;
```

# OpenMP: Reduction operands/initial-values

- Many different associative operands can be used with reduction:
- Initial values are the ones that make sense mathematically.

Operator	Initial value
<b>+</b>	<b>0</b>
<b>*</b>	<b>1</b>
<b>-</b>	<b>0</b>
<b>min</b>	Largest pos. number
<b>max</b>	Most neg. number

C/C++ only	
Operator	Initial value
<b>&amp;</b>	<b>~0</b>
<b> </b>	<b>0</b>
<b>^</b>	<b>0</b>
<b>&amp;&amp;</b>	<b>1</b>
<b>  </b>	<b>0</b>

Fortran Only	
Operator	Initial value
<b>.AND.</b>	<b>.true.</b>
<b>.OR.</b>	<b>.false.</b>
<b>.NEQV.</b>	<b>.false.</b>
<b>.IEOR.</b>	<b>0</b>
<b>.IOR.</b>	<b>0</b>
<b>.IAND.</b>	<b>All bits on</b>
<b>.EQV.</b>	<b>.true.</b>

OpenMP includes user defined reductions and array-sections as reduction variables (we just don't cover those topics here)

# Exercise: PI with loops

- Go back to the serial pi program and parallelize it with a loop construct
- Your goal is to minimize the number of changes made to the serial program.

```
#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();
```

# Example: PI with a loop and a reduction

```
#include <omp.h>
```

```
static long num_steps = 100000;    double step;
```

```
void main ()
```

```
{  int i;        double x, pi, sum = 0.0;
```

```
    step = 1.0/(double) num_steps;
```

```
    #pragma omp parallel
```

```
    {
```

```
        double x;
```

```
        #pragma omp for reduction(+:sum)
```

```
            for (i=0;i< num_steps; i++){
```

```
                x = (i+0.5)*step;
```

```
                sum = sum + 4.0/(1.0+x*x),
```

```
            }
```

```
    }
```

```
    pi = step * sum;
```

```
}
```

Create a team of threads ...  
without a parallel construct, you'll  
never have more than one thread

Create a scalar local to each thread to hold  
value of x at the center of each interval

Break up loop iterations  
and assign them to  
threads ... setting up a  
reduction into sum.  
Note ... the loop index is  
local to a thread by default.

# Example: PI with a loop and a reduction

```
#include <omp.h>
static long num_steps = 100000;      double step;
void main ()
{
    double pi, sum = 0.0;
    step = 1.0/(double) num_steps;

    #pragma omp parallel for reduction(+:sum)
    for (int i=0;i< num_steps; i++){
        double x = (i+0.5)*step;
        sum = sum + 4.0/(1.0+x*x);
    }
    pi = step * sum;
}
```

Using modern C style, we put declarations close to where they are used ... which lets me use the parallel for construct.

# Results\*: PI with a loop and a reduction

- Original Serial pi program with 100000000 steps ran in 1.83 seconds.

Example: Pi with a		threads	1 <sup>st</sup> SPMD	1 <sup>st</sup> SPMD padded	SPMD critical	PI Loop
<pre>#include &lt;omp.h&gt; static long num_steps = 100000000; void main () {     int i;     double x, pi, sum;     step = 1.0/(double) num_steps;     #pragma omp parallel     {         double x;         #pragma omp for reduction(+:sum)         for (i=0; i&lt; num_steps; i++){             x = (i+0.5)*step;             sum = sum + 4.0/(1.0+x*x);         }         pi = step * sum;     } }</pre>		1	1.86	1.86	1.87	1.91
		2	1.03	1.01	1.00	1.02
		3	1.08	0.69	0.68	0.80
		4	0.97	0.53	0.53	0.68

\*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.

# The nowait clause

- Barriers are really expensive. You need to understand when they are implied and how to skip them when it's safe to do so.

```
double A[big], B[big], C[big];
```

```
#pragma omp parallel  
{
```

```
    int id=omp_get_thread_num();  
    A[id] = big_calc1(id);
```

```
#pragma omp barrier
```

```
#pragma omp for
```


```
    for(i=0;i<N;i++){C[i]=big_calc3(i,A);} 
```

```
#pragma omp for nowait
```

```
    for(i=0;i<N;i++){ B[i]=big_calc2(C, i); }  
    A[id] = big_calc4(id);
```

```
}
```

implicit barrier at the end of a for  
worksharing construct

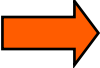


implicit barrier at the end  
of a parallel region



no implicit barrier  
due to nowait



- Introduction to OpenMP
- Creating Threads
- Synchronization
- Parallel Loops
-  • Data Environment
- Memory Model
- Irregular Parallelism and Tasks
- Recap
- Beyond the Common Core:
  - Worksharing Revisited
  - Synchronization Revisited: Options for Mutual exclusion
  - Thread Affinity and Data Locality
  - Thread Private Data
  - Memory Models and Point-to-Point Synchronization
  - Programming your GPU with OpenMP

# Data Environment:

## Default storage attributes

- Shared memory programming model:
  - Most variables are shared by default
- Global variables are SHARED among threads
  - Fortran: COMMON blocks, SAVE variables, MODULE variables
  - C: File scope variables, static
  - Both: dynamically allocated memory (ALLOCATE, malloc, new)
- But not everything is shared...
  - Stack variables in subprograms(Fortran) or functions(C) called from parallel regions are PRIVATE
  - Automatic variables within a statement block are PRIVATE.

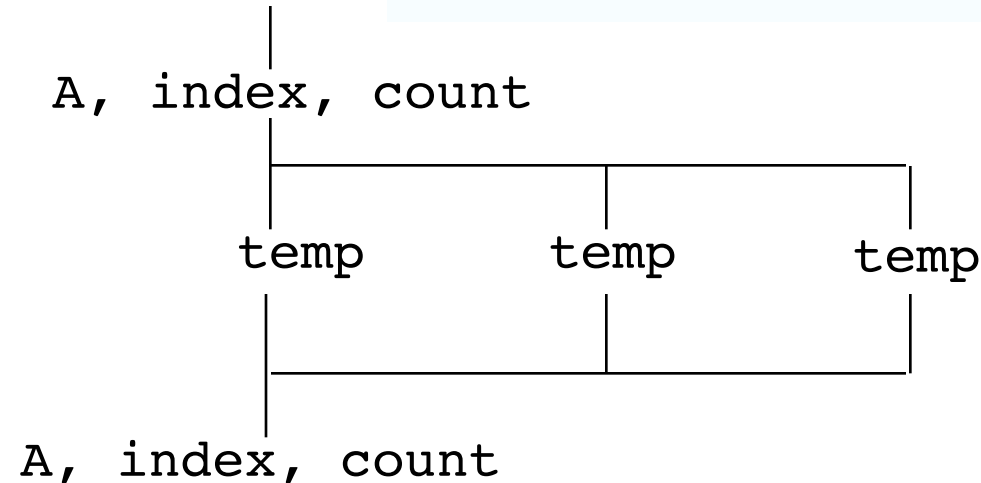
# Data Sharing: Examples

```
double A[10];
int main() {
    int index[10];
    #pragma omp parallel
        work(index);
    printf("%d\n", index[0]);
}
```

A, index and count are shared by all threads.

temp is local to each thread

```
extern double A[10];
void work(int *index) {
    double temp[10];
    static int count;
    ...
}
```



# Data Sharing: Changing storage attributes

- One can selectively change storage attributes for constructs using the following clauses\* (note: *list* is a comma-separated list of variables)
  - shared(list)
  - private(list)
  - firstprivate(list)
- These can be used on parallel and for constructs ... other than shared which can only be used on a parallel construct
- Force the programmer to explicitly define storage attributes
  - default (none)

default() can only be used  
on parallel constructs

# 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

# Data Sharing: Private and the original variable

- The original variable's value is unspecified if it is referenced outside of the construct
  - Implementations may reference the original variable or a copy ..... a dangerous programming practice!
  - For example, consider what would happen if the compiler inlined work()?

```
int tmp;  
void danger() {  
    tmp = 0;  
    #pragma omp parallel private(tmp)  
    work();  
    printf("%d\n", tmp);  
}
```

tmp has unspecified value

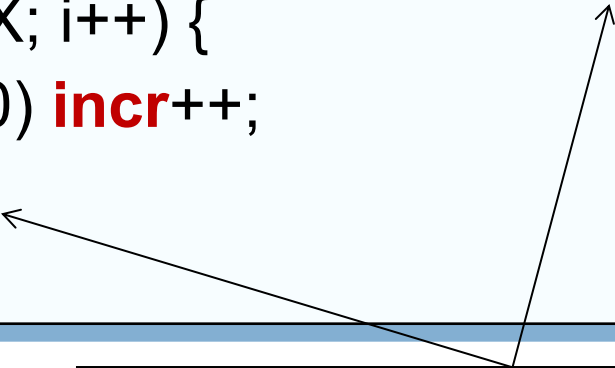
```
extern int tmp;  
void work() {  
    tmp = 5;  
}
```

unspecified which  
copy of tmp

# Firstprivate clause

- Variables initialized from a shared variable
- C++ objects are copy-constructed

```
incr = 0;  
#pragma omp parallel for firstprivate(incr)  
for (i = 0; i <= MAX; i++) {  
    if ((i%2)==0) incr++;  
    A[i] = incr;  
}
```



Each thread gets its own copy of  
incr with an initial value of 0

# Data sharing:

## A data environment test

- Consider this example of PRIVATE and FIRSTPRIVATE

```
variables: A = 1, B = 1, C = 1  
#pragma omp parallel private(B) firstprivate(C)
```

- Are A,B,C private to each thread or shared inside the parallel region?
- What are their initial values inside and values after the parallel region?

Inside this parallel region ...

- “A” is shared by all threads; equals 1
- “B” and “C” are private to each thread.
  - B’s initial value is undefined
  - C’s initial value equals 1

Following the parallel region ...

- B and C revert to their original values of 1
- A is either 1 or the value it was set to inside the parallel region

# Data Sharing: Default clause

- **default(none)**: Forces you to define the storage attributes for variables that appear inside the static extent of the construct ... if you fail the compiler will complain. Good programming practice!
- You can put the default clause on parallel and parallel + workshare constructs.

The static extent is the code in the compilation unit that contains the construct.

```
#include <omp.h>
int main()
{
    int i, j=5;    double x=1.0, y=42.0;
    #pragma omp parallel for default(none) reduction(*:x)
    for (i=0;i<N;i++){
        for(j=0; j<3; j++)
            x+= foobar(i, j, y);
    }
    printf(" x is %f\n",(float)x);
}
```

The compiler would complain about j and y, which is important since you don't want j to be shared

The full OpenMP specification has other versions of the default clause, but they are not used very often so we skip them in the common core

# 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 (hint ... the problem is with the data environment).
- Once you have a working version, try to optimize the program.
  - Try different schedules on the parallel loop.
  - Try different mechanisms to support mutual exclusion ... do the efficiencies change?

# The Mandelbrot Set Area Program

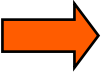
```
#include <omp.h>
# define NPOINTS 1000
# define MXITR 1000
struct d_complex{
    double r;    double i;
};
void testpoint(struct d_complex);
struct d_complex c;
int numoutside = 0;

int main(){
    int i, j;
    double area, error, eps = 1.0e-5;
#pragma omp parallel for private(c, j) firstprivate(eps)
    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);
        }
    }
    area=2.0*2.5*1.125*(double)(NPOINTS*NPOINTS-
numoutside)/(double)(NPOINTS*NPOINTS);
    error=area/(double)NPOINTS;
}
```

```
void testpoint(struct d_complex c){
struct d_complex z;
    int iter;
    double temp;

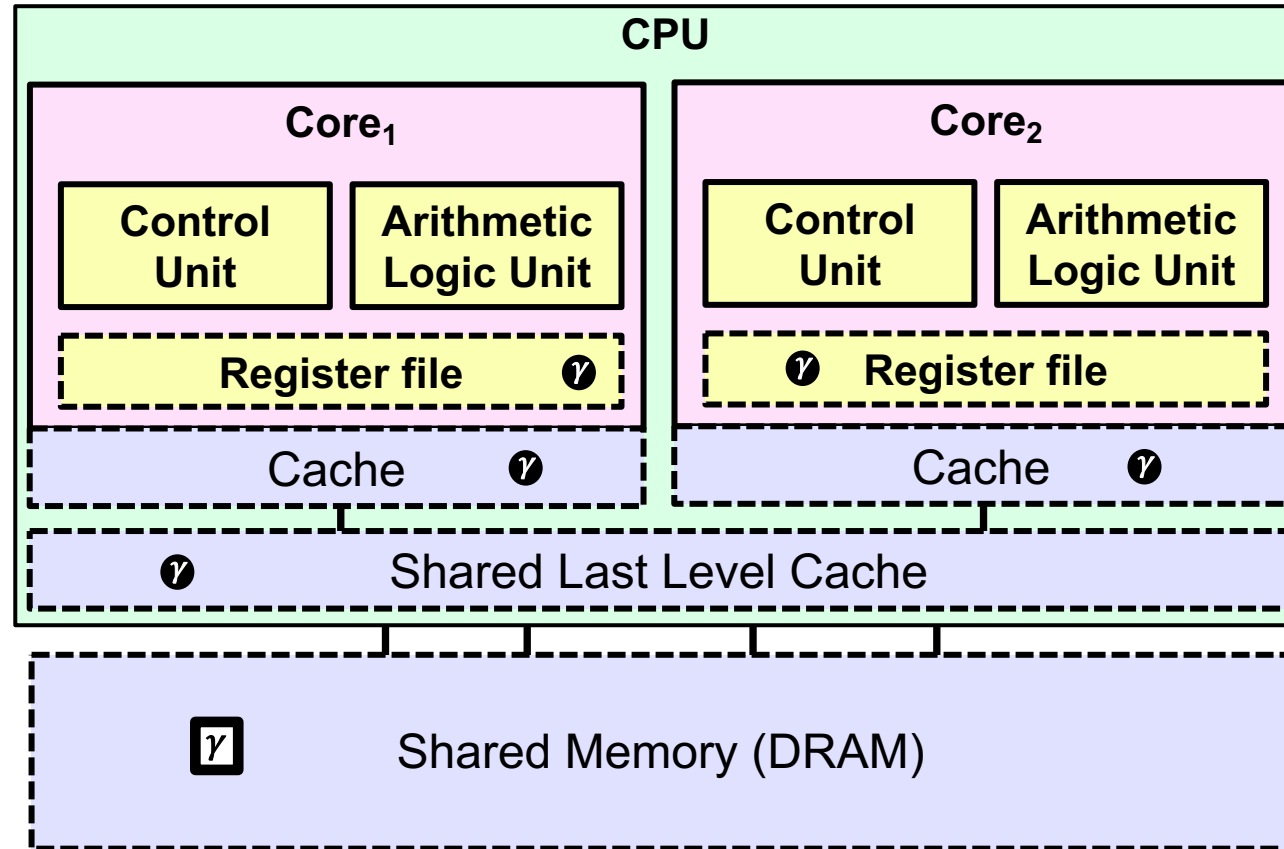
    z=c;
    for (iter=0; iter<MXITR; iter++){
        temp = (z.r*z.r)-(z.i*z.i)+c.r;
        z.i = z.r*z.i*2+c.i;
        z.r = temp;
        if ((z.r*z.r+z.i*z.i)>4.0) {
#pragma omp critical
            numoutside++;
            break;
        }
    }
}
```

- eps was not initialized
- Protect updates of numoutside
- Which value of c does testpoint() see? Global or private?

- Introduction to OpenMP
- Creating Threads
- Synchronization
- Parallel Loops
- Data Environment
-  • Memory Model
- Irregular Parallelism and Tasks
- Recap
- Beyond the Common Core:
  - Worksharing Revisited
  - Synchronization Revisited: Options for Mutual exclusion
  - Thread Affinity and Data Locality
  - Thread Private Data
  - Memory Models and Point-to-Point Synchronization
  - Programming your GPU with OpenMP

# Memory Models ...

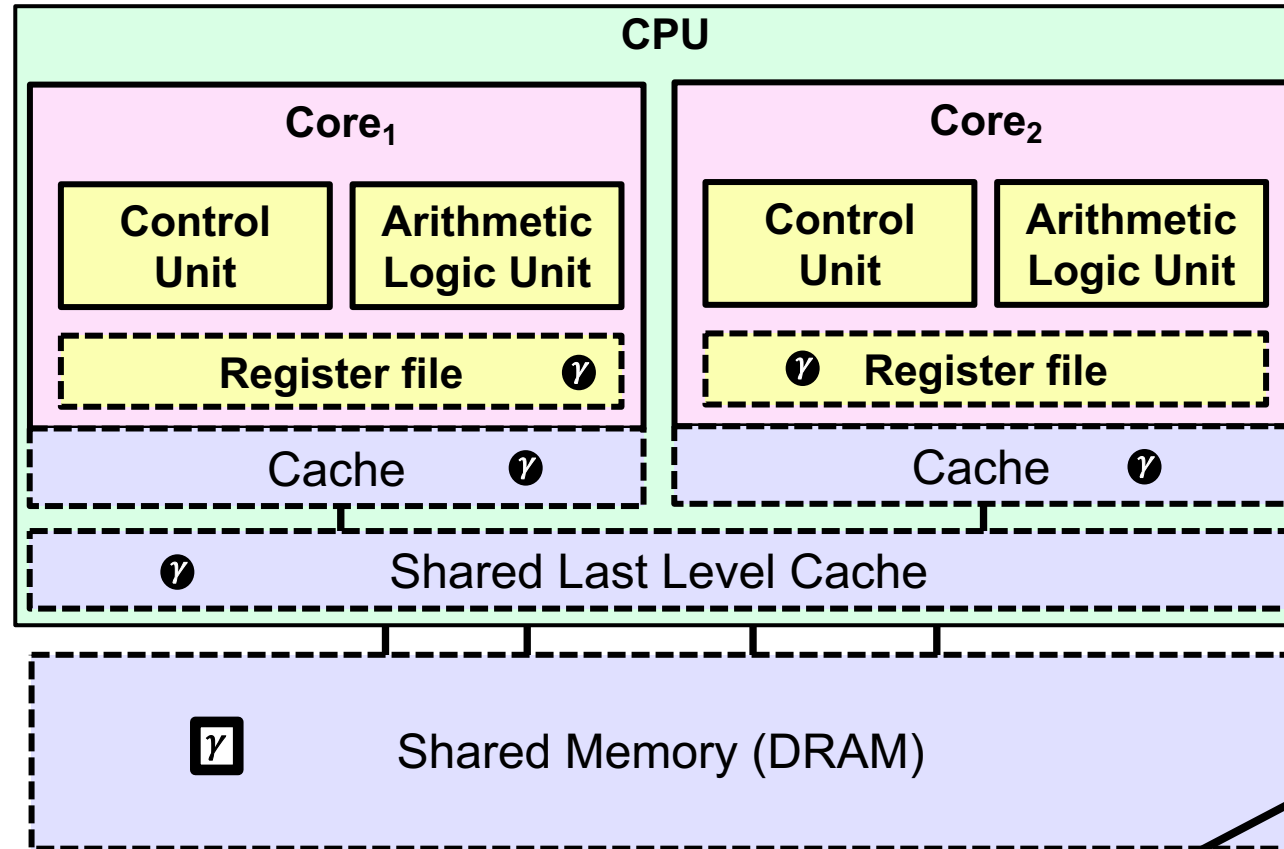
- Programming models for Multithreading support shared memory.
- All threads share an address space ... but consider the variable  $\gamma$



- Multiple copies of a variable (such as  $\gamma$ ) may be present at various levels of cache, or in registers and they may ALL have different values.
- So which value of  $\gamma$  is the one a thread should see at any point in a computation?

# Memory Models ...

- Programming models for Multithreading support shared memory.
- All threads share an address space ... but consider the variable  $\gamma$



A memory consistency model (or “memory model” for short) provides the rules needed to answer this question.

- Multiple copies of a variable (such as  $\gamma$ ) may be present at various levels of cache, or in registers and they may ALL have different values.
- So which value of  $\gamma$  is the one a thread should see at any point in a computation?

# OpenMP and Relaxed Consistency

- Most (if not all) multithreading programming models (including OpenMP) supports a **relaxed-consistency** memory model
  - Threads can maintain a **temporary view** of shared memory that is not consistent with that of other threads
  - These temporary views are made consistent only at certain points in the program
  - The operation that enforces consistency is called the **flush operation**\*

\*Note: in OpenMP 5.0 the name for the flush described here was changed to a "strong flush". This was done so we could distinguish the traditional OpenMP flush (the strong flush) from the new synchronizing flushes (acquire flush and release flush).

# Flush Operation

- Defines a sequence point at which a thread is guaranteed to see a consistent view of memory\*
  - Previous read/writes by this thread have completed and are visible to other threads
  - No subsequent read/writes by this thread have occurred
- A flush operation is analogous to a **fence** in other shared memory APIs

\* This applies to the set of shared variables visible to a thread at the point the flush is encountered. We call this “**the flush set**”

# Flush Example

- Flush forces data to be updated in memory so other threads see the most recent value\*

```
double A;  
A = compute();  
#pragma omp flush(A)  
    // flush to memory to make sure other  
    // threads can pick up the right value
```

Note: OpenMP's flush is analogous to a fence in other shared memory APIs

\* If you pass a list of variables to the flush directive, then that list is **“the flush set”**

# What is the BIG DEAL with Flush?

- Compilers routinely reorder instructions implementing a program
  - Can better exploit the functional units, keep the machine busy, hide memory latencies, etc.
- Compilers generally cannot move instructions:
  - Past a barrier
  - Past a flush on all variables
- But it can move them past a flush with a list of variables so long as those variables are not accessed
- Keeping track of consistency when flushes are used can be confusing ... especially if “flush(list)” is used.

Warning: the flush operation (a strong flush) does not actually synchronize different threads. It just ensures that a thread's variables are made consistent with main memory


# Flush and Synchronization

- A flush operation is implied by OpenMP synchronizations, e.g.,
  - at entry/exit of parallel regions
  - at implicit and explicit barriers
  - at entry/exit of critical regions
  - ....(but not on entry to worksharing regions)

## WARNING:

If you find your self 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.

- Introduction to OpenMP
- Creating Threads
- Synchronization
- Parallel Loops
- Data Environment
- Memory Model
-  • Irregular Parallelism and Tasks
- Recap
- Beyond the Common Core:
  - Worksharing Revisited
  - Synchronization Revisited: Options for Mutual exclusion
  - Thread Affinity and Data Locality
  - Thread Private Data
  - Memory Models and Point-to-Point Synchronization
  - Programming your GPU with OpenMP

# Irregular Parallelism

- Let's call a problem “irregular” when one or both of the following hold:
  - Data Structures are sparse
  - 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

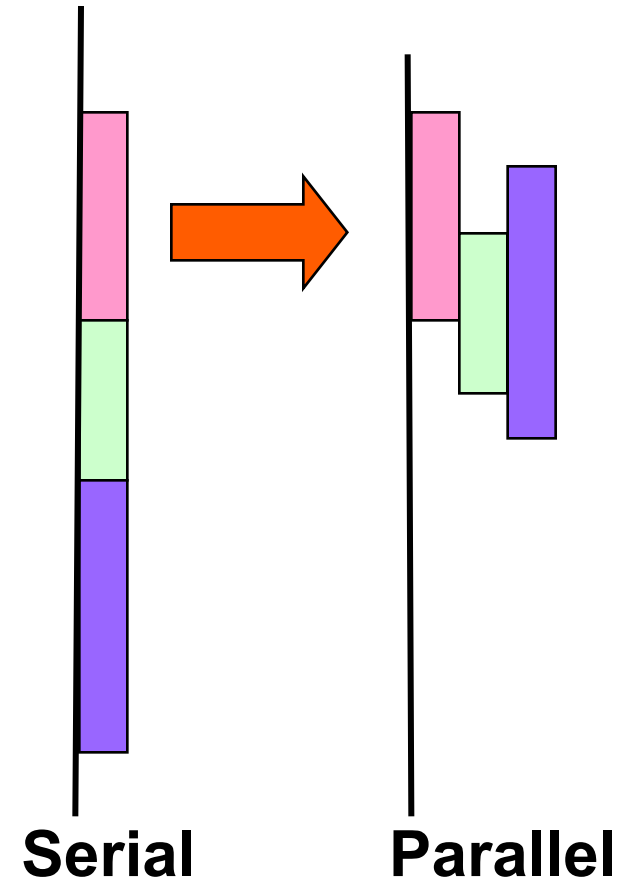
With so much code to add and three passes through the data, this is really ugly.  
There has got to be a better way to do this

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

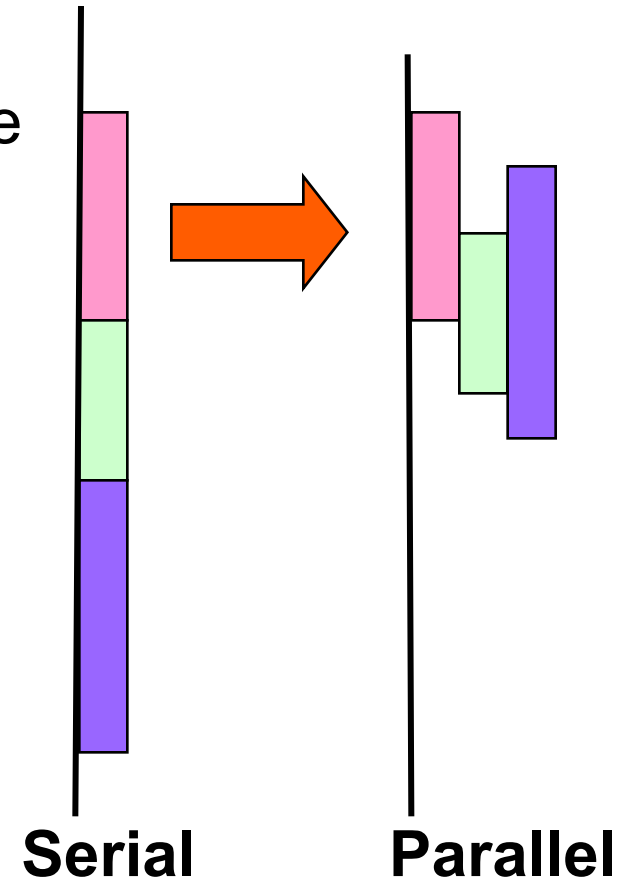
# What are Tasks?

- 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

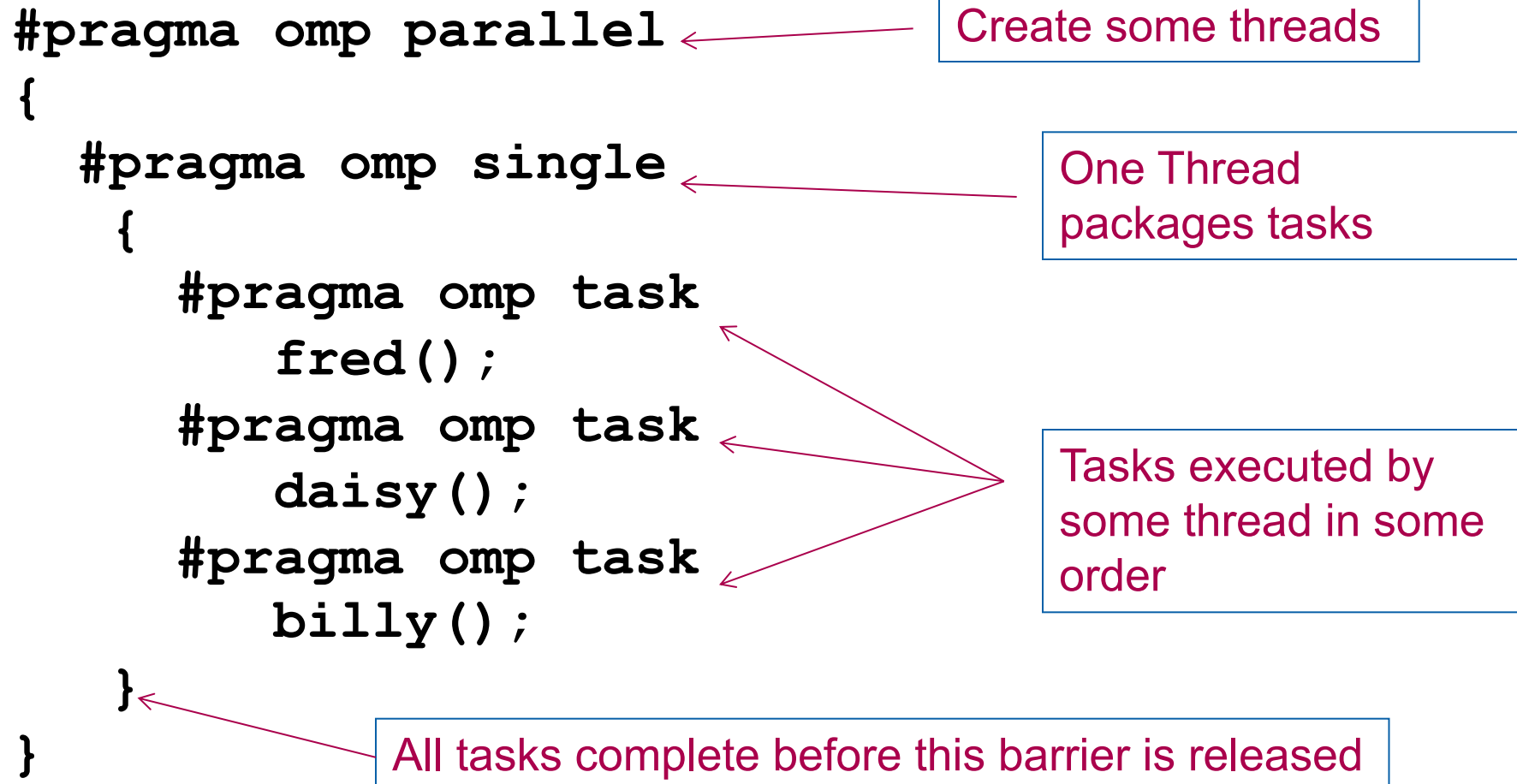
- The **single** construct denotes a block of code that is executed by only one thread (not necessarily the primary\* thread).
- A barrier is implied at the end of the single block (can remove the barrier with a *nowait* clause).

```
#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

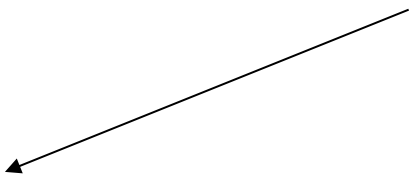
- Variables can be shared, private or firstprivate with respect to task
- These concepts are a little bit different compared with threads:
  - If a variable is **shared** on a task construct, the references to it inside the construct are to the storage with that name at the point where the task was encountered
  - 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  $O(n^2)$  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);
```

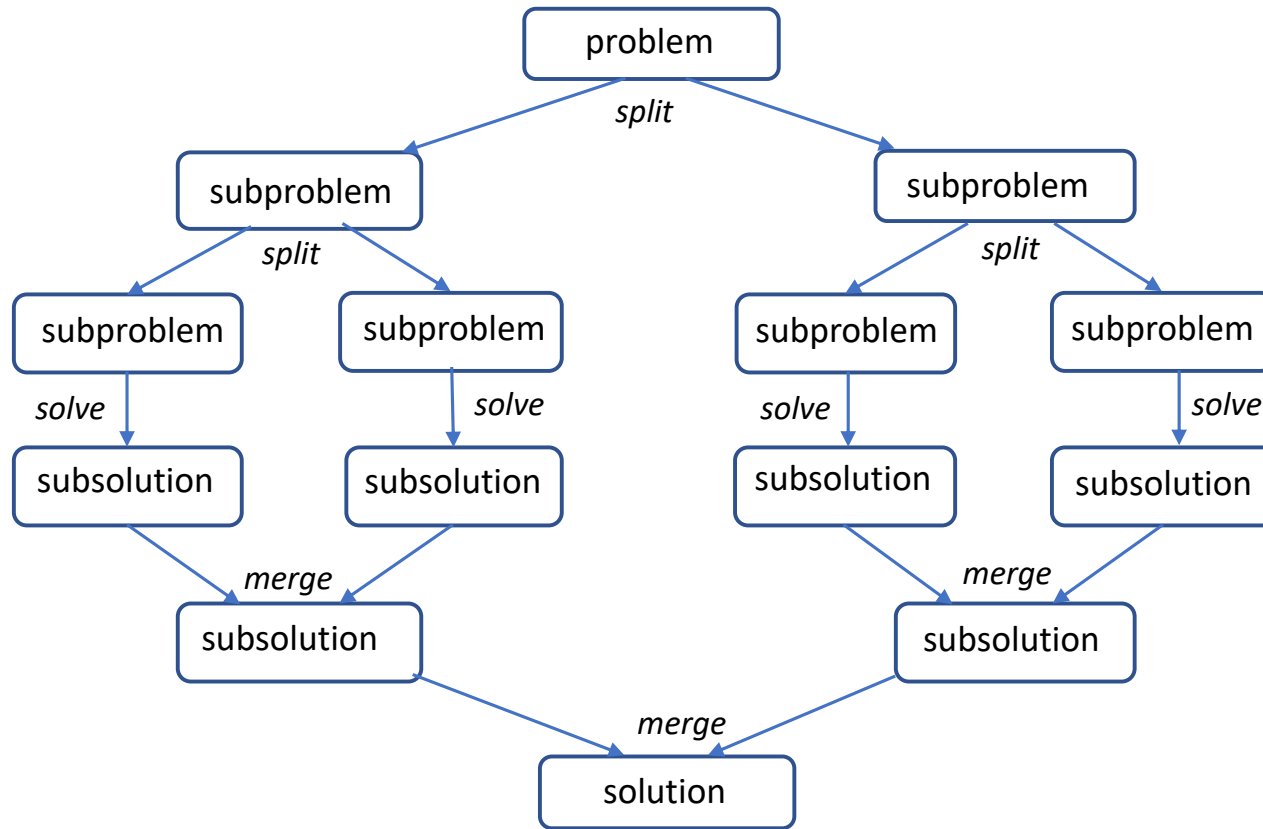
```
  }
```

```
}
```

- 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

# Exercise: PI with tasks

- Go back to the original pi.c program
  - Parallelize this program using OpenMP tasks

```
#pragma omp parallel
#pragma omp task
#pragma omp taskwait
#pragma omp single
double omp_get_wtime()
int omp_get_thread_num();
int omp_get_num_threads();
```

- Hint: first create a recursive pi program and verify that it works. Think about the computation you want to do at the leaves. If you go all the way down to one iteration per leaf-node, won't you just swamp the system with tasks?

# 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;
}
```

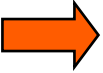
# Results\*: Pi with tasks

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

- Introduction to OpenMP
- Creating Threads
- Synchronization
- Parallel Loops
- Data Environment
- Memory Model
- Irregular Parallelism and Tasks
-  • Recap
- Beyond the Common Core:
  - Worksharing Revisited
  - Synchronization Revisited: Options for Mutual exclusion
  - Thread Affinity and Data Locality
  - Thread Private Data
  - Memory models and point-to-point Synchronization
  - Programming your GPU with OpenMP

# 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.

# There is Much More to OpenMP than the Common Core

- Synchronization mechanisms
  - locks, synchronizing flushes and several forms of atomic
- Data environment
  - lastprivate, threadprivate, default(private|shared)
- Fine grained task control
  - dependencies, tied vs. untied tasks, task groups, task loops ...
- Vectorization constructs
  - simd, uniform, simdlen, inbranch vs. nobranch, ....
- Map work onto an attached device (such as a GPU)
  - target, teams distribute parallel for, target data ...
- ... and much more. The OpenMP 5.0 specification is over 618 pages!!!

Don't become overwhelmed. Master the common core and move on to other constructs when you encounter problems that require them.

# OpenMP Organizations

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

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

- OpenMP User’s Group (cOMPunity) URL:

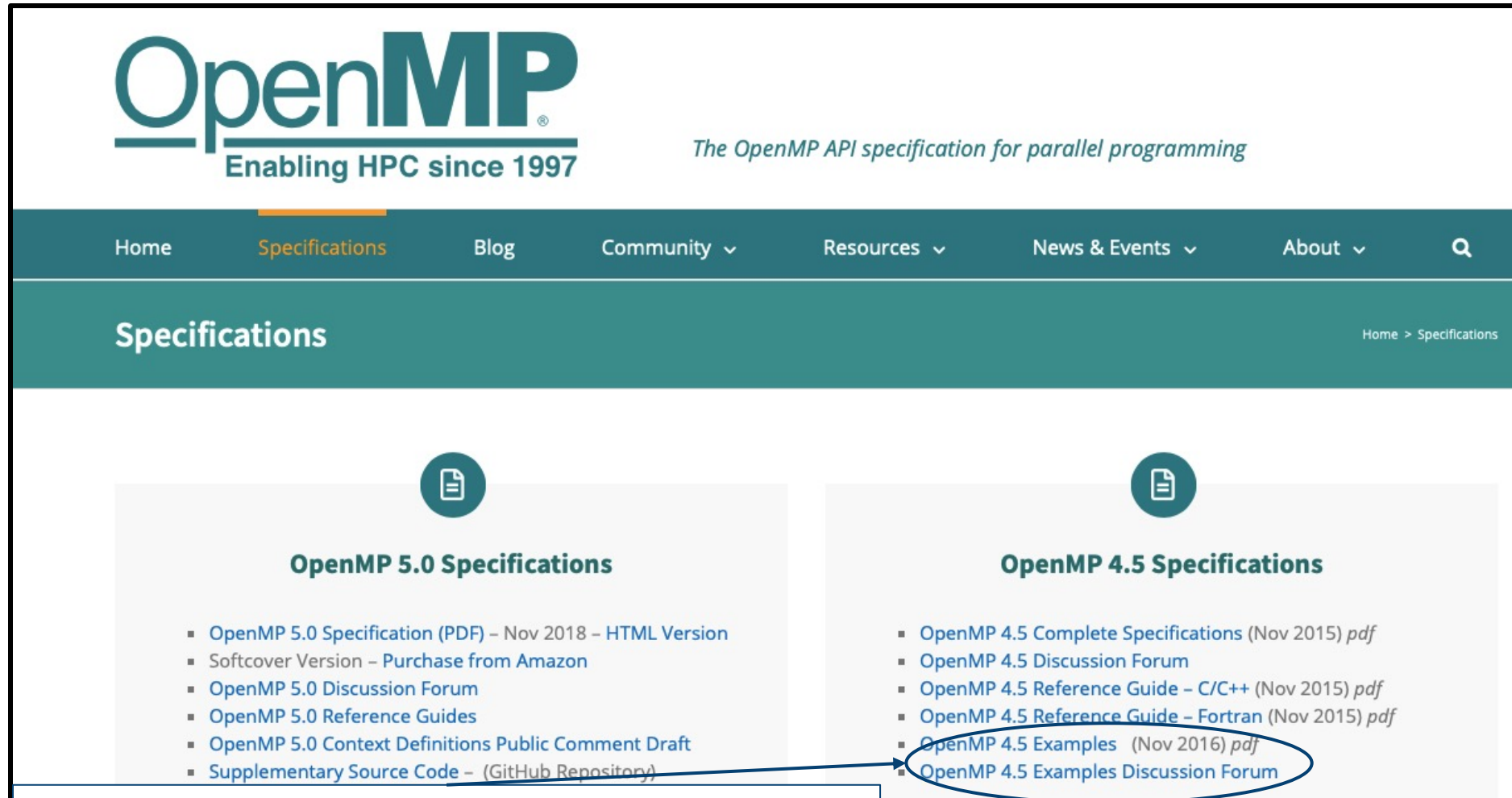
[www.compunity.org](http://www.compunity.org)

Get involved, join the ARB and cOMPunity.

Help define the future of OpenMP

# Resources

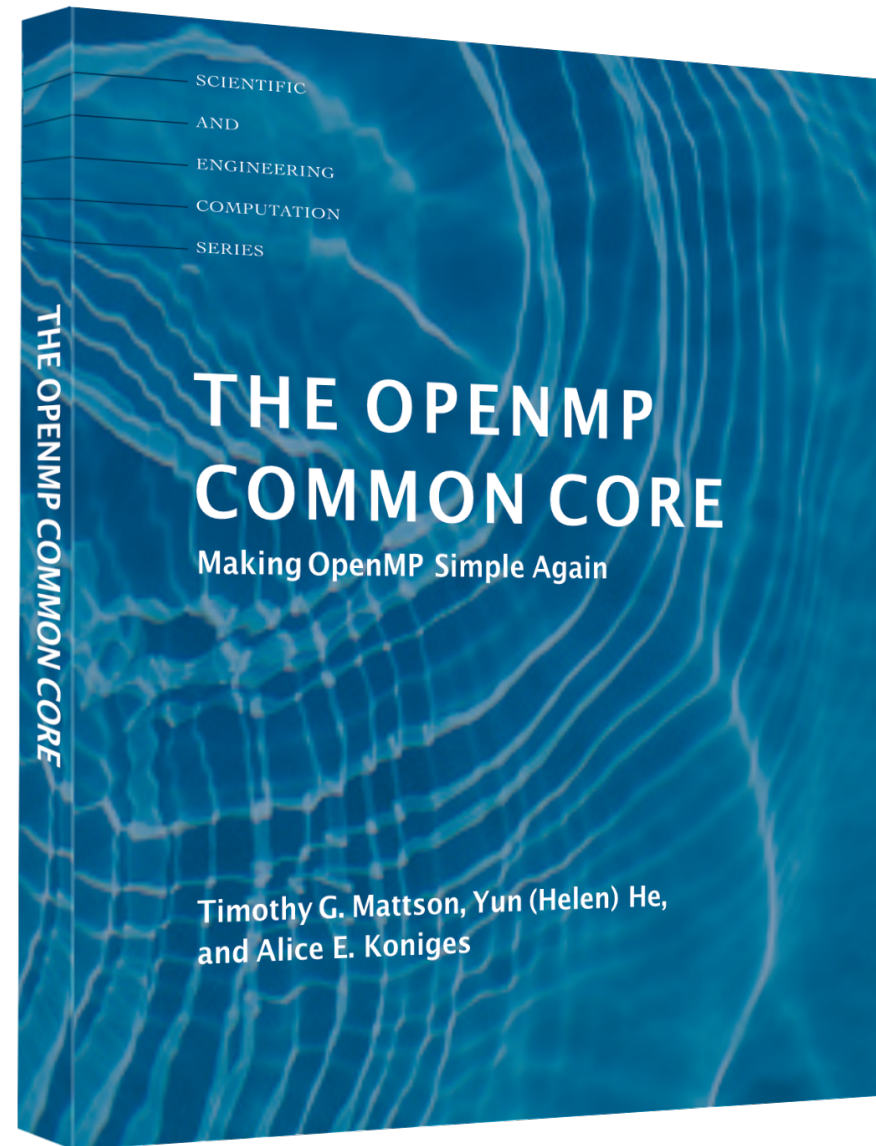
- [www.openmp.org](http://www.openmp.org) has a wealth of helpful resources



Including a comprehensive collection of examples of code using the OpenMP constructs

# 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

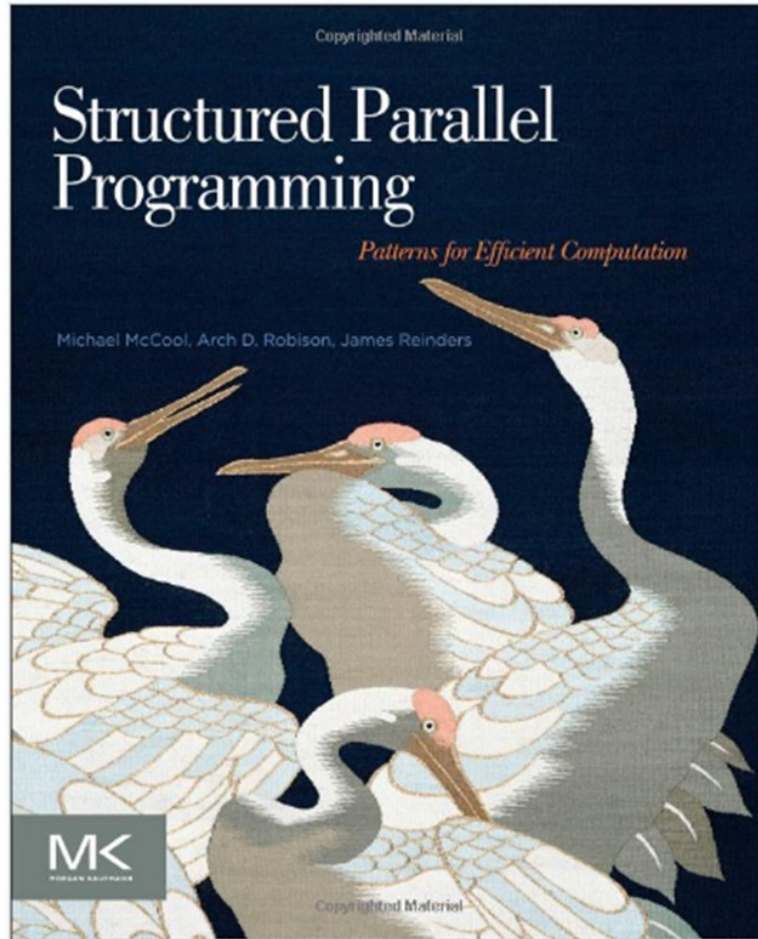


# Books about OpenMP

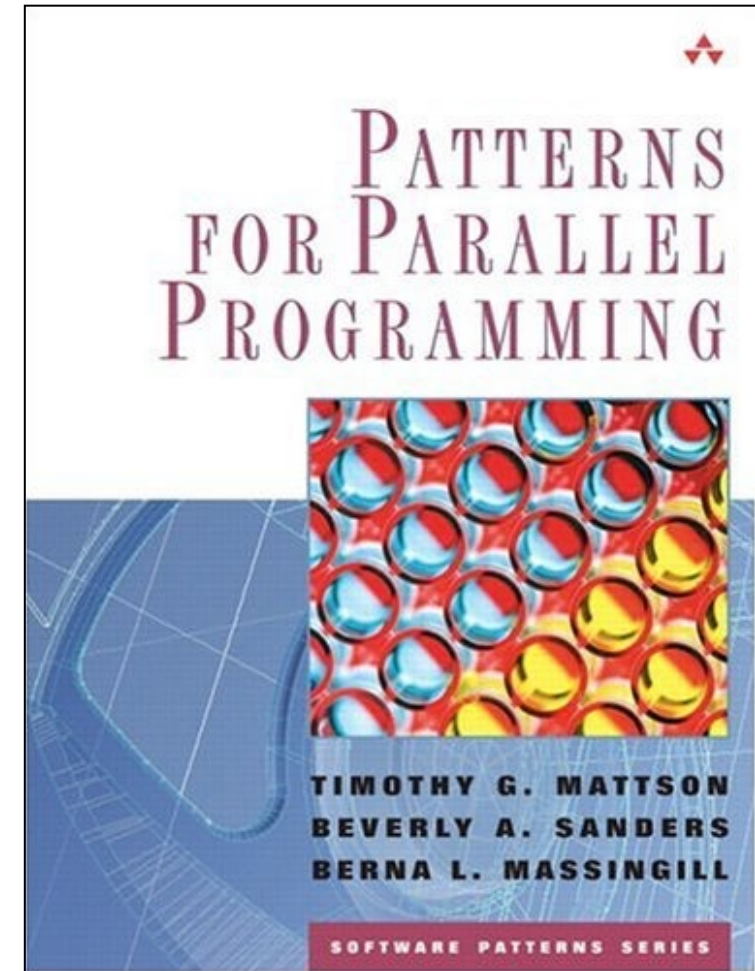
A great book that covers  
OpenMP features beyond  
OpenMP 2.5



# Background references



A great book that explores key patterns with Cilk, TBB, OpenCL, and OpenMP (by McCool, Robison, and Reinders)



- A book about how to “think parallel” with examples in OpenMP, MPI and java

- Introduction to OpenMP
- Creating Threads
- Synchronization
- Parallel Loops
- Data Environment
- Memory Model
- Irregular Parallelism and Tasks
- Recap
- Beyond the Common Core:
  - ➔ – Worksharing Revisited
    - Synchronization Revisited: Options for Mutual exclusion
    - Thread Affinity and Data Locality
    - Thread Private Data
    - Memory models and point-to-point Synchronization
    - Programming your GPU with OpenMP

# 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, nonmontonic 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

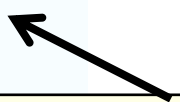
Least work at runtime :  
scheduling done at compile-time

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);
    }
}
mtime = 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

- Introduction to OpenMP
- Creating Threads
- Synchronization
- Parallel Loops
- Data Environment
- Memory Model
- Irregular Parallelism and Tasks
- Recap
- Beyond the Common Core:
  - Worksharing Revisited
  -  – Synchronization Revisited: Options for Mutual exclusion
  - Thread Affinity and Data Locality
  - Thread Private Data
  - Memory models and point-to-point Synchronization
  - Programming your GPU with OpenMP

# 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

- locks (both simple and nested)

Covered in this section

# 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

# The OpenMP 3.1 Atomics (1 of 2)

- Atomic was expanded to cover the full range of common scenarios where you need to protect a memory operation so it occurs atomically:

**# pragma omp atomic [read | write | update | capture]**

- Atomic can protect loads

**# pragma omp atomic read**

**v = x;**

- Atomic can protect stores

**# pragma omp atomic write**

**x = expr;**

- Atomic can protect updates to a storage location (this is the default behavior ... i.e. when you don't provide a clause)

**# pragma omp atomic update**

**x++; or ++x; or x--; or --x; or**

**x binop= expr; or x = x binop expr;**

This is the  
original OpenMP  
atomic

# The OpenMP 3.1 Atomics (2 of 2)

- Atomic can protect the assignment of a value (its capture) AND an associated update operation:

```
# pragma omp atomic capture  
statement or structured block
```

- Where the statement is one of the following forms:

```
v = x++;      v = ++x;      v = x--;    v = -x;    v = x binop expr;
```

- Where the structured block is one of the following forms:

<b>{v = x; x binop = expr;}</b>	<b>{x binop = expr; v = x;}</b>
<b>{v=x; x=x binop expr;}</b>	<b>{X = x binop expr; v = x;}</b>
<b>{v = x; x++;}</b>	<b>{v=x; ++x;}</b>
<b>{++x; v=x;}</b>	<b>{x++; v = x;}</b>
<b>{v = x; x--;}</b>	<b>{v= x; --x;}</b>
<b>{--x; v = x;}</b>	<b>{x--; v = x;}</b>

The capture semantics in atomic were added to map onto common hardware supported atomic operations and to support modern lock free algorithms

# Synchronization: Lock Routines

- Simple Lock routines:

- A simple lock is available if it is unset.
  - `omp_init_lock()`, `omp_set_lock()`,  
`omp_unset_lock()`, `omp_test_lock()`, `omp_destroy_lock()`

A lock implies a memory fence (a “flush”) of all thread visible variables

- Nested Locks

- A nested lock is available if it is unset or if it is set but owned by the thread executing the nested lock function
  - **`omp_init_nest_lock()`, `omp_set_nest_lock()`, `omp_unset_nest_lock()`,  
`omp_test_nest_lock()`, `omp_destroy_nest_lock()`**

Note: a thread always accesses the most recent copy of the lock, so you don't need to use a flush on the lock variable.

Locks with hints were added in OpenMP 4.5 to suggest a lock strategy based on intended use (e.g. contended, uncontended, speculative, unspeculative)

# Synchronization: Simple Locks Example

- Count odds and evens in an input array(x) of N random values.

```
int i, ix, even_count = 0, odd_count = 0;
```

```
omp_lock_t odd_lck, even_lck;
```

```
omp_init_lock(&odd_lck);
```

```
omp_init_lock(&even_lck);
```

One lock per case ... even and odd

```
#pragma omp parallel for private(ix) shared(even_count, odd_count)
```

```
for(i=0; i<N; i++){
```

```
    ix = (int) x[i]; //truncate to int
```

```
    if(((int) x[i])%2 == 0) {
```

```
        omp_set_lock(&even_lck);
```

```
        even_count++;
```

```
        omp_unset_lock(&even_lck);
```

```
    }
```

```
    else{
```

```
        omp_set_lock(&odd_lck);
```

```
        odd_count++;
```

```
        omp_unset_lock(&odd_lck);
```

```
    }
```

```
}
```

```
omp_destroy_lock(&odd_lck);
```

```
omp_destroy_lock(&even_lck);
```

```
}
```

Enforce mutual exclusion updates,  
but in parallel for each case.

Free-up storage when done.

# Exercise


- In the file hist.c, we provide a program that generates a large array of random numbers and then generates a histogram of values.
- This is a "quick and informal" way to test a random number generator ... if all goes well the bins of the histogram should be the same size.
- Parallelize the filling of the histogram You must assure that your program is race free and gets the same result as the sequential program.
- Using everything we've covered today, **manage updates to shared data in two different ways**. Try to minimize the time to generate the histogram.
- Time ONLY the assignment to the histogram. Can you beat the sequential time?

# 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



# Histogram program: one lock per histogram bin

- Example: conflicts are rare, but to play it safe, we must assure mutual exclusion for updates to histogram elements.

```
#pragma omp parallel for
for(i=0;i<NBUCKETS; i++){
    omp_init_lock(&hist_locks[i]);    hist[i] = 0;
}
#pragma omp parallel for
for(i=0;i<NVALS;i++){
    ival = (int) x[i];
    omp_set_lock(&hist_locks[ival]);
    hist[ival]++;
    omp_unset_lock(&hist_locks[ival]);
}

#pragma omp parallel for
for(i=0;i<NBUCKETS; i++)
    omp_destroy_lock(&hist_locks[i]);
```

One lock per element of hist

Enforce mutual exclusion on update to hist array

Free-up storage when done.

# Histogram program: reduction with an array

- We can give each thread a copy of the histogram, they can fill them in parallel, and then combine them when done

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

Easy to write and correct, Uses a lot of memory on the stack, but its fast ... sometimes faster than the serial method.

sequential	0.0019 secs
critical	0.079 secs
Locks per bin	0.029 secs
Reduction, replicated histogram array	0.00097 secs

1000000 random values in X sorted into 50 bins. 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 for the above loop only (we do not time set-up for locks, destruction of locks or anything else)


# Sometimes when working with multiple interacting locks, you have to pay attention to the locking orders

Lock Example from Gafort (SpecOMP'2001)

- Genetic algorithm in Fortran
- Most “interesting” loop: shuffle the population.
  - Original loop is not parallel; performs pair-wise swap of an array element with another, randomly selected element. There are 40,000 elements.
  - Parallelization idea:
    - Perform the swaps in parallel
    - Need to prevent simultaneous access to same array element: use one lock per array element → 40,000 locks.

# Parallel Loop In shuffle.f of Gafort


Exclusive access  
to array  
elements.

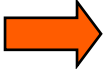


Ordered locking  
prevents  
deadlock.



```
!$OMP PARALLEL PRIVATE(rand, ioother, itemp, temp, my_cpu_id)
    my_cpu_id = 1
!$  my_cpu_id = omp_get_thread_num() + 1
!$OMP DO
    DO j=1,npopsiz-1
        CALL ran3(1,rand,my_cpu_id,0)
        ioother=j+1+DINT(DBLE(npopsiz-j)*rand)
!$      IF (j < ioother) THEN
!$        CALL omp_set_lock(lck(j))
!$        CALL omp_set_lock(lck(ioother))
!$      ELSE
!$        CALL omp_set_lock(lck(ioother))
!$        CALL omp_set_lock(lck(j))
!$      END IF
        itemp(1:nchrome)=iparent(1:nchrome,ioother)
        iparent(1:nchrome,ioother)=iparent(1:nchrome,j)
        iparent(1:nchrome,j)=itemp(1:nchrome)
        temp=fitness(ioother)
        fitness(ioother)=fitness(j)
        fitness(j)=temp
!$      IF (j < ioother) THEN
!$        CALL omp_unset_lock(lck(ioother))
!$        CALL omp_unset_lock(lck(j))
!$      ELSE
!$        CALL omp_unset_lock(lck(j))
!$        CALL omp_unset_lock(lck(ioother))
!$      END IF
    END DO
!$OMP END DO
!$OMP END PARALLEL
```

- Introduction to OpenMP
- Creating Threads
- Synchronization
- Parallel Loops
- Data Environment
- Memory Model
- Irregular Parallelism and Tasks
- Recap
- Beyond the Common Core:
  - Worksharing Revisited
  - Synchronization Revisited: Options for Mutual exclusion
  -  – Thread Affinity and Data Locality
  - Thread Private Data
  - Memory models and point-to-point Synchronization
  - Programming your GPU with OpenMP

- Introduction to OpenMP
- Creating Threads
- Synchronization
- Parallel Loops
- Data Environment
- Memory Model
- Irregular Parallelism and Tasks
- Recap
- Beyond the Common Core:
  - Worksharing Revisited
  - Synchronization Revisited: Options for Mutual exclusion
  - Thread Affinity and Data Locality
  -  – Thread Private Data
  - Memory models and point-to-point Synchronization
  - Programming your GPU with OpenMP

# Data Sharing: Threadprivate

- Makes global data private to a thread
  - Fortran: **COMMON** blocks
  - C: File scope and static variables, static class members
- Different from making them **PRIVATE**
  - with **PRIVATE** global variables are masked.
  - **THREADPRIVATE** preserves global scope within each thread
- Threadprivate variables can be initialized using **COPYIN** or at time of definition (using language-defined initialization capabilities)

# A Threadprivate Example (C)

Use threadprivate to create a counter for each thread.

```
int counter = 0;
#pragma omp threadprivate(counter)

int increment_counter()
{
    counter++;
    return (counter);
}
```

# Data Copying: Copyin

You initialize threadprivate data using a copyin clause.

```
parameter (N=1000)
common/buf/A(N)
!$OMP THREADPRIVATE(/buf/)

!$ Initialize the A array
  call init_data(N,A)

!$OMP PARALLEL COPYIN(A)

... Now each thread sees threadprivate array A initialized
... to the global value set in the subroutine init_data()

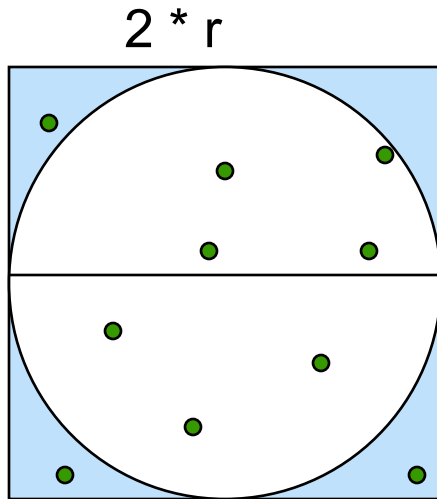
!$OMP END PARALLEL

end
```

# Exercise: Monte Carlo Calculations

Using random numbers to solve tough problems

- Sample a problem domain to estimate areas, compute probabilities, find optimal values, etc.
- Example: Computing  $\pi$  with a digital dart board:



N= 10	$\pi = 2.8$
N=100	$\pi = 3.16$
N= 1000	$\pi = 3.148$

- Throw darts at the circle/square.
- Chance of falling in circle is proportional to ratio of areas:  
$$A_c = r^2 * \pi$$
$$A_s = (2*r) * (2*r) = 4 * r^2$$
$$P = A_c/A_s = \pi / 4$$
- Compute  $\pi$  by randomly choosing points;  $\pi$  is four times the fraction that falls in the circle

# Exercise: Monte Carlo pi (cont)

- We provide three files for this exercise
  - pi\_mc.c: the Monte Carlo method pi program
  - random.c: a simple random number generator
  - random.h: include file for random number generator
- Create a parallel version of this program.
- Run it multiple times with varying numbers of threads.
- Is the program working correctly? Is there anything wrong?

# Parallel Programmers love Monte Carlo algorithms

```
#include "omp.h"
static long num_trials = 10000;
int main ()
{
    long i;    long Ncirc = 0;    double pi, x, y;
    double r = 1.0; // radius of circle. Side of square is 2*r
    seed(0,-r, r); // The circle and square are centered at the origin
#pragma omp parallel for private (x, y) reduction (+:Ncirc)
    for(i=0;i<num_trials; i++)
    {
        x = random();    y = random();
        if ( x*x + y*y <= r*r)  Ncirc++;
    }

    pi = 4.0 * ((double)Ncirc/((double)num_trials);
    printf("\n %d trials, pi is %f \n",num_trials, pi);
}
```

Embarrassingly parallel: the parallelism is so easy its embarrassing.

Add two lines and you have a parallel program.

# Random Numbers: Linear Congruential Generator (LCG)

- LCG: Easy to write, cheap to compute, portable, OK quality

```
random_next = (MULTIPLIER * random_last + ADDEND)% PMOD;  
random_last = random_next;
```

- If you pick the multiplier and addend correctly, LCG has a period of PMOD.
- Picking good LCG parameters is complicated, so look it up (Numerical Recipes is a good source). I used the following:
  - ◆ MULTIPLIER = 1366
  - ◆ ADDEND = 150889
  - ◆ PMOD = 714025

# LCG code

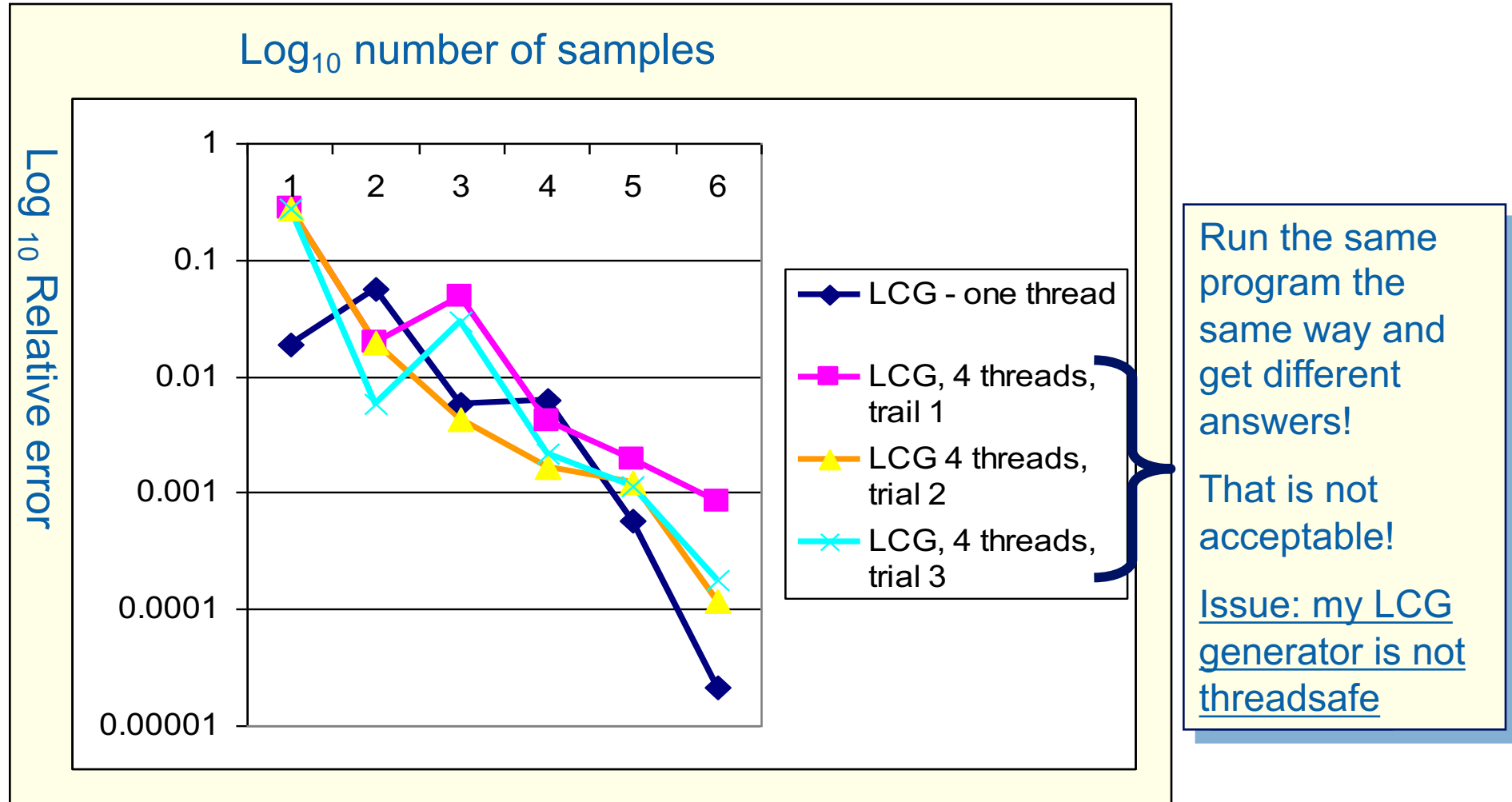
```
static long MULTIPLIER = 1366;
static long ADDEND     = 150889;
static long PMOD       = 714025;
long random_last = 0;
double random ()
{
    long random_next;

    random_next = (MULTIPLIER * random_last + ADDEND)% PMOD;
    random_last = random_next;

    return ((double)random_next/((double)PMOD));
}
```

Seed the pseudo random  
sequence by setting  
random\_last

# Running the PI\_MC program with LCG generator



Program written using the Intel C/C++ compiler (10.0.659.2005) in Microsoft Visual studio 2005 (8.0.50727.42) and running on a dual-core laptop (Intel T2400 @ 1.83 Ghz with 2 GB RAM) running Microsoft Windows XP.

## Exercise: Monte Carlo pi (cont)

- Create a threadsafe version of the monte carlo pi program
- Do not change the interfaces to functions in random.c
  - This is an exercise in modular software ... why should a user of your parallel random number generator have to know any details of the generator or make any changes to how the generator is called?
  - The random number generator must be thread-safe
- Verify that the program is thread safe by running multiple times for a fixed number of threads.
- Any concerns with the program behavior?

# LCG code: threadsafe version

```
static long MULTIPLIER = 1366;
static long ADDEND     = 150889;
static long PMOD       = 714025;
long random_last = 0;
#pragma omp threadprivate(random_last)
double random ()
{
    long random_next;

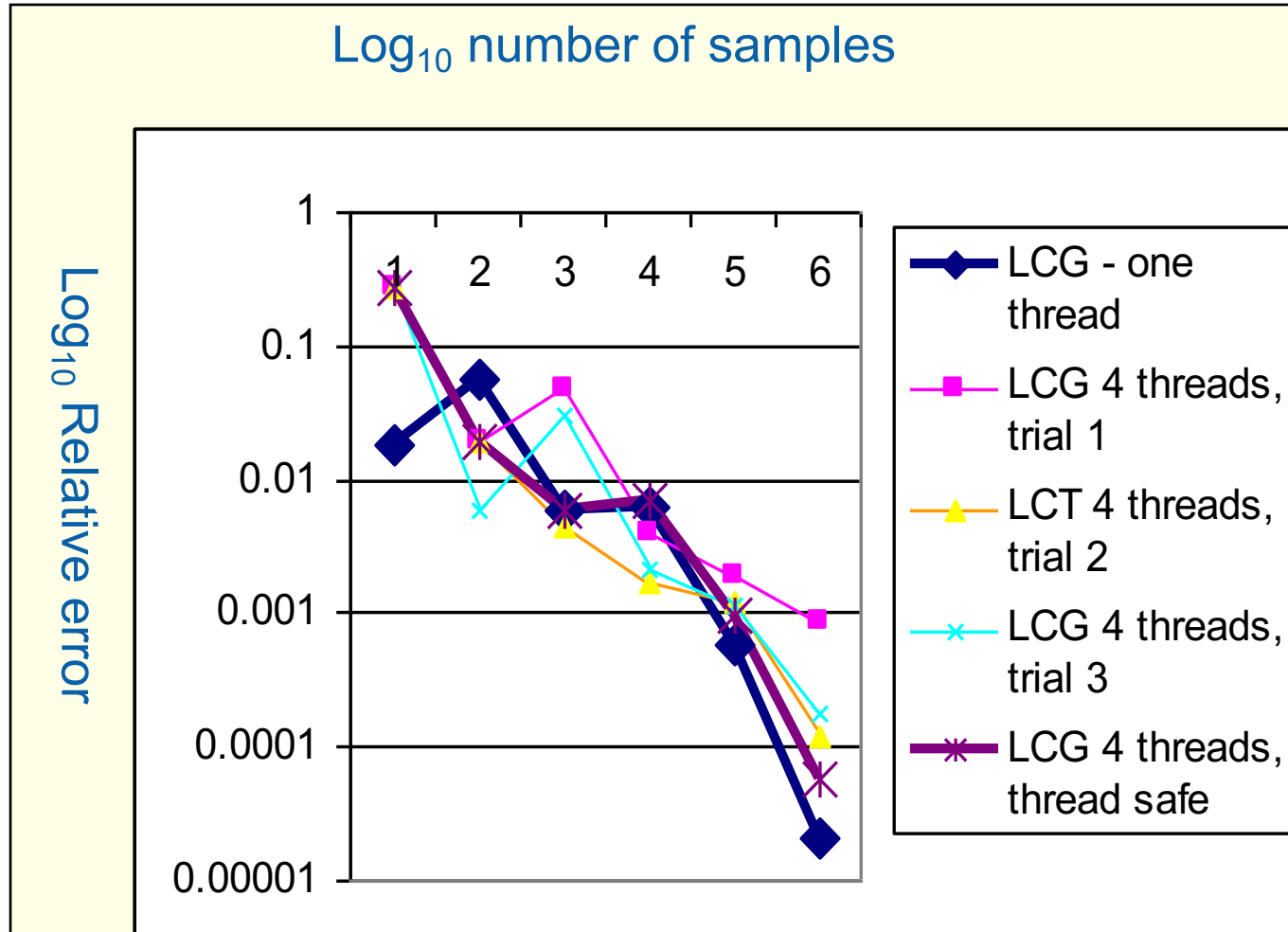
    random_next = (MULTIPLIER * random_last + ADDEND)% PMOD;
    random_last = random_next;

    return ((double)random_next/(double)PMOD);
}
```

random\_last carries state between random number computations,

To make the generator threadsafe, make random\_last threadprivate so each thread has its own copy.

# Thread Safe Random Number Generators



Thread safe version gives the same answer each time you run the program.

But for large number of samples, its quality is lower than the one thread result!

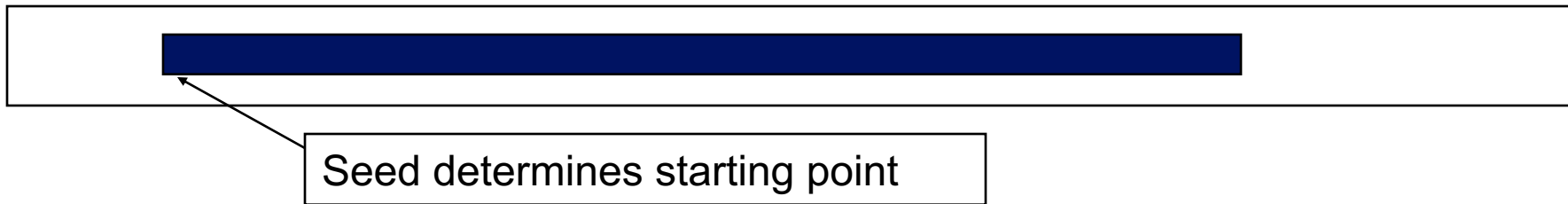
Why?

# Pseudo Random Sequences

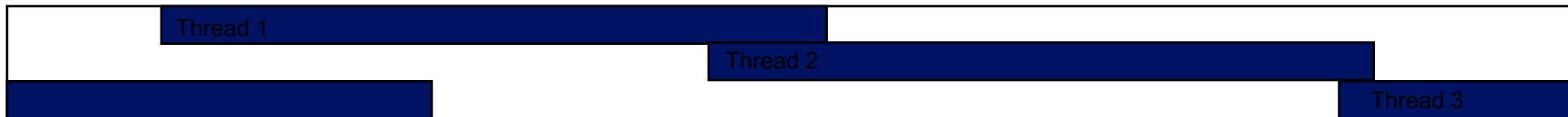
- Random number Generators (RNGs) define a sequence of pseudo-random numbers of length equal to the period of the RNG



- In a typical problem, you grab a subsequence of the RNG range



- Grab arbitrary seeds and you may generate overlapping sequences
  - ◆ E.g. three sequences ... last one wraps at the end of the RNG period.



- Overlapping sequences = over-sampling and bad statistics ... lower quality or even wrong answers!

# Parallel random number generators

- Multiple threads cooperate to generate and use random numbers.

- Solutions:

- Replicate and Pray
- Give each thread a separate, independent generator
- Have one thread generate all the numbers.
- Leapfrog ... deal out sequence values “round robin” as if dealing a deck of cards.
- Block method ... pick your seed so each threads gets a distinct contiguous block.

- Other than “replicate and pray”, these are difficult to implement. Be smart ... get a math library that does it right.

If done right, can generate the same sequence regardless of the number of threads ...

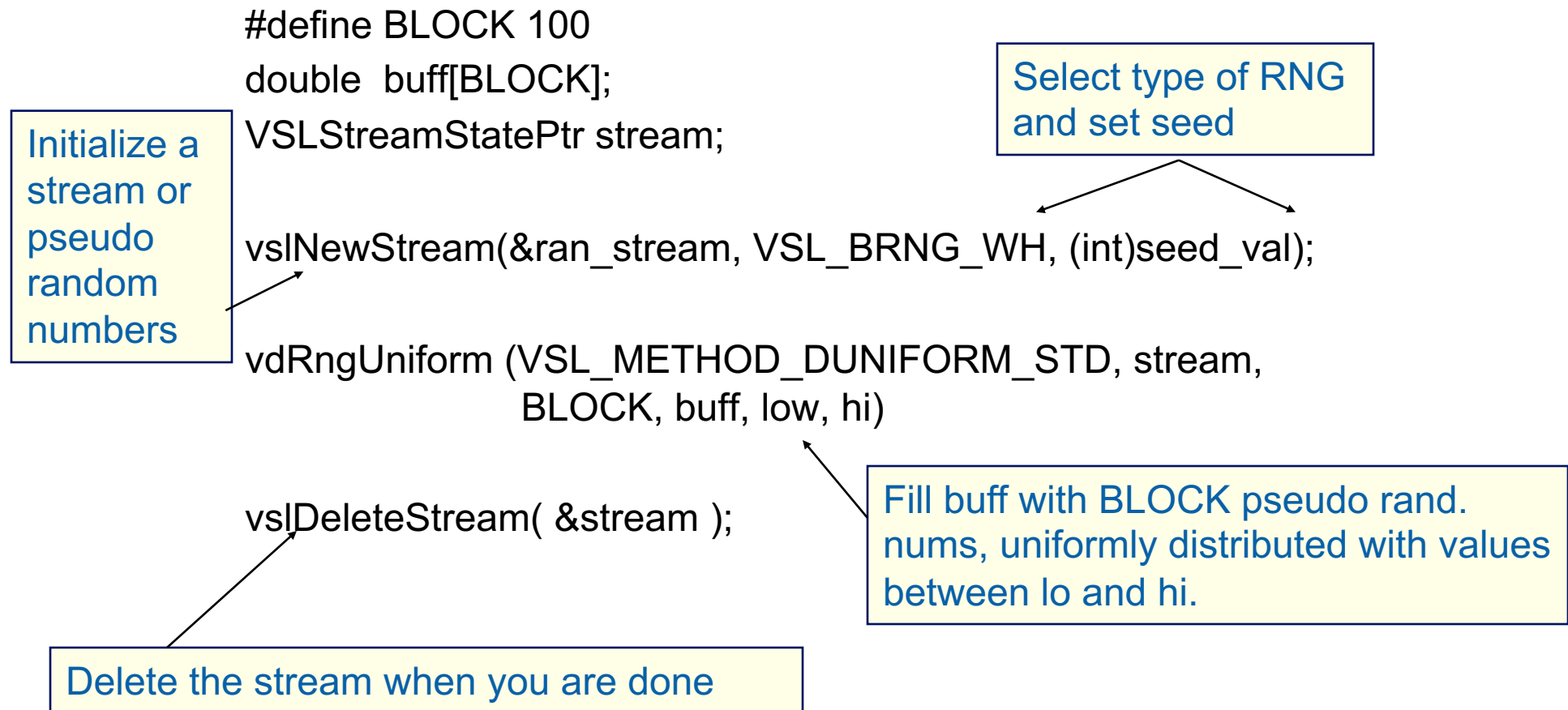
Nice for debugging, but not really needed scientifically.

Intel's Math kernel Library supports a wide range of parallel random number generators.

For an open alternative, the state of the art is the Scalable Parallel Random Number Generators Library (SPRNG): <http://www.sprng.org/> from Michael Mascagni's group at Florida State University.

# MKL Random Number Generators (RNG)

- MKL includes several families of RNGs in its vector statistics library.
- Specialized to efficiently generate vectors of random numbers

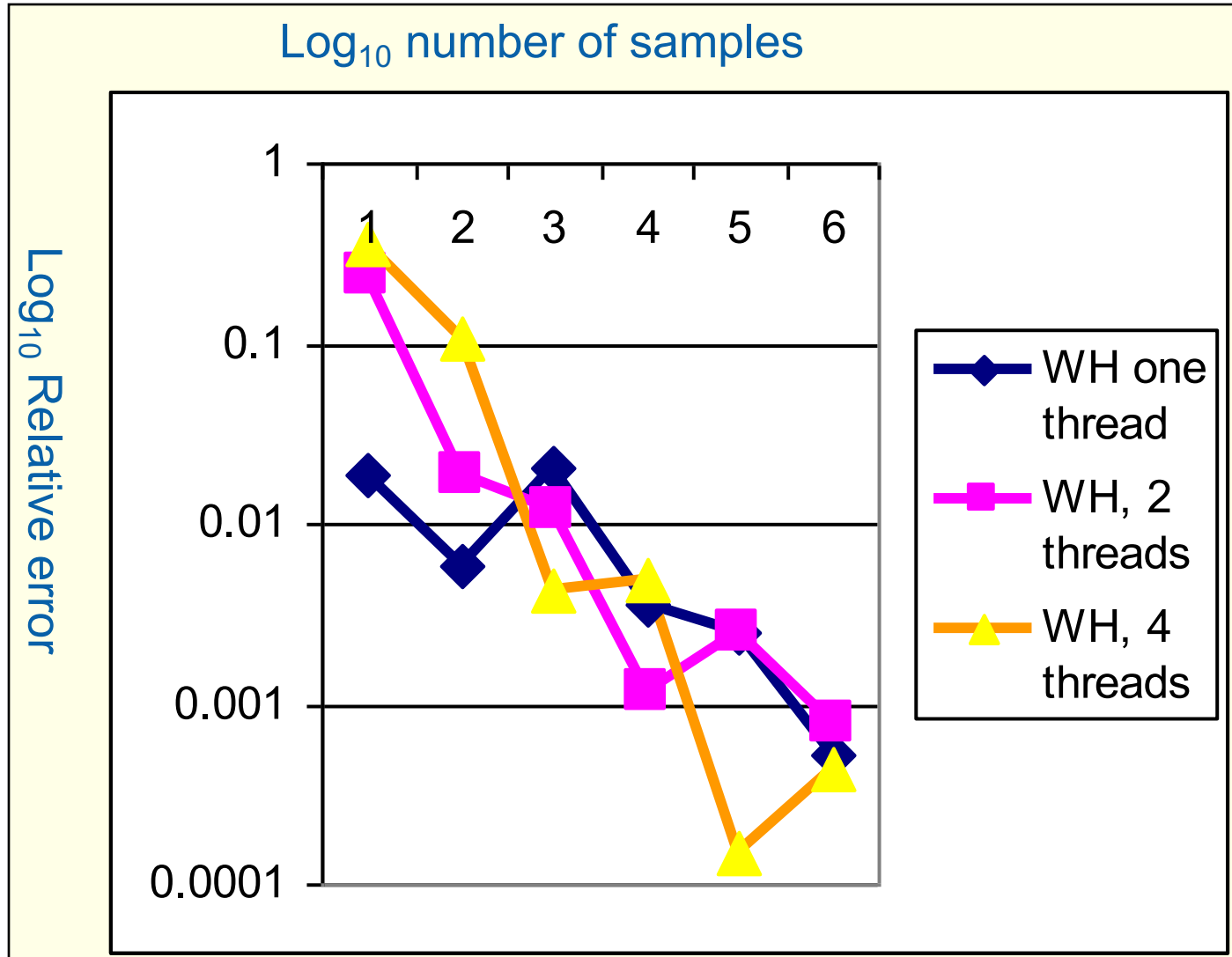


# Wichmann-Hill Generators (WH)

- WH is a family of 273 parameter sets each defining a non-overlapping and independent RNG.
- Easy to use, just make each stream threadprivate and initiate RNG stream so each thread gets a unique WG RNG.

```
VSLStreamStatePtr stream;  
#pragma omp threadprivate(stream)  
  
...  
vslNewStream(&ran_stream, VSL_BRNG_WH+Thrd_ID, (int)seed);
```

# Independent Generator for each thread



Notice that once you get beyond the high error, small sample count range, adding threads doesn't decrease quality of random sampling.

# Leap Frog Method

- Interleave samples in the sequence of pseudo random numbers:
  - Thread  $i$  starts at the  $i^{\text{th}}$  number in the sequence
  - Stride through sequence, stride length = number of threads.
- Result ... the same sequence of values regardless of the number of threads.

```
#pragma omp single
{
    nthreads = omp_get_num_threads();
    iseed = PMOD/MULTIPLIER;    // just pick a seed
    pseed[0] = iseed;
    mult_n = MULTIPLIER;
    for (i = 1; i < nthreads; ++i)
    {
        iseed = (unsigned long long)((MULTIPLIER * iseed) % PMOD);
        pseed[i] = iseed;
        mult_n = (mult_n * MULTIPLIER) % PMOD;
    }
}

random_last = (unsigned long long) pseed[id];
```

One thread  
computes offsets  
and strided  
multiplier

LCG with Addend = 0 just  
to keep things simple

Each thread stores offset starting  
point into its threadprivate "last  
random" value


# Same sequence with many threads.

- We can use the leapfrog method to generate the same answer for any number of threads

Steps	One thread	2 threads	4 threads
1000	3.156	3.156	3.156
10000	3.1168	3.1168	3.1168
100000	3.13964	3.13964	3.13964
1000000	3.140348	3.140348	3.140348
10000000	3.141658	3.141658	3.141658

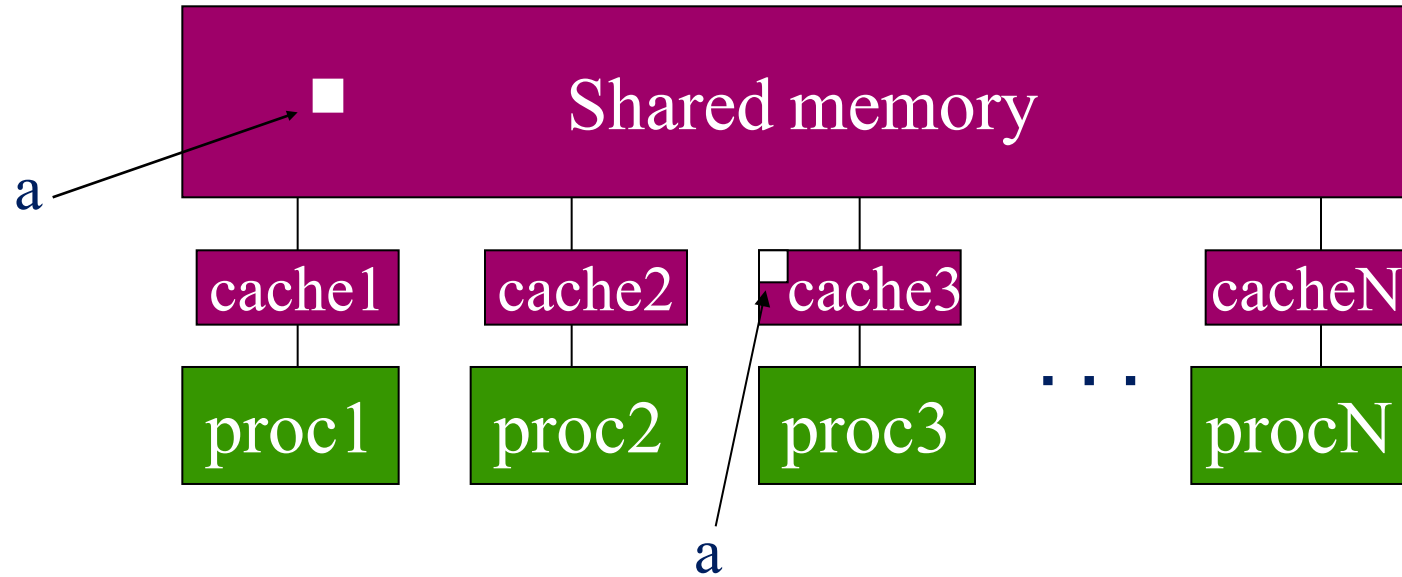
Used the MKL library with two generator streams per computation: one for the x values (WH) and one for the y values (WH+1). Also used the leapfrog method to deal out iterations among threads.



- Introduction to OpenMP
- Creating Threads
- Synchronization
- Parallel Loops
- Data Environment
- Memory Model
- Irregular Parallelism and Tasks
- Recap
- Beyond the Common Core:
  - Worksharing Revisited
  - Synchronization Revisited: Options for Mutual exclusion
  - Thread Affinity and Data Locality
  - Thread Private Data
  -  – Memory models and point-to-point Synchronization
  - Programming your GPU with OpenMP

# OpenMP Memory Model

- OpenMP supports a shared memory model
- All threads share an address space, where variable can be stored or retrieved:



- Threads maintain their own temporary view of memory as well ... the details of which are not defined in OpenMP but this temporary view typically resides in caches, registers, write-buffers, etc.

# Flush Operation

- Defines a sequence point at which a thread enforces a consistent view of memory.
- For variables visible to other threads and associated with the flush operation (the **flush-set**)
  - The compiler can't move loads/stores of the flush-set around a flush:
    - All previous read/writes of the flush-set by this thread have completed
    - No subsequent read/writes of the flush-set by this thread have occurred
  - Variables in the flush set are moved from temporary storage to shared memory.
  - Reads of variables in the flush set following the flush are loaded from shared memory.

IMPORTANT POINT: The flush makes the calling threads temporary view match the view in shared memory. Flush by itself does not force synchronization.

# Memory Consistency: Flush Example

- Flush forces data to be updated in memory so other threads see the most recent value

```
double A;  
A = compute();  
#pragma omp flush(A)  
  
// flush to memory to make sure other  
// threads can pick up the right value
```

Flush without a list: flush set is all thread visible variables

Flush with a list: flush set is the list of variables

Note: OpenMP's flush is analogous to a fence in other shared memory APIs

# Flush and Synchronization

- A flush operation is implied by OpenMP synchronizations, e.g.,
  - at entry/exit of parallel regions
  - at implicit and explicit barriers
  - at entry/exit of critical regions
  - whenever a lock is set or unset
- ....
- (but not at entry to worksharing regions or entry/exit of primary\* regions)

\*the term “master” has been deprecated in OpenMP 5.1 and replaced with the term “primary”.

# Example: prod\_cons.c

- Parallelize a producer/consumer program
  - One thread produces values that another thread consumes.

```
int main()
{
    double *A, sum, runtime;    int flag = 0;

    A = (double *) malloc(N*sizeof(double));

    runtime = omp_get_wtime();

    fill_rand(N, A);           // Producer: fill an array of data

    sum = Sum_array(N, A); // Consumer: sum the array

    runtime = omp_get_wtime() - runtime;

    printf(" In %lf secs, The sum is %lf \n",runtime,sum);
}
```

- Often used with a stream of produced values to implement “pipeline parallelism”
- The key is to implement pairwise synchronization between threads

# Pairwise Synchronization in OpenMP

- OpenMP lacks synchronization constructs that work between pairs of threads.
- When needed, you have to build it yourself.
- Pairwise synchronization
  - Use a shared flag variable
  - Reader spins waiting for the new flag value
  - Use flushes to force updates to and from memory

# Exercise: Producer/Consumer

```
int main()
{
    double *A, sum, runtime;    int numthreads, flag = 0;
    A = (double *)malloc(N*sizeof(double));
    #pragma omp parallel sections
    {
        #pragma omp section
        {
            fill_rand(N, A);

            flag = 1;
        }
        #pragma omp section
        {
            while (flag == 0){

            }

            sum = Sum_array(N, A);
        }
    }
}
```

Put the flushes in the right places to make this program race-free.

Do you need any other synchronization constructs to make this work?

# Solution (try 1): Producer/Consumer

```
int main()
{
    double *A, sum, runtime;    int numthreads, flag = 0;
    A = (double *)malloc(N*sizeof(double));
    #pragma omp parallel sections
    {
        #pragma omp section
        {
            fill_rand(N, A);
            #pragma omp flush
            flag = 1;
            #pragma omp flush (flag)
        }
        #pragma omp section
        {
            #pragma omp flush (flag)
            while (flag == 0){
                #pragma omp flush (flag)
            }
            #pragma omp flush
            sum = Sum_array(N, A);
        }
    }
}
```

Use flag to Signal when the  
“produced” value is ready

Flush forces refresh to memory;  
guarantees that the other thread  
sees the new value of A

Flush needed on both “reader” and “writer”  
sides of the communication

Notice you must put the flush inside the  
while loop to make sure the updated flag  
variable is seen

This program works with the x86 memory model (loads and stores use relaxed  
atomics), but it technically has a race ... on the store and later load of flag

# The OpenMP 3.1 Atomics (1 of 2)

- Atomic was expanded to cover the full range of common scenarios where you need to protect a memory operation so it occurs atomically:

**# pragma omp atomic [read | write | update | capture]**

- Atomic can protect loads

**# pragma omp atomic read**

**v = x;**

- Atomic can protect stores

**# pragma omp atomic write**

**x = expr;**

- Atomic can protect updates to a storage location (this is the default behavior ... i.e. when you don't provide a clause)

**# pragma omp atomic update**

**x++; or ++x; or x--; or --x; or**

**x binop= expr; or x = x binop expr;**

This is the  
original OpenMP  
atomic

# The OpenMP 3.1 Atomics (2 of 2)

- Atomic can protect the assignment of a value (its capture) AND an associated update operation:

# pragma omp atomic capture  
statement or structured block

- Where the statement is one of the following forms:

**`v = x++;      v = ++x;      v = x--;      v = -x;      v = x binop expr;`**

- Where the structured block is one of the following forms:

**`{v = x; x binop = expr;}`**

**`{x binop = expr; v = x;}`**

**`{v=x; x=x binop expr;}`**

**`{X = x binop expr; v = x;}`**

**`{v = x; x++;}`**

**`{v=x; ++x;}`**

**`{++x; v=x;}`**

**`{x++; v = x;}`**

**`{v = x; x--;}`**

**`{v= x; --x;}`**

**`{--x; v = x;}`**

**`{x--; v = x;}`**

The capture semantics in atomic were added to map onto common hardware supported atomic operations and to support modern lock free algorithms

# Atomics and Synchronization Flags

```
int main()
{
    double *A, sum, runtime;
    int numthreads, flag = 0, flg_tmp;
    A = (double *)malloc(N*sizeof(double));
    #pragma omp parallel sections
    {
        #pragma omp section
        {
            fill_rand(N, A);
            #pragma omp flush
            #pragma omp atomic write
            flag = 1;
            #pragma omp flush (flag)
        }
        #pragma omp section
        {
            while (1){
                #pragma omp flush(flag)
                #pragma omp atomic read
                flg_tmp= flag;
                if (flg_tmp==1) break;
            }
            #pragma omp flush
            sum = Sum_array(N, A);
        }
    }
}
```

This program is truly race free ... the reads and writes of flag are protected so the two threads cannot conflict

Still painful and error prone due to all of the flushes that are required

# OpenMP 4.0 Atomic: Sequential consistency



- Sequential consistency:
  - The order of loads and stores in a race-free program appear in some interleaved order and all threads in the team see this same order.
- OpenMP 4.0 added an optional clause to atomics
  - `#pragma omp atomic [read | write | update | capture] [seq_cst]`
- In more pragmatic terms:
  - If the `seq_cst` clause is included, OpenMP adds a flush without an argument list to the atomic operation so you don't need to.
- In terms of the C++'11 memory model:
  - Use of the `seq_cst` clause makes atomics follow the sequentially consistent memory order.
  - Leaving off the `seq_cst` clause makes the atomics relaxed.

Advice to programmers: save yourself a world of hurt ... let OpenMP take care of your flushes for you whenever possible ... use `seq_cst`

# Atomics and Synchronization Flags (4.0)


```
int main()
{
    double *A, sum, runtime;
    int numthreads, flag = 0, flg_tmp;
    A = (double *)malloc(N*sizeof(double));
    #pragma omp parallel sections
    {
        #pragma omp section
        {
            fill_rand(N, A);

            #pragma omp atomic write seq_cst
            flag = 1;
        }
        #pragma omp section
        {
            while (1){

                #pragma omp atomic read seq_cst
                flg_tmp= flag;
                if (flg_tmp==1) break;
            }

            sum = Sum_array(N, A);
        }
    }
}
```

This program is truly race free ... the reads and writes of flag are protected so the two threads cannot conflict – and you do not use any explicit flush constructs (OpenMP does them for you)

- Introduction to OpenMP
- Creating Threads
- Synchronization
- Parallel Loops
- Data Environment
- Memory Model
- Irregular Parallelism and Tasks
- Recap
- Beyond the Common Core:
  - Worksharing Revisited
  - Synchronization Revisited: Options for Mutual exclusion
  - Thread Affinity and Data Locality
  - Thread Private Data
  - Memory models and point-to-point Synchronization
  -  – Programming your GPU with OpenMP

# How do we execute code on a GPU: The SIMT model (Single Instruction Multiple Thread)

1. Turn source code into a scalar work-item

```
extern void reduce( __local float*, __global float*);

__kernel void pi( const int niters, float step_size,
                 __local float* l_sums, __global float* p_sums)
{
    int n_wrk_items = get_local_size(0);
    int loc_id      = get_local_id(0);
    int grp_id      = get_group_id(0);
    float x, accum = 0.0f;  int i, istart, iend;

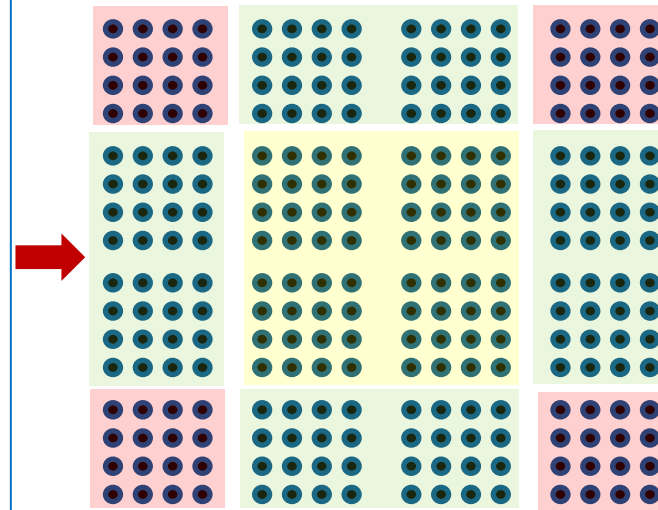
    istart = (grp_id * n_wrk_items + loc_id) * niters;
    iend   = istart + niters;

    for(i= istart; i<iend; i++){
        x = (i+0.5f)*step_size;  accum += 4.0f/(1.0f+x*x); }

    l_sums[loc_id] = accum;
    barrier(CLK_LOCAL_MEM_FENCE);
    reduce(l_sums, p_sums);
}
```

This is OpenCL kernel code ... the sort of code the OpenMP compiler generates on your behalf

2. Map work-items onto an N dim index space.



3. Map data structures onto the same index space

4. Run on hardware designed around the same SIMT execution model



# How do we execute code on a GPU:

## OpenCL and CUDA nomenclature

Turn source code into a scalar **work-item** (a CUDA **thread**)

```
extern void reduce( __local float*, __global float*);

__kernel void pi( const int niters, float step_size,
                 __local float* l_sums, __global float* p_sums)
{
    int n_wrk_items = get_local_size(0);
    int loc_id      = get_local_id(0);
    int grp_id      = get_group_id(0);
    float x, accum = 0.0f;  int i,istart,iend;

    istart = (grp_id * n_wrk_items + loc_id) * niters;
    iend   = istart+niters;

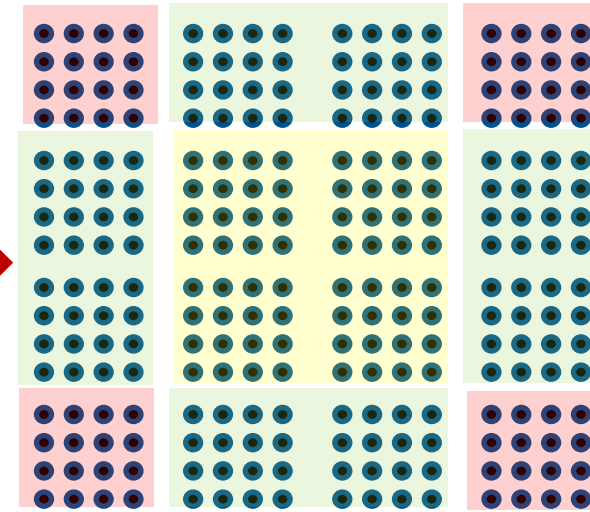
    for(i= istart; i<iend; i++){
        x = (i+0.5f)*step_size;  accum += 4.0f/(1.0f+x*x); }

    l_sums[loc_id] = accum;
    barrier(CLK_LOCAL_MEM_FENCE);
    reduce(l_sums, p_sums);
}
```

This code defines a **kernel**

Submit a kernel to an OpenCL **command queue** or a CUDA **stream**

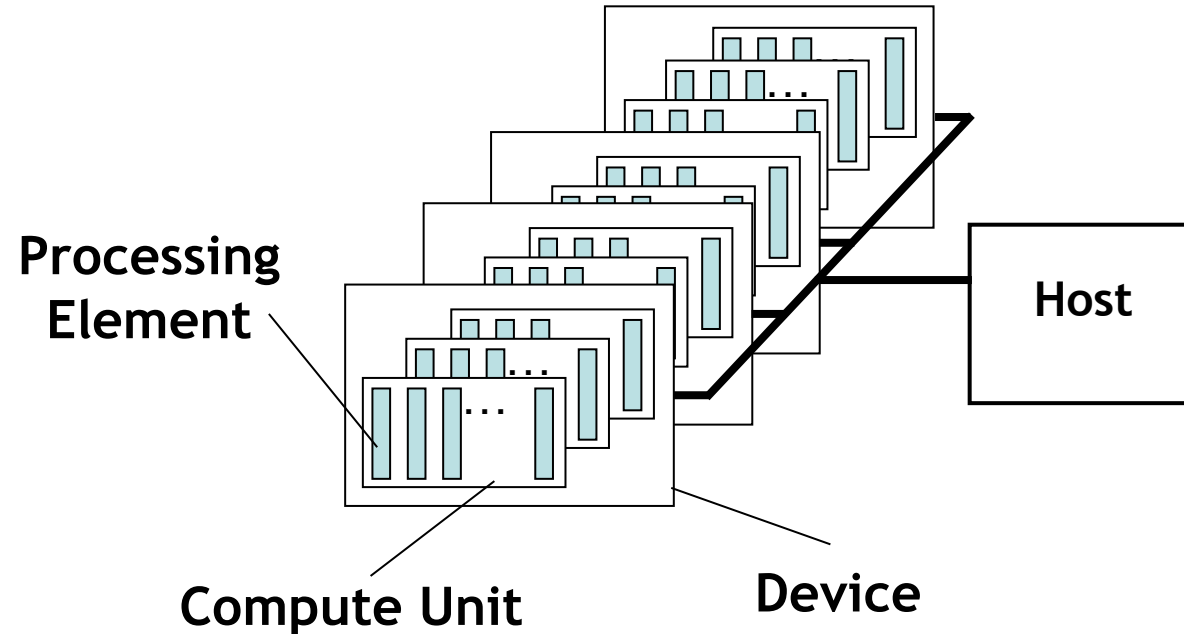
Organize work-items into **work-groups** and map onto an N dim index space. CUDA calls a work-group a **thread-block**



OpenCL index space is called an **NDRange**.  
CUDA calls this a **Grid**

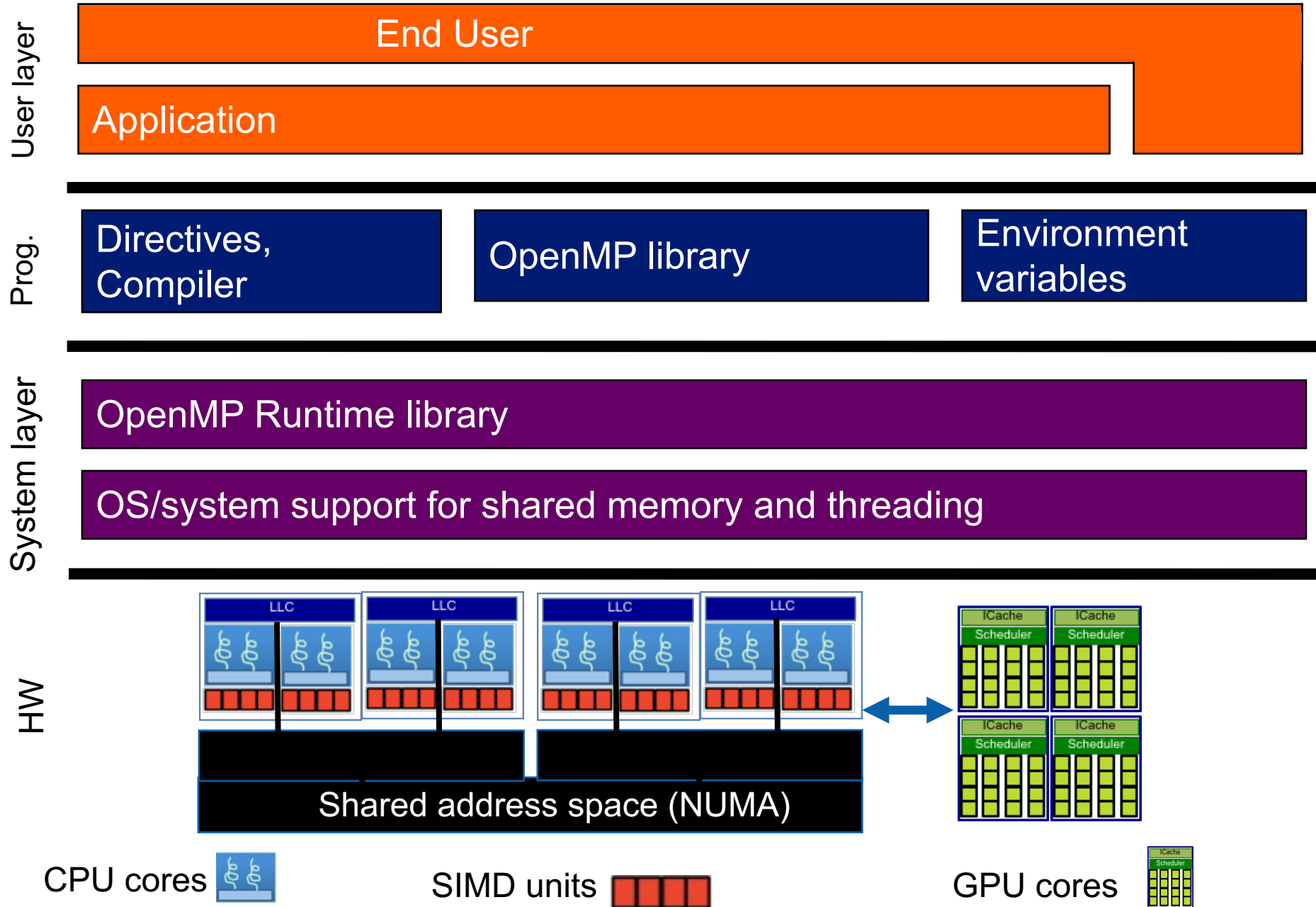
It's called SIMT, but GPUs are really vector-architectures with a block of work-items executing together (a **subgroup** in OpenCL or a **warp** with CUDA)

# A Generic Host/Device Platform Model



- One **Host** and one or more **Devices**
  - Each Device is composed of one or more Compute Units
  - Each Compute Unit is divided into one or more **Processing Elements**
- Memory divided into **host memory** and **device memory**

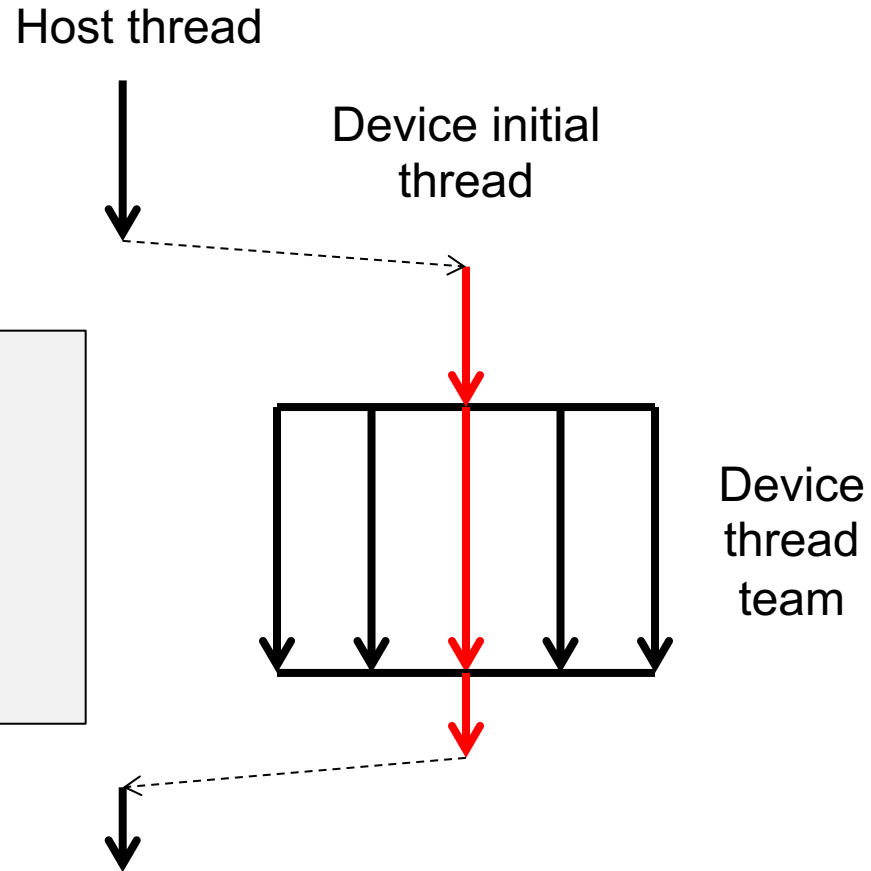
# OpenMP Basic Definitions: Solution stack



# Accelerated workshare v1.0

- Simply add a target construct

```
#pragma omp target  
#pragma omp parallel for  
for (i=0;i<N;i++)  
...
```



- Transfer control of execution to a **SINGLE** device thread
- Only one team of threads workshares the loop

# The target data environment

- Remember: distinct memory spaces on host and device.
- OpenMP uses a combination of *implicit* and *explicit* memory movement.
- Data may move between the host and the device in well defined places:
  - Firstly, at the beginning and end of a **target** region:

```
#pragma omp target
```

```
{ // Data may move here
```

```
...
```

```
} // and here
```

- We'll discuss the other places later...

# Default Data Mapping:

## implicit movement with a target region

- Scalar variables:
  - Examples:
    - `int N; double x;`
  - OpenMP implicitly maps scalar variables as **firstprivate**
    - A new value per work-item initialized with the original value (in OpenCL nomenclature, the firstprivate goes in private memory).
  - The variable *is not* copied back to the host at the end of the target region.
  - OpenMP target regions for GPUs execute with CUDA/OpenCL, and a firstprivate scalar can be launched as a parameter to a kernel function without the overhead of setting up a variable in device memory.

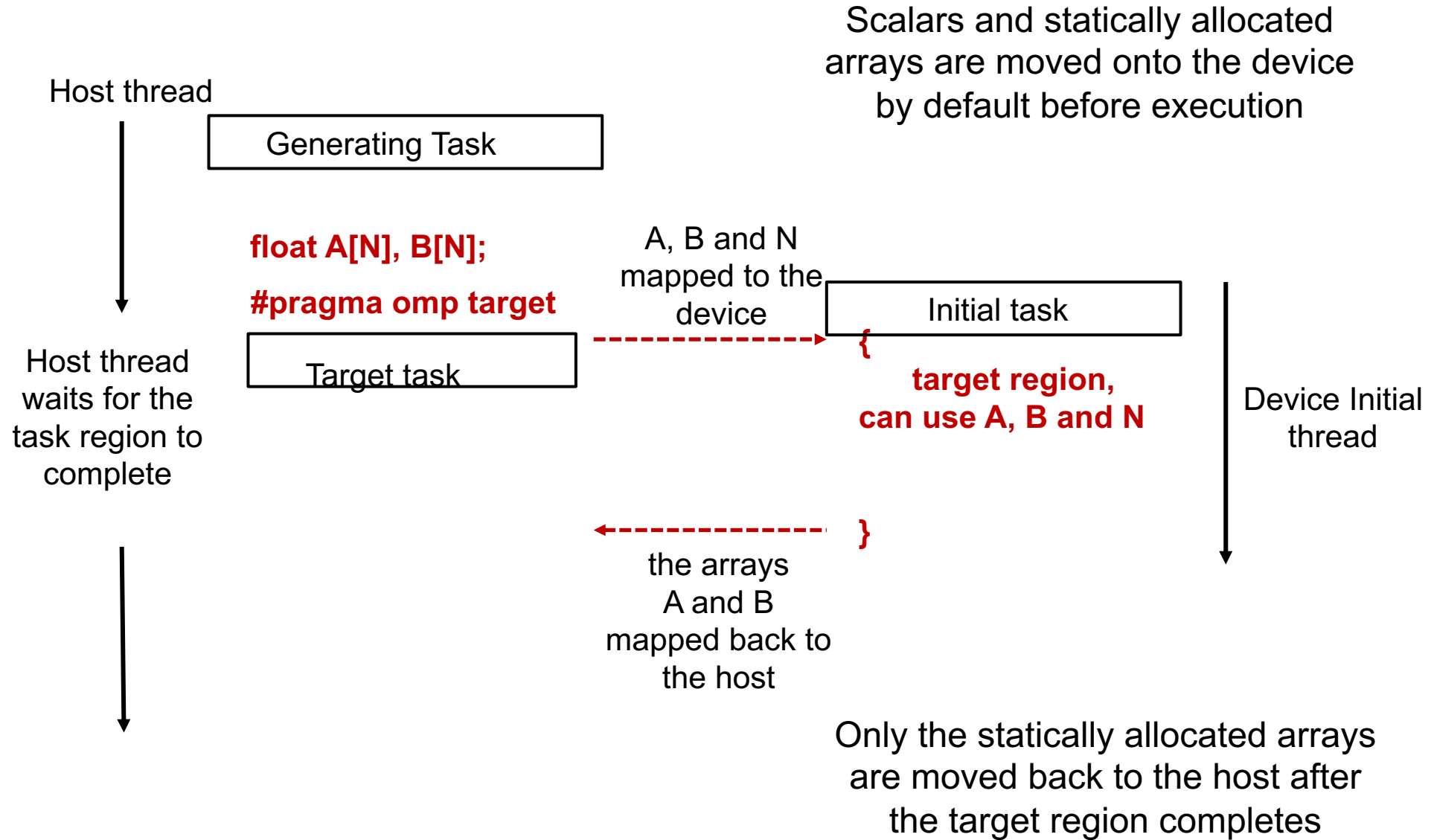
# Default Data Mapping: implicit movement with a target region

- Non-scalar variables:
  - Must have a *complete type*.
  - Example: fixed sized (stack) array:
    - `double A[1000];`
  - Copied to the device at the start of the **target** region, *and* copied back at the end. In OpenCL nomenclature, these are placed in device global memory.
  - A new value is created in the target region and initialized with the original data, but it is shared between threads on the device. Data is copied back to the host at the end of the target region.
  - OpenMP calls this mapping **tofrom**

# Default Data Mapping: implicit movement with a target region

- Pointers and their data:
  - *Example: arrays allocated on the heap*
    - `double *A = malloc(sizeof(double)*1000);`
  - The pointer value will be mapped.
  - But the data it points to ***will not*** be mapped by default.

# The target data environment



# Default Data Sharing: example

```
int main(void) {  
    int N = 1024;  
    double A[N], B[N];
```

1. Variables created in host memory.

```
#pragma omp target
```

2. Scalar **N** and stack arrays **A** and **B** are copied to device memory. Execution transferred to device.

```
{
```

```
    for (int ii = 0; ii < N; ++ii) {
```

3. **ii** is **private** on the device as it's declared within the target region

```
        A[ii] = A[ii] + B[ii];
```

4. Execution on the device.

```
    }
```

```
} // end of target region
```

5. stack arrays **A** and **B** are copied *from* device memory back to the host. Host resumes execution.

```
}
```

# Explicit Data Sharing

- Previously, we described the rules for *implicit* data movement.
- We *explicitly* control the movement of data using the **map** clause.
- Data allocated on the heap needs to explicitly copied to/from the device:

```
int main(void) {  
    int ii=0, N = 1024;  
    int* A = malloc(sizeof(int)*N);  
  
    #pragma omp target  
    {  
        // N, ii and A all exist here  
        // The data that A points to (*A , A[ii]) DOES NOT exist here!  
    }  
}
```

# Controlling data movement

```
int i, a[N], b[N], c[N];  
#pragma omp target map(to:a,b) map(tofrom:c)
```

Data movement  
defined from the  
*host* perspective.

- The various forms of the map clause
  - **map(to:list)**: On entering the region, variables in the list are initialized on the device using the original values from the host (host to device copy).
  - **map(from:list)**: At the end of the target region, the values from variables in the list are copied into the original variables (device to host copy). On entering the region, initial value of the variable is not initialized.
  - **map(tofrom:list)**: the effect of both a map-to and a map-from (host to device copy at start of region, device to host copy at end)
  - **map(alloc:list)**: On entering the region, data is allocated and uninitialized on the device.
  - **map(list)**: equivalent to **map(tofrom:list)**.
- For pointers you must use array section notation ..
  - **map(to:a[0:N])**. Notation is **A[lower-bound : length]**

# Moving arrays with the map clause

```
int main(void) {  
    int N = 1024;  
    int* A = malloc(sizeof(int)*N);  
  
    #pragma omp target map(A[0:N])  
    {  
        // N, ii and A all exist here  
        // The data that A points to DOES exist here!  
    }  
}
```

Default mapping  
**map(tofrom: A[0:N])**

Copy at start and end of  
**target** region.

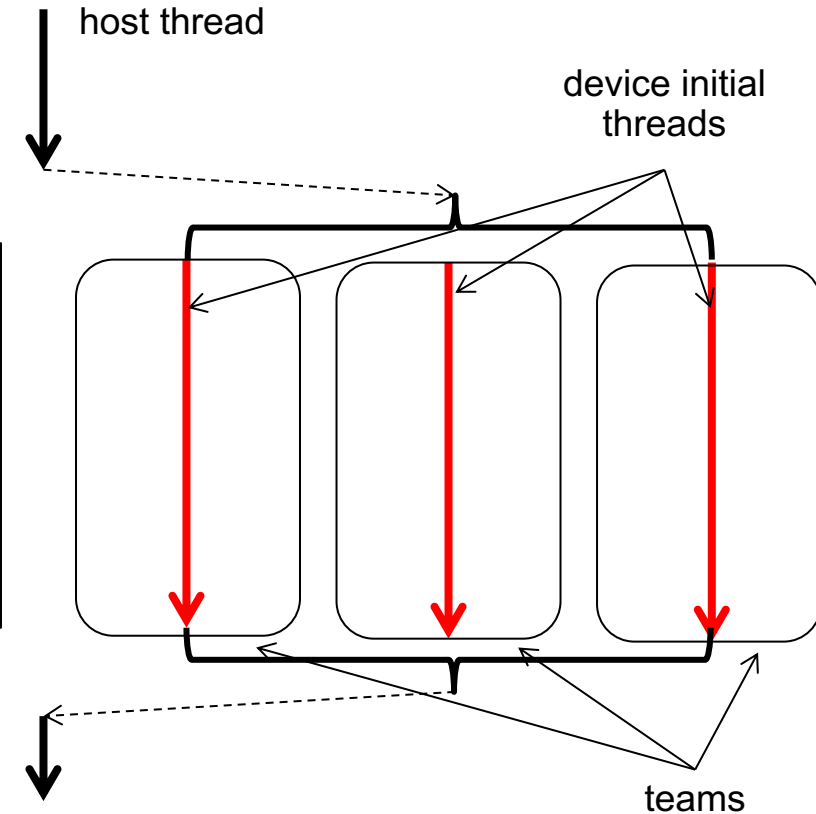
# teams and distribute constructs

- The **teams** construct
  - Similar to the **parallel** construct
  - It starts a league of thread teams
  - Each team in the league starts as one initial thread – a team of one
  - Threads in different teams cannot synchronize with each other
  - The construct must be “perfectly” nested in a **target** construct
- The **distribute** construct
  - Similar to the **for** construct
  - Loop iterations are workshared across the initial threads in a league
  - No implicit barrier at the end of the construct
  - **dist\_schedule(kind[, chunk\_size])**
    - If specified, scheduling kind must be static
    - Chunks are distributed in round-robin fashion in chunks of size **chunk\_size**
    - If no chunk size specified, chunks are of (almost) equal size; each team receives at least one chunk

# Accelerated workshare v2.0

- teams construct
- distribute construct

```
#pragma omp target  
#pragma omp teams  
#pragma omp distribute  
for (i=0;i<N;i++)  
  ...
```

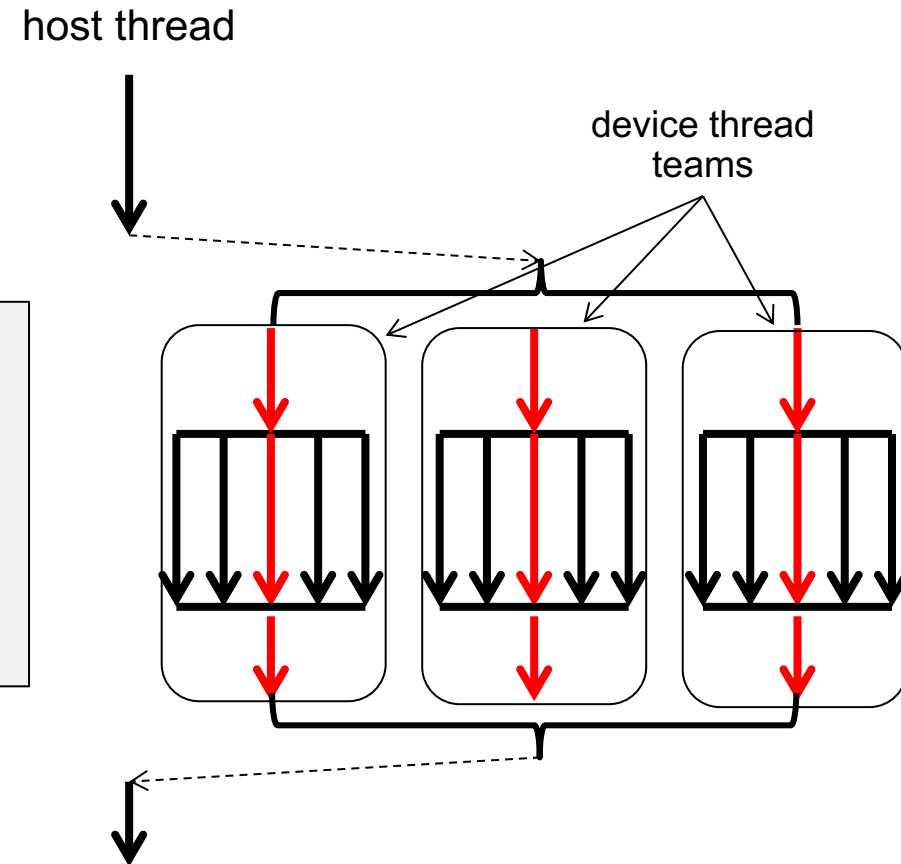


- Transfer execution control to **MULTIPLE** device initial threads
- Workshare loop iterations across the initial threads.

# Accelerate workshare v3.0

- teams distribute
- parallel for simd

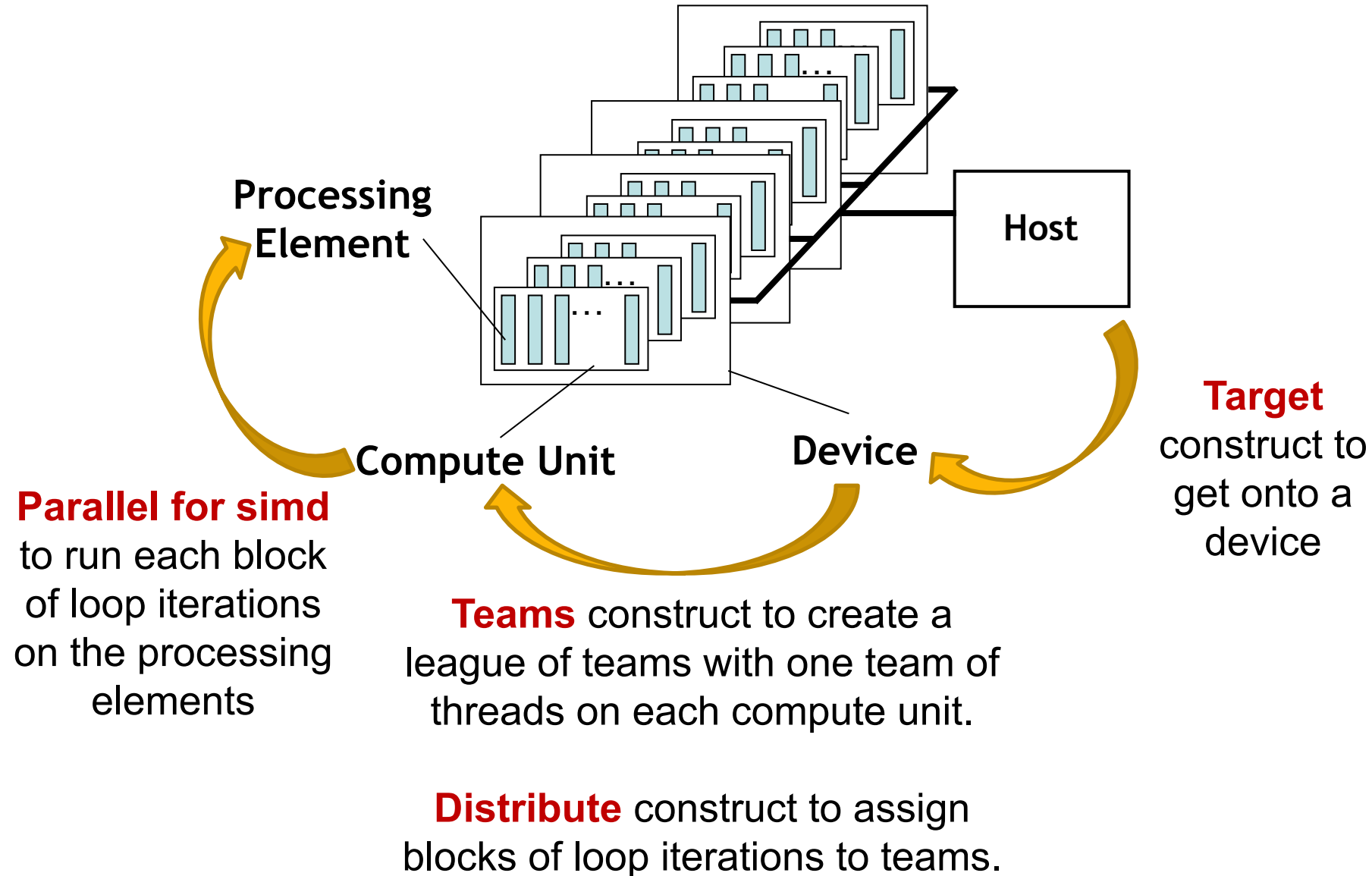
```
#pragma omp target  
#pragma omp teams distribute  
for (i=0;i<N;i++)  
#pragma omp parallel for simd  
for (j=0;j<M;i++)  
...
```



- Transfer execution control to **MULTIPLE** device initial threads
  - Workshare loop iterations across the initial threads (teams distribute)
- Each initial thread becomes the primary\* thread in a thread team
  - Workshare loop iterations across the threads in a team (parallel for)

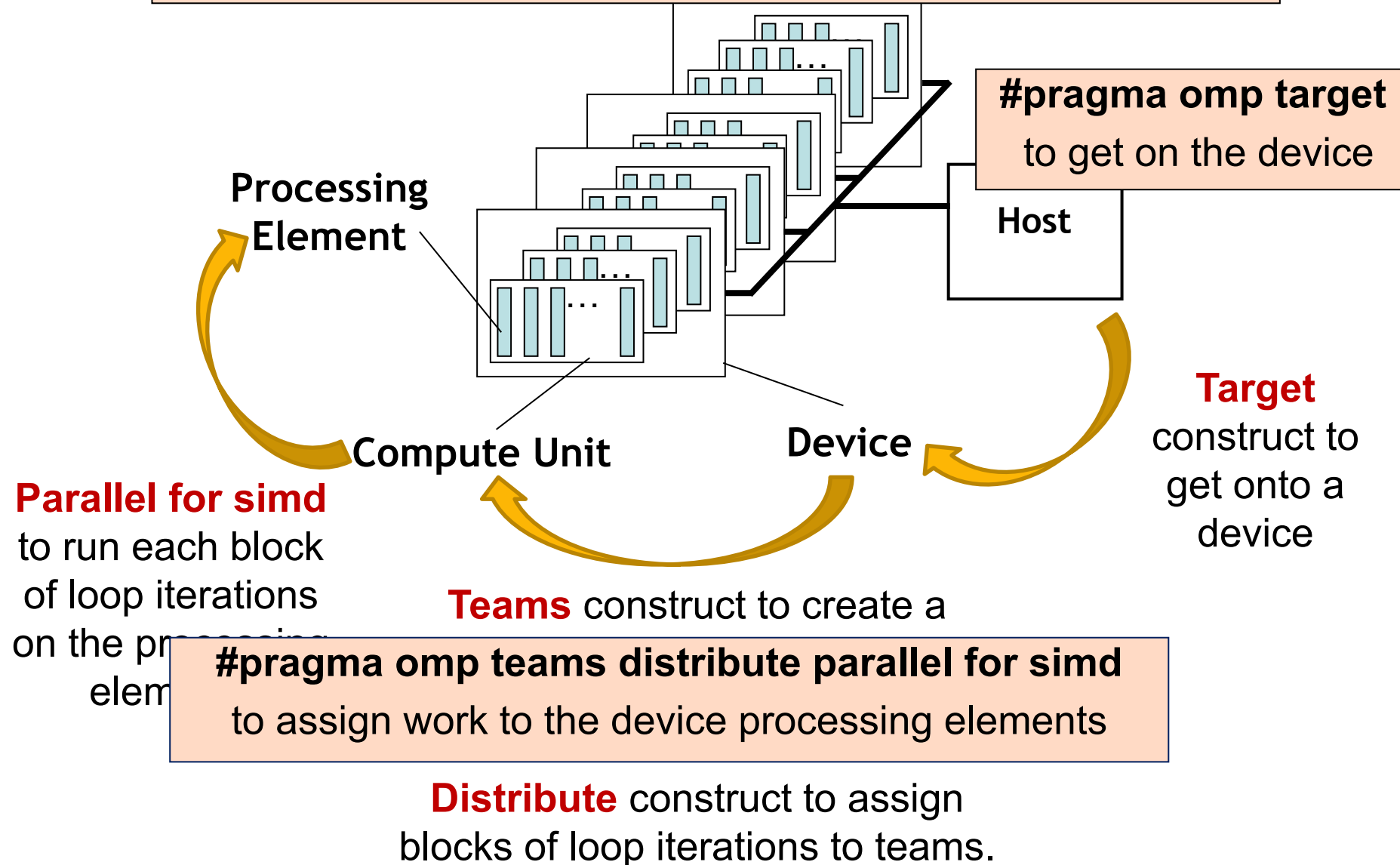
\*the term “master” has been deprecated in OpenMP 5.1 and replaced with the term “primary”.

# Our host/device Platform Model and OpenMP



# Our host/device Platform Model and OpenMP

Typical usage ... let the compiler do what's best for the device:



# Our running example: Jacobi solver

- An iterative method to solve a system of linear equations
  - Given a matrix  $A$  and a vector  $b$  find the vector  $x$  such that  $Ax=b$
- The basic algorithm:
  - Write  $A$  as a lower triangular ( $L$ ), upper triangular ( $U$ ) and diagonal matrix
$$Ax = (L+D+U)x = b$$
  - Carry out multiplications and rearrange
$$Dx = b - (L+U)x \rightarrow x = (b - (L+U)x) / D$$
  - Iteratively compute a new  $x$  using the  $x$  from the previous iteration
$$X_{\text{new}} = (b - (L+U)x_{\text{old}}) / D$$
- Advantage: we can easily test if the answer is correct by multiplying our final  $x$  by  $A$  and comparing to  $b$
- Disadvantage: It takes many iterations and only works for diagonally dominant matrices

# Jacobi Solver

Iteratively update  $x_{\text{new}}$  until the value stabilizes (i.e. change less than a preset TOL)

```
<<< allocate and initialize the matrix A >>>
<<< and vectors x1, x2 and b >>>

while((conv > TOL) && (iters<MAX_ITERS))
{
    iters++;

    for (i=0; i<Ndim; i++){
        xnew[i] = (TYPE) 0.0;
        for (j=0; j<Ndim;j++){
            if(i!=j)
                xnew[i] += A[i*Ndim + j]*xold[j];
        }
        xnew[i] = (b[i]-xnew[i])/A[i*Ndim+i];
    }
}
```

```
// test convergence
conv = 0.0;
for (i=0; i<Ndim; i++){
    tmp = xnew[i]-xold[i];
    conv += tmp*tmp;
}
conv = sqrt((double)conv);

// swap pointers for next
// iteration
TYPE* tmp = xold;
xold = xnew;
xnew = tmp;

} // end while loop
```

# Jacobi Solver (Par Targ, 1/2)

```
while((conv > TOL) && (iters<MAX_ITERS))
{
    iters++;
    #pragma omp target map(tofrom:xnew[0:Ndim],xold[0:Ndim]) \
        map(to:A[0:Ndim*Ndim], b[0:Ndim])
    #pragma omp teams distribute parallel for simd private(i,j)
    for (i=0; i<Ndim; i++){
        xnew[i] = (TYPE) 0.0;
        for (j=0; j<Ndim;j++){
            if(i!=j)
                xnew[i]+= A[i*Ndim + j]*xold[j];
        }
        xnew[i] = (b[i]-xnew[i])/A[i*Ndim+i];
    }
}
```

# Jacobi Solver (Par Targ, 2/2)

```
//
// test convergence
//
conv = 0.0;
#pragma omp target map(to:xnew[0:Ndim],xold[0:Ndim]) \
                    map(tofrom:conv)
#pragma omp teams distribute parallel for simd \
                    private(i,tmp) reduction(+:conv)
for (i=0; i<Ndim; i++){
    tmp = xnew[i]-xold[i];
    conv += tmp*tmp;
}
conv = sqrt((double)conv);
TYPE* tmp = xold;
xold = xnew;
xnew = tmp;
} // end while loop
```

This worked but the performance was awful. Why?

System	Implementation	Ndim = 4096
NVIDIA® K20X™ GPU	Target dir per loop	131.94 secs

Cray® XC40™ Supercomputer running Cray® Compiling Environment 8.5.3.  
Intel® Xeon ® CPU E5-2697 v2 @ 2.70GHz with 32 GB DDR3. NVIDIA® Tesla® K20X, 6GB.

# Data movement dominates!!!

```
while((conv > TOLERANCE) && (iters<MAX_ITERS))
```

```
{ iters++;
```

```
  xnew = iters % s ? x2 : x1;
```

```
  xold  = iters % s ? x1 : x2;
```

Typically over 4000 iterations!

```
#pragma omp target map(tofrom:xnew[0:Ndim],xold[0:Ndim]) \
  map(to:A[0:Ndim*Ndim], b[0:Ndim] )
#pragma omp teams distribute parallel for simd private(i,j)
for (i=0; i<Ndim; i++){
  xnew[i] = (TYPE) 0.0;
  for (j=0; j<Ndim;j++){
    if(i!=j)
      xnew[i]+= A[i*Ndim + j]*xold[j];
  }
  xnew[i] = (b[i]-xnew[i])/A[i*Ndim+i];
}
```

For each iteration, **copy to** device  
 $(3*Ndim+Ndim^2)*sizeof(TYPE)$  bytes

```
// test convergence
```

```
conv = 0.0;
```

```
#pragma omp target map(to:xnew[0:Ndim],xold[0:Ndim]) \
  map(tofrom:conv)
#pragma omp teams distribute parallel for private(i,tmp) reduction(+:conv)
for (i=0; i<Ndim; i++){
  tmp = xnew[i]-xold[i];
  conv += tmp*tmp;
}
```

For each iteration, **copy from** device  
 $2*Ndim*sizeof(TYPE)$  bytes

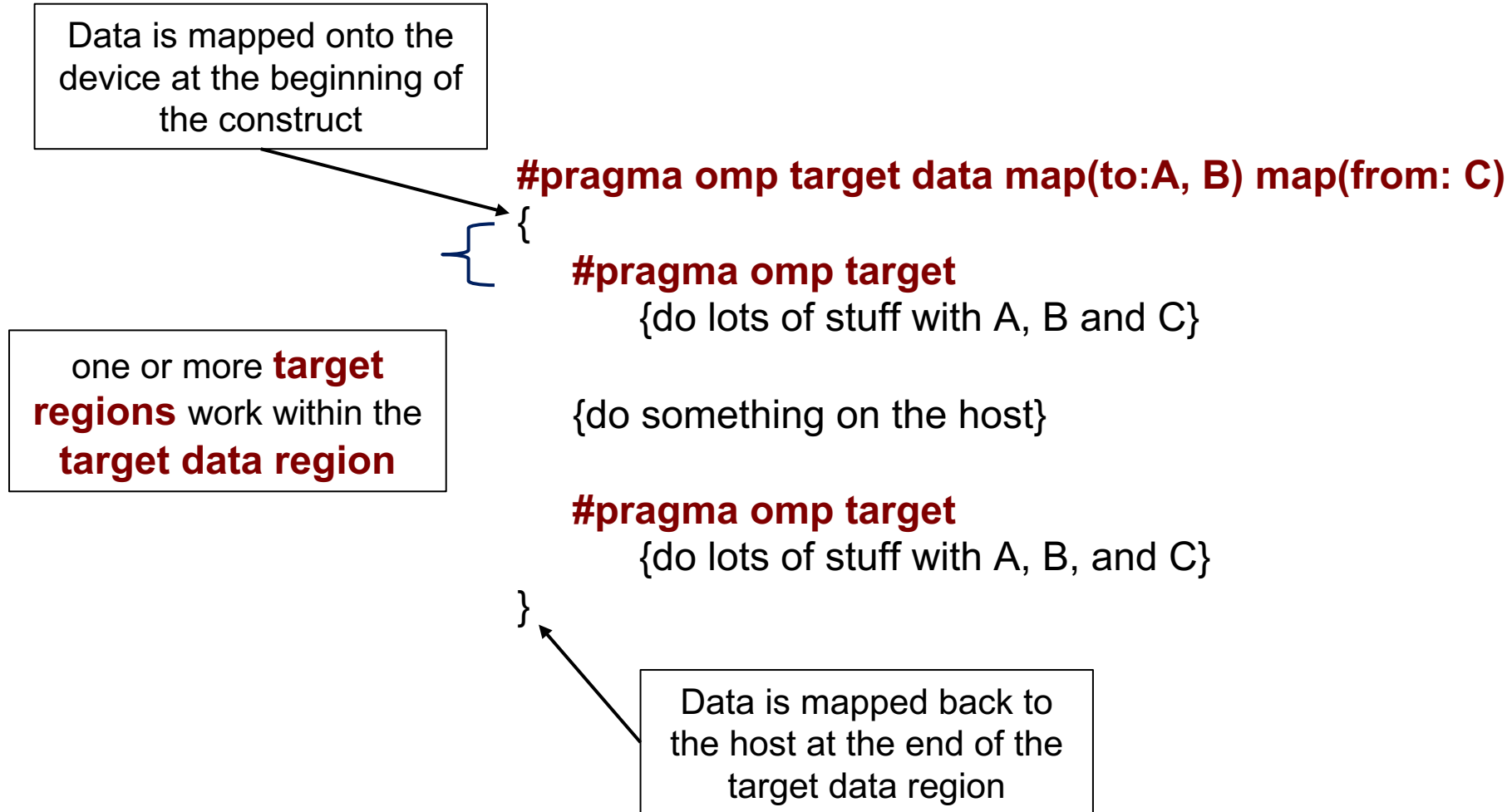
For each iteration, **copy to**  
device  
 $2*Ndim*sizeof(TYPE)$  bytes

```
conv = sqrt((double)conv);
```

```
}
```

# Target data directive

- The **target data** construct creates a target data region  
... use **map** clauses for explicit data management



# Jacobi Solver (Par Target Data, 1/2)

```
#pragma omp target data map(tofrom:x1[0:Ndim],x2[0:Ndim]) \  
                        map(to:A[0:Ndim*Ndim], b[0:Ndim] ,Ndim)
```

```
while((conv > TOL) && (iters<MAX_ITERS))  
    { iters++;
```

```
#pragma omp target
```

```
#pragma omp teams distribute parallel for simd private(j) firstprivate(xnew,xold)
```

```
    for (i=0; i<Ndim; i++){  
        xnew[i] = (TYPE) 0.0;  
        for (j=0; j<Ndim;j++){  
            if(i!=j)  
                xnew[i]+= A[i*Ndim + j]*xold[j];  
        }  
        xnew[i] = (b[i]-xnew[i])/A[i*Ndim+i];  
    }
```

# Jacobi Solver (Par Target Data, 2/2)

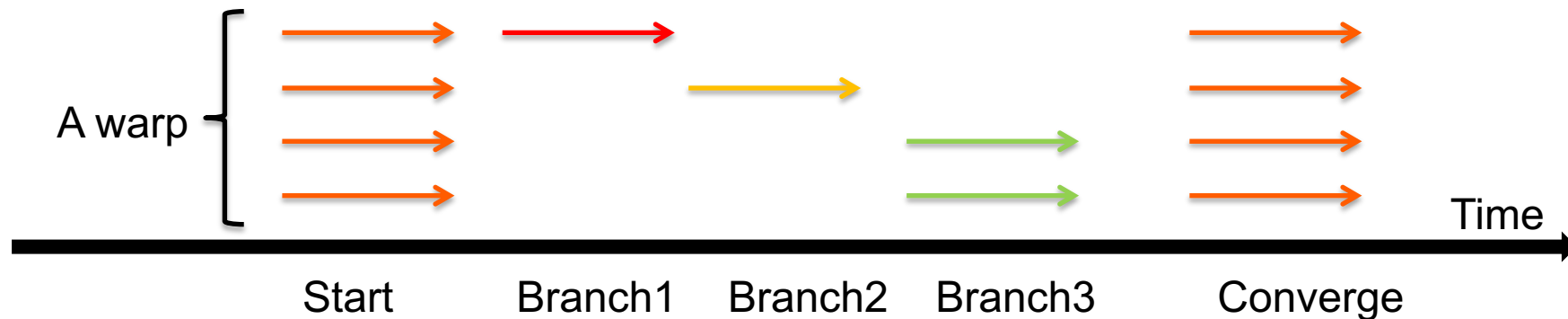
```
// test convergence
conv = 0.0;
#pragma omp target map(tofrom: conv)
#pragma omp teams distribute parallel for simd \
private(tmp) firstprivate(xnew,xold) reduction(+:conv)
    for (i=0; i<Ndim; i++){
        tmp = xnew[i]-xold[i];
        conv += tmp*tmp;
    }
// end target region
conv = sqrt((double)conv);

TYPE* tmp = xold;
xold = xnew;
xnew = tmp;
} // end while loop
```

System	Implementation	Ndim = 4096
NVIDIA® K20X™ GPU	Target dir per loop	131.94 secs
	Above plus target data region	18.37 secs

# Single Instruction Multiple Data

- Individual work-items of a warp start together at the same program address
- Each work-item has its own instruction address counter and register state
  - Each work-item is free to branch and execute independently
  - Supports the SPMD pattern.
- Branch behavior
  - Each branch will be executed serially
  - Work-items not following the current branch will be disabled



# Branching

## Conditional execution

```
// Only evaluate expression  
// if condition is met  
if (a > b)  
{  
    acc += (a - b*c);  
}
```

## Selection and masking

```
// Always evaluate expression  
// and mask result  
temp = (a - b*c);  
mask = (a > b ? 1.f : 0.f);  
acc += (mask * temp);
```

# Coalescence

- Coalesce - to combine into one
- Coalesced memory accesses are key for high bandwidth
- Simply, it means, if thread  $i$  accesses memory location  $n$  then thread  $i+1$  accesses memory location  $n+1$
- In practice, it's not quite as strict...

```
for (int id = 0; id < size; id++)
{
    // ideal
    float val1 = memA[id];

    // still pretty good
    const int c = 3;
    float val2 = memA[id + c];

    // stride size is not so good
    float val3 = memA[c*id];

    // terrible
    const int loc =
        some_strange_func(id);

    float val4 = memA[loc];
}
```

# Jacobi Solver (Target Data/branchless/coalesced mem, 1/2)

```
#pragma omp target data map(tofrom:x1[0:Ndim],x2[0:Ndim]) \
    map(to:A[0:Ndim*Ndim], b[0:Ndim] ,Ndim)
while((conv > TOL) && (iters<MAX_ITERS))
{
    iters++;
    #pragma omp target
        #pragma omp teams distribute parallel for simd private(j)
    for (i=0; i<Ndim; i++){
        xnew[i] = (TYPE) 0.0;
        for (j=0; j<Ndim;j++){
            xnew[i]+= (A[j*Ndim + i]*xold[j])*((TYPE) (i != j));
        }
        xnew[i] = (b[i]-xnew[i])/A[i*Ndim+i];
    }
}
```

We replaced the original code with a poor memory access pattern

$xnew[i] += (A[i*Ndim + j]*xold[j])$

With the more efficient

$xnew[i] += (A[j*Ndim + i]*xold[j])$

# Jacobi Solver (Target Data/branchless/coalesced mem, 2/2)

```
//
// test convergence
conv = 0.0;
#pragma omp target map(tofrom: conv)
#pragma omp teams distribute parallel for simd \
    private(tmp) reduction(+:conv)
for (i=0; i<Ndim; i++){
    tmp = xnew[i]-xold[i];
    conv += tmp*tmp;
}
conv = sqrt((double)conv);
TYPE* tmp = xold;
xold = xnew;
xnew = tmp;
} // end while loop
```

System	Implementation	Ndim = 4096
NVIDIA® K20X™ GPU	Target dir per loop	131.94 secs
	Above plus target data region	18.37 secs
	Above plus reduced branching	13.74 secs
	Above plus improved mem access	7.64 secs