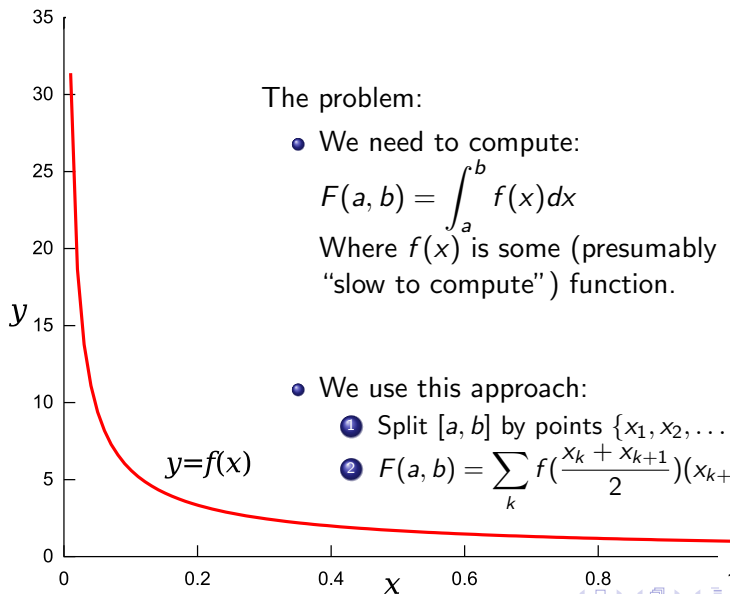


MPI profiling with Alinea MAP

Alexander Gaenko

February 5, 2016

Problem: Calculation of a definite integral.



The problem:

- We need to compute:

$$F(a, b) = \int_a^b f(x) dx$$

Where $f(x)$ is some (presumably “slow to compute”) function.

- We use this approach:

① Split $[a, b]$ by points $\{x_1, x_2, \dots, x_k\}$

②
$$F(a, b) = \sum_k f\left(\frac{x_k + x_{k+1}}{2}\right)(x_{k+1} - x_k)$$

Problem size behavior: how to measure.

Do we even need to parallelize?

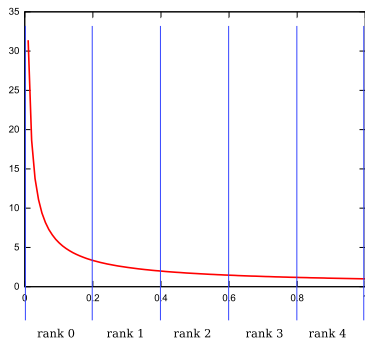
```
1 # Compilation:
2 $ gcc -O3 -o integral_seq.x integral_seq.cxx\
3     -L./mylib -lmymath
4 # Timed runs:
5 $ time -p ./integral_seq.x 1000000
6 Result=3.775045 Exact=3.775063 Difference=-0.000019
7 real 2.12
8 $ time -p ./integral_seq.x 2000000
9 Result=3.775058 Exact=3.775063 Difference=-0.000005
10 real 4.28
11 $ time -p ./integral_seq.x 8000000
12 Result=3.775062 Exact=3.775063 Difference=-0.000001
13 real 17.61
```

The time grows linearly with the problems size. Acceptable accuracy at 8M points. Can we speed it up?

Parallelization: domain decomposition.

Approach:

- Split $[a, b]$ into several domains;
- Compute integrals independently.



$$F(a, b) = \int_a^b f(x) dx$$

- 1 Assign a process to each domain $[x_k, x_{k+1}]$
- 2 Let each process compute $F(x_k, x_{k+1})$
- 3
$$F(a, b) = \sum_k F(x_k, x_{k+1})$$

A sketch of the parallel code

```
1  MPI_Init(&argc, &argv);
2
3  int rank, nprocs;
4  MPI_Comm_rank(MPI_COMM_WORLD, &rank);
5  MPI_Comm_size(MPI_COMM_WORLD, &nprocs);
6
7  // ...Get total number of steps, broadcast...
8
9  // Each rank figures out its integration limits and number of steps
10 unsigned long my_stepbase, my_nsteps;
11 get_steps(nsteps_all, nprocs, rank, &my_stepbase, &my_nsteps);
12
13 const double per_step=(global_b-global_a)/nsteps_all;
14 const double x1=global_a + my_stepbase*per_step;
15 const double x2=x1 + my_nsteps*per_step;
16
17 // Compute my own part of the integral
18 double my_y=integral(integrand, my_nsteps, x1, x2);
19
20 // Sum all numbers on master
21 double y=0;
22 MPI_Reduce(&my_y, &y, 1, MPI_DOUBLE, MPI_SUM, 0, MPI_COMM_WORLD);
23
24 // ... print results ...
25 MPI_Barrier(MPI_COMM_WORLD);
26 // Here we could start another computation.
27 MPI_Finalize();
```

Parallel performance: how to measure.

Now let's see how much we achieved...

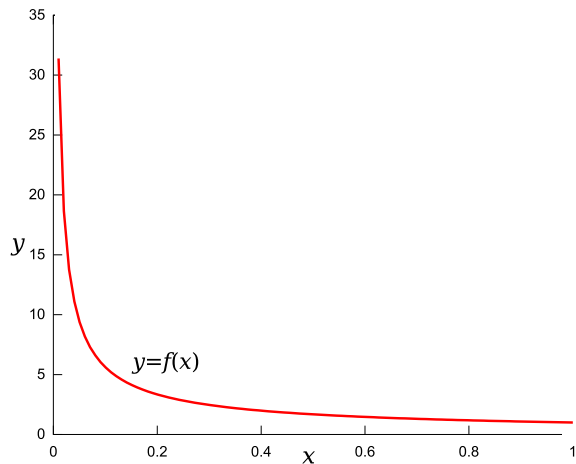
- **Strong scaling**: as we add processes, how do we fare?
- **Weak scaling**: as we add *both* processes and work?

```
1 $ mpicc -O3 -o integral_par.x integral_par.cxx \  
2     -L./mylib -lmymath  
3 $ time -p mpirun -np 1 ./integral_par.x 8000000  
4 Result=3.775062 Exact=3.775063 Difference=-0.000001  
5 real 17.23  
6 user 17.08  
7 sys 0.02  
8 $ time -p mpirun -np 2 ./integral_par.x 8000000  
9 Result=3.775062 Exact=3.775063 Difference=-0.000001  
10 real 17.24  
11 user 31.98  
12 sys 0.05
```

Parallel performance: results

Is there a performance problem?

THIS GRAPH IS WRONG



How does performance analysis work?

How to collect data?

- **Instrumentation:**
 - Insert timers & counters in the code
 - Requires source or binary processing
- **Sampling:**
 - Interrupt & check the program at regular intervals
 - Introduces statistical error

What kind of data?

- **Profile:**
 - Summary information only
 - Relatively small file
- **Trace:**
 - Detailed recording during the run
 - Potentially huge file
 - Profile can be restored

Allinea MAP does **tracing** by **sampling**.

Prepare for profiling

To prepare for profiling, one needs:

- Compile with full optimization
- Generate debugging symbols
- Link with system libs dynamically
 - Usually the default
 - Notable exception: Cray
- On Flux: load ddt module

```
1 $ mpicc -g -O3 -o integral_par.x \  
2     integral_par.cxx -L ./mylib -lmymath  
3 $ module add ddt
```

Running Map: simple way (demo)

- 1 Get interactive access to a compute node
- 2 Change to your working directory
- 3 Optionally, set *sampling interval*
- 4 Run as you would, prefixed by `map`

```
1 $ qsub -V -I -X -q flux -l qos=flux,nproc=12 \  
2     -l walltime=10:0:0 -A account_flux  
3 $ cd $PBS_O_WORKDIR  
4 $ export ALLINEA_SAMPLER_INTERVAL=5  
5 $ map mpirun -np 12 ./integral_par.x 10000
```

Caution:

- Too small interval: large overhead!
- Too large interval: not enough samples!

Running Map: other options

What if you can not or would not run GUI?

- Have slow or non-existing X connection to compute nodes
- Do not want to run interactively

Use `-profile` option.

```
1 #PBS -V
2 #PBS -q flux -l qos=flux -A account_flux
3 #PBS -l nproc=12,walltime=10:0:0
4 cd $PBS_O_WORKDIR
5 export ALLINEA_SAMPLER_INTERVAL=5
6 map -profile mpirun -np 12 ./integral_par.x 10000
```

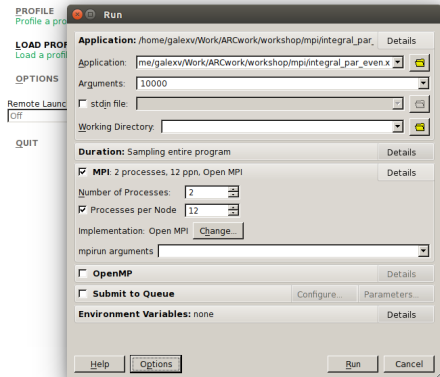
This will create a `*.map` file. Then run from the login node:

```
1 $ map integral_par_even_12p_*.map
```

If you are submitting to a Flux queue...

- 1 Run map from the login node:

```
$ map ./integral_par.x 8000000
```
- 2 Set number of processes
- 3 Check **Submit to queue**
- 4 Click **Configure...**
- 5 Load a proper *submission template file* (see next page)
- 6 Click **OK**
- 7 Click **Run**



Submission template for Flux

```
1 #PBS -V
2 #PBS -l walltime=WALL_CLOCK_LIMIT_TAG
3 #PBS -l nodes=NUM_NODES_TAG:ppn=PROCS_PER_NODE_TAG
4 #PBS -q QUEUE_TAG -l qos=flux -A account_flux
5 #PBS -o PROGRAM_TAG-allinea.stdout
6 #PBS -e PROGRAM_TAG-allinea.stderr
7
8 cd $PBS_O_WORKDIR
9 AUTO_LAUNCH_TAG
```

Time for a live demo!

SCREENSHOT HERE

- Most of the time is spent in MPI
- As the run progresses, *even more* time is spent in MPI
- Problem: some processes spend more time in calculating the integrand $f(x)$!
- It's called “Load Imbalance”

Possible solutions:

- Distribute work unevenly (but how?)
- Implement dynamic load balancing

Dynamic load balancing

Manager-Workers approach:

Manager

- 1 Listen to all workers
- 2 Worker sent READY ?
send GO with a job chunk
- 3 Worker sent DONE?
add result to the sum
- 4 No more job chunks?
send STOP to the worker
- 5 No more workers?
we are done

Worker

- 1 Send READY to the Manager
- 2 Listen to the Manager
- 3 Manager sent GO ?
 - Get job chunk
 - Do the calculation
 - Send DONE with result to the Manager
 - Go to (1)
- 4 Manager sent STOP?
exit.

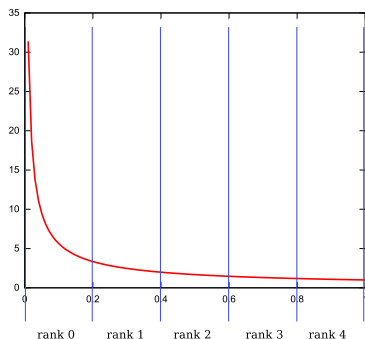
Strong scaling graph.

Map demo and screenshot.

The problem and the first approach

$$I(a, b) = \int_a^b x^{-3/4} dx$$

Analytic solution to cross-check: $I(a, b) = 4(b^{1/4} - a^{1/4})$

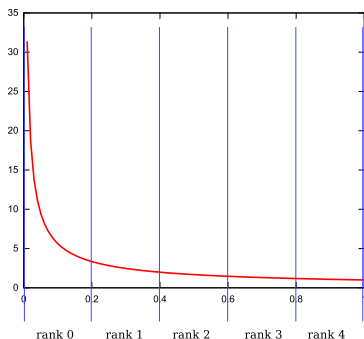


- Split $[a, b]$ equally between workers
- $I(a, b) = \sum_k I(x_k, x_{k+1})$

The problem and the first approach

$$I(a, b) = \int_a^b x^{-3/4} dx$$

Analytic solution to cross-check: $I(a, b) = 4(b^{1/4} - a^{1/4})$



- Split $[a, b]$ equally between workers
- $I(a, b) = \sum_k I(x_k, x_{k+1})$
- Caveat: integrand takes more time to compute at low x !

Dynamic load balancing

Try various chunk sizes, various number of processes.
Where are the bottlenecks?

(Alternative: give smaller chunks to low X. Scalability.
Profile.)