

Optimization of a Shallow Water Simulation

Ravi Patel (rgp62), Saurabh Netravalkar (sn575), Greg Granito (gdg38)

Profiling

Original Code

The ampxe profiler was used with hotspots to observe the running times of different functions. The functions which take longer to run and are a potential bottleneck are:

1. limit_derivs (CPU time: 1.155s)
2. compute_step (CPU time: 0.5688)
3. compute_fg_speeds (CPU time: 0.191)

Naïve Parallelization

Loop unrolling was performed and OpenMP parallel for pragmas were naively added in front of the looping constructs of each of the three functions.

Observations

1. The parallel code runs slower for the 200^2 cell problem than the serial un-tuned code. The reason is the inherent read after write and write after read dependencies among the loop bodies. Hence, the waiting time for each thread combined with the added overhead of communication costs blows up the time of execution.
2. The run time for loop unrolling alone for the bottlenecks above were found to be higher:
 - a. limit_derivs (CPU time: 1.480s)
 - b. compute_step (CPU time: 0.717s)
 - c. compute_fg_speeds (CPU time: 0.397s)
3. The run time for loop unrolling and naïve parallelization (24 threads) in the above functions were also found to be higher. About 20s of additional time was spent on overhead for ampxe.
 - a. limit_derivs (CPU time: 1.925s)
 - b. compute_step (CPU time: 1.315s)
 - c. compute_fg_speeds (CPU time: 0.703s)
4. Weak and strong scaling analysis was performed on the naïve parallelized code:

Strong Scaling (for 400^2 threads)

Problem size (cells)	Total time
1	19.3s
8	6.2s
24	7.4s

Weak Scaling (for 200^2 cells/core)

Threads	Total time
1	2.5s
4	7.5s
16	62.5s

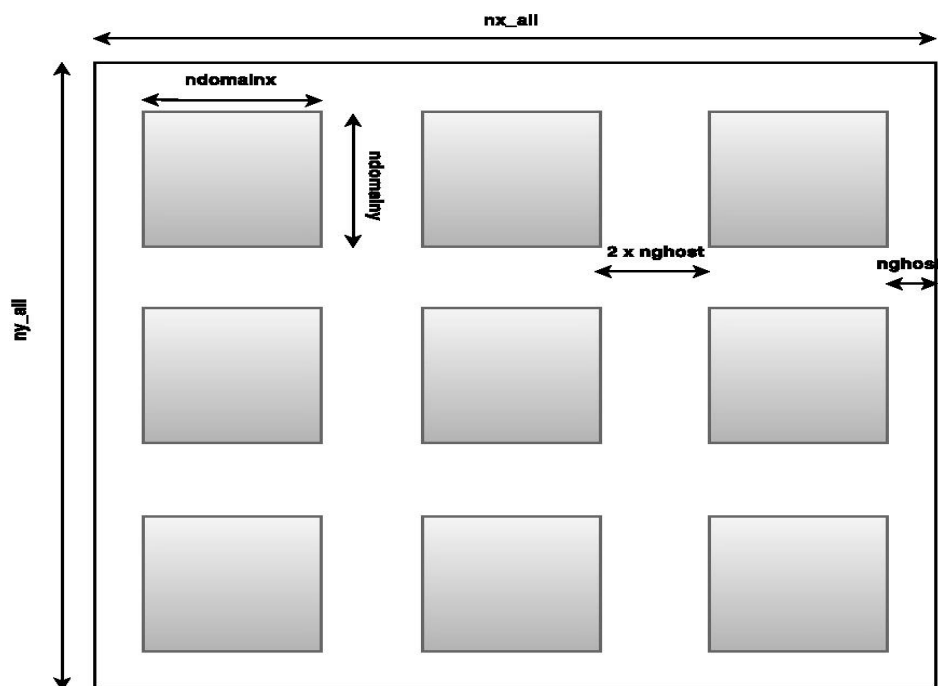
- From strong scaling, performance improved using 8 cores over 1 (39% efficiency), but 24 cores result in worse performance.
- From weak scaling, this code shows 33% efficiency using 4 cores over 1, but extremely poor efficiency using 16 cores.
- Naïve parallelization results in some improvements for low thread counts, but scales poorly for higher thread counts.

High Level Strategy for using Domain decomposition

- Remove data dependencies by tiling the grid and padding the tiled grids with ghost cells as well as the necessary cell layers of adjacent tiles.
- For time t_s in blocks of a step:
 - For Each Processor p_i , operate in parallel on a block b_i of the board:
 - Copy into local board of each processor along with ghost cell info
 - Advance by making the processors compute some steps independently
 - Copy out to a new board
 - Barrier synchronization
 - Swap the board pointers, i.e. new board becomes the old board and vice-versa

Further Potential Optimizations

- Sizing blocks to fit in cache considering the sharing into account
- Align memory access for the blocks
- Compute wave speeds every few steps and choose 'dt', the time step more conservatively



Parallelization

Direct Domain Decomposition

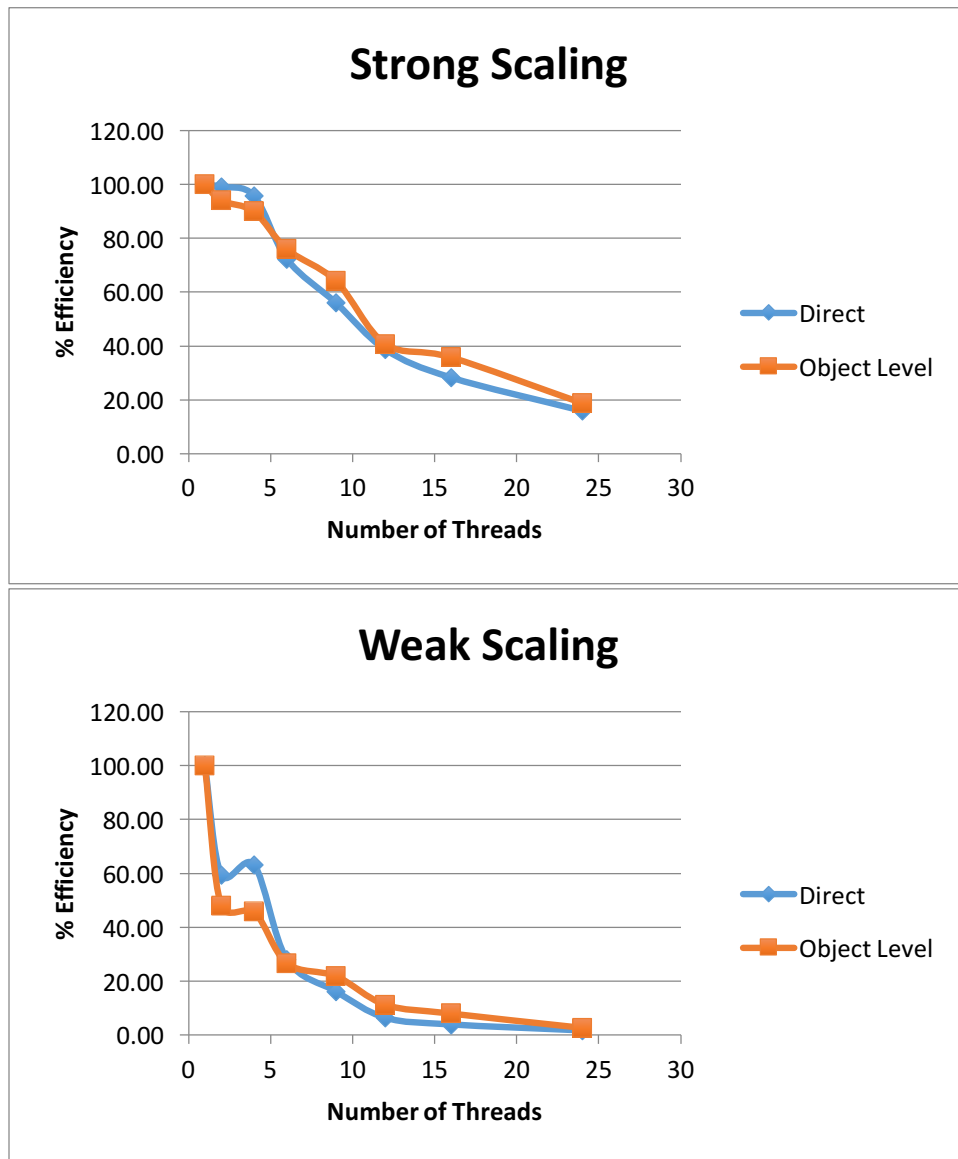
Workflow: The workflow as implemented, for parallel execution of the shallow water wave simulation using domain decomposition is as follows:

1. The original board is divided into square tiles/blocks
2. Each block is padded on all 4 sides with layers of ghost cells as shown in the figure.
3. Each processor is given its own tile padded with ghost cells and works in parallel by time stepping on its own tile.
4. A single processor initially performs the ghost cell updates for the entire board, i.e. updating ghost cells values by copying out rows from left and right and then top and bottom, by wrapping around the board boundaries. Now, the processors are ready to compute their own time step advances in parallel.
5. Each processor also computes the local wave speeds 'cx' and 'cy' within their own tile and write the values to a shared array, indexing the array by their own unique thread identifiers which range from 0 to number of threads - 1.
6. At the end of 'n' time steps, ('n' = number of layers of ghost cells - 3, since a minimum of 3 ghost cells are required to compute one entire time step in parallel without requiring data from adjacent tiles), the processors perform a barrier sync.
7. A single processor traverses through the array of local tile wave speeds and computes the maximum speed among them in both x and y directions, and uses this value to perform the time stepping increment.
8. Once the time steps are completed, the true data values are copied back from the ghost cell padded board into the original board.

Object Level Domain Decomposition

We also created an implementation of domain decomposition where we divided our original $n \times n$ grid into tiles for each thread at the driver level. A shared Sim object was generated to contain the full simulation. Sim objects were also created for each thread. The simulations were run separately on each thread. Between iterations, the threads copy results to and ghost cells from the main simulation as well as min reduce to obtain synchronize the time step.

Weak and strong scaling analysis was performed on the two domain decomposition approaches.



From this analysis, direct domain decomposition performs better than object level for low thread counts but the object level performs better at higher thread counts. Both approaches achieve better performance than naïve parallelization.

References

<https://github.com/amirajdhawan/water> : Held meetings with him to understand their group's approach, brainstormed and reasoned about potential further optimizations.