

CME 213, ME 339—Spring 2021

Eric Darve, ICME, Stanford



“Controlling complexity is the essence of computer programming.” (Brian Kernigan)



C++ threads are great for low-level multicore programming
But too general and complicated for engineering applications



Two common scenarios

1. For loop: partition the loop into chunks
Have each thread process one chunk.
2. Hand-off a block of code to a separate thread

OpenMP simplifies the programming significantly.

In many cases, adding one line is sufficient to make it run in parallel.

OpenMP is the standard approach in scientific computing for multicore processors

What is OpenMP?

Application Programming Interface (API)

Jointly defined by a group of major computer hardware and software vendors

Portable, scalable model for developers of shared memory parallel applications

Supports C/C++ and Fortran on a wide variety of computers



OpenMP website

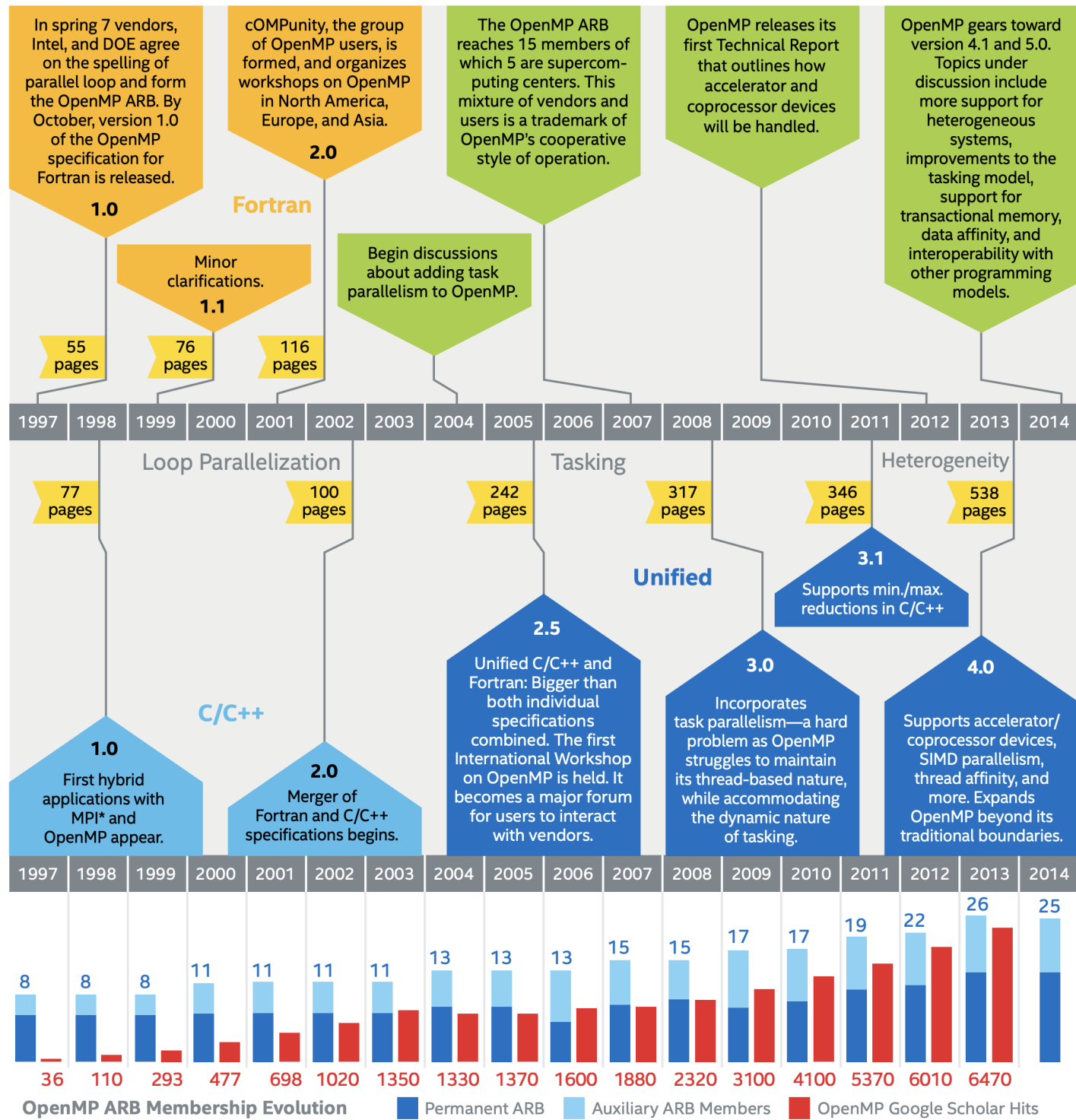
<https://openmp.org>

Wikipedia

<https://en.wikipedia.org/wiki/OpenMP>

LLNL tutorial

<https://computing.llnl.gov/tutorials/openMP/>



[Parallel Universe Magazine Issue 18](#)

[Download PDF of paper](#)

Compiling your code

Header file:

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

Compiler	Flag
gcc	
g++	
g77	-fopenmp
gfortran	
icc	
icpc	-openmp
ifort	

Installation on macOS

Option 1: libomp

```
$ brew install libomp
```

Use the system compiler `/usr/bin/g++`

Compile with options

```
-I/usr/local/include -Xpreprocessor -fopenmp  
-L/usr/local/lib -lomp
```

Option 2: gcc

```
$ brew install gcc
```

```
Compiler: /usr/local/bin/g++-9
```

```
Flag: -fopenmp
```


Which version of openMP do you have?

This determines the set of features available

```
$ echo | cpp -I/usr/local/include -Xpreprocessor -fopenmp -dM | grep OPENMP  
#define _OPENMP 201511
```

Mapping

Year	OpenMP version
200505	2.5
200805	3.0
201107	3.1
201307	4.0
201511	4.5
201811	5.0

Parallel regions

```
#pragma omp parallel
```

This is the basic building block

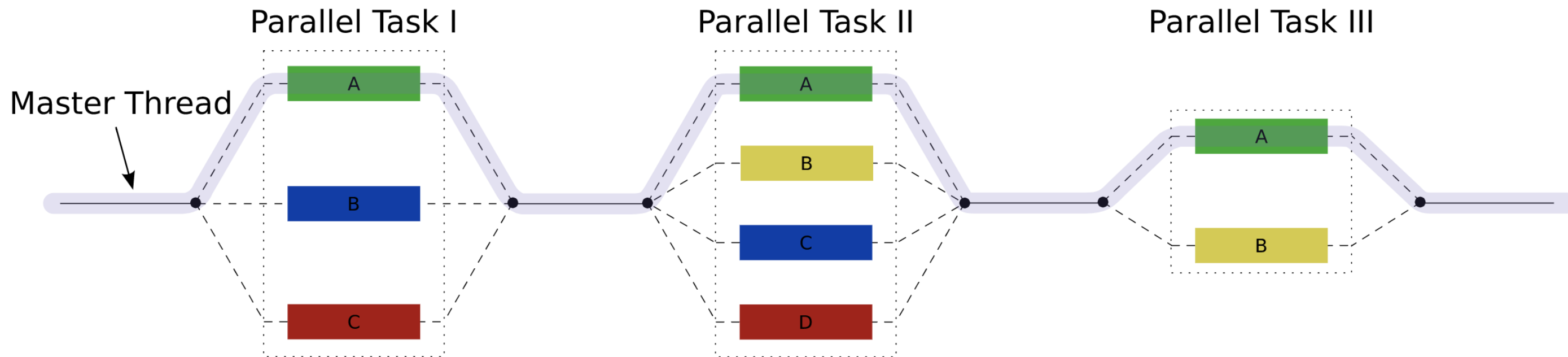
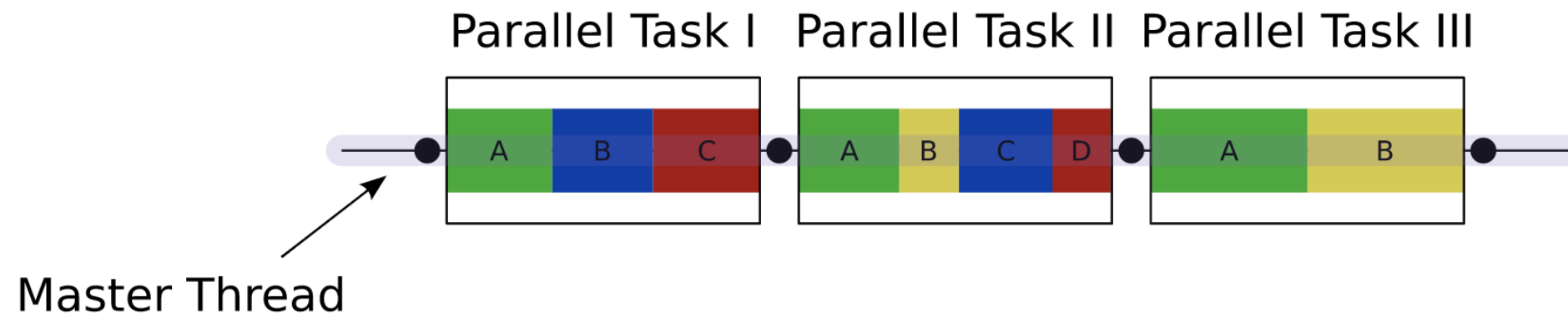
This is not how OpenMP is used in most cases

Serves simply as an introduction



```
#pragma omp parallel
```

The block of code that follows is executed
by all threads in the team



This is called the
fork-join model

Computing π

```
3. 1415926535897932384626433832795028841971693993751058209749445
923078164062862089986280348253421170679821480865132823066470938
44609550582231725359408128481117450284102701938521105559644
6229489549303819644288109756659334461284756482337867831652
7120190914564856692346034861045432664821339360726024914127
3724587006606315588174881520920962829254091715364367892590
3600113305305488204665213841469519415116094330572703657595
9195309218611738193261179310511854807446237996274956735188575
2724891227938183011949129833673362440656643086021394946395224
73719070217986094370277053921717629317675238467481846766940513
20005681271452635608277857713427577896091736371787214684409012
24953430146549585371050792279689258923542019956112129021960864
03441815981362977477130996051870721134999999837297804995105973
17328160963185950244594553469083026425223082533446850352619311
88171010003137838752886587533208381420617177669147303598253490
42875546873115956286388235378759375195778185778053217122680661
30019278766111959092164201989380952572010654858632788659361533
81827968230301952035301852968995773622599413891249721775283479
13151557485724245415069595082953311686172785588907509838175463
74649393192550604009277016711390098488240128583616035637076601
04710181942955596198946767837449448255379774726847104047534646
2080466842590694912933136770289891521047521620569660240580381
5019351125338243003558764024749647326391419927260426992279678
2354781636009341721641219924586315030286182974555706749838505
4945885869269956909272107975093029553211653449872027559602364
806654991198818347977535663698074265425278625518184175746728
90977727938000816470600161452491921732172147723501414419735
685481613611573525521334757418494684385233239073941433345477
6241686251898356948556209921922218427255025425688767179049460
16534668049886272327917860857843838279679766814541009538837863
609506800642251252051173929848960841284886269456042419652850222
106611863067442786220391949450471237137869609563643719172874677
```

[hello_world_openmp.cpp](#)


```
#pragma omp parallel num_threads(nthreads)
{
    long tid = omp_get_thread_num();
    // Only thread 0 does this
    if (tid == 0)
    {
        int n_threads = omp_get_num_threads();
        printf("[info] Number of threads = %d\n", n_threads);
    }
    // Print the thread ID
    printf("Hello World from thread = %ld\n", tid);

    // Compute digits of pi
    DoWork(tid, ndigits[tid], etime[tid]);
}
// All threads join the master thread and terminate
```

Choose your compiler in [Makefile](#)

```
$ make
```

```
$ ./hello_world_openmp
```

Sample output

```
Let's compute pi = 3.1415926535897932384626433832795028841...
[info] Number of threads = 8
Hello World from thread = 0
Thread 0 approximated Pi as 3141592653589793
Hello World from thread = 1
Thread 1 approximated Pi as 314159265358979323846264
...
Thread 7 approximated Pi as 31415926535897932384626433832795028841...
Thread 0 computed 16 digits of pi in 67 musecs ( 4.188 musec per digit)
Thread 1 computed 24 digits of pi in 16 musecs ( 0.667 musec per digit)
...
Thread 7 computed 72 digits of pi in 68 musecs ( 0.944 musec per digit)
```

Formula used to approximate π

$$\frac{\pi}{2} = 1 + \frac{1}{3} \left(1 + \frac{2}{5} \left(1 + \frac{3}{7} \left(1 + \frac{4}{9} \left(1 + \dots \right) \right) \right) \right)$$

Common use case: for loop

This example cover 99% of the needs for scientific computing



for loop_omp.cpp

```
#pragma omp parallel for  
  for (int i = 0; i < n; ++i)  
    z[i] = x[i] + y[i];
```

Exercise

[matrix_prod_openmp.cpp](#)

Parallelize the matrix-matrix product

Experiment with different options

```
./matrix_prod_openmp -p PROC
```

PROC number of threads to use

for loops can be scheduled in different ways by the library

```
#pragma omp for schedule(kind,chunk_size)
```

kind: static, dynamic, guided

static

chunk size is fixed (chunk_size)

round-robin assignment

0, 1, 2, 3, 0, 1, 2, 3, ...

Pro: low overhead

Con: assumes that running time is the same for all chunks

dynamic

Each thread executes a chunk

Then, requests another chunk until none remain

Pro: low overhead, adapts to threads that run at different speeds

Con: last thread may terminate long after the others

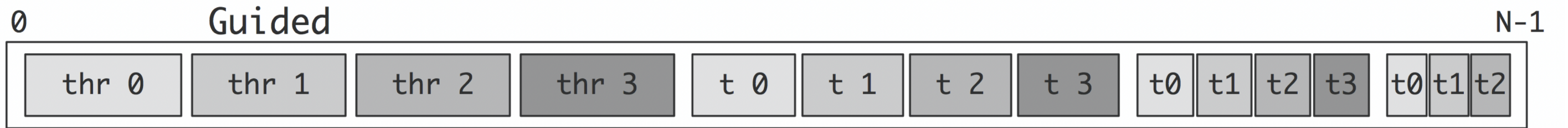
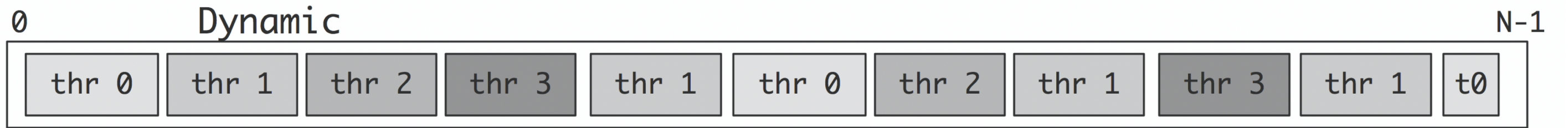
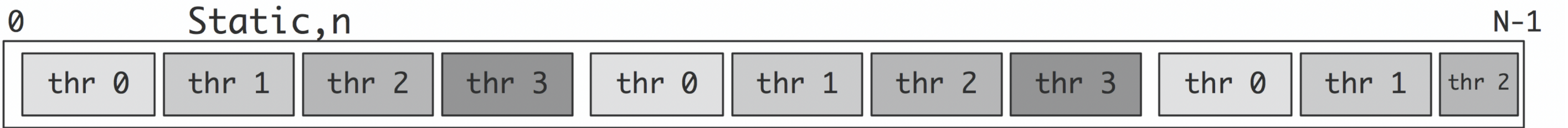
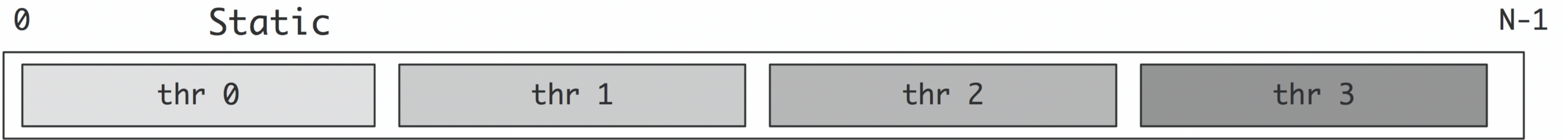
guided

Chunk size is different for each chunk

Each successive chunk is smaller than the last

Pro: all threads tend to finish at the same time

Con: high overhead for scheduling



iteration number →

nowait

```
#pragma omp parallel
{
    #pragma omp for nowait
    for (i=1; i<n; i++)
        b[i] = (a[i] + a[i-1]) / 2.0;
    #pragma omp for nowait
    for (i=0; i<m; i++)
        y[i] = sqrt(z[i]);
}
```

collapse

```
#pragma omp for collapse(2) private(i, k, j)
  for (k=kl; k<=ku; k+=ks)
    for (j=jl; j<=ju; j+=js)
      for (i=il; i<=iu; i+=is)
        bar(a,i,j,k);
```

OpenMP clause

Recall in C++ threads:

Variables passed as argument to a thread are **shared**

Variables **inside the function** that a thread is executing are **private** to that thread

OpenMP makes some reasonable default choices
But they can be changed using `shared` and `private`

shared_private_openmp.cpp

Variables declared before the block are shared

```
int shared_int = -1;
#pragma omp parallel
{
    printf("Thread ID %2d | shared_int = %d\n", omp_get_thread_num(),
          shared_int);
}
```

private clause

```
int is_private = -2;

#pragma omp parallel private(is_private)
{
    const int rand_tid = rand();
    is_private = rand_tid;
    printf("Thread ID %2d | is_private = %d\n", omp_get_thread_num(),
          is_private);
    assert(is_private == rand_tid);
}
```

Data sharing attribute clause

Most common:

- `shared(list)`
- `private(list)`

Less common: `firstprivate`, `lastprivate`, `linear`

`firstprivate`: private variable; initialized using value when construct is encountered

`lastprivate`: set equal to the private version of whichever thread executes the final iteration

`linear`: see p. 25 of [specifications](#); variable is private and has a linear relationship with respect to the iteration space of a loop associated with the construct



<https://www.openmp.org/spec-html/5.1/openmp.html>

[PDF](#)