# Introduction to OpenMP
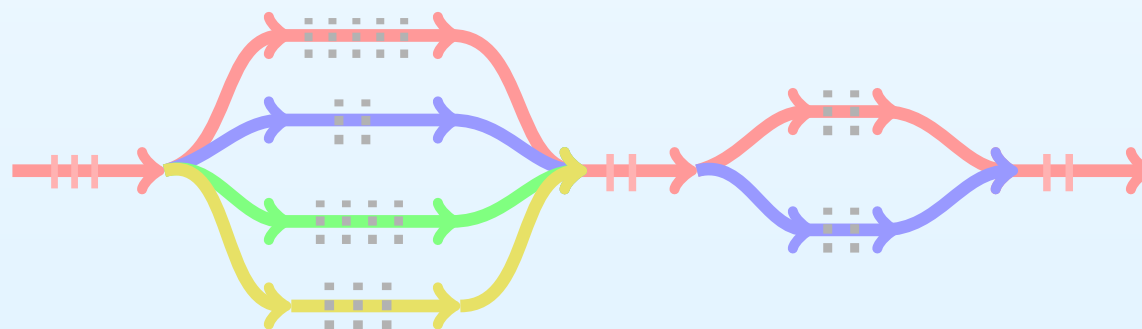
Xiaoxu Guan

*High Performance Computing, LSU*

*April 6, 2016*

- Overview of Parallel Computing

- Overview of Parallel Computing
- Parallel Programming on **Shared**-Memory and **Distributed**-Memory Machines

# Overview

- Overview of Parallel Computing
- Parallel Programming on **Shared**-Memory and **Distributed**-Memory Machines
- Introduction to OpenMP
  - Prerequisite for Parallel Computing;
  - Constructs for **Parallel Execution**;
  - **Data Communications**;
  - **Synchronization**;

# Overview

- Overview of Parallel Computing
- Parallel Programming on **Shared**-Memory and **Distributed**-Memory Machines
- Introduction to OpenMP
  - Prerequisite for Parallel Computing;
  - Constructs for **Parallel Execution**;
  - **Data Communications**;
  - **Synchronization**;
- OpenMP Programming: Directives/Pragams, Environment Variables, and Run-time Libraries
  - Variables Peculiar to OpenMP Programming;
  - Loop Level Parallelism;
  - Non-Loop Level Parallelism;

# Overview

- Overview of Parallel Computing
- Parallel Programming on **Shared**-Memory and **Distributed**-Memory Machines
- Introduction to OpenMP
  - Prerequisite for Parallel Computing;
  - Constructs for **Parallel Execution**;
  - **Data Communications**;
  - **Synchronization**;
- OpenMP Programming: Directives/Pragams, Environment Variables, and Run-time Libraries
  - Variables Peculiar to OpenMP Programming;
  - Loop Level Parallelism;
  - Non-Loop Level Parallelism;
- Summary and Further Reading

- Why parallel or concurrency computing?
- Goes beyond the single-core capability (memory and flops per unit time), and therefore increase performance;
- Reduces wall-clock time, and saves energy;
- Finish those impossible tasks in my lifetime;
- Handles larger and larger-scale problems;
- **There is no free lunch, however!**
- Different techniques other than serial coding are needed;
- Effective parallel algorithms in terms of performance;
- Increasing flops per unit time is one of our endless goals in the HPC community;
- Think in parallel;
- Start parallel programming as soon as possible;
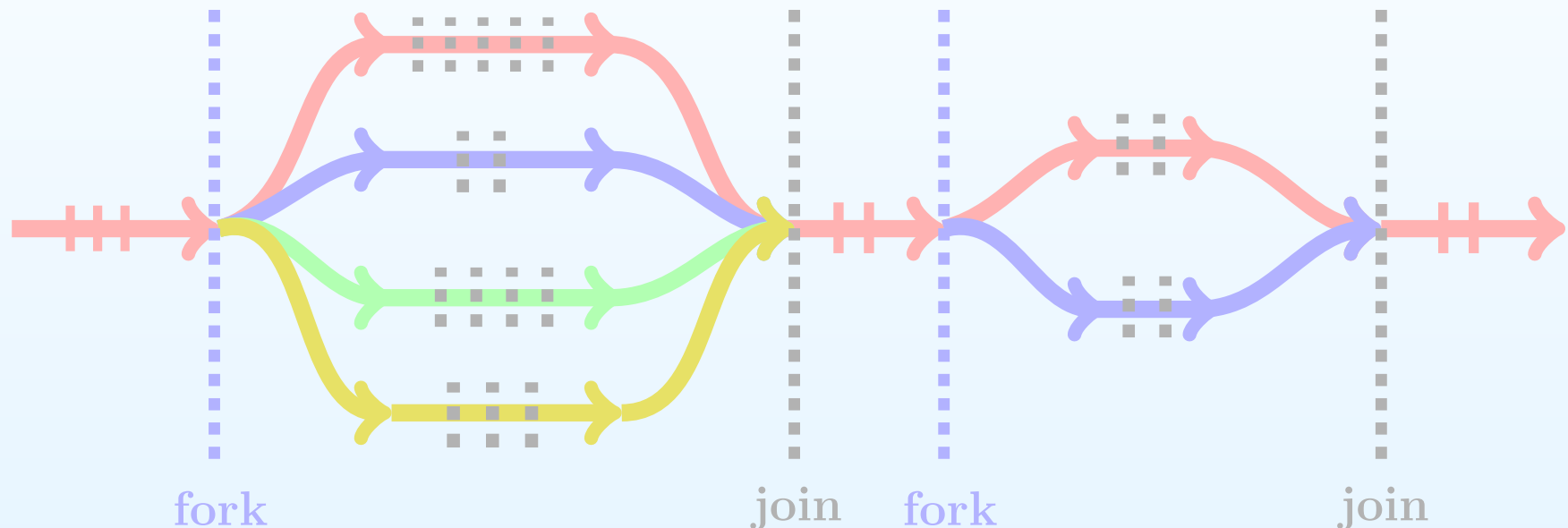
# Parallel programming

- Parallel programming environment;
  - Essential language extensions to the existing language (Fortran 95);
  - New constructs for directives/pragmas to existing serial programs (**OpenMP** and HPF);
  - Run-time libraries that support data communication and synchronization (**MPI** and Pthreads);
- **OpenMP** stands for **Open M**ulti-**P**rocessing (API);
- **OpenMP** is one of the directives/pragmas approaches that support parallelism on **shared** memory systems;
- **OpenMP** is supported by Fortran, and C/C++;
- **OpenMP** allows us to start from a serial code and provides an incremental approach to express parallelism;

# The "Three Knights" in OpenMP

**(1) Directives/pragmas** need to express parallelism;
**(2) Run-time libraries** can dynamically control or change code execution at run-time;
**(3) Environment variables** specify the run-time options;

- How does OpenMP achieve parallel computing?
  - Specify parallel execution – parallel constructs allowing parallel execution;
  - Data communication – data constructs for communication among threads;
  - Synchronization – synchronization constructs;
- OpenMP directives/pragmas:

  Fortran: `!$omp`, `c$omp`, or `*$omp [clauses]`

  C/C++: `#pragma [clauses]`

# Parallel execution

- Constructs for parallel execution: OpenMP starts with a single thread, but it supports the directives/pragmas to spawn multiple threads in a fork-join model;



fork          join    fork        join

- OpenMP `do` and `parallel` directives;
- OpenMP also allows you to change the number of threads at run-time;

# Data communication

- When multiple threads were spawned, each thread was assigned to a unique thread ID from $0$ to $N-1$. Here $N$ is the total number of threads;

- The key point is that there are three types of variables: `private`, `shared`, and `reduction` variables;

- At run-time, there is always a common region in global memory that allows all threads to access, and this memory region is used to store all `shared` variables;

- Each thread was also assigned a private memory region to store all `private` variables; Thread `a` cannot access the private variables stored in the private memory space for thread `b`;

- Data Communications are achieved through `read` and `write` operations among the threads;

LSU
CENTER FOR COMPUTATION
& TECHNOLOGY

**Information Technology Services**
LSU HPC Training Series, Spring 2016

LONI

p. 7/44

# Synchronization

- In OpenMP, synchronization is used to **(1)** control the access to **shared** variables and **(2)** coordinate the workflow;
- Event and mutual exclusion synchronization;
- **Event synchronization** includes **barrier** directives, which are either explicit or implicit; a thread has to wait until all threads reach the same point;
- **Mutual exclusion** is supported through **critical**, **atomic**, **single**, and **master** directives. All these are used to control how many threads, which thread, or when a thread can execute a specified code block or modify shared variables;
- Be careful with synchronization!

# Compile OpenMP code

**LSU** INFORMATION TECHNOLOGY SERVICES

- Compiler options that enable OpenMP directive/pragmas:

| Compiler | Fortran | C | C++ |
| --- | --- | --- | --- |
| Intel | ifort **-openmp** | icc **-openmp** | icpc **-openmp** |
| PGI | pgf90 **-mp** | pgcc **-mp** | pgCC **-mp** |
| GCC | gfortran **-fopenmp** | gcc **-fopenmp** | g++ **-fopenmp** |

- If the above flags are left out, OpenMP code is compiled as serial code (except Intel compilers but `-openmp-stubs` needed);
- Load modules on the HPC or LONI machines:

`$ module load [package name]`

`$ soft add [+package name] (resoft)` # intel, pgi, or gcc.

- Set up an environment variable:

`$ export OMP_NUM_THREADS=[number of threads]`

# Loop level parallelism

# First OpenMP "Hello World!" in Fortran and C

Fortran (hello.f90)

```fortran
1      program hello_world
2      implicit none
3
4      integer ::  id, omp_get_thread_num
5
6   !$omp parallel
7      id = omp_get_thread_num()
8      write(*,'(1x,a,i3)') "Hello World! from", id
9   !$omp end parallel
10
11     end program hello_world
```

```
$ export OMP_NUM_THREADS=20
    # for instance, on SuperMIC in bash shell
$ ifort -o hello hello.f90 -openmp
```

# First OpenMP "Hello World!" in Fortran and C

C (hello.c)

```c
 1    #include <stdio.h>
 2    #include <stdlib.h>
 3    #include <omp.h>
 4
 5    int main( ) {
 6    int id;
 7
 8    #pragma omp parallel {
 9    id = omp_get_thread_num();
10    printf("Hello World! from %3d\n", id);
11                                    }
12    }
```

```
$ export OMP_NUM_THREADS=20
    # for instance, on SuperMIC in bash shell
$ icc -o hello hello.c -openmp
```

**Information Technology Services**
LSU HPC Training Series, Spring 2016

LSU
CENTER FOR COMPUTATION & TECHNOLOGY

LONI

p. 12/44

# Loop-level parallelism

- Loop-level parallelism is one of the fine-grained approaches supported by OpenMP;

- `parallel do` directives in Fortran and `parallel for` pragmas in C/C++;

```fortran
1  !$omp parallel do [clauses]                Fortran
2     do i = imin, imax, istep
3        loop body ...
4     end do
5  [!$omp end parallel do]
```

```c
1  #pragma omp parallel for [clauses]         C/C++
2     for (i = imin; i < imax; increment_expr)
3     {
4        loop body ...;
5     }
```

# Loop-level parallelism

```fortran
1  !$omp parallel [clauses]            Fortran
2  !$omp do [clauses]
3      do i = imin, imax, istep
4         loop body ...
5      end do
6  !$omp end do
7  !$omp end parallel
```

```c
1  #pragma omp parallel [clauses]      C/C++
2    {
3  #pragma omp for [clauses]
4      for (i = imin; i < imax; increment_expr)
5      {
6         loop body ...;
7      }
8    }
```

# Loop-level parallelism

- How about nested multiple loops? Where do we add `parallel for`, right above outer loop or inner loop?

```
1  for (i = imin; i < imax; increment_i)     C/C++
2  {                                          (inner loop)
3  #pragma omp parallel for
4      for (j = jmin; j < jmax; increment_j)
5      { loop body ...; }
6  }
```

```
1  #pragma omp parallel for                   C/C++ (outer loop)
2  for (i = imin; i < imax; increment_i)
3  {
4      for (j = jmin; j < jmax; increment_j)
5      { loop body ...; }
6  }
```

# More words on parallel loops

- OpenMP only supports Fortran `do` loops and C/C++ `for` loops that the number of loop iterations is known for at run-time;
- However, it doesn't support other loops, including `do-while` and `repeat-until` loops in Fortran and `while` loops and `do-while` loops in C/C++. In these cases, the trip count of loop is unknown before entering the loop;
- Loop body has to follow `parallel do` or `parallel for` immediately, and nothing in between them!
- There is an implicit `barrier` at the end of `parallel do` or `for` loops;
- All loops must have a single entry point and single exit point. We are **not** allowed to jump into a loop or branch out of a loop;

# How to control loops?

- Once we entered the parallel region, for some variables multiple threads need to use the same named variables, but they store different values at different memory locations; these variables are called **private** variables;

- This leads to the fact that all private variables are **undefined** or **uninitialized** before entry and after exit from parallel regions;

- The **shared** variables are also necessary to allow data communication between threads;

- **Default** scopes for variables: By default all the variables are considered to be **shared** in parallel regions, unless they are explicitly declared as **private**, **reduction**, or **other** types;

- Remember Fortran and C/C++ may have different settings regarding default rules;

- Let's see how we can do it, for instance, in parallel loops;
- OpenMP provides a means to change the default rules;
- Clauses `default(none)`, `default(private)`, and `default(shared)` in Fortran;
- But only `default(none)` and `default(shared)` in C/C++;

```fortran
1  ALLOCATE( da(1:nsize), db(1:nsize) )
2  !$omp parallel do default(none),     &
3  !$omp private(i,temp),               &
4  !$omp shared(imin,imax,istep,scale,da,db)
5     do i = imin, imax, istep
6        temp = scale * da(i)
7        da(i) = temp + db(i)
8     end do
9  !$omp end parallel do
```

Fortran

# How to control loops?

- OpenMP `reduction` operations;
- The reduction variable is very special that it has both characters of private and shared variables;
- Compiler needs to know what type of operation is associated with the `reduction` variable; `operation = +, *, max, min`, etc;
- `reduction(operation : variables_list)`

```fortran
1 ALLOCATE( da(1:nsize) )                    Fortran
2       prod = 1.0d0
3 !$omp parallel do default(none), private(i),  &
4 !$omp reduction(* : prod)
5    do i = imin, imax, istep
6       prod = prod * da(i)
7    end do              What happens if we compile it?
8 !$omp end parallel do
```

- Two special "private" variables: `firstprivate` and `lastprivate`; they are used to initialize and finalize some `private` variables;

- `firstprivate`: upon entering a `parallel do/for`, the private variable for each **slave** thread has a copy of the **master** thread's value;

- `lastprivate`: upon exiting a `parallel do/for`, no matter which thread executed the **last** iteration (sequential), the private variable was copied back to the **master** thread;

- Why do we need them? **(1)** all private variables are **undefined** outside of a parallel region, **(2)** they provide a simply way to exchange data to some extent through these special **private** variables;

- In a parallel region, a given variable can only be one of `private`, `shared`, or `reduction`, but it can be both of `firstprivate` and `lastprivate`;

```
 1  double ashift = shift ;                        C/C++
 2  #pragma omp parallel for default(none),           \
 3               firstprivate(ashift), shared(a),   \
 4               private(i)
 5   {
 6     for (i = imin; i <= imax; ++i)
 7        {
 8        ashift = ashift + (double) i ;
 9        a[i] = a[i] + ashift ;
10        }
11   }
```

# How to control loops?

- Exception of the default rules: Fortran and C/C++ behave differently;
- The index in a parallel loop is always **private**. The index in a sequential loop is also **private** in Fortran, but is **shared** in C by default!
- Is the following code correct?
- Has the loop j been parallelized?

```
1  #pragma omp parallel for                        C/C++
2  for (i = imin; i <= imax; ++i)
3    {
4      for (j = jmin; j <= jmax; ++j)
5        a[i][j] = (double) (i + j) ;
6    }
```

- Do we have the same issues in the Fortran version?

- Exception of the default rules. Fortran and C/C++ behave differently;
- The index in a parallel loop is always **private**. The index in a sequential loop is also **private** in Fortran, but is **shared** in C by default!
- Is the following code correct?
- Has the loop j been parallelized?

```
1  #pragma omp parallel for private(i,j)        C/C++
2  for (i = imin; i <= imax; ++i)
3    {
4      for (j = jmin; j <= jmax; ++j)
5        a[i][j] = (double) (i + j) ;
6    }
```

- Do we have the same issues in the Fortran version?

- Parallelize multiple nested loops;
- The `collapse(n)` ($n \geqslant 1$) for nested parallel loops;
- Each thread takes a chunk of the `i` loop and a chunk of the `j` loop at the same time;
- No statements in between;

```
1  #pragma omp parallel for private(i,j), \    C/C++
2            collapse(2)
3  for (i = imin; i <= imax; ++i)
4     {
5        for (j = jmin; j <= jmax; ++j)
6         a[i][j] = (double) (i + j) ;
7     }
```

# Restrictions on parallel loops

- Not all loops are parallelizable. What can we do?
- Think parallely and change your algorithms;
- We have to maintain the correctness of the results;
- One of the common mistakes is `data race`;

```
1  #pragma omp parallel for                        C/C++
2    {
3      for (i = imin; i <= imax; ++i)
4       r[i] = r[i] + r[i-1] ;
5    }
```

- **Data race** means that in a parallel region, the same memory location is referred by **two** or **more** statements, and at least one of them is a **write** operation;
- Data race requires more attention and might lead to incorrect results!

# Restrictions on parallel loops

- A closer look at at the data race: let's run it on `2` threads and assume that `r[0]=a`; `r[1]=b`; `r[2]=c`; and `imin=1`; `imax=2`;
- Note `r[1]` is referred twice, and thus we have two scenarios:

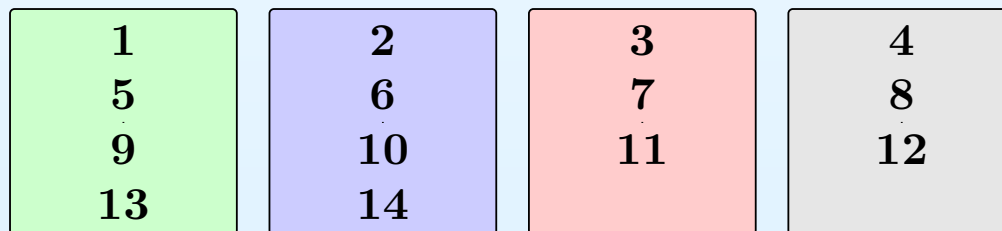| if `thread` `0` finished first | | if `thread` `1` finished first | |
|---|---|---|---|
| `thread` `0` | `thread` `1` | `thread` `1` | `thread` `0` |
| `i = 1` | `i = 2` | `i = 2` | `i = 1` |
| `r[0]=a` | | `r[1]=b` | `r[0]=a` |
| `r[1]=a+b` | `r[2]=a+b+c` | `r[2]=b+c` | `r[1]=b+a` |
| | time | | time |

- OpenMP standard does not guarantee which thread finishes first or later;

# How to control loops again?

- OpenMP supports three loop schedulings as clauses:
  `static`, `dynamic`, and `guided` in the code, plus `run-time`
  scheduling;

- `schedule( type[, chunk_size] )`

For `static`, if `chunk_size` is given, loop iterations are divided
into multiple blocks and each block contains `chun_size`
iterations. The iterations will be assigned to threads in a
round-robin fashion. If `chunk_size` is not present, the loop
iterations will be (nearly) evenly divided and assigned to each
thread.

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

14 iterations
on 4 threads
in round-robin fashion

# How to control loops again?

- For `dynamic`, if `chunk_size` is given, the partition is almost the same as those of `static`. The difference is that with `static`, the mapping between loop iterations and threads are done during **compilation**, while for `dynamic`, it will be done at **run-time** (therefore, more potentially overhead); if `chunk_size` is not present, then it was set to `1`.
- The `guided` scheduling means the `chunk_size` assigned to threads decreases exponentially;
- Run-time scheduling: set the environment variable `OMP_SCHEDULE`;
- `$ export OMP_SCHEDULE=10`, for instance;
- Each scheduling has its own pros and cons, so be careful with `chunk_size` and potential overhead;

# Non-loop-level parallelism

# Parallel regions

- In addition to `parallel do` or `for`, most importantly OpenMP supports the parallelism beyond loop levels;

```
1  !$omp parallel [clauses]          Fortran
2      code block
3  !$omp end parallel
```

```
1  #pragma omp parallel [clauses]    C/C++
2      { code block ; }
```

- Each thread in the parallel team executes the same block of code, but with different data;

- In `parallel` directives, **clauses** include:
`private(list)`, `shared(list)`, `reduction(operation : list)`, `default(none | private | shared)`, `if(logical operation)`, `copyin(list)`;

# Any differences?

```fortran
1 !$omp parallel
2    id = omp_get_thread_num()
3    write(*,*) "Hello World!  from ", id
4 !$omp end parallel
```
Fortran

```fortran
1 !$omp parallel
2 do k = 1, 5
3    id = omp_get_thread_num()
4    write(*,*) "Hello World!  from ", id, k
5 end do
6 !$omp end parallel
```
Fortran

```fortran
1 !$omp parallel do
2 do k = 1, 5
3    id = omp_get_thread_num()
4    write(*,*) "Hello World!  from ", id, k
5 end do
6 !$omp end parallel do
```
Fortran

# Global variables in OpenMP

- In addition to **automatic** or **static** variables in Fortran and C/C++, we also need **global** variables;
- `Command blocks` or `modules` in Fortran, while `globals` in C/C++, and we might have issues with private variables;
- **Global**/**local** variables between different code units for a given thread;
- **Private**/**shared** variables between multiple threads in a given code unit;
- The default data scoping rule is only apply to its **lexical** region, and all rest are **shared**; How can we make **private** variables "propagate" to **other** code units?
- OpenMP introduced the `threadprivate` directive to solve data scoping issues;
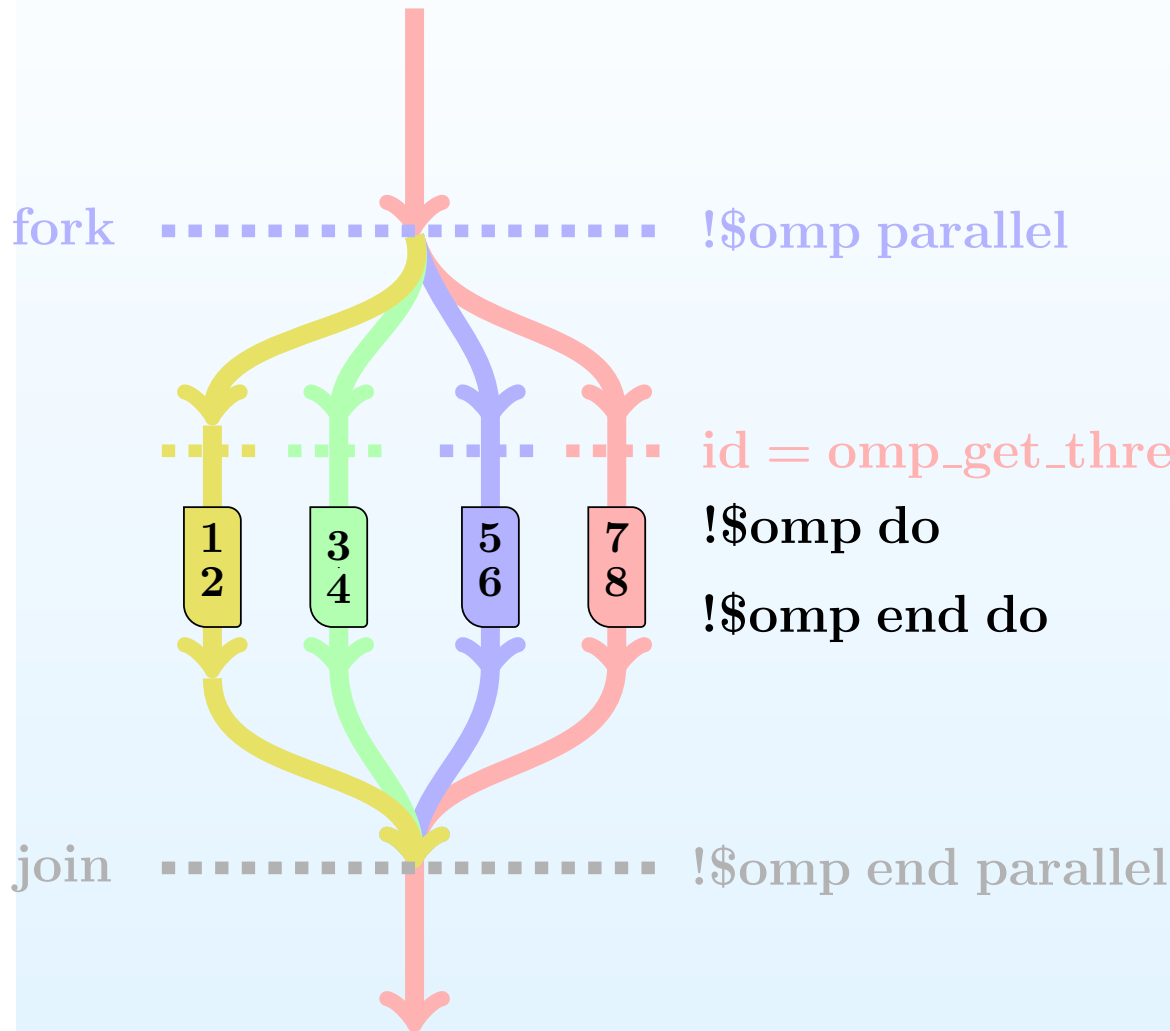
# Global variables in OpenMP

- `!$omp threadprivate (list_common_variables)` in Fortran;
- `#pragma omp threadprivate (list_variables)` in C/C++;
- We have **global** but **private** variables;
- The `threadprivate` variables are special private variables; thus thread `a` cannot access the `threaprivate` variables stored on thread `b`;
- The `threadprivate` variables persist from one parallel region to another, because they are globals;
- Furthermore, OpenMP supports the `copyin (list)` clause to initialize global variables on slave threads to be the values on the master thread;
- `#pragma omp parallel copyin (a,b,c) { code block; }`
- Sounds familiar with the Intel `Xeon Phi` programming?

# Work-sharing directives

```fortran
 1  program mapping                          Fortran
 2  implicit none
 3  integer ::  i,id,nothread,  &
 4           omp_get_thread_num, omp_get_num_threads
 5
 6  !$omp parallel private (k,id), shared(nothread)
 7     id = omp_get_thread_num()
 8     nothread = omp_get_num_threads()
 9  !$omp do
10   do k = 1, 40
11   write(*,'(1x,2(a,i4))') "id = ",id, " k = ",k
12   end do
13  !$omp end do [nowait]
14  !$omp end parallel
15  end program mapping
```

# Work-sharing directives

The point is that `!$omp do` directive does not spawn threads. Instead, only `!$omp parallel` spawns multiple threads!
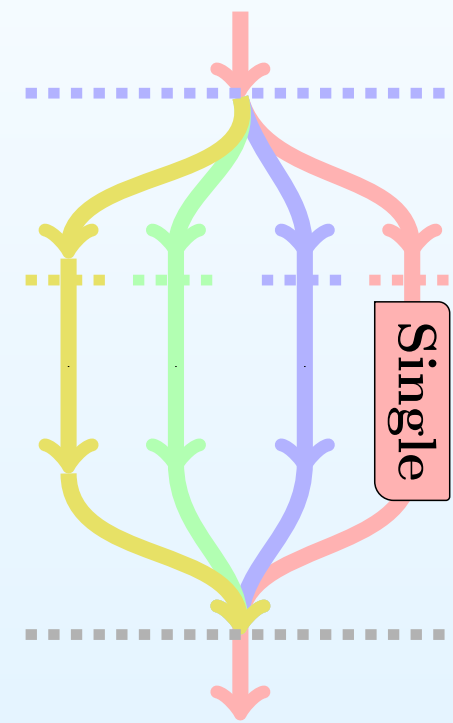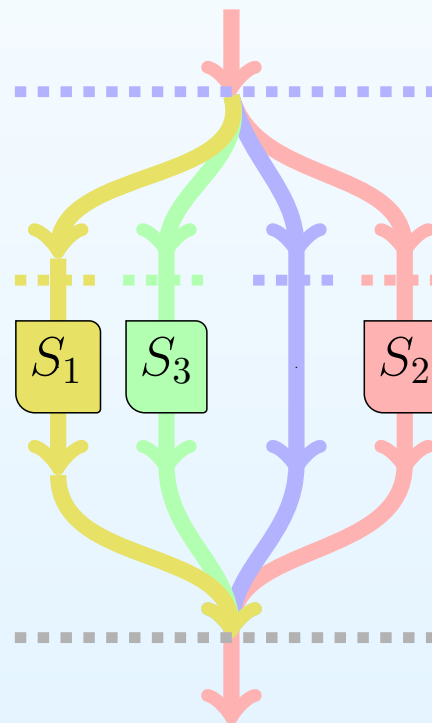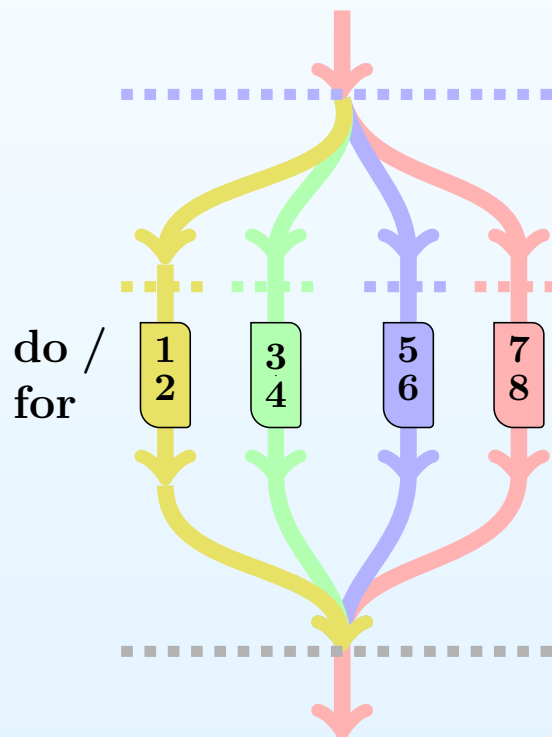
fork ............ !$omp parallel

id = omp_get_thread_num()

| 1 2 | 3 4 | 5 6 | 7 8 |

!$omp do

!$omp end do

join ............ !$omp end parallel

`!$omp do` needs to be embedded in an existing parallel region.

# Work-sharing directives

- Work-sharing constructs do **not** spawn multiple threads; they need to be **embedded** in a parallel region; if not, only one thread will run work-sharing constructs;

- There is an **implicit** barrier at the **end** of a work-sharing construct, but no implicit barrier upon the entry to it;

- **Three** work-sharing constructs:

```
!$omp do                    #pragma for
!$omp section               #pragma section
!$omp single                #pragma single
```

- A given thread may work on zero, one or more `omp sections`; but only one thread runs `omp single` at a given time;

- Be sure there are no data dependencies between `sections`;

- All threads must encounter the same workflow (though it may or may not execute the same code block at run-time);

# Work-sharing directives



`!$omp do`
`#pragma for`

`!$omp section`
`#pragma section`

`!$omp single`
`#pragma single`

do / for: 1 2, 3 4, 5 6, 7 8

$S_1$  $S_3$  $S_2$

Single

# Work-sharing directives

```c
 1  #include <omp.h>
 2  #define nsize 500
 3  main() { int i, j, k ;
 4  double a[nsize], b[nsize], c[nsize] ;
 5  for (k = 0; k <= nsize, ++k) {
 6  a[k] = (double) k; b[k] = a[k]; c[k] = 0.5*a[k];}
 7
 8  #pragma omp parallel {
 9      #pragma omp sections {
10      #pragma omp section { code block_1; }
11      #pragma omp section { code block_2; }
12      #pragma omp section { code block_3; }
13                              }
14                  }
15          }
```

# Synchronization

# Synchronization

- OpenMP provides the constructs for **mutual exclusion**:

  `critical`, `atomic`, `master`, `barrier`, and `run-time` routines;
  `!$omp critical [name] code block`
  `!$omp end critical [name]` in Fortran;
  `#pragma omp critical [name] {code block;}` in C/C++;

- [name] is an optional; But in Fortran, name here should be unique (cannot be the same as those of `do` loops or `if/endif` blocks, etc);

- At a given time, `critical` only allows **one** thread to run it, and all other threads also need to go through `critical` section, but have to wait to enter `critical` section;

- Don't jump into or branch out of a critical section;

- It is useful in a parallel region;

- It might have a tremendous impact on code performance;

- The other way to think of `reduction` variable (say addition):

```fortran
 1        tsum = 0.0d0 ; nsize = 10000         Fortran
 2  !$omp parallel private(temp), shared(tsum,nsize)
 3        temp = 0.0d0
 4  !$omp do
 5      do i = 1, nsize
 6          temp = temp + array(i)
 7      end do
 8  !$omp end do
 9
10  !$omp critical
11          tsum = tsum + temp
12  !$omp end critical
13
14  !$omp end parallel
```

# Synchronization

- Using `atomic` to protect a shared variable:

```c
1  #include <omp.h>
2  #include <stdio.h>
3  #include <stdlib.h>
4  #define nsize 1000
5  int main () {
6  int i; double x = 0.0, answer;
7  #pragma omp parallel for private(i) shared(x) {
8  for (i = 0; i < nsize; ++i) {
9      #pragma omp atomic
10     x += (double) i; }           }
11 answer = (double) 0.5*(nsize-1)*nsize;
12 printf("%f\n", x);
13 printf("correct answer is %f\n", answer);
14 }
```

# OpenMP run-time libraries

- `integer omp_get_num_threads()`
  `int omp_get_num_threads(void)`
  **#** No. of threads in the current collaborating parallel region;

- `integer omp_get_thread_num()`
  `int omp_get_thread_num(void)`
  **#** Return the thread IDs in a parallel team;

- `integer omp_get_num_procs()`
  `int omp_get_num_procs(void)`
  **#** Get the number of "processors" available to the code;

- `call omp_set_num_threads(num_threads)`
  `omp_set_num_threads(num_threads)`
  **#** Set number of threads to be `num_threads` for the following parallel regions;

- `omp_get_wtime()` **#** Measure elapsed wall-clock time (in seconds) relative to an arbitrary reference time;

- OpenMP loop-level, non-loop level parallelism, synchronization, and run-time libraries;
- How to protect **shared** variables; pay attention to them and synchronization; data races;
- Global and local variables in OpenMP programming (**global private** variables);
- Develop a defensive programming style;

**Parallel Programming in OpenMP**, R. Chandra et al. (Morgan Kaufmann Publishers, 2001).

# Questions?

`sys-help@loni.org`

LSU
CENTER FOR COMPUTATION & TECHNOLOGY