

PHY1610 - High Performance Scientific Computing with OpenMP, part 2

Ramses van Zon, Marcelo Ponce

March 25, 2021

Section 1

Reductions

- Dot product of two vectors
- Start from a serial implementation, then will add OpenMP
- Program tells answer, correct answer, time.

$$n = \vec{x} \cdot \vec{y} = \sum_i x_i y_i$$

```
// ndot_main.cc
#include <iostream>
#include <rarray>
#include "ticktock.h"
double ndot(const rarray<double,1>& x,
            const rarray<double,1>& y);
int main()
{
    const int n = 20*1000*1000;
    rarray<double,1> x(n), y(n);
    for (int i=0; i<n; i++)
        x[i]=y[i]=i;
    double nn = n;
    double ans = (nn-1)*nn*(2*nn-1)/6;
    TickTock tt;
    tt.tick();
    double dot = ndot(x,y);
    std::cout << "Dot product: " << dot << "\n"
              << "Exact answer: " << ans << "\n";
    tt.tock("Took");
}
```

```
// serial_ndot.cc
#include <rarray>
#include <algorithm>
double ndot(const rarray<double,1>& x,
            const rarray<double,1>& y)
{
    const int n = std::min(x.size(), y.size());
    double tot=0;
    for (int i=0; i<n; i++)
        tot += x[i] * y[i];
    return tot;
}
```

```
// ndot_main.cc
#include <iostream>
#include <rarray>
#include "ticktock.h"
double ndot(const rarray<double,1>& x,
            const rarray<double,1>& y);
int main()
{
    const int n = 20*1000*1000;
    rarray<double,1> x(n), y(n);
    for (int i=0; i<n; i++)
        x[i]=y[i]=i;
    double nn = n;
    double ans = (nn-1)*nn*(2*nn-1)/6;
    TickTock tt;
    tt.tick();
    double dot = ndot(x,y);
    std::cout << "Dot product: " << dot << "\n"
              << "Exact answer: " << ans << "\n";
    tt.tock("Took");
}
```

```
// serial_ndot.cc
#include <rarray>
#include <algorithm>
double ndot(const rarray<double,1>& x,
            const rarray<double,1>& y)
{
    const int n = std::min(x.size(), y.size());
    double tot=0;
    for (int i=0; i<n; i++)
        tot += x[i] * y[i];
    return tot;
}
```

```
$ make serial_ndot
$ ./serial_ndot
Dot product: 2.66667e+21
Exact answer: 2.66667e+21
Took 0.1055 sec
$
```

- We could clearly parallelize the loop.
- We could make `tot` shared, then all threads can add to it.

- We could clearly parallelize the loop.
- We could make `tot` shared, then all threads can add to it.

```
// omp_ndot_race.cc
#include <rarray>
#include <algorithm>
double ndot(const rarray<double,1>& x,
            const rarray<double,1>& y) {
    const int n = std::min(x.size(), y.size());
    double tot=0;
    #pragma omp parallel for default(none) shared(tot,x,y)
    for (int i=0; i<n; i++)
        tot += x[i] * y[i];
    return tot;
}
```

- We could clearly parallelize the loop.
- We could make `tot` shared, then all threads can add to it.

```
// omp_ndot_race.cc
#include <rarray>
#include <algorithm>
double ndot(const rarray<double,1>& x,
            const rarray<double,1>& y) {
    const int n = std::min(x.size(), y.size());
    double tot=0;
    #pragma omp parallel for default(none) shared(tot,x,y)
    for (int i=0; i<n; i++)
        tot += x[i] * y[i];
    return tot;
}
```

```
$ make omp_ndot_race
$ export OMP_NUM_THREADS=16
$ ./omp_ndot_race
Dot product: 2.64925e+20
Exact answer: 2.66667e+21
Took 0.5431 sec
$ ./omp_ndot_race
Dot product: 2.62621e+20
Exact answer: 2.66667e+21
Took 0.5383 sec
```


- We could clearly parallelize the loop.
- We could make `tot` shared, then all threads can add to it.

```
// omp_ndot_race.cc
#include <rarray>
#include <algorithm>
double ndot(const rarray<double,1>& x,
            const rarray<double,1>& y) {
    const int n = std::min(x.size(), y.size());
    double tot=0;
    #pragma omp parallel for default(none) shared(tot,x,y)
    for (int i=0; i<n; i++)
        tot += x[i] * y[i];
    return tot;
}
```

```
$ make omp_ndot_race
$ export OMP_NUM_THREADS=16
$ ./omp_ndot_race
Dot product: 2.64925e+20
Exact answer: 2.66667e+21
Took 0.5431 sec
$ ./omp_ndot_race
Dot product: 2.62621e+20
Exact answer: 2.66667e+21
Took 0.5383 sec
```

Wrong answer!

- We could clearly parallelize the loop.
- We could make tot shared, then all threads can add to it.

```
// omp_ndot_race.cc
#include <rarray>
#include <algorithm>
double ndot(const rarray<double,1>& x,
            const rarray<double,1>& y) {
    const int n = std::min(x.size(), y.size());
    double tot=0;
    #pragma omp parallel for default(none) shared(tot,x,y)
    for (int i=0; i<n; i++)
        tot += x[i] * y[i];
    return tot;
}
```

```
$ make omp_ndot_race
$ export OMP_NUM_THREADS=16
$ ./omp_ndot_race
Dot product: 2.64925e+20
Exact answer: 2.66667e+21
Took 0.5431 sec
$ ./omp_ndot_race
Dot product: 2.62621e+20
Exact answer: 2.66667e+21
Took 0.5383 sec
```

Wrong answer!

Answer varies!

- We could clearly parallelize the loop.
- We could make `tot` shared, then all threads can add to it.

```
// omp_ndot_race.cc
#include <rarray>
#include <algorithm>
double ndot(const rarray<double,1>& x,
            const rarray<double,1>& y) {
    const int n = std::min(x.size(), y.size());
    double tot=0;
    #pragma omp parallel for default(none) shared(tot,x,y)
    for (int i=0; i<n; i++)
        tot += x[i] * y[i];
    return tot;
}
```

```
$ make omp_ndot_race
$ export OMP_NUM_THREADS=16
$ ./omp_ndot_race
Dot product: 2.64925e+20
Exact answer: 2.66667e+21
Took 0.5431 sec
$ ./omp_ndot_race
Dot product: 2.62621e+20
Exact answer: 2.66667e+21
Took 0.5383 sec
```

Wrong answer!

Answer varies!

Slower computation!

- Can be very subtle, and only appear intermittently.
- Your program can have a bug but not display any symptoms for small runs!
- Primarily a problem with shared memory.
- Classical parallel bug.
- Multiple writers to some shared resource.

Say, initially, $tot=0$, and one threads want to add 1 to it while a second thread want to add 2 at the same time.

- The correct answer for tot is, clearly, three.
- However, we may see any of the answers 1, 2, or 3.

How does this issue arise?

Say, initially, $tot=0$, and one threads want to add 1 to it while a second thread want to add 2 at the same time.

- The correct answer for tot is, clearly, three.
- However, we may see any of the answers 1, 2, or 3.

How does this issue arise?

Non-atomic adding and updating

Say, initially, $tot=0$, and one threads want to add 1 to it while a second thread want to add 2 at the same time.

- The correct answer for tot is, clearly, three.
- However, we may see any of the answers 1, 2, or 3.

How does this issue arise?

Non-atomic adding and updating

Thread 0: add 1	Thread 1: add 2
-----------------	-----------------

Say, initially, $tot=0$, and one threads want to add 1 to it while a second thread want to add 2 at the same time.

- The correct answer for tot is, clearly, three.
- However, we may see any of the answers 1, 2, or 3.

How does this issue arise?

Non-atomic adding and updating

Thread 0: add 1	Thread 1: add 2
read $tot=0$ to $reg0$.

Say, initially, $tot=0$, and one threads want to add 1 to it while a second thread want to add 2 at the same time.

- The correct answer for tot is, clearly, three.
- However, we may see any of the answers 1, 2, or 3.

How does this issue arise?

Non-atomic adding and updating

Thread 0: add 1

read $tot=0$ to $reg0$

$reg0 = reg0+1$

Thread 1: add 2

.

read $tot=0$ to $reg1$

Say, initially, $tot=0$, and one threads want to add 1 to it while a second thread want to add 2 at the same time.

- The correct answer for tot is, clearly, three.
- However, we may see any of the answers 1, 2, or 3.

How does this issue arise?

Non-atomic adding and updating

Thread 0: add 1	Thread 1: add 2
read $tot=0$ to $reg0$.
$reg0 = reg0 + 1$	read $tot=0$ to $reg1$
store $reg0(=1)$ in tot	$reg1 = reg1 + 2$

Say, initially, $tot=0$, and one threads want to add 1 to it while a second thread want to add 2 at the same time.

- The correct answer for tot is, clearly, three.
- However, we may see any of the answers 1, 2, or 3.

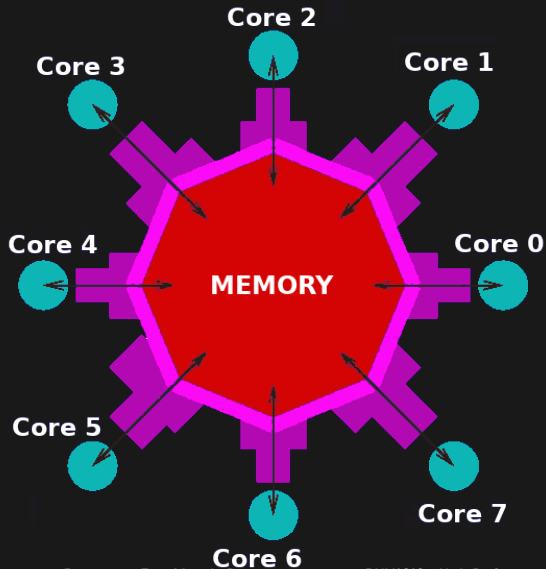
How does this issue arise?

Non-atomic adding and updating

Thread 0: add 1	Thread 1: add 2
read $tot=0$ to $reg0$.
$reg0 = reg0 + 1$	read $tot=0$ to $reg1$
store $reg0(=1)$ in tot	$reg1 = reg1 + 2$
.	store $reg1(=2)$ in tot

You might think the parallel version should at least still be faster, though it may be wrong. But even that's not the case.

- Here, multiple cores repeatedly try to read, access and store the same variable in memory.
- This means the shared variable that is updated in a register, cannot stay in register: It has to be copied back to main memory, so the other threads see it correctly.
- The other threads then have to re-read the variable.
- This write-back would not be necessary if the variable was shared but not written to.



- Memory is layered: between registers and shared main memory there are further layers called **caches**.
- Caches are faster but more expensive and therefore smaller. They are like private memory for each core.
- Main memory is the slowest part of the memory.
- Caches are automatically kept coherent between cores.

Section 2

Fixing the race condition

Our code get it wrong because different threads are updating the tot variable at the same time.

The `critical` construct:

- Defines a critical region.
- Only one thread can be operating within this region at a time.
- Keeps modifications to shared resources safe.

Our code get it wrong because different threads are updating the tot variable at the same time.

The critical construct:

- Defines a critical region.
- Only one thread can be operating within this region at a time.
- Keeps modifications to shared resources safe.

```
// omp_ndot_critical.cc
#include <rarray>
#include <algorithm>
double ndot(const rarray<double,1>& x,
            const rarray<double,1>& y)
{
    const int n = std::min(x.size(), y.size());
    double tot=0;
    #pragma omp parallel for default(none) shared(tot,x,y)
    for (int i=0; i<n; i++)
        #pragma omp critical
        tot += x[i] * y[i];
    return tot;
}
```



```
// omp_ndot_critical.cc
#include <rarray>
#include <algorithm>
double ndot(const rarray<double,1>& x,
            const rarray<double,1>& y)
{
    const int n = std::min(x.size(), y.size());
    double tot=0;
    #pragma omp parallel for default(none) shared(tot,x,y)
    for (int i=0; i<n; i++)
        #pragma omp critical
        tot += x[i] * y[i];
    return tot;
}
```

```
// omp_ndot_critical.cc
#include <rarray>
#include <algorithm>
double ndot(const rarray<double,1>& x,
            const rarray<double,1>& y)
{
    const int n = std::min(x.size(), y.size());
    double tot=0;
    #pragma omp parallel for default(none) shared(tot,x,y)
    for (int i=0; i<n; i++)
        #pragma omp critical
        tot += x[i] * y[i];
    return tot;
}
```

```
$ make omp_ndot_critical
$ export OMP_NUM_THREADS=16
$ ./omp_ndot_critical
Dot product:  2.66667e+21
Exact answer: 2.66667e+21
Took 4.6697 sec
```

Correct, but 44× slower than serial version!

- Most hardware has support for atomic instructions (indivisible so cannot get interrupted)
- Small subset, but load/add/store usually in it.
- Not as general as critical
- Much lower overhead.
- `#pragma omp atomic [read|write|update|capture]`

```
// omp_ndot_atomic.cc
#include <rarray>
#include <algorithm>
double ndot(const rarray<double,1>& x,
            const rarray<double,1>& y)
{
    const int n = std::min(x.size(), y.size());
    double tot=0;
    #pragma omp parallel for default(none) shared(tot,x,y)
    for (int i=0; i<n; i++)
        #pragma omp atomic update
        tot += x[i] * y[i];
    return tot;
}
```

```
// omp_ndot_atomic.cc
#include <rarray>
#include <algorithm>
double ndot(const rarray<double,1>& x,
            const rarray<double,1>& y)
{
    double tot=0;
    #pragma omp parallel for default(none) shared(tot,n,x,y)
    for (int i=0; i<n; i++)
        #pragma omp atomic update
        tot += x[i] * y[i];
    return tot;
}
```

```
// omp_ndot_atomic.cc
#include <rarray>
#include <algorithm>
double ndot(const rarray<double,1>& x,
            const rarray<double,1>& y)
{
    double tot=0;
    #pragma omp parallel for default(none) shared(tot,n,x,y)
    for (int i=0; i<n; i++)
        #pragma omp atomic update
        tot += x[i] * y[i];
    return tot;
}
```

```
$ make omp_ndot_atomic
$ export OMP_NUM_THREADS=16
$ ./omp_ndot_atomic
Dot product: 2.66667e+21
Exact answer: 2.66667e+21
Took 2.177 sec
```

About twice faster than critical, but still not great.

The issue we have not resolved is that we're still updating `tot`, which causes copies to main memory at every iteration.

What if we accumulated `tot` for each core, and sum them up later?

```
double ndot(const rarray<double,1>& x,
            const rarray<double,1>& y)
{
    const int n = std::min(x.size(), y.size());
    double tot=0;
    #pragma omp parallel default(none) shared(tot,x,y)
    {
        double localtot=0;
        #pragma omp for
        for (int i=0; i<n; i++)
            localtot += x[i] * y[i];
        #pragma omp atomic update
        tot += localtot;
    }
    return tot;
}
```

The issue we have not resolved is that we're still updating `tot`, which causes copies to main memory at every iteration.

What if we accumulated `tot` for each core, and sum them up later?

```
double ndot(const rarray<double,1>& x,
            const rarray<double,1>& y)
{
    const int n = std::min(x.size(), y.size());
    double tot=0;
    #pragma omp parallel default(none) shared(tot,x,y)
    {
        double localtot=0;
        #pragma omp for
        for (int i=0; i<n; i++)
            localtot += x[i] * y[i];
        #pragma omp atomic update
        tot += localtot;
    }
    return tot;
}
```

```
$ export OMP_NUM_THREADS=16
$ ./omp_ndot_local
Dot product: 2.66667e+21
Exact answer: 2.66667e+21
Took 0.01715 sec
```

Correct answer, 6x faster!

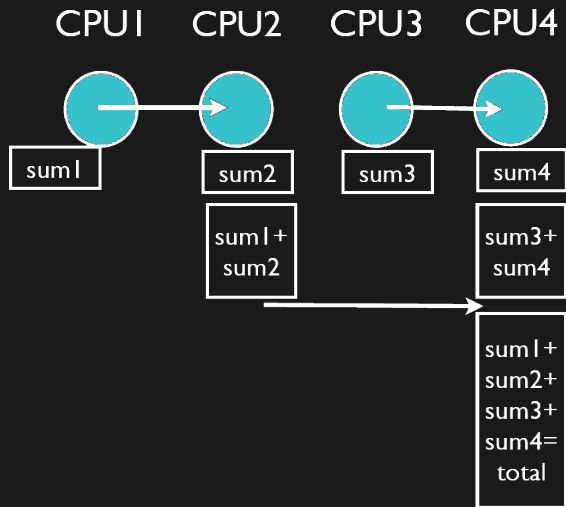
- What we did is quite common, taking a bunch of data and summing it to one value: **reduction**

- What we did is quite common, taking a bunch of data and summing it to one value: **reduction**
- OpenMP supports this using **reduction variables**.

- What we did is quite common, taking a bunch of data and summing it to one value: **reduction**
- OpenMP supports this using **reduction variables**.
- When declaring a variables as reduction variables, private copies are made (much as for private variables), which are combined at the end of a parallel region through some operation (+, *, min, max).

- What we did is quite common, taking a bunch of data and summing it to one value: **reduction**
- OpenMP supports this using **reduction variables**.
- When declaring a variables as reduction variables, private copies are made (much as for private variables), which are combined at the end of a parallel region through some operation (+, *, min, max).
- `omp_ndot_reduction.cc`

- What we did is quite common, taking a bunch of data and summing it to one value: **reduction**
- OpenMP supports this using **reduction variables**.
- When declaring a variables as reduction variables, private copies are made (much as for private variables), which are combined at the end of a parallel region through some operation (+, *, min, max).
- `omp_ndot_reduction.cc`



```
// omp_ndot_reduction.cc
#include <rarray>
#include <algorithm>
double ndot(const rarray<double,1>& x,
            const rarray<double,1>& y)
{
    const int n = std::min(x.size(), y.size());
    double tot=0;
    #pragma omp for default(none) shared(x,y) reduction(+:tot)
    for (int i=0; i<n; i++)
        tot += x[i] * y[i];
    return tot;
}
```

```
// omp_ndot_reduction.cc
#include <rarray>
#include <algorithm>
double ndot(const rarray<double,1>& x,
            const rarray<double,1>& y)
{
    const int n = std::min(x.size(), y.size());
    double tot=0;
    #pragma omp for default(none) shared(x,y) reduction(+:tot)
    for (int i=0; i<n; i++)
        tot += x[i] * y[i];
    return tot;
}
```

```
$ make omp_ndot_reduction
$ export OMP_NUM_THREADS=8
$ ./omp_ndot_reduction
Dot product: 2.66667e+21
Exact answer: 2.66667e+21
Took 0.01691 sec
$
```

```
// omp_ndot_reduction.cc
#include <rarray>
#include <algorithm>
double ndot(const rarray<double,1>& x,
            const rarray<double,1>& y)
{
    const int n = std::min(x.size(), y.size());
    double tot=0;
    #pragma omp for default(none) shared(x,y) reduction(+:tot)
    for (int i=0; i<n; i++)
        tot += x[i] * y[i];
    return tot;
}
```

```
$ make omp_ndot_reduction
$ export OMP_NUM_THREADS=8
$ ./omp_ndot_reduction
Dot product: 2.66667e+21
Exact answer: 2.66667e+21
Took 0.01691 sec
$
```

Correct, same timing as local sums, but simpler code.

As we saw in the random-number lecture, a **stream** of pseudo random numbers can be generated by:

- having some state, initialized using a seed
- advancing a state through an algorithm
- computing the random number from that state (adjusted for particular probability distribution)

This means that the n th random number that is generated depends on having the n compute.

It would seem that we cannot parallelize this to generate multiple random numbers at the same time.

Keeping 1 stream

- Produce in one thread and distribute to threads.
- Skip ahead: works only for some algorithms
- Interleave

Low performance, may need to know number of threads and number of draws.

Keeping 1 stream

- Produce in one thread and distribute to threads.
- Skip ahead: works only for some algorithms
- Interleave

Low performance, may need to know number of threads and number of draws.

Separate streams per thread

- Different seeds.
- Different parameters for same kind of rng.
- Different kinds of rng.

Why we can do this is that we only care that the streams look random and mutually independent.

Harder to ensure consistency with serial code.

Even more important to have good RNGs.

```
// serial computation of pi
#include <iostream>
#include <random>
#include <climits>
using namespace std;

double uniform() {
    static mt19937 rng(13);
    static uniform_real_distribution<double> dist(0,1);
    return dist(rng);
}

int main(int argc, char* argv[]) {
    long tries = (argc<2)?10000:atoi(argv[1]);
    long count = 0;

    for (long i = 0; i < tries; i++) {
        double x = uniform();
        double y = uniform();
        if (x*x + y*y < 1)
            count++;
    }
    double pi = (4.0*count)/tries;
    cout << "Pi is approximately " << pi << endl;
}
```

```
// serial computation of pi
#include <iostream>
#include <random>
#include <climits>
using namespace std;

double uniform() {
    static mt19937 rng(13);
    static uniform_real_distribution<double> dist(0,1);
    return dist(rng);
}

int main(int argc, char* argv[]) {
    long tries = (argc<2)?10000:atoi(argv[1]);
    long count = 0;

    for (long i = 0; i < tries; i++) {
        double x = uniform();
        double y = uniform();
        if (x*x + y*y < 1)
            count++;
    }

    double pi = (4.0*count)/tries;
    cout << "Pi is approximately " << pi << endl;
}
```

```
// parallel computation of pi
#include <iostream>
#include <random>
#include <climits>
using namespace std;
mt19937 srng(13);
uniform_int_distribution<unsigned> sdist(0,UINT_MAX);
double uniform() {
    thread_local mt19937 rng(sdist(srng));
    thread_local uniform_real_distribution<double> dist(0,1);
    return dist(rng);
}

int main(int argc, char* argv[]) {
    long tries = (argc<2)?10000:atoi(argv[1]);
    long count = 0;
    #pragma omp parallel for shared(tries) reduction(+:count)
    for (long i = 0; i < tries; i++) {
        double x = uniform();
        double y = uniform();
        if (x*x + y*y < 1)
            count++;
    }

    double pi = (4.0*count)/tries;
    cout << "Pi is approximately " << pi << endl;
}
```

static (C)

This keyword has (too?) many different meanings in C++:

- Before variables declared inside a function, the variable is initialized upon the first call of the function, and reused in subsequent calls.
- Before variables declared inside a class, the variable is shared among all instances of the class.
- Before functions inside a class, it means these should be called with the class name and can only use static class variables.
- Before variables or functions declared globally inside a file, the variable is only known inside that file.

static (C)

This keyword has (too?) many different meanings in C++:

- Before variables declared inside a function, the variable is initialized upon the first call of the function, and reused in subsequent calls.
- Before variables declared inside a class, the variable is shared among all instances of the class.
- Before functions inside a class, it means these should be called with the class name and can only use static class variables.
- Before variables or functions declared globally inside a file, the variable is only known inside that file.

thread_local (C++11)

- Before variables declared inside a function, the variable is initialized per thread upon the first call of the function. Each thread will see its own copy, and reuse it in subsequent calls. This is analogous to openmp's 'private' for static variables.
- Before declarations of global variables, each thread gets an independent copy.

```
// parallel computation of pi
#include <iostream>
#include <random>
#include <climits>
using namespace std;
mt19937 srng(13);
uniform_int_distribution<unsigned> sdist(0,UINT_MAX);
double uniform() {
    thread_local mt19937 rng(sdist(srng));
    thread_local uniform_real_distribution<double> dist(0,1);
    return dist(rng);
}
int main(int argc, char* argv[]) {
    long tries = (argc<2)?10000:atoi(argv[1]);
    long count = 0;
    #pragma omp parallel for shared(tries) reduction(+:count)
    for (long i = 0; i < tries; i++) {
        double x = uniform();
        double y = uniform();
        if (x*x + y*y < 1)
            count++;
    }
    double pi = (4.0*count)/tries;
    cout << "Pi is approximately " << pi << endl;
}
```

```
// parallel computation of pi
#include <iostream>
#include <random>
#include <climits>
using namespace std;
mt19937 srng(13);
uniform_int_distribution<unsigned> sdist(0,UINT_MAX);
double uniform() {
    thread_local mt19937 rng(sdist(srng));
    thread_local uniform_real_distribution<double> dist(0,1);
    return dist(rng);
}
int main(int argc, char* argv[]) {
    long tries = (argc<2)?10000:atoi(argv[1]);
    long count = 0;
    #pragma omp parallel for shared(tries) reduction(+:count)
    for (long i = 0; i < tries; i++) {
        double x = uniform();
        double y = uniform();
        if (x*x + y*y < 1)
            count++;
    }
    double pi = (4.0*count)/tries;
    cout << "Pi is approximately " << pi << endl;
}
```

```
// parallel computation of pi
#include <iostream>
#include <random>
#include <climits>
using namespace std;
mt19937 srng(13);
uniform_int_distribution<unsigned> sdist(0,UINT_MAX);
double uniform() {
    thread_local int uninitialized = true;
    thread_local mt19937 rng;
    thread_local uniform_real_distribution<double> dist(0,1);
    if (uninitialized) {
        #pragma omp critical
        rng.seed(sdist(srng));
        uninitialized = false;
    }
    return dist(rng);
}
int main(int argc, char* argv[]) {
    long tries = (argc<2)?10000:atoi(argv[1]);
    long count = 0;
    #pragma omp parallel for shared(tries) reduction(+:count)
    for (long i = 0; i < tries; i++) {
        double x = uniform();
        double y = uniform();
        if (x*x + y*y < 1)
            count++;
    }
}
```


Section 3

Load Balancing

- So far every iteration of the loop had the same amount of work.
- Not always the case.
- Sometimes cannot predict beforehand how unbalanced the problem is

OpenMP has work sharing constructs that allow you do statically or dynamically balance the load.

- The Mandelbrot set is a boundary in the (complex) plane between a region from which point can escape and one from which they can't.
- Based on a mapping in complex plane:

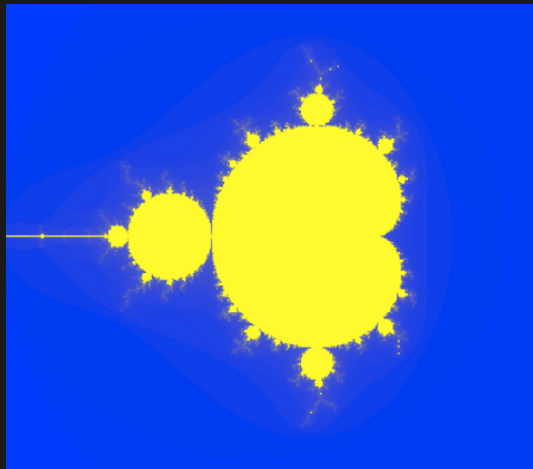
$$b_{n+1} = b_n^2 + a$$

- Mandelbrot set is boundary between diverging points ($b_0 = 0 \Rightarrow \|b_\infty\| = \infty$) and converging points ($\|b_\infty\| < \infty$).

Note: if $\|b_n\| > 2$, point diverges.

- Calculation:
 - ▶ iterate for each point a in square, see if $\|b_n\| > 2$.
 - ▶ $n < n_{\max}$, then blue, else yellow.
- On the outside points diverge quickly.

Inside points: lots of work.



```
// iterations for each point
int how_many_iter(std::complex<double> a, int maxiter);
// compute iterations for each point in a rectangle
rarray<int,2> make_mandel_map(double xmin, double xmax, double ymin,
                           double ymax, int npix, int maxiter)

// display specific stuff
char display_map(const rarray<int,2>&,float,double,double&,double&,double&,double&);
void my_pgctab(float,float,float,float,float,float,int);
// driver routine
int main();
```

Compile and run:

```
$ make mandel mandel-parallel
$ ./mandel
    5.06 sec
...
$ export OMP_NUM_THREADS=16
$ ./mandel-parallel
    1.366 sec
```

```
rarray<int,2> make_mandel_map(double xmin, double xmax,  
                             double ymin, double ymax,  
                             int npix, int maxiter) {  
    rarray<int,2> mymap(npix,npix);  
  
    for (int i=0; i<npix; i++)  
        for (int j=0; j<npix; j++) {  
            double x = ((double)i)/((double)npix)*(xmax-xmin)+xmin;  
            double y = ((double)j)/((double)npix)*(ymax-ymin)+ymin;  
            std::complex<double> a(x,y);  
            mymap[i][j] = how_many_iter(a,maxiter);  
        }  
    return mymap;  
}
```

```
rarray<int,2> make_mandel_map(double xmin, double xmax,
                           double ymin, double ymax,
                           int npix, int maxiter) {
    rarray<int,2> mymap(npix,npix);

    for (int i=0; i<npix; i++)
        for (int j=0; j<npix; j++) {
            double x = ((double)i)/((double)npix)*(xmax-xmin)+xmin;
            double y = ((double)j)/((double)npix)*(ymax-ymin)+ymin;
            std::complex<double> a(x,y);
            mymap[i][j] = how_many_iter(a,maxiter);
        }
    return mymap;
}
```

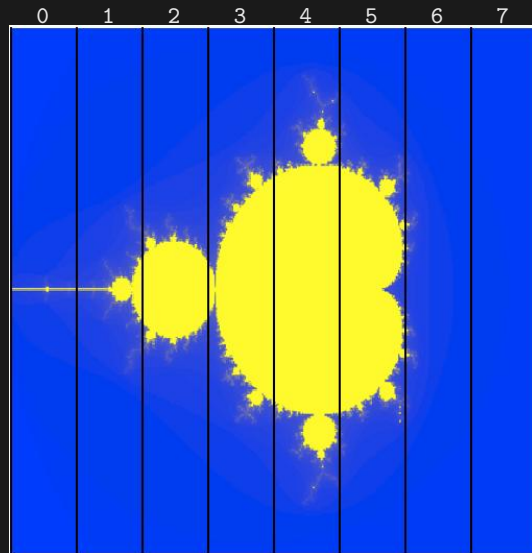
```
int how_many_iter(std::complex<double> a, int maxiter) {
    std::complex<double> b = a;
    for (int i=0; i<maxiter; i++) {
        if (std::norm(b) > 4) return i;
        b = b*b + a;
    }
    return maxiter;
}
```

```
rarray<int,2> make_mandel_map(double xmin, double xmax,
                           double ymin, double ymax,
                           int npix, int maxiter) {
    rarray<int,2> mymap(npix,npix);
    #pragma omp parallel for default(none) shared(mymap,xmin,xmax,ymin,ymax,npix,maxiter)
    for (int i=0; i<npix; i++)
        for (int j=0; j<npix; j++) {
            double x = ((double)i)/((double)npix)*(xmax-xmin)+xmin;
            double y = ((double)j)/((double)npix)*(ymax-ymin)+ymin;
            std::complex<double> a(x,y);
            mymap[i][j] = how_many_iter(a,maxiter);
        }
    return mymap;
}
```

```
int how_many_iter(std::complex<double> a, int maxiter) {
    std::complex<double> b = a;
    for (int i=0; i<maxiter; i++) {
        if (std::norm(b) > 4) return i;
        b = b*b + a;
    }
    return maxiter;
}
```

- Default work sharing breaks N iterations into $N/nthreads$ chunks and assigns them to threads.
- But threads 0, 1, 6 and 7 will be done and sitting idle while threads 2, 3, 4 and 5 work on the rest
- Inefficient use of resources.

Serial	5.060s
Nthreads=16	1.336s
Speedup	3.8x
Efficiency	24%

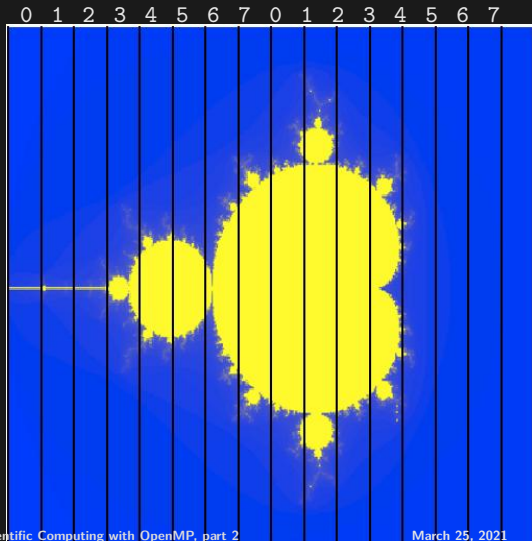


- Default: each thread gets a big consecutive chunk of the loop. Often better to give each thread many smaller interleaved chunks.
- Can add `schedule` clause to `omp for` to change work sharing.
- We can decide either at compile-time (static schedule) or run-time (dynamic schedule) how work will be split.
- `#pragma omp parallel for schedule(static, m)` gives `m` consecutive loop elements to each thread instead of a big chunk.
- With `'schedule(dynamic, m)`, each thread will work through `m` loop elements, then go to the OpenMP run-time system and ask for more.
- Load balancing (possibly) better with dynamic, but larger overhead than with static.

```
#pragma omp parallel for schedule(static,25)
```

- Can change the chunk size different from $\sim N/nthreads$
- In this case, more columns – work distributed a bit better.
- Now, for instance, thread 7 gets both a big work chunk and a little one.

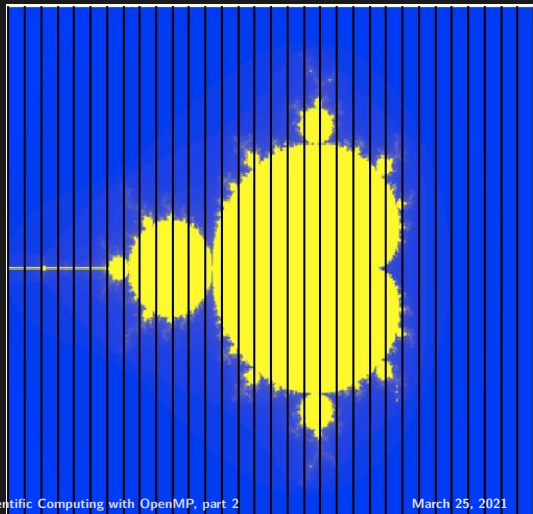
Serial	5.060s
Nthreads=16	0.7693s
Speedup	6.6x
Efficiency	41%



```
#pragma omp parallel for schedule(dynamic)
```

- Break up into many pieces and hand them to threads when they are ready.
- Dynamic scheduling.
- Increases overhead, decreases idling threads.
- Can also choose chunk size.

Serial	5.060s
Nthreads=16	0.7686s
Speedup	6.6x
Efficiency	41%



- `schedule(static)` or `schedule(dynamic)` are good starting points.
- To get best performance in badly imbalanced problems, may have to play with chunk size; depends on your problem, hardware, and compiler.

<code>static,1</code>	<code>dynamic,1</code>
0.4347s	0.4121s
11.6x	12.3x
72%	77%

There are many more features to OpenMP not discussed here.

- Collapsed loops
- Tasks
- Tasks with dependencies
- Nested OpenMP parallelism
- Locks
- SIMD
- Thread affinities
- Compute devices (e.g. NVIDIA/AMD graphics cards, Intel Xeon Phi)