

HPC Project

Parallel and Distributed Computing

Gabriele Pintus

May 27, 2024

University of Trieste

The aim of this project is to ...

Assignment 1:

Latency of CoOps

Assignment 1

Implement a C code which measures the latency of COOPs.

Assignment 2:

The Mandelbrot Set

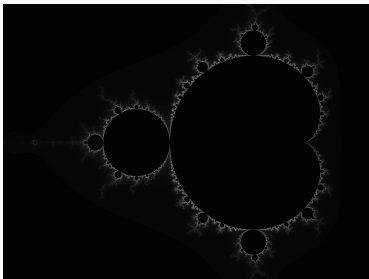
Assignment 2

Implement a hybrid C code which computes the Mandelbrot set using:

- **OpenMP**: Open Multi-Processing
- **MPI**: Message Passing Interface

Mandelbrot Set

Is the set of all complex numbers c for which the sequence defined by $z_{n+1} = z_n^2 + c$ does not diverge.



Arithmetic Intensity

Defined as the ratio of the number of arithmetic operations to the number of memory operations.

Arithmetic operations:

- 2 mult and 1 add: $z.x * z.x + z.y * z.y$
- 2 mult, 1 sub and 1 add: $temp = z.x * z.x - z.y * z.y + c.x$
- 2 mult and 1 add: $z.y = 2 * z.x * z.y + c.y$
- 1 add: $++n$

Memory operations:

- 1 write: `int n = 0`
- 2 writes: `Complex z = {0,0}`
- 2 reads: $z.x * z.x + z.y * z.y$
- 3 reads and 1 write: $temp = z.x * z.x - z.y * z.y + c.x$
- 3 reads and 1 write: $z.y = 2 * z.x * z.y + c.y$
- 1 write: `z.x = temp`
- 1 write: `++n`

The arithmetic intensity is then

$$\frac{11 \cdot N}{10 \cdot N + 3} = \frac{11}{10 + \frac{3}{N}} \xrightarrow{N \rightarrow \infty} \frac{11}{10} = 1.1 \text{ flops/byte} > 1$$

Compute-bound when the number of iterations is large.

We are not considering optimizations like:

- Intense use of registers
- Vectorization
- Cache
- Loop unrolling

Therefore, the AI could be higher.

General purpose API for parallel programming in C, C++ and Fortran.
Allow to parallelize loops, sections, tasks and more with minimal effort.

```
void mandelbrot_set(uint8_t *image) {  
    ...  
    #pragma omp parallel for schedule(dynamic, CHUNK_SIZE)  
    for (int i = 0; i < total_size; ++i) {  
        int x = i / HEIGHT;  
        int y = i % HEIGHT;  
        Complex c = {X_MIN + x * dx, Y_MIN + y * dy};  
  
        // Compute mandelbrot set in the point  
        int iter = mandelbrot(c);  
        iter = iter % MAX_ITER;  
  
        // store the value  
        image[i] = iter;  
    }  
}
```

- `#pragma omp parallel for:`
parallelizes the loop
- `schedule(dynamic, CHUNK_SIZE):`
divides the iterations in chunks of size `CHUNK_SIZE` and assigns them to threads as they become available.

Code Optimizations

The following C function `mandelbrot` computes the Mandelbrot set in a given point `c`.

```
int mandelbrot(const Complex c) {
    int n = 0;
    Complex z = {0, 0};
    while ((z.x * z.x + z.y * z.y) < 4 && n < MAX_ITER){
        double temp = z.x * z.x - z.y * z.y + c.x;
        z.y = 2 * z.x * z.y + c.y;
        z.x = temp;
        ++n;
    }
    return n;
}
```

Code Optimizations

Let's consider a couple of optimizations:

- `const Complex c`: the `const` keyword makes the variable `c` read-only. However, this also enables the compiler to pass a pointer to the variable instead of copying it.
- `Magnitude < MAX_ITER`: for most of the points the magnitude of `z` will be greater than 2 before reaching the maximum number of iterations. Once the first condition is false, we don't need to check the second one.
- `++n`: the pre-increment operator is more efficient than the post-increment one since it doesn't need to create a temporary variable.

Compiler's Optimizations

The **gcc** compiler has an optimization flag `-O{0,1,2,3}` which enables several optimizations. We compare the basic code with the level 3 optimized one by giving a look at the assembly code.

- **Less memory accesses:** the optimized code uses more intensively the CPU registers, avoiding memory accesses. This help reducing the *waiting for contention* time.

20 vs 70 commands for memory accesses, 71% less.

- **Optimized memory:** the optimized code uses the `.p2align` directive to align the memory accesses to the cache line size. This reduces the number of cache misses.

There could still be a bit of room for manual optimization.

- **movapd:** the optimized code uses more often the `movapd` instruction, which can moves two doubles at once.
- **Parallelization:** more wisely, the optimized code uses the `gomp_parallel_loop_nonmonotonic_dynamic` directive. Chunk size can vary at runtime, allowing a better load balancing.

Results

With the `perf stat` command when can measure several metrics.

Metric	Basic	Optimized
Wall-clock time	40.53 s	16.09 s
User time	316.67 s	126.71 s
System time	199.23 ms	70.34 ms
Instructions	881 billions	466 billions
Instructions per cycle	1.11	1.29
Cache references	44.743.584	43.485.500
Cache misses	0.98%	0.34%

Intel core i7-8550U, 1.8 GHz, 4 cores, 8 threads, 16 GB RAM.
OMP threads: 8, image size: 2048×2048 .

Message Passing Interface is a standard for parallel and distributed computing. Interaction between processes must be explicitly defined.

Idea: we want to distribute the load in a dynamic way as the OpenMP code does.

Solution: we define a `MPI_CHUNK_SIZE` and we build a *work queue*. Each process will get a chunk from the queue and will compute it. A new chunk will be assigned to a process as it becomes available.

This is a simple, yet effective, way to implement the dynamic scheduler, which allows to better distribute the load among the processes.

MPI - Work Queue

The queue data structure, implemented in the `sys/queue.h` library, is built by the root process.

```
// Create the work queue
work_queue = (WorkQueue *)malloc(sizeof(WorkQueue));
TAILQ_INIT(work_queue);

// Add work items to the queue
for (uint32_t i = 0; i < n_chunks; ++i){
    WorkItem *item = (WorkItem *)malloc(sizeof(WorkItem));
    item->start_idx = i * MPI_CHUNK_SIZE;
    item->end_idx = (i + 1) * MPI_CHUNK_SIZE;
    if (item->end_idx > total_size)
        item->end_idx = total_size;
    TAILQ_INSERT_TAIL(work_queue, item, entries);
}
```

MPI - Communication

We use blocking communication to send the chunks to the processes.

- `MPI_Send`: sends a message to a process
- `MPI_Recv`: receives a message from a process

A further optimization could be to use non-blocking communication to overlap computation and communication and reduce the waiting time.

The root process sends only 2×4 bytes, but receive

$$2 \times 4 + \text{MPI_CHUNK_SIZE} \times \text{image_t}$$

bytes, where `image_t` is either 1 or 2 bytes.

Hence, the root process is the bottleneck of the communication.

One could decide to use two “weak” processes, one for sending the chunks and one for receiving them. Moreover, using non-blocking communication reduce the waiting time.

The worker processes can receive a new chunk immediately after finishing the previous one. They don't need to wait for:

- The root process to end its own computation
- The root process to copy the previous chunk
- The communication to be completed

We want to measure the scalability of the code in two ways:

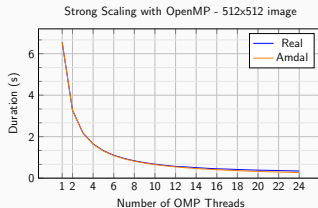
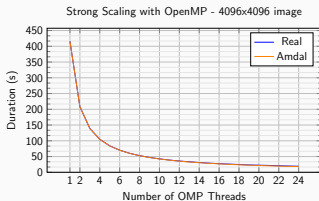
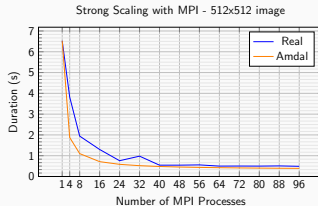
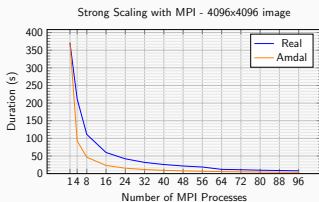
- **Strong scaling:** we fix the problem size and we increase the number of processes. We expect the execution time to decrease.
- **Weak scaling:** we increase the problem size proportionally to the number of processes. We expect the execution time to remain constant.

In both cases we perform the same test once by fixing the number of MPI processes to 1 and varying the number of OpenMP threads, and once by fixing the number of OpenMP threads to 1 and varying the number of MPI processes.

Every test is repeated for different problem sizes: 512×512 and 4096×4096 .

Strong Scaling

We compare the experimental data with the Amdal's law given $p = 1$.



The OpenMP code scales better, touching the lower bound.

Weak Scaling

Theoretically, the execution time should remain constant.

