

CME 213, ME 339 Spring 2023

Introduction to parallel computing using MPI, openMP, and CUDA

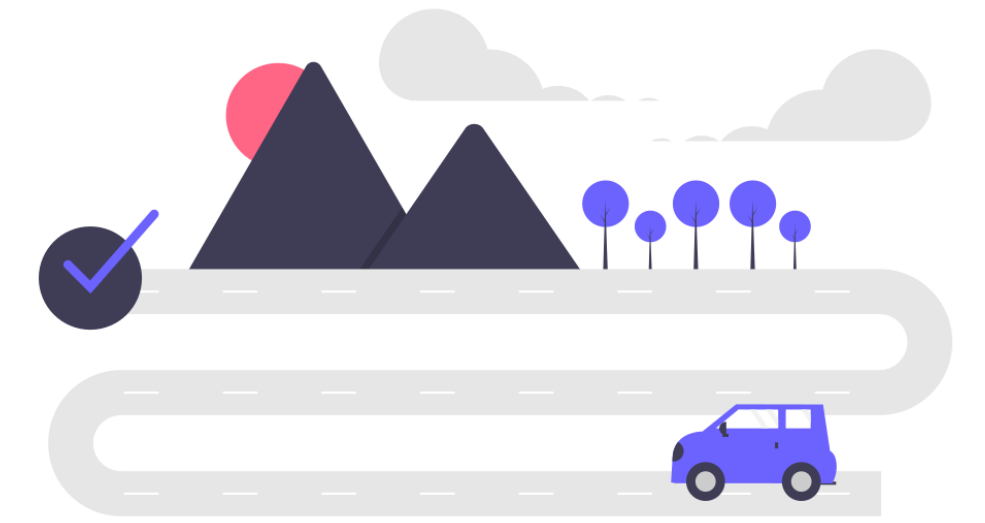
Stanford University

Eric Darve, ICME, Stanford

“If debugging is the process of removing bugs, then programming must be the process of putting them in.” (Edsger W. Dijkstra)



Recap



- OpenMP: loop scheduling: static, dynamic, guided
- Optimization: nowait, collapse
- Data affinity, NUMA, first-touch policy: a thread that reads from/writes to a memory location (computation) should be the first one to write to that location (initialization).
- Data sharing: shared, private, first/lastprivate, threadprivate, reduction, copyin/copyprivate.

Data sharing, visually

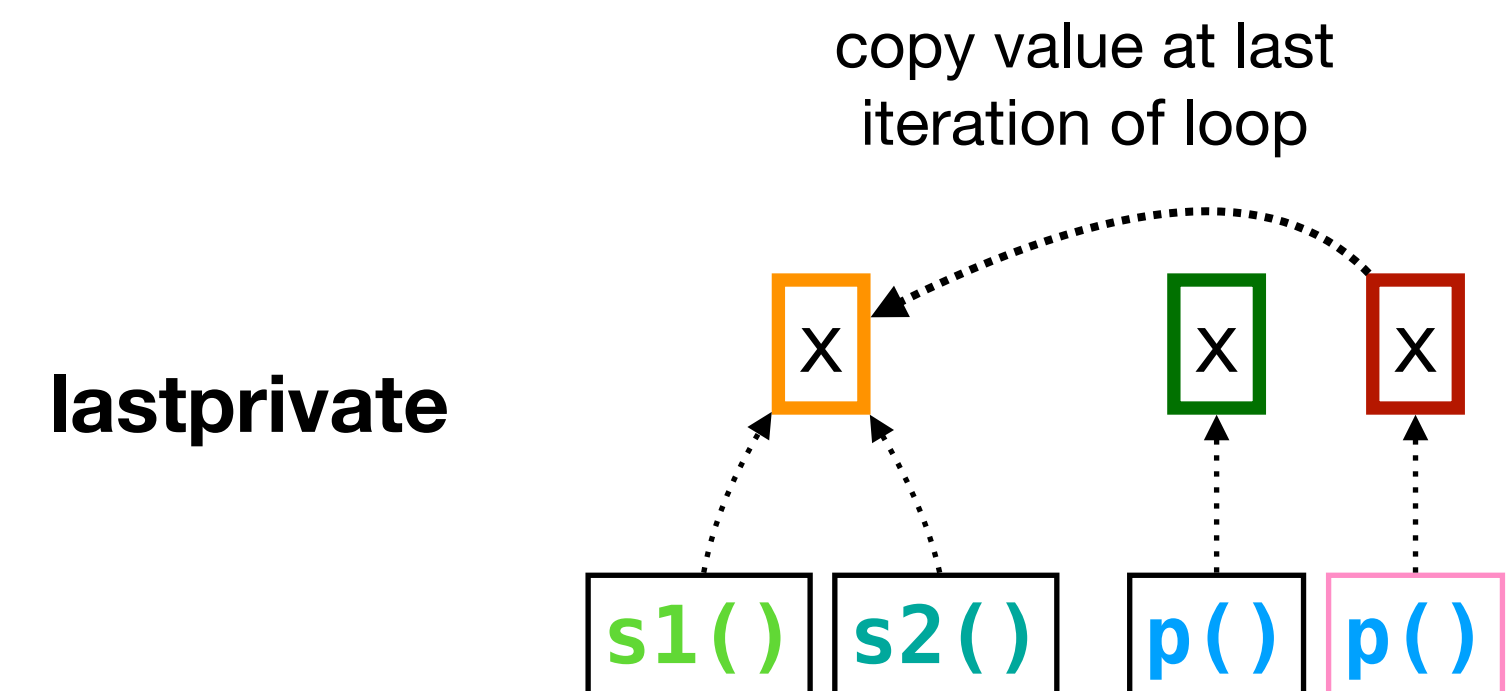
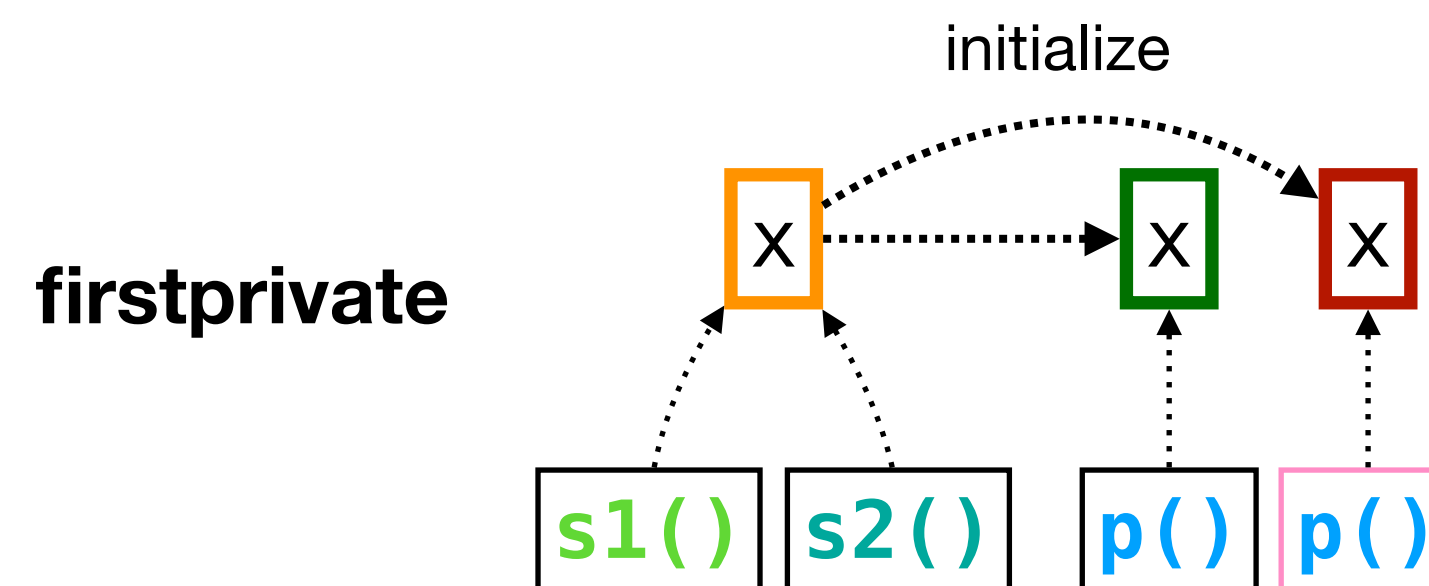
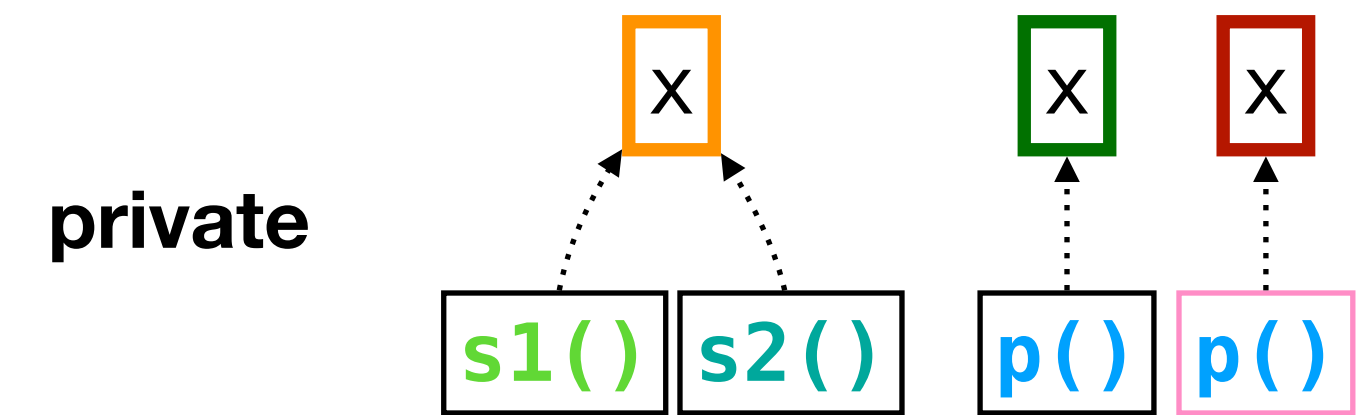
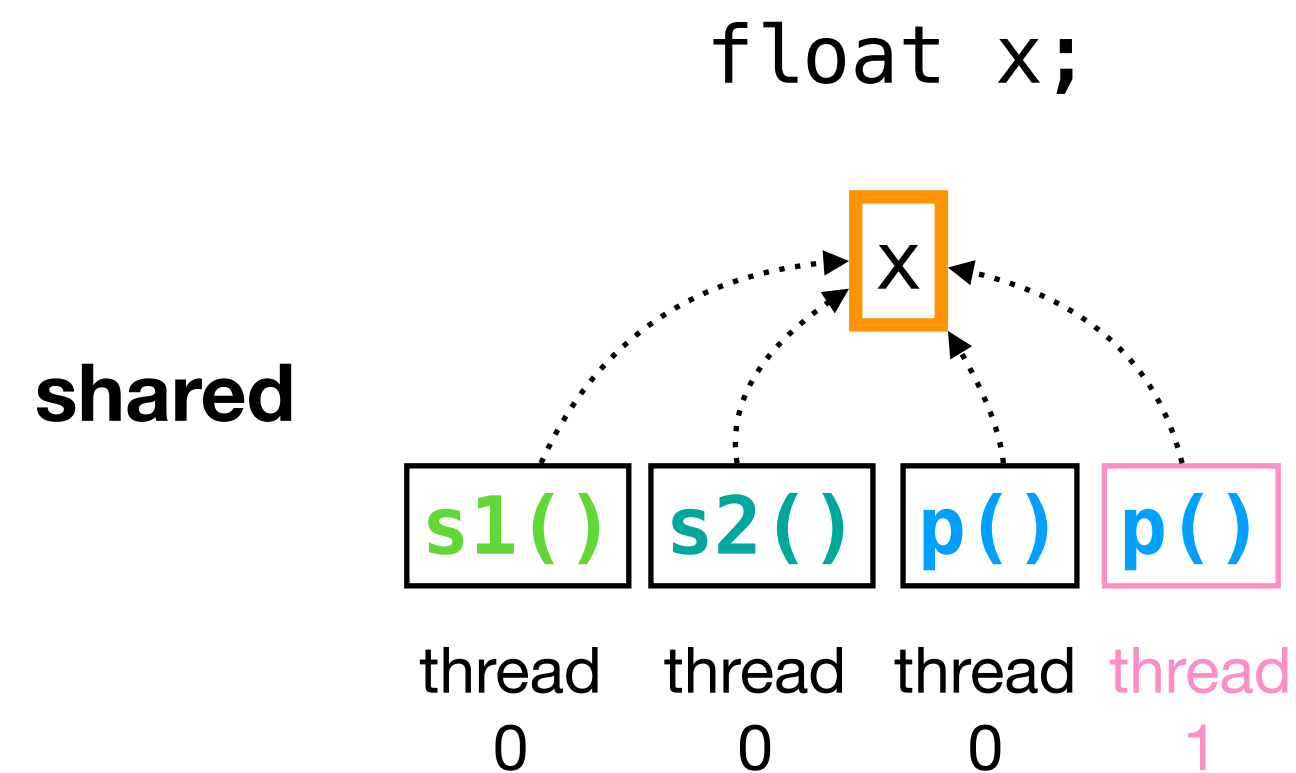
openMP code

Sequential region `s1()`

Parallel region `p()`

Sequential region `s2()`

2 threads

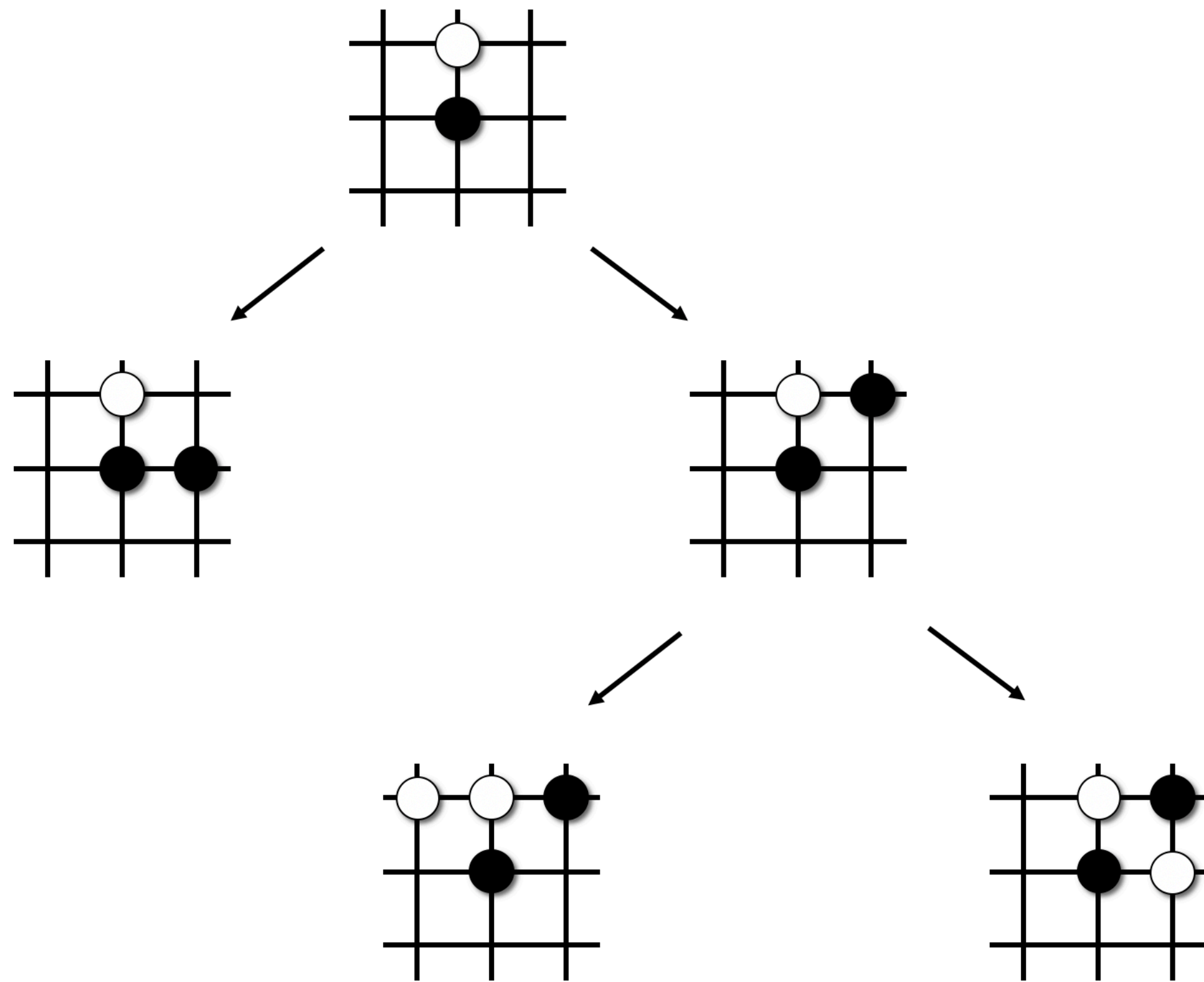


OpenMP tasks

#pragma omp task

- Many situations require a more flexible way of expressing parallelism
- Example: tree traversal
- Let's play a game of go.
- How would you describe the different ways to play the game?

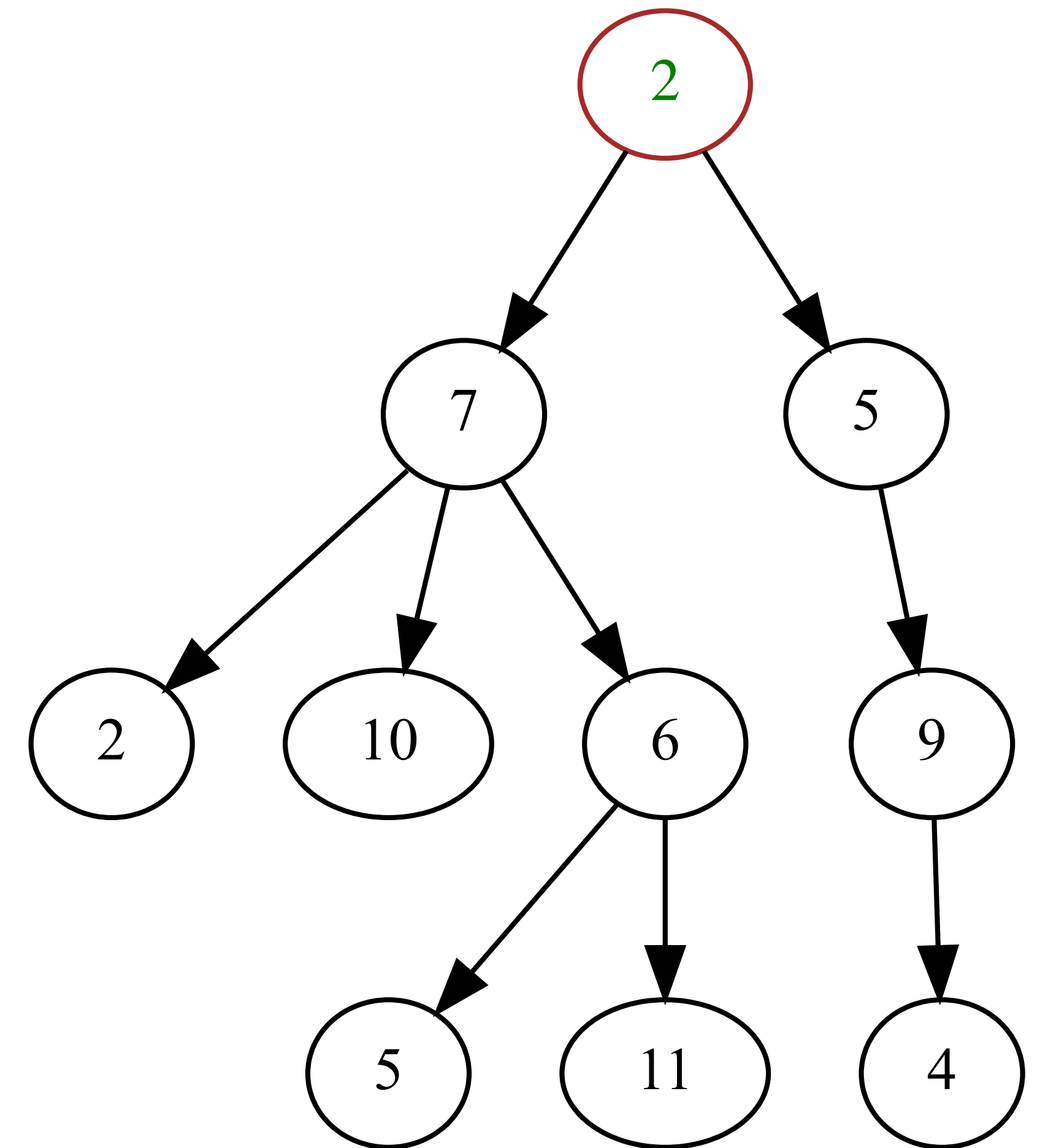
Playing Go



- Playing a game of go can be described using a tree structure.
- Traversing the tree and doing some computations is a key step in many algorithms.

Tree traversal

- Go through each node and execute some operation
- But, the tree is not full, e.g., number of child nodes varies
- As a result, the calculation is very unstructured.
- Work is discovered dynamically as the tree is traversed.
- A for loop is insufficient for this type of calculations.



omp tasks

- See `tree.cpp`
- Create tasks on the fly as we encounter new work to do.

```
|  if (curr_node->left)
|  #pragma omp task
|      Traverse(curr_node->left);
|
|  if (curr_node->right)
|  #pragma omp task
|      Traverse(curr_node->right);
```


Creating the parallel region

- We create a parallel region to have a pool of threads ready to do work.
- But the block of code is only executed by a single thread.

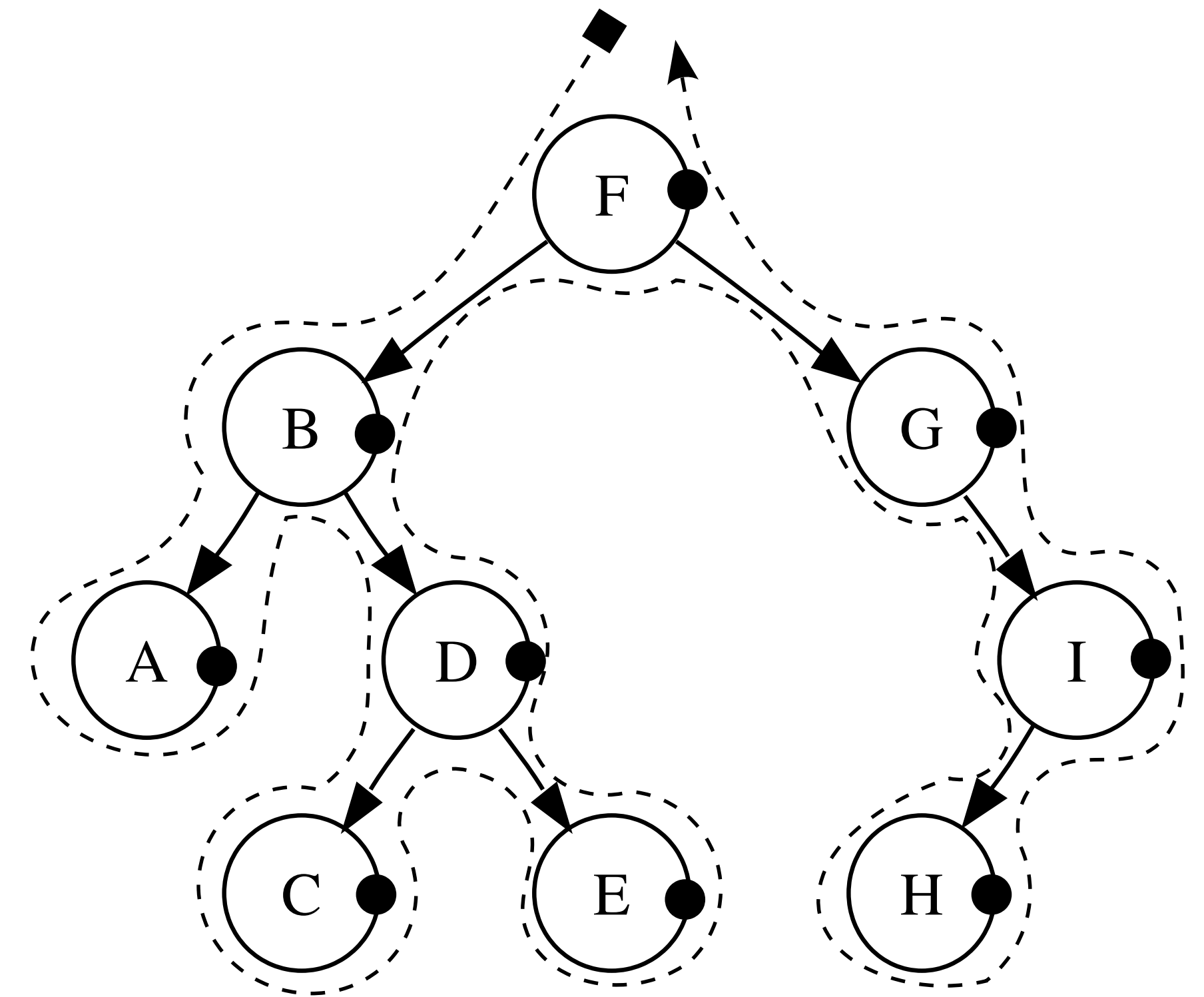
```
#pragma omp parallel
#pragma omp single
    // Only a single thread should execute this
    Traverse(root);
```

Task execution

- The encountering thread may immediately execute the task, or defer its execution.
- Any thread in the team may be assigned the task.
- This is managed by the OpenMP library.

Post-order traversal

- Let's now assume that we traverse the tree and based on the result of the children node, we do some calculation.
- This requires a synchronization. We have to wait for the children tasks to complete before we can proceed.



taskwait

- See `tree_postorder.cpp`
- `taskwait` is required to make sure left and right are up to date.

```
int PostOrderTraverse(struct Node *curr_node) {
    int left = 0, right = 0;

    if (curr_node->left)
#pragma omp task shared(left)
        left = PostOrderTraverse(curr_node->left);
    // Default attribute for task constructs is firstprivate
    if (curr_node->right)
#pragma omp task shared(right)
        right = PostOrderTraverse(curr_node->right);

#pragma omp taskwait
    curr_node->data = left + right; // Number of children nodes

    Visit(curr_node);
    return 1 + left + right;
}
```

Synchronization constructs

`taskwait` specifies a wait on the completion of child tasks of the current task.

`taskgroup` specifies a wait on the completion of child tasks of the current task and their descendent tasks.

`barrier` specifies an explicit barrier at the point at which the construct appears. All threads of the team that is executing the binding parallel region must execute the barrier region and complete execution of all explicit tasks bound to this parallel region before any are allowed to continue execution beyond the barrier.

Data sharing for tasks

- tasks are different from parallel for regions.
- In for loops, thread wait at the end of the loop.
- For tasks, the main thread just continues the execution and a worker thread takes care of the task.
- Default shared attribute is usually not the expected attribute.
- For tasks, the default attribute rules are different.

Default for tasks

- In a task generating construct, if no default clause is present, a variable that in the enclosing context is determined to be shared by all implicit tasks bound to the current team is shared.
- In a task generating construct, if no default clause is present, a variable for which the data-sharing attribute is not determined by the rule above is `firstprivate`.

firstprivate

- The `firstprivate` clause declares one or more list items to be private to a task, and initializes each of them with the value that the corresponding original item has when the construct is encountered.
- In this case, however, we want the variables `left` and `right` to be shared.
- A `shared` clause must be used.

```
int left = 0, right = 0;

if (curr_node->left)
#pragma omp task shared(left)
    left = PostOrderTraverse(curr_node->left);
// Default attribute for task constructs is firstprivate
if (curr_node->right)
#pragma omp task shared(right)
    right = PostOrderTraverse(curr_node->right);

#pragma omp taskwait
curr_node->data = left + right; // Number of children nodes

Visit(curr_node);
```


car race!



See `car_race.cpp`

Can you explain what each directive is doing?

What are the possible outputs for each case?

Is the output deterministic or undetermined?

What parts are deterministic, and what parts are undetermined?

Output of code

```
darve@icme-gpu1:~/2023/openMP$ srun -c 16 -p CME ./car_race
Run 1
A race car is fun to watch
A race car is fun to watch
A race A race car is fun to watch
A race car is fun to watch
car is fun to watch
A race car is fun to watch
A race car is fun to watch
A race car is fun to watch

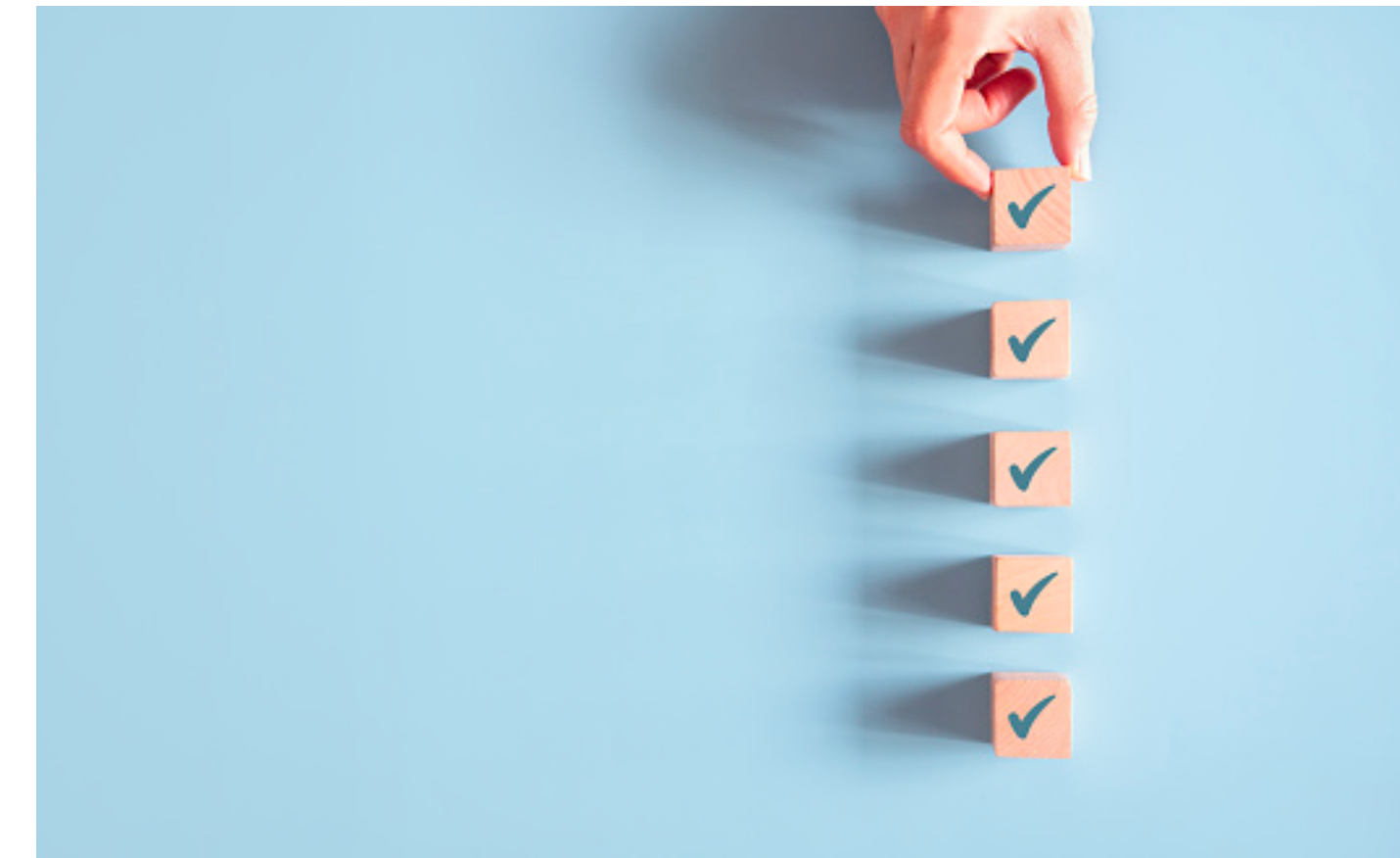
Run 2
A race car is fun to watch
A race car is fun to watch

Run 3
A race car is fun to watch
A A race A race car is fun to watch
race car car is fun to watch
is fun to watch
A A race A A race car race car is fun to watch
is fun to watch
car is fun to watch
race car is fun to watch

Run 4
A is fun to watch
race car A race is fun to watch
car

Run 5
A race car is fun to watch
darve@icme-gpu1:~/2023/openMP$ █
```

Processing entries in a list using tasks



- We want to read and modify entries in a linked list using OpenMP tasks.
- Each entry can be updated independently of the others.
- See `list.cpp`

```

#pragma omp parallel
#pragma omp single
{
    Node *curr_node = head;

    while (curr_node) {
        printf("Main thread. %p\n", (void *)curr_node);
#pragma omp task
        {
            // curr_node is firstprivate by default
            Wait();
            int tid = omp_get_thread_num();
            Visit(curr_node);
            printf("Task @%2d: node %p data %d\n", tid, (void *)curr_node,
                curr_node->data);
        }
        curr_node = curr_node->next;
    }
}

```

- The variable `curr_node` is `firstprivate` by default.
- This is the correct data sharing attribute.

More advanced openMP features

Task priority and data-dependencies

Recent additions to tasks

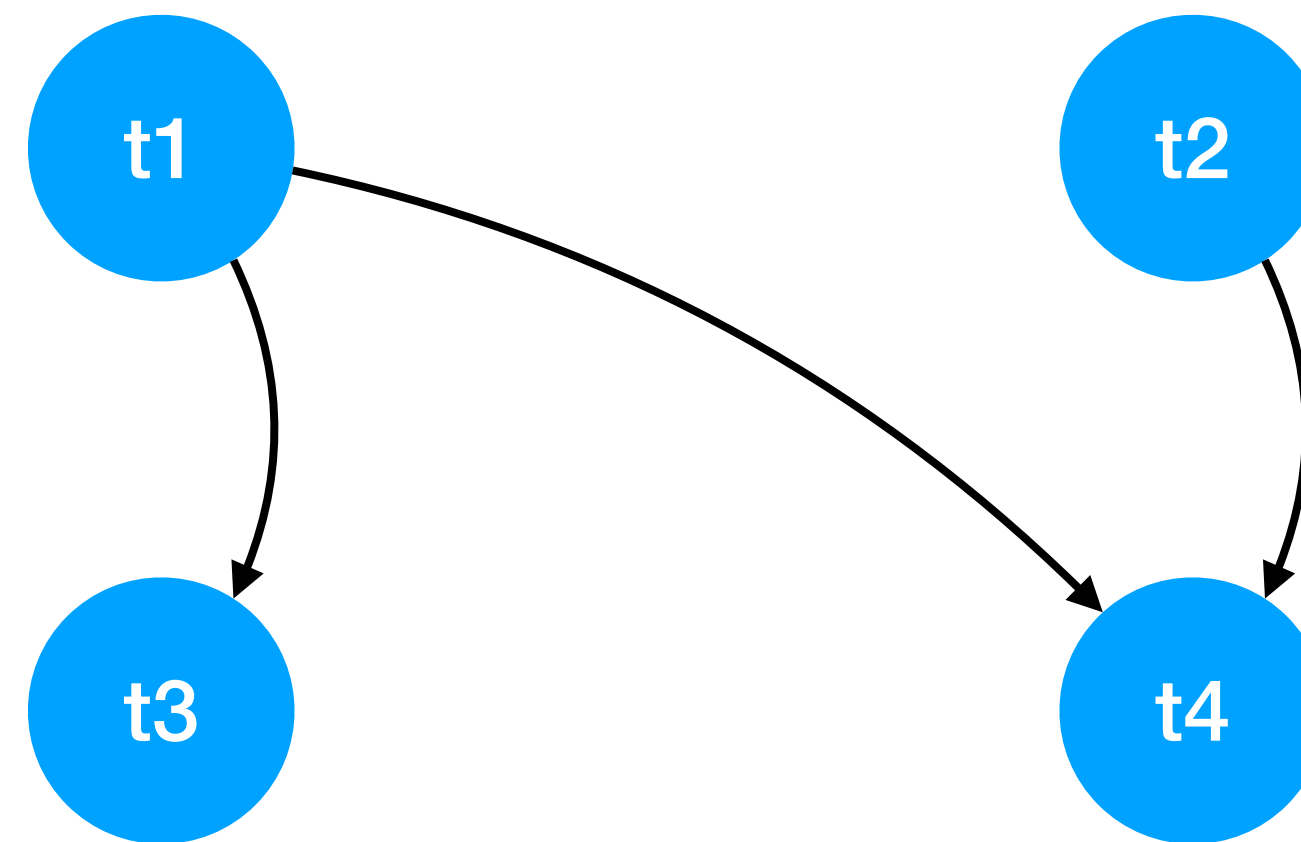
- Task priority.
- You may wish to assign higher priority to certain tasks.
- This is important in situations where for example, a task takes much longer than other tasks. In that case, we want OpenMP to assign a thread to that task as soon as possible.
- The priority clause specifies a hint for the task execution order. Among all tasks ready to be executed, higher priority tasks (those with a higher numerical priority-value) are recommended to execute before lower priority ones.

```
|   for (i = 0; i < N; i++)  
#pragma omp task priority(i)  
|       compute_array(&array[i * M], M);
```

What is the expected order of execution of the tasks?

taskwait and depend

- Although `taskwait` allows to synchronize tasks, it is not suitable when a fine-grained synchronization is required.
- For example: `t3` can start as soon as `t1` finishes; `t4` needs to wait for `t1` and `t2`.



depend

- This type of fine-grained dependency can be expressed using the depend clause.
- `depend(dep-type: x)`
- `dep-type` is one of
 - `in, out, inout, mutexinoutset`
 - `x` variable or array section.

List of dependencies

dep-type	waits on	waits on	waits on	mutually exclusive
in		out/inout	mutexinoutset	
out/inout	in	out/inout	mutexinoutset	
mutexinoutset	in	out/inout		mutexinoutset

Example

```
|  int x = 1;  
|  #pragma omp parallel  
|  #pragma omp single  
|  {  
|    #pragma omp task shared(x) depend(in : x)  
|    printf("x = %d\n", x);  
|    #pragma omp task shared(x) depend(out : x)  
|    x = 2;  
|  }
```

Always prints $x = 1$

```

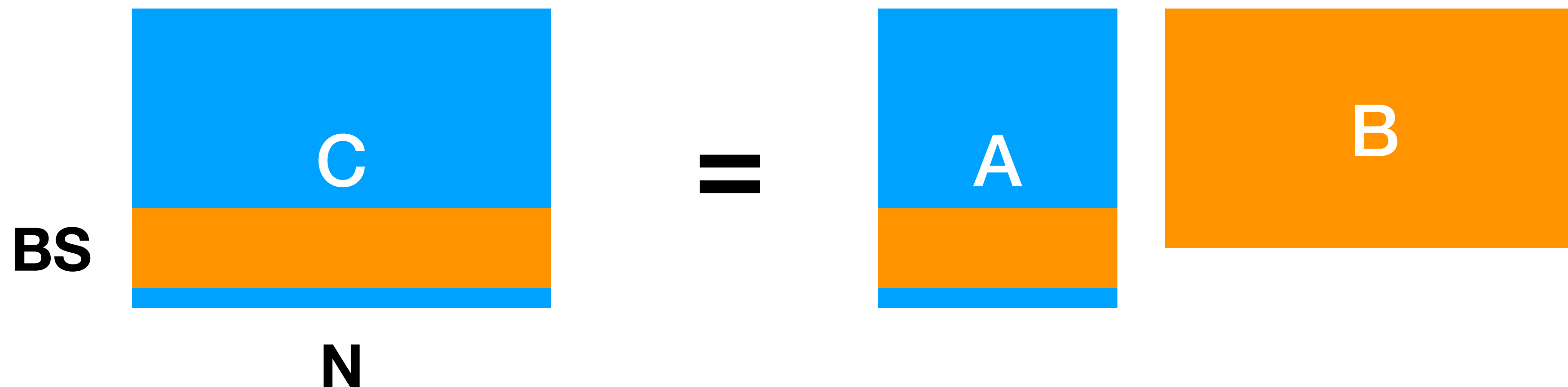
int d;
#pragma omp parallel
#pragma omp single
{
    int a, b, c;
#pragma omp task depend(out : c) shared(c)
    c = 1; /* Task T1 */
#pragma omp task depend(out : a) shared(a)
    a = 2; /* Task T2 */
#pragma omp task depend(out : b) shared(b)
    b = 3; /* Task T3 */
#pragma omp task depend(in : a) depend(mutexinoutset : c) shared(a, c)
    c += a; /* Task T4 */
#pragma omp task depend(in : b) depend(mutexinoutset : c) shared(b, c)
    c += b; /* Task T5 */
#pragma omp task depend(in : c) shared(c)
    d = c; /* Task T6 */
}
printf("%d\n", d);

```

1. Does t4 wait on t2?
2. Does t5 wait on t3?
3. Does t5 wait on t4?
4. Does t5 wait on t1?
5. Can t4 and t5 execute at the same time?
6. Can t6 execute at the same time as t4?
7. What is the output of this program?

Matrix-matrix product

```
#pragma omp parallel
#pragma omp single
| for (int i = 0; i < N; i += BS) {
// i is firstprivate by default
#pragma omp task depend(in: A[i * N:BS * N], B) depend(inout: C[i * N:BS * N])
|   for (int ii = i; ii < i + BS; ii++)
|     for (int j = 0; j < N; j++)
|       for (int k = 0; k < N; k++)
|         C[ii * N + j] += A[ii * N + k] * B[k * N + j];
| }
}
```



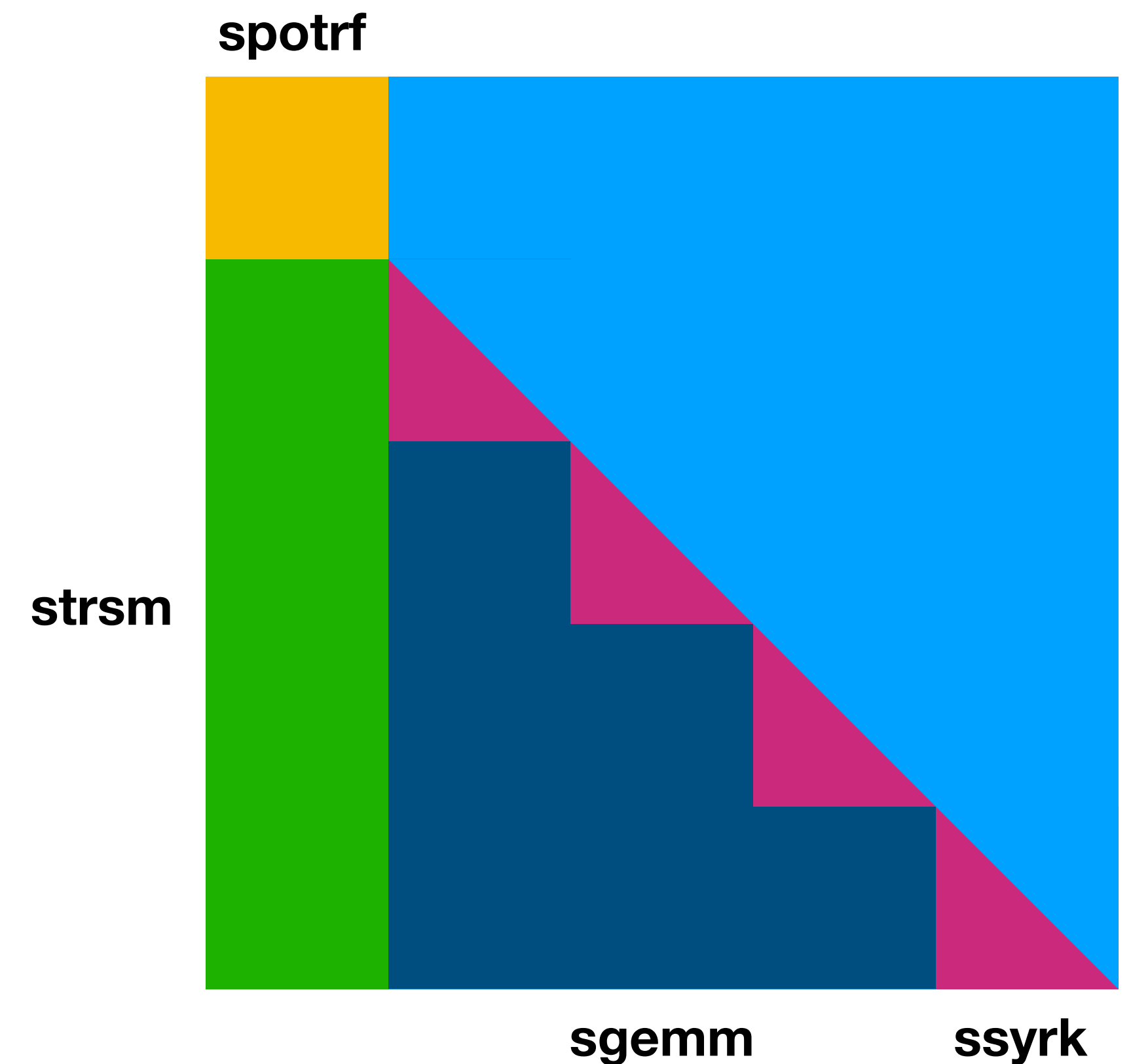
depend syntax

- `depend(in: A[i * N:BS * N])`
- Specifies the entries in A for which there is an in dependency.
- Syntax:

`A[lower-bound : length : stride]`

Cholesky algorithm

```
for (int k = 0; k < NB; k++) {  
#pragma omp task depend(inout : A[k][k])  
    spotrf(A[k][k]);  
    for (int i = k + 1; i < NB; i++)  
#pragma omp task depend(in : A[k][k]) depend(inout : A[k][i])  
        strsm(A[k][k], A[k][i]);  
    // update trailing submatrix  
    for (int i = k + 1; i < NB; i++) {  
        for (int j = k + 1; j < i; j++)  
#pragma omp task depend(in : A[k][i], A[k][j]) depend(inout : A[j][i])  
            sgemm(A[k][i], A[k][j], A[j][i]);  
#pragma omp task depend(in : A[k][i]) depend(inout : A[i][i])  
        // diagonal block update  
        ssyrk(A[k][i], A[i][i]);  
    }  
}
```



OpenMP synchronization constructs

- Some race conditions have a simple pattern and can be resolved using openMP directives.
- Learning these patterns is important to make sure the code is correct and to get good performance.

Reduction

This is the most common data race condition.

```
#pragma omp parallel for reduction(+ : sum)
for (int i = 0; i < size; i++) {
    sum += a[i];
}
```

Reduction optimization

- Although += is a race condition, there is a strategy to generate parallel code.
- Have each thread compute a local sum, private to each thread.
- **Then do an efficient reduction over partial results.**
- OpenMP does not analyze the content of the loop. It relies on the user to correctly specify the reduction operator.
- The reduction operator is used to:
 - **Correctly initialize sum**
 - **Correctly combine the partial results after the for loop is complete.**

Initializers and combiners

Identifier	Initializer	Combiner
+	<code>omp_priv = 0</code>	<code>omp_out += omp_in</code>
-	<code>omp_priv = 0</code>	<code>omp_out += omp_in</code>
*	<code>omp_priv = 1</code>	<code>omp_out *= omp_in</code>
&	<code>omp_priv = ~0</code>	<code>omp_out &= omp_in</code>
	<code>omp_priv = 0</code>	<code>omp_out = omp_in</code>
^	<code>omp_priv = 0</code>	<code>omp_out ^= omp_in</code>
&&	<code>omp_priv = 1</code>	<code>omp_out = omp_in && omp_out</code>
	<code>omp_priv = 0</code>	<code>omp_out = omp_in omp_out</code>
max	<code>omp_priv = <i>Least representable number in the reduction list item type</i></code>	<code>omp_out = omp_in > omp_out ? omp_in : omp_out</code>
min	<code>omp_priv = <i>Largest representable number in the reduction list item type</i></code>	<code>omp_out = omp_in < omp_out ? omp_in : omp_out</code>

Example: entropy calculation

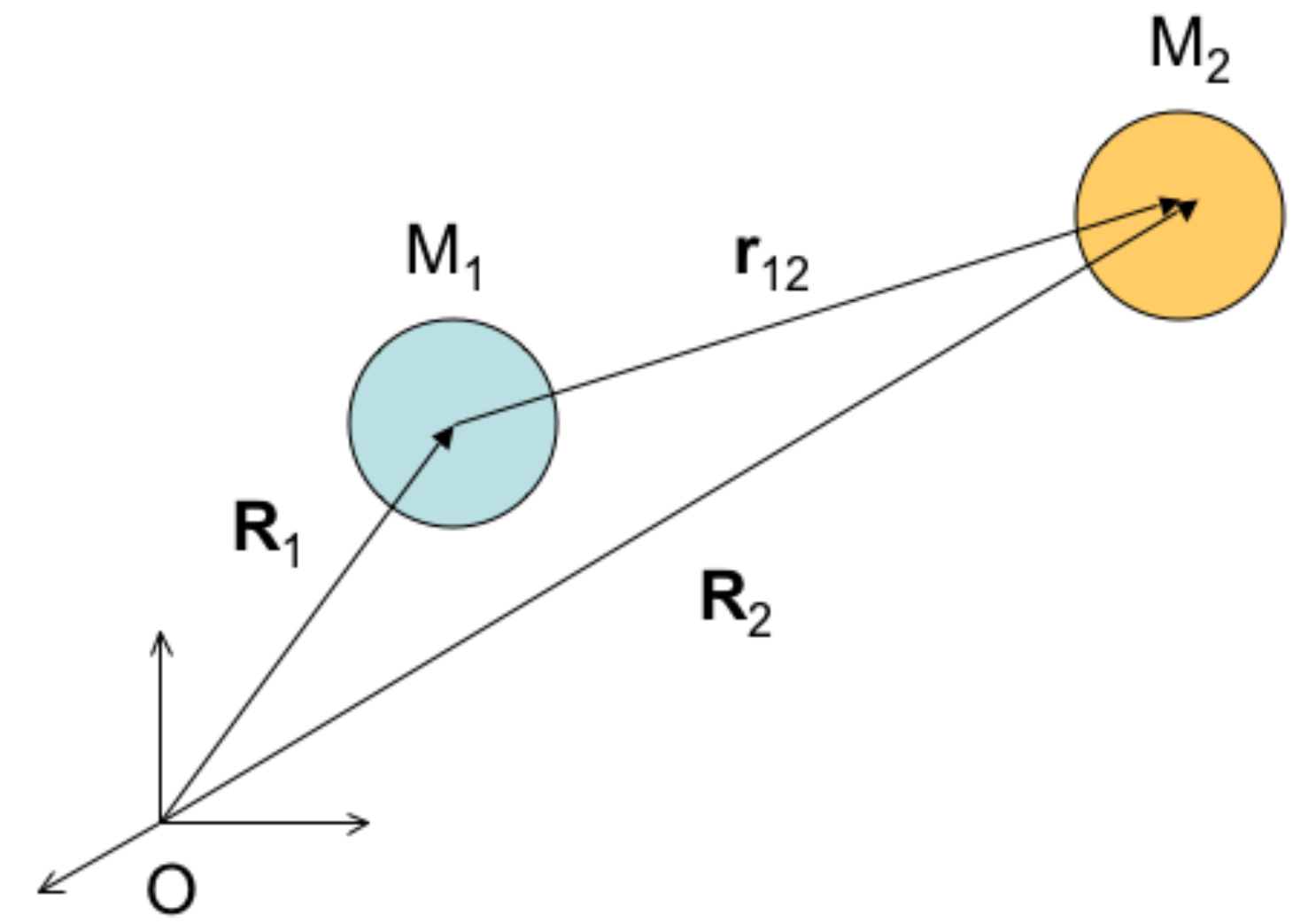


See `entropy.cpp`

Atomic operations

- There are cases where a similar `+=` operation need to be computed but `reduction` does not apply.
- In that case, we need to use an `atomic` clause.
- Atomic operations should be: `+`, `*`, `-`, `/`, `&`, `^`, `|`, `<<`, or `>>`
- But it is not as efficient as the `reduction` clause.
- The `atomic` construct ensures that a specific storage location is accessed atomically, rather than exposing it to the possibility of multiple, simultaneous reading and writing threads that may result in indeterminate values.

atomic example



N-body calculation.

Summation of forces leads to a race condition that can be resolved with an `atomic` clause.

See `atomic.cpp`

atomic is required both for i and j
in order to guarantee a correct result.

```
#pragma omp parallel for
    for (int i = 0; i < n; ++i)
        for (int j = i + 1; j < n; ++j) {
            const float x_ = x[i] - x[j];
            const float f_ = force(x_);
#pragma omp atomic
            f[i] += f_;
#pragma omp atomic
            f[j] -= f_;
        }
```

critical

- Similar to **mutex**, there might be regions of code that can be executed by **a single thread at a time** only.
- `critical` restricts execution of the associated structured block to a single thread at a time.
- See `critical.cpp`

Example: insertion in set

```
set<int> m;
#pragma omp parallel for
for (int i = 2; i <= n; ++i) {
    bool is_prime = is_prime_test(i);
#pragma omp critical
    if (is_prime) m.insert(i);
    // Save this prime; this requires a critical clause.
    // Only a single thread can execute this line at a time.
}
```

Insertion leads to a race condition that needs to be protected using `critical`.

Performance considerations: highlight

- Use **scheduling** clause for for loops.
- Prefer: reduction > atomic > critical.
- Use **nowait** whenever possible.
- Avoid creating parallel regions one after the other if they can safely be merged.
- **Memory access** is key for performance. **Data placement** matters for CC-NUMA processors.
- See Mastering OpenMP Performance.

Other topics (not covered):

`affinity, target, simd, locks`

5.0 OpenMP specification

5.2 is much less readable.