

Programming OpenMP

Christian Terboven

Michael Klemm





Agenda (in total 7 Sessions)



- Session 1: OpenMP Introduction
- Session 2: Tasking

Session 3: Optimization for NUMA and SIMD

- → Review of Session 2 / homework assignments
- → OpenMP and NUMA architectures
- → Task Affinity
- →SIMD

→Homework assignments ☺

- Session 4: What Could Possibly Go Wrong Using OpenMP
- Session 5: Introduction to Offloading with OpenMP
- Session 6: Advanced Offloading Topics
- Session 7: Selected / Remaining Topics



Programming OpenMP

Review

Christian Terboven Michael Klemm





Questions?

4



Fibonacci

Fibonacci illustrated



```
int main(int argc,
 1
 2
               char* argv[])
 3
    {
         [...]
 4
 5
         #pragma omp parallel
 6
 7
             #pragma omp single
 8
 9
                 fib(input);
            }
10
11
         }
12
         [...]
13.}
```

```
int fib(int n)
14
                      {
15
        if (n < 2) return n;
16
        int x, y;
17
        #pragma omp task shared(x)
18
19
            x = fib(n - 1);
20
        }
21
        #pragma omp task shared(y)
22
            y = fib(n - 2);
23
24
25
        #pragma omp taskwait
26
             return x+y;
27.}
```

- Only one Task / Thread enters fib() from main(), it is responsible for creating the two initial work tasks
- Taskwait is required, as otherwise x and y would get lost



T1 enters fib(4)



Task Queue

7



- T1 enters fib(4)
- T1 creates tasks for fib(3) and fib(2)







- T1 enters fib(4)
- T1 creates tasks for fib(3) and fib(2)
- T1 and T2 execute tasks from the queue



Task Queue



- T1 enters fib(4)
- T1 creates tasks for fib(3) and fib(2)
- T1 and T2 execute tasks from the queue
- T1 and T2 create 4 new tasks







- T1 enters fib(4)
- T1 creates tasks for fib(3) and fib(2)
- T1 and T2 execute tasks from the queue
- T1 and T2 create 4 new tasks
- T1 T4 execute tasks



Task Queue



- T1 enters fib(4)
- T1 creates tasks for fib(3) and fib(2)
- T1 and T2 execute tasks from the queue
- T1 and T2 create 4 new tasks
- T1 T4 execute tasks
- • •





For / Work-distribution

Example solution: For w/ Tasking wo/ Red.

```
#pragma omp parallel firstprivate(presult)
#pragma omp single
ł
        for (int i = 0; i < dimension; i++)</pre>
#pragma omp task shared(presult)
{
                result += do_some_computation(i);
}
} // end omp single
#pragma omp critical
{
        result += presult;
}
} // end omp parallel
```







```
#pragma omp parallel reduction(task,+:result)
{
#pragma omp single
{
#pragma omp taskloop in_reduction(+:result)
for (int i = 0; i < dimension; i++)
{
result += do_some_computation(i);
}
} // end omp single
} // end omp parallel</pre>
```



QuickSort

Example solution: Quick Sort



Example solution: Quick Sort





Programming OpenMP

Review

Christian Terboven Michael Klemm





Questions?



Fibonacci

Fibonacci illustrated



```
int main(int argc,
 1
 2
               char* argv[])
 3
    {
         [...]
 4
 5
         #pragma omp parallel
 6
 7
             #pragma omp single
 8
 9
                 fib(input);
            }
10
11
         }
12
         [...]
13.}
```

```
int fib(int n)
14
                      {
15
        if (n < 2) return n;
16
        int x, y;
17
        #pragma omp task shared(x)
18
19
            x = fib(n - 1);
20
        }
21
        #pragma omp task shared(y)
22
            y = fib(n - 2);
23
24
25
        #pragma omp taskwait
26
             return x+y;
27.}
```

- Only one Task / Thread enters fib() from main(), it is responsible for creating the two initial work tasks
- Taskwait is required, as otherwise x and y would get lost



T1 enters fib(4)



Task Queue



- T1 enters fib(4)
- T1 creates tasks for fib(3) and fib(2)







- T1 enters fib(4)
- T1 creates tasks for fib(3) and fib(2)
- T1 and T2 execute tasks from the queue







- T1 enters fib(4)
- T1 creates tasks for fib(3) and fib(2)
- T1 and T2 execute tasks from the queue
- T1 and T2 create 4 new tasks







- T1 enters fib(4)
- T1 creates tasks for fib(3) and fib(2)
- T1 and T2 execute tasks from the queue
- T1 and T2 create 4 new tasks
- T1 T4 execute tasks



Task Queue



- T1 enters fib(4)
- T1 creates tasks for fib(3) and fib(2)
- T1 and T2 execute tasks from the queue
- T1 and T2 create 4 new tasks
- T1 T4 execute tasks
- • •





For / Work-distribution

Example solution: For w/ Tasking wo/ Red.

```
#pragma omp parallel firstprivate(presult)
#pragma omp single
ł
        for (int i = 0; i < dimension; i++)</pre>
#pragma omp task shared(presult)
{
                result += do_some_computation(i);
}
} // end omp single
#pragma omp critical
{
        result += presult;
}
} // end omp parallel
```







```
#pragma omp parallel reduction(task,+:result)
{
#pragma omp single
{
#pragma omp taskloop in_reduction(+:result)
for (int i = 0; i < dimension; i++)
{
result += do_some_computation(i);
}
} // end omp single
} // end omp parallel</pre>
```



QuickSort



Programming OpenMP

Christian Terboven

Michael Klemm





Review: NUMA concept
Non-uniform Memory





double* A;

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

```
for (int i = 0; i < N; i++) {
    A[i] = 0.0;
}</pre>
```



Non-uniform Memory



 Serial code: all array elements are allocated in the memory of the NUMA node closest to the core executing the initializer thread (first touch)

```
double* A;
A = (double*)
    malloc(N * sizeof(double));
for (int i = 0; i < N; i++) {
    A[i] = 0.0;
}
```



About Data Distribution



Important aspect on cc-NUMA systems
 If not optimal, longer memory access times and hotspots

Placement comes from the Operating System
 This is therefore Operating System dependent

 Windows, Linux and Solaris all use the "First Touch" placement policy by default

 \rightarrow May be possible to override default (check the docs)

Non-uniform Memory



 Serial code: all array elements are allocated in the memory of the NUMA node closest to the core executing the initializer thread (first touch)

```
double* A;
A = (double*)
    malloc(N * sizeof(double));
for (int i = 0; i < N; i++) {
    A[i] = 0.0;
}
```







 First Touch w/ parallel code: all array elements are allocated in the memory of the NUMA node that contains the core that executes the thread that initializes the partition

```
double* A;
A = (double*)
  malloc(N * sizeof(double));
omp_set_num_threads(2);
#pragma omp parallel for
for (int i = 0; i < N; i++) {
  A[i] = 0.0;
}
```





Serial vs. Parallel Initialization

Stream example on 2 socket sytem with Xeon X5675 processors, 12 OpenMP threads:

	сору	scale	add	triad
ser_init	18.8 GB/s	18.5 GB/s	18.1 GB/s	18.2 GB/s
par_init	41.3 GB/s	39.3 GB/s	40.3 GB/s	40.4 GB/s



OpenMP Tutorial Members of the OpenMP Language Committee

33



Thread Binding and Memory Placement

Get Info on the System Topology



- Before you design a strategy for thread binding, you should have a basic understanding of the system topology. Please use one of the following options on a target machine:
 - →Intel MPI's cpuinfo tool

→ cpuinfo

→Delivers information about the number of sockets (= packages) and the mapping of processor ids to cpu cores that the OS uses.

hwlocs' hwloc-ls tool

 \rightarrow hwloc-ls

→ Displays a (graphical) representation of the system topology, separated into NUMA nodes, along with the mapping of processor ids to cpu cores that the OS uses and additional info on caches.

Decide for Binding Strategy



- Selecting the "right" binding strategy depends not only on the topology, but also on application characteristics.
 - →Putting threads far apart, i.e., on different sockets
 - \rightarrow May improve aggregated memory bandwidth available to application
 - \rightarrow May improve the combined cache size available to your application
 - →May decrease performance of synchronization constructs
 - →Putting threads close together, i.e., on two adjacent cores that possibly share some caches
 - →May improve performance of synchronization constructs
 - \rightarrow May decrease the available memory bandwidth and cache size

Places + Binding Policies (1/2)



Define OpenMP Places

→ set of OpenMP threads running on one or more processors

- → can be defined by the user, i.e. OMP_PLACES=cores
- Define a set of OpenMP Thread Affinity Policies

 SPREAD: spread OpenMP threads evenly among the places, partition the place list
 - → CLOSE: pack OpenMP threads near primary thread
 - → PRIMARY: collocate OpenMP thread with primary thread
- Goals
 - → user has a way to specify where to execute OpenMP threads
 - → locality between OpenMP threads / less false sharing / memory bandwidth

OMP_PLACES env. variable



Assume the following machine:



 \rightarrow 2 sockets, 4 cores per socket, 4 hyper-threads per core

- Abstract names for OMP_PLACES:
 - \rightarrow threads: Each place corresponds to a single hardware thread.
 - \rightarrow cores: Each place corresponds to a single core (having one or more hardware threads).
 - \rightarrow sockets: Each place corresponds to a single socket (consisting of one or more cores).
 - \rightarrow II_caches (5.1): Each place corresponds to a set of cores that share the last level cache.
 - → numa_domains (5.1): Each places corresponds to a set of cores for which their closest memory is: the same memory; and at a similar distance from the cores.

OpenMP 4.0: Places + Policies



Example's Objective:

→ separate cores for outer loop and near cores for inner loop

Outer Parallel Region: proc_bind(spread), Inner: proc_bind(close)
 > spread creates partition, compact binds threads within respective partition

```
OMP_PLACES=(0,1,2,3), (4,5,6,7), ... = (0-4):4:8 = cores
#pragma omp parallel proc_bind(spread) num_threads(4)
#pragma omp parallel proc_bind(close) num_threads(4)
```

Example

 \rightarrow initial

→ spread 4



→ close 4

OpenMP 4.0: Places + Policies



Example's Objective:

→ separate cores for outer loop and near cores for inner loop

Outer Parallel Region: proc_bind(spread), Inner: proc_bind(close)
 > spread creates partition, compact binds threads within respective partition

```
OMP_PLACES=(0,1,2,3), (4,5,6,7), ... = (0-4):4:8 = cores
#pragma omp parallel proc_bind(spread) num_threads(4)
#pragma omp parallel proc_bind(close) num_threads(4)
```

Example

 \rightarrow initial

→ spread 4

→ close 4





More Examples (1/3)

Assume the following machine:



Parallel Region with two threads, one per socket
 >OMP_PLACES=sockets

+#pragma omp parallel num_threads(2) proc_bind(spread)



More Examples (2/3)

Assume the following machine:

p0 p1 p2 p3 p4 p5 p6 p7

- Parallel Region with four threads, one per core, but only on the first socket
 - →OMP_PLACES=cores
 - +#pragma omp parallel num_threads(4) proc_bind(close)



More Examples (3/3)

 Spread a nested loop first across two sockets, then among the cores within each socket, only one thread per core

→OMP_PLACES=cores

 \rightarrow #pragma omp parallel num_threads(2) proc_bind(spread)

+#pragma omp parallel num_threads(4) proc_bind(close)

Places API routines allow to

 -> query information about binding...

 \rightarrow query information about the place partition...



Places API: Example

Simple routine printing the processor ids of the place the calling thread is bound to:



Places API: Example

Simple routine printing the processor ids of the place the calling thread is bound to:

```
void print binding info() {
     int my place = omp get place num();
     int place num procs = omp get place num procs (my place);
     printf("Place consists of %d processors: ", place num procs);
     int *place processors = malloc(sizeof(int) * place_num_procs);
     omp get place proc ids (my place, place processors)
     for (int i = 0; i < place num procs - 1; i++) {
             printf("%d ", place processors[i]);
     printf("\n");
     free(place processors);
```

OpenMP 5.0 way to do this



• Set OMP_DISPLAY_AFFINITY=TRUE

→Instructs the runtime to display formatted affinity information

 \rightarrow Example output for two threads on two physical cores:

→Output can be formatted with OMP_AFFINITY_FORMAT env var or
corresponding routine

→Formatted affinity information can be printed with
omp_display_affinity(const char* format)

OpenMP 5.0 way to do this



• Set OMP_DISPLAY_AFFINITY=TRUE

→Instructs the runtime to display formatted affinity information

 \rightarrow Example output for two threads on two physical cores:

nesting_level= 1, thread_num= 0, thread_affinity= 0,1
nesting_level= 1, thread_num= 1, thread_affinity= 2,3

corresponding routine

→Formatted affinity information can be printed with
omp display affinity(const char* format)

Affinity format specification



- t omp_get_team_num()
- T omp_get_num_teams()
- L omp_get_level()
- n omp_get_thread_num()
- N omp_get_num_threads()
- Example:

 \rightarrow Possible output:

omp_get_ancestor_thread_num() at level-1
beatname

H hostname

а

P process identifier

native thread identifier

A thread affinity: list of processors (cores)

Affinity format specification



- T omp_get_num_teams()
- L omp_get_level()
- n omp_get_thread_num()
- N omp_get_num_threads()
- Example:

- a omp_get_ancestor_thread_num() at level-1
 H hostname
 P process identifier
 i native thread identifier
- A thread affinity: list of processors (cores)

OMP_AFFINITY_FORMAT="Affinity: %0.3L %.8n %.15{A} %.12H" →Possible output:

Affinity:	001	0	0-1,16-17	host003
Affinity:	001	1	2-3,18-19	host003

Fine-grained control of Memory Affinity



- Explicit NUMA-aware memory allocation:
 - \rightarrow By carefully touching data by the thread which later uses it
 - →By changing the default memory allocation strategy
 - Linux: numactl command
 - \rightarrow By explicit migration of memory pages
 - →Linux: move_pages()
- Example: using numactl to distribute pages roundrobin:
 - >numactl -interleave=all ./a.out



Managing Memory Spaces

Different kinds of memory

- Traditional DDR-based memory
- High-bandwidth memory
- Non-volatile memory

- - -

<u>OpenMP</u>

Cascade Lake (Leonide at INRIA)

CPU: Intel(R) Xeon(R) Gold 6230 CPU @ 2.10GHz Freq Govenor: performance			
node 0 cpus: 0 2 4 6 8 10 12 14 16 18			
20 22 24 26 28 30 32 34 36 38			
node 0 size: 191936 MB			
node 0 free: 178709 MB			
node 1 cpus: 1 3 5 7 9 11 13 15 17 19 21 23			
25 27 29 31 33 35 37 39			
node 1 size: 192016 MB			
node 1 free: 179268 MB			
node 2 cpus:			
node 2 size: 759808 MB			
node 2 free: 759794 MB			
node 3 cpus:			
node 3 size: 761856 MB			
node 3 tree: /61851 MB			
node distances:			
10000 U I Z 3			
$ \begin{array}{cccccccccccccccccccccccccccccccccccc$			
$2 \cdot 17 28 10 28$			
3. 28 17 28 10			
DRAM + Optane			

<u>OpenMP</u>

Memory Management

- Allocator := an OpenMP object that fulfills requests to allocate and deallocate storage for program variables
- **OpenMP allocators are of type** omp_allocator_handle_t
- Default allocator for Host
 >via OMP_ALLOCATOR env. var. or corresponding API
- OpenMP 5.0 supports a set of memory allocators

OpenMP allocators



Selection of a certain kind of memory

Allocator name	Storage selection intent
omp_default_mem_alloc	use default storage
omp_large_cap_mem_alloc	use storage with large capacity
omp_const_mem_alloc	use storage optimized for read-only variables
omp_high_bw_mem_alloc	use storage with high bandwidth
omp_low_lat_mem_alloc	use storage with low latency
omp_cgroup_mem_alloc	use storage close to all threads in the contention group of the thread requesting the allocation
omp_pteam_mem_alloc	use storage that is close to all threads in the same parallel region of the thread requesting the allocation
omp_thread_local_mem_allo c	use storage that is close to the thread requesting the allocation

Using OpenMP Allocators



• New clause on all constructs with data sharing clauses:

>allocate([allocator:] list)

Allocation:

> omp_alloc(size_t size, omp_allocator_handle_t allocator)

Deallocation:

>omp_free(void *ptr, const omp_allocator_handle_t allocator)

> allocator argument is optional

 allocate directive: standalone directive for allocation, or declaration of allocation stmt.

OpenMP allocator traits / 1



Allocator traits control the behavior of the allocator

sync_hint	contended, uncontended, serialized, private default: contended	
alignment	positive integer value that is a power of two default: 1 byte	
access	all, cgroup, pteam, thread default: all	
pool_size	positive integer value	
fallback	<pre>default_mem_fb, null_fb, abort_fb, allocator_fb default: default_mem_fb</pre>	
fallback fb_data	default_mem_fb, null_fb, abort_fb, allocator_fb default: default_mem_fb an allocator handle	
fallback fb_data pinned	default_mem_fb, null_fb, abort_fb, allocator_fb default: default_mem_fb an allocator handle true, false default: false	

OpenMP Tutorian Members of the OpenMP Language Committee

OpenMP Allocator Traits / 2



fallback: describes the behavior if the allocation cannot be fulfilled
 default_mem_fb: return system's default memory

 \rightarrow Other options: null, abort, or use different allocator

pinned: request pinned memory, i.e. for GPUs

OpenMP Allocator Traits / 3



- partition: partitioning of allocated memory of physical storage resources (think of NUMA)
 - >environment: use system's default behavior
 - >nearest: most closest memory
 - >blocked: partitioning into approx. same size with at most one block per storage resource
 - →interleaved: partitioning in a round-robin fashion across the storage
 resources

OpenMP Allocator Traits / 4



Construction of allocators with traits via

>omp_allocator_handle_t omp_init_allocator(

omp_memspace_handle_t memspace,

int ntraits, const omp_alloctrait_t traits[]);

→ Selection of memory space mandatory

 \rightarrow Empty traits set: use defaults

Allocators have to be destroyed with *_destroy_*

Custom allocator can be made default with omp_set_default_allocator(omp_allocator_handle_t allocator)

OpenMP Memory Spaces



Storage resources with explicit support in OpenMP:

	omp_default_mem_space	System's default memory resource
	omp_large_cap_mem_spa ce	Storage with larg(er) capacity
	omp_const_mem_space	Storage optimized for variables with constant value
→Exac	omp_high_bw_mem_spac e	Storage with high bandwidth
→Pre-c	omp_low_lat_mem_space	Storage with low latency

Memory Management Status



- LLVM OpenMP runtime internally already uses libmemkind (libnuma, numactl)
 - → Support for various kinds of memory: DDR, HBW and Persistent Memory (Optane)
 - → Library loaded at initialization (checks for availability)
 - \rightarrow If requested memory space for allocator is not available \rightarrow fallback to DDR

Memory Management implementation in LLVM still not complete

- → Some allocator traits not implemented yet
- Some partition values not implemented yet (environment, interleaved, nearest, blocked)
- Semantics of omp_high_bw_mem_space and omp_large_cap_mem_space unclear. Which memory should be used?
 - \rightarrow Explicitly target HBM \rightarrow currently implemented in LLVM

LLVM has custom implementation of aligned memory allocation

→ Allocation covers → {Allocator Information + Requested Size + Buffer based on alignment}



Programming OpenMP

Christian Terboven

Michael Klemm





Improving Tasking Performance: Task Affinity
Motivation



Techniques for process binding & thread pinning available

→ OpenMP thread level: OMP_PLACES & OMP_PROC_BIND

→OS functionality: taskset -c

OpenMP Tasking:

■ In general: Tasks may be executed by any thread in the team →Missing task-to-data affinity may have detrimental effect on performance

OpenMP 5.0:

affinity clause to express affinity to data

affinity clause



- New clause: #pragma omp task affinity (list)
 - \rightarrow Hint to the runtime to execute task closely to physical data location

→Clear separation between dependencies and affinity

Expectations:

→Improve data locality / reduce remote memory accesses

→ Decrease runtime variability

Still expect task stealing
 In particular, if a thread is under-utilized

Code Example



Excerpt from task-parallel STREAM

```
1
    #pragma omp task \
2
        shared(a, b, c, scalar) \
3
        firstprivate(tmp_idx_start, tmp_idx_end) \
        affinity( a[tmp_idx_start] )
4
5
    {
6
       int i:
7
       for(i = tmp_idx_start; i <= tmp_idx_end; i++)</pre>
            a[i] = b[i] + scalar + c[i];
8
   oops have been blocked manually (see tmp idx start/end)
```

→Assumption: initialization and computation have same blocking and same affinity



P



A map is introduced to store location information of data that was previously used

MP



MP



OpenMP

Program runtime Median of 10 runs





Program runtime Median of 10 runs





Program runtime Median of 10 runs



LIKWID: reduction of remote data volume from 69% to 13%



⁶⁴ OpenMP Tutorial Members of the OpenMP Language Committee



Program runtime Median of 10 runs



Distribution of single task execution times



LIKWID: reduction of remote data volume from 69% to 13%

⁶⁴ OpenMP Tutorial Members of the OpenMP Language Committee





- Requirement for this feature: thread affinity enabled
- The affinity clause helps, if
 - →tasks access data heavily
 - →single task creator scenario, or task not created with data affinity
 - \rightarrow high load imbalance among the tasks

Different from thread binding: task stealing is absolutely allowed



Programming OpenMP

Christian Terboven Michael Klemm



Topics



- Exploiting SIMD parallelism with OpenMP
- Using SIMD directives with loops
- Creating SIMD functions

SIMD on x86 Architectures



Width of SIMD registers has been growing in the past:



More Powerful SIMD Units



SIMD instructions become more powerful



More Powerful SIMD Units



SIMD instructions become more powerful



vadd dest{k1}, source2, source3



SIMD instructions become more powerful

More Powerful SIMD Units



More Powerful SIMD Units

OpenMP.

SIMD instructions become more powerful





Compilers offer auto-vectorization as an optimization pass
 Usually, part of the general loop optimization passes



- Compilers offer auto-vectorization as an optimization pass
 Usually, part of the general loop optimization passes
 - →Code analysis detects code properties that inhibit SIMD vectorization



- Compilers offer auto-vectorization as an optimization pass
 Usually, part of the general loop optimization passes
 - →Code analysis detects code properties that inhibit SIMD vectorization
 - \rightarrow Heuristics determine if SIMD execution might be beneficial



- Compilers offer auto-vectorization as an optimization pass
 Usually, part of the general loop optimization passes
 - →Code analysis detects code properties that inhibit SIMD vectorization
 - \rightarrow Heuristics determine if SIMD execution might be beneficial
 - \rightarrow If all goes well, the compiler will generate SIMD instructions



- Compilers offer auto-vectorization as an optimization pass
 Usually, part of the general loop optimization passes
 - →Code analysis detects code properties that inhibit SIMD vectorization
 - →Heuristics determine if SIMD execution might be beneficial
 - \rightarrow If all goes well, the compiler will generate SIMD instructions

 Example: clang/LLVM GCC Intel Compiler -ftree-vectorize -vec (enabled w/ -O2)
 -Rpass=loop-.* -ftree-loop-vectorize -qopt-report=vec
 -mprefer-vector-width=<width> -fopt-info-vec-all



Compilers offer auto-vectorization as an optimization pass
 Usually, part of the general loop optimization passes

 \rightarrow Code analysis detects code properties that inhibit SIMD vectorization

Heuristics determine if SIMD execution might be beneficial

 \rightarrow If all goes well, the compiler will generate SIMD instructions

Example: clang/LLVM GCC Intel Compiler
 -fvectorize -ftree-vectorize -vec (enabled w/ -O2)
 -Rpass=loop-.* -ftree-loop-vectorize -qopt-report=vec
 -mprefer-vector-width=<width> -fopt-info-vec-all

Why Auto-vectorizers Fail

<u>OpenM</u>

Data dependencies

- Other potential reasons
 - →Alignment
 - →Function calls in loop block
 - →Complex control flow / conditional branches
 - →Loop not "countable"
 - \rightarrow e.g., upper bound not a runtime constant
 - →Mixed data types
 - →Non-unit stride between elements
 - →Loop body too complex (register pressure)
 - →Vectorization seems inefficient
- Many more ... but less likely to occur

Data Dependencies



- Suppose two statements S1 and S2
- S2 depends on S1, iff S1 must execute before S2
 →Control-flow dependence
 - →Data dependence
 - → Dependencies can be carried over between loop iterations
- Important flavors of data dependencies





Dependencies may occur across loop iterations
 >Loop-carried dependency

The following code contains such a dependency:

```
void lcd_ex(float* a, float* b, size_t n, float c1, float c2)
{
    size_t i;
    for (i = 0; i < n; i++) {
        a[i] = c1 * a[i + 17] + c2 * b[i];
    }
}</pre>
```

 Some iterations of the loop have to complete before the next iteration can run
 Simple trick: Can you reverse the loop w/o getting wrong results?



Dependencies may occur across loop iterations

 ->Loop-carried dependency

The following code contains such a dependency:

```
void lcd_ex(float* a, float* b, size_t n, float c1, float c2)
{
    size_t i;
    for (i = 0; i < n; i++) {
        a[i] = c1 * a[i + 17] + c2 * b[i];
    }
}</pre>
```

 Some iterations of the loop have to complete before the next iteration can run
 Simple trick: Can you reverse the loop w/o getting wro

Loop-carried dependency for a[i] and a[i+17]; distance is 17.

 \rightarrow Simple trick: Can you reverse the loop w/o getting wrong results?



```
void lcd_ex(float* a, float* b, size_t n, float c1, float c2) {
    for (int i = 0; i < n; i++) {
        a[i] = c1 * a[i + 17] + c2 * b[i];
    }
}</pre>
```





```
void lcd_ex(float* a, float* b, size_t n, float c1, float c2) {
    for (int i = 0; i < n; i++) {
        a[i] = c1 * a[i + 17] + c2 * b[i];
    }
}</pre>
```





```
void lcd_ex(float* a, float* b, size_t n, float c1, float c2) {
    for (int i = 0; i < n; i++) {
        a[i] = c1 * a[i + 17] + c2 * b[i];
}</pre>
```









Can we parallelize or vectorize the loop?



(except for very specific loop schedules)



Can we parallelize or vectorize the loop?





(except for very specific loop schedules)



Can we parallelize or vectorize the loop?





(except for very specific loop schedules)
Loop-carried Dependencies



Can we parallelize or vectorize the loop?





(except for very specific loop schedules)

Loop-carried Dependencies



Can we parallelize or vectorize the loop?





(except for very specific loop schedules)

Loop-carried Dependencies



Can we parallelize or vectorize the loop?





(except for very specific loop schedules)

 \rightarrow Vectorization: yes

(iff vector length is shorter than any distance of any dependency)

In a Time Before OpenMP 4.0



```
→Compiler pragmas (e.g., #pragma vector)
```

→Low-level constructs (e.g., _mm_add_pd())

```
#pragma omp parallel for
#pragma vector always
#pragma ivdep
for (int i = 0; i < N; i++) {
    a[i] = b[i] + ...;
```

In a Time Before OpenMP 4.0

Support required vendor-specific extensions
 Programming models (e.g., Intel® Cilk Plus)

```
→Compiler pragmas (e.g., #pragma vector)
```

→Low-level constructs (e.g., _mm_add_pd())

```
#pragma omp parallel for
#pragma vector always
#pragma ivdep
for (int i = 0; i < N; i++) {
    a[i] = b[i] + ...;
```



SIMD Loop Construct



Vectorize a loop nest

→Cut loop into chunks that fit a SIMD vector register

 \rightarrow No parallelization of the loop body

Syntax (C/C++)
#pragma omp simd [clause[[,] clause],...]
for-loops

Syntax (Fortran)

!\$omp simd [clause[[,] clause],...]
do-loops
[!\$omp end simd]









```
float sprod(float *a, float *b, int n) {
  float sum = 0.0f;
#pragma omp simd reduction(+:sum)
  for (int k=0; k<n; k++)
    sum += a[k] * b[k];
  return sum;
}</pre>
```







private(var-list):

Uninitialized vectors for variables in var-list



private(var-list):

Uninitialized vectors for variables in var-list



firstprivate(var-list):
 Initialized vectors for variables in var-list

$$\mathbf{X}: \begin{array}{c} 4\\ 2 \end{array} \longrightarrow \begin{array}{c} 4\\ 2 \end{array} \end{array}$$



private(var-list):

Uninitialized vectors for variables in var-list



firstprivate(var-list):
 Initialized vectors for variables in var-list

$$\mathbf{X}: \begin{array}{c} 4\\ 2 \end{array} \longrightarrow \begin{array}{c} 4\\ 2 \end{array} \begin{array}{c} 4\\ 2 \end{array}$$

reduction(op:var-list):

Create private variables for *var-list* and apply reduction operator *op* at the end of the construct



safelen (length)

 \rightarrow Maximum number of iterations that can run concurrently without breaking a dependence

 \rightarrow In practice, maximum vector length



safelen (length)

 \rightarrow Maximum number of iterations that can run concurrently without breaking a dependence

 \rightarrow In practice, maximum vector length

linear (list[:linear-step])

 \rightarrow The variable's value is in relationship with the iteration number

 $\rightarrow x_i = x_{orig} + i * linear-step$

OpenMP?

safelen (length)

 \rightarrow Maximum number of iterations that can run concurrently without breaking a dependence

 \rightarrow In practice, maximum vector length

linear (list[:linear-step])

 \rightarrow The variable's value is in relationship with the iteration number

 $\rightarrow x_i = x_{orig} + i * linear-step$

aligned (list[:alignment])

→ Specifies that the list items have a given alignment

 \rightarrow Default is alignment for the architecture

OpenMP?

safelen (length)

 \rightarrow Maximum number of iterations that can run concurrently without breaking a dependence

 \rightarrow In practice, maximum vector length

linear (list[:linear-step])

 \rightarrow The variable's value is in relationship with the iteration number

 $\rightarrow x_i = x_{orig} + i * linear-step$

aligned (list[:alignment])

→ Specifies that the list items have a given alignment

→Default is alignment for the architecture

collapse (n)

SIMD Worksharing Construct



Parallelize and vectorize a loop nest

→Distribute a loop's iteration space across a thread team

→ Subdivide loop chunks to fit a SIMD vector register

Syntax (C/C++)
#pragma omp for simd [clause[[,] clause],...]
for-loops

Syntax (Fortran)

!\$omp do simd [clause[[,] clause],...]
do-loops
[!\$omp end do simd [nowait]]









```
float sprod(float *a, float *b, int n) {
  float sum = 0.0f;
#pragma omp for simd reduction(+:sum)
  for (int k=0; k<n; k++)
    sum += a[k] * b[k];
  return sum;
}</pre>
```





```
float sprod(float *a, float *b, int n) {
  float sum = 0.0f;
#pragma omp for simd reduction(+:sum)
  for (int k=0; k<n; k++)
    sum += a[k] * b[k];
  return sum;
}</pre>
```



84



```
float sprod(float *a, float *b, int n) {
  float sum = 0.0f;
#pragma omp for simd reduction(+:sum)
  for (int k=0; k<n; k++)
    sum += a[k] * b[k];
  return sum;
}</pre>
```



Be Careful What You Wish For...



- You should choose chunk sizes that are multiples of the SIMD length
 - → Remainder loops are not triggered
 - \rightarrow Likely better performance

Be Careful What You Wish For...



- You should choose chunk sizes that are multiples of the SIMD length
 - → Remainder loops are not triggered
 - → Likely better performance



Be Careful What You Wish For...



- You should choose chunk sizes that are multiples of the SIMD length
 - → Remainder loops are not triggered
 - → Likely better performance
- In the above example ...
 - \rightarrow and AVX2, the code will only execute the remainder loop!
 - \rightarrow and SSE, the code will have one iteration in the SIMD loop plus one in the remainder loop!

OpenMP 4.5 Simplifies SIMD Chunks

- Chooses chunk sizes that are multiples of the SIMD length
 First and last chunk may be slightly different to fix alignment and to handle loops that are not exact multiples of SIMD width
 - →Remainder loops are not triggered
 - →Likely better performance

OpenMP 4.5 Simplifies SIMD Chunks

Chooses chunk sizes that are multiples of the SIMD length

 \rightarrow First and last chunk may be slightly different to fix alignment and to handle loops that are

not exact multiples of SIMD width

→Remainder loops are not triggered

→Likely better performance



```
float min(float a, float b) {
   return a < b ? a : b;
float distsq(float x, float y) {
   return (x - y) * (x - y);
void example() {
#pragma omp parallel for simd
   for (i=0; i<N; i++) {
        d[i] = min(distsq(a[i], b[i]), c[i]);
```



 Declare one or more functions to be compiled for calls from a SIMD-parallel loop

Syntax (C/C++):

#pragma omp declare simd [clause[[,] clause],...]

[#pragma omp declare simd [clause[[,] clause],...]]

[...]

function-definition-or-declaration

Syntax (Fortran):

!\$omp declare simd (proc-name-list)



```
#pragma omp declare simd
float min(float a, float b) {
   return a < b ? a : b;
#pragma omp declare simd
float distsq(float x, float y) {
   return (x - y) * (x - y);
void example() {
#pragma omp parallel for simd
   for (i=0; i<N; i++) {
        d[i] = min(distsq(a[i], b[i]), c[i]);
```



OpenMP Tutorial Members of the OpenMP Language Committee



simdlen (length)

 \rightarrow generate function to support a given vector length



simdlen (length)

→ generate function to support a given vector length

uniform (argument-list)

 \rightarrow argument has a constant value between the iterations of a given loop



simdlen (length)

→ generate function to support a given vector length

uniform (argument-list)

 \rightarrow argument has a constant value between the iterations of a given loop

inbranch

 \rightarrow function always called from inside an if statement

notinbranch

 \rightarrow function never called from inside an if statement



simdlen (length)

→ generate function to support a given vector length

uniform (argument-list)

 \rightarrow argument has a constant value between the iterations of a given loop

inbranch

→ function always called from inside an if statement

notinbranch

 \rightarrow function never called from inside an if statement

- linear (argument-list[:linear-step])
- aligned (argument-list[:alignment])
inbranch & notinbranch





inbranch & notinbranch



inbranch & notinbranch



SIMD Constructs & Performance



M.Klemm, A.Duran, X.Tian, H.Saito, D.Caballero, and X.Martorell. Extending OpenMP with Vector Constructs for Modern Multicore SIMD Architectures. In Proc. of the Intl. Workshop on OpenMP, pages 59-72, Rome, Italy, June 2012. LNCS 7312.



Programming OpenMP Hands-on Exercises

Christian Terboven

Michael Klemm Yun (Helen) He





Exercises



- We have implemented a series of small hands-on examples that you can use and play with.
 Download: git clone https://github.com/NERSC/openmp-series-2024
 - Subfolder: Session-3-NUMA_SIMD/exercises, with instructions in Exercises_OMP_2024.pdf
 - → Build: make (or follow README files)
 - \rightarrow You can then find the compiled executable to run with sample Slurm commands
 - \rightarrow We use the GCC compiler mostly
- Each hands-on exercise has a folder "solution"
 It shows the OpenMP directive that we have added
 - \rightarrow You can use it to cheat \odot , or to check if you came up with the same solution

Exercises: Overview



Exercise no.	Exercise name	OpenMP Topic	Day / Order (proposal)
1	PI	Apply OpenMP SIMD	Third day
2	xthi	Review for NUMA	Third day
3	Stream	Optimize / review for NUMA	Third day
4	Jacobi	Optimize / review for NUMA	Third day