

Implementation of cohort parallelisation in the Seapodym code

Alexander Pletzer, Chris Scott,
Inna Senina, Lucas Bonnin and Romain Forestier

September 11, 2025

This report describes NeSI/REANNZ’s consultancy project entitled “SEAPODYM cohort parallelisation” whose work was performed from May to September 2025.

1 Statement of the problem

SEAPODYM is a quantitative spatio-temporal model of population dynamics that solves partial differential equations with initial and boundary conditions in C++. The model is parameterized using a maximum likelihood estimation approach that integrates georeferenced datasets obtained from industrial fishing and scientific campaigns.

The objective of this work is to parallelise the SEAPODYM code over cohorts. A cohort is a group of fish that are born at the same time. Each cohort can be integrated forward in time independently of other cohorts, until the cohort reaches a certain age and dies. At the next time step, a new cohort is born and integrated forward in time until it dies. Except at the first time step, the initial condition of each cohort depends on the density of the other cohorts at the previous time step. Thus, the integration of each cohort must take into account the state of the other cohorts and this requires careful consideration when parallelising the code.

An example of cohort time integration is shown below:

$$\begin{array}{ccc} 0 & 1 & 2 \\ 3 & 1 & 2 \\ 3 & 4 & 2 \\ 3 & 4 & 5 \\ 6 & 4 & 6 \\ 6 & 7 & 6 \end{array} \tag{1}$$

Here, the vertical axis represents time increasing from top to bottom ($N_t = 6$). The horizontal axis represents the tasks than can be performed concurrently. Each row represents different age groups ($N_a = 3$). Each cohort is identified by an integer ($0, 1, \dots 7$).

Cohort 1 is one unit time older than cohort 2. Therefore cohort 1 dies one time step before cohort 2. The number of time steps a cohort is integrated is $i = i_{beg} \dots i_{end} - 1$ where $i_{beg} \geq 0$ and $i_{end} \leq N_a$. In our example, we have When a cohort dies, it is replaced by a new cohort at the next time step.

Table 1: Start and end integration indices for each cohort

cohort Id	i_{beg}	i_{end}
0	2	3
1	1	3
2	0	3
3	0	3
4	0	3
5	0	3
6	0	2
7	0	1

To instantiate the new cohort, data will need to be communicated from the cohorts at the previous time step. Figure 1 shows the dependency of the new cohorts on the older cohorts at different steps.

2 Task farming with dependencies on other tasks' steps

We opted for a task farming approach. This involves creating a pool of tasks that can be executed concurrently and a manager that assigns the tasks to the workers. Each task consists of advancing a cohort in time.

In classical task farming, the manager starts by assigning tasks to the workers. This involves sending a message to each worker, along with some input parameters. Each worker then executes the task and reports the results back to the manager, who then assigns a new task to the worker. When no more tasks are available, the manager sends a message to the workers to shut down. Task farming is particularly suitable when tasks take various amount of time to execute.

Clearly, the dependencies between tasks and sub-tasks shown in Fig. 1 prevents us from using such a naive task farming approach. In our case, tasks

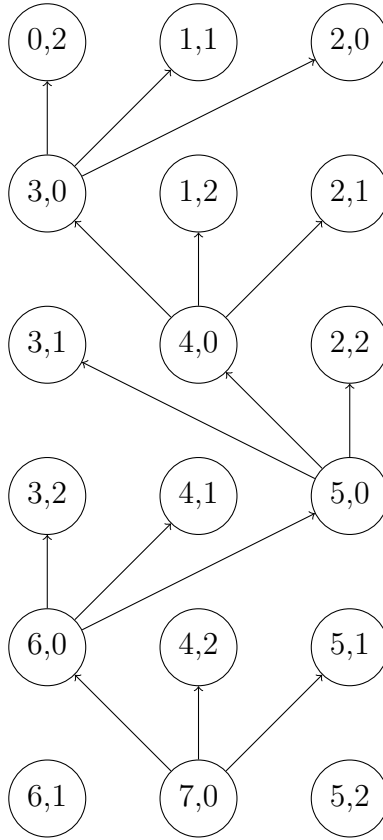


Figure 1: Dependency of cohort tasks on other cohort, step tuples. The first digit is the cohort Id and the second the step.

can only be started when other tasks' particular steps have completed (as indicated by the arrows in Fig. 1). This requires the manager to keep track of task-step dependencies.

We start by describing a worker's task. A worker waits for a task to be assigned by the manager. A task has a identification number `taskid` that fully describes the task to accomplish. The worker reports back the status of the execution of each task to the manager. The corresponding pseudo-code is shown in Algorithm 2.

Algorithm 1 A worker's pseudo-code.

```

1: while true do
    taskid = GETTASKIDFROMMANAGER;
2:   if taskid < 0 then Break
3:   end if
        ▷ Perform the task, stepping from stepBeg to stepEnd - 1
4:   TASKFUNC(taskid, stepBeg, stepEnd)
        ▷ Notify the manager that this worker is available again
5:   SENDSIGNALTOMANAGER(DONE)
6: end while

```

The task function defined in Algorithm 2 is responsible for advancing the cohort associated with `taskid` from `stepBeg` to `stepEnd - 1`. Note the call

Algorithm 2 The task function executed by the worker.

```

1: function MYFUNCTION(taskid, stepBeg, stepEnd)
2:   data = GETDATAFROMMANAGER(taskid)
3:   cohort = CREATECOHORT(taskid)           ▷ Advance a cohort
4:   for step = stepBeg to stepEnd - 1 do
5:     STEPFORWARD(cohort)
6:     SENDSTEPCOMPLETEMESSAGE(TOMANAGER(taskid, step))
7:   end for
8: end function

```

to `sendStepCompleteMessageToManager` which informs the manager that a specific step for a task has been completed. This is crucial for managing task dependencies.

The manager code orchestrates the work. It maintains a list of all active workers, a task queue, a list of completed tasks and stores the results. Since workers will be sending various types of messages ("worker is available" or "result of task-step X"), the manager needs to be able to distinguish between

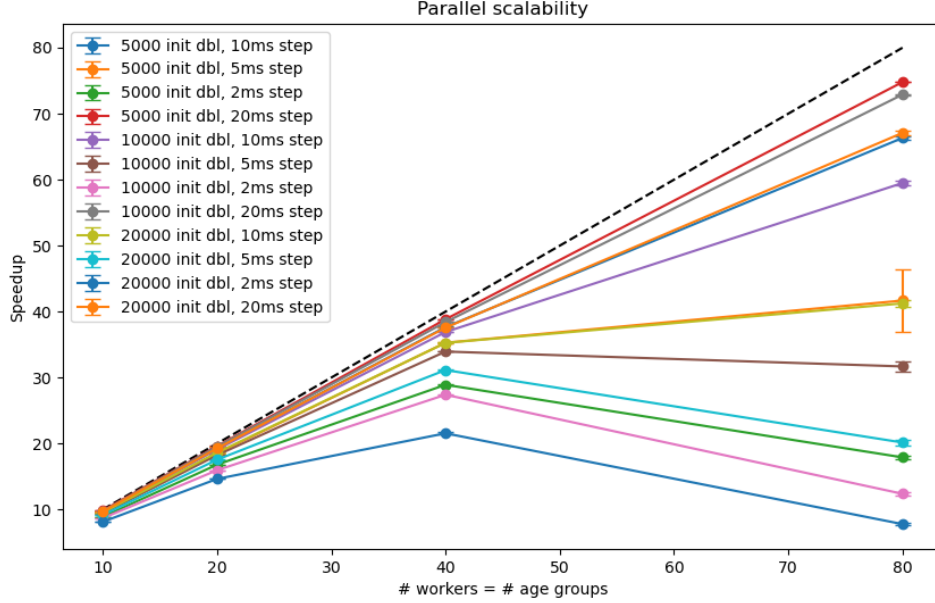


Figure 2: Parallel scalability for different step execution times and number of doubles initially fetched from the manager. The number of age groups matches the number of workers. These results assume no initialisation setup time.

them. This is done using message tags. The manager’s pseudo-code is shown in Algorithm 3.

3 Scalability

Figure 2 shows the parallel scalability of the cohort parallelisation approach described above as the number of workers is increased. The speedup numbers were obtained for an idealized scenario, assuming a fixed cost of communication (to initialize a new cohort) and a constant cost of executing a step on the NeSI/REANNZ AMD Milan node cluster. In all cases the number of workers was chosen to match the number of age groups N_a .

Perfect parallel scaling is shown as the black dashed line. Good scalability depends on the time taken to execute each step and the amount of data fetched from the manager at the start of each task. Parallel efficiency $> 50\%$ is achieved for time step times $\geq 10\text{ms}$ and initial data sizes ≤ 10000 doubles.

Algorithm 3 The manager's pseudo-code.

```
1: while not taskQueue.empty() or not assigned.empty() do
    ▷ Look for messages "task-step complete"
2:   while true do
    ▷ Probe for messages from workers
3:     msg = getMessageFromAnyWorker
4:     if msg.type != "task-step complete" then
5:       break
6:     end if
    ▷ Store the result
7:     results.insert(msg.output)
8:     taskid = output[0]
9:     step = output[1] completed.inserttaskid, step
10:    if step == stepEnd - 1 then assigned.erase(taskid);
11:    end if
12:  end while
    ▷ Assign ready tasks to any available worker
13:  for taskid in taskQueue do
14:    taskDependencies = getTaskStepDepenenciestaskid bool ready =
    true;
15:    for dep in taskDependencies do
16:      if completed.find(dep) == completed.end() then
17:        ready = false
18:        break
19:      end if
20:    end for
21:    if ready and !activeWorkers.empty() then
22:      worker = activeWorkers.begin()
23:      activeWorkers.erase(worker)
24:      SENDTASKTOWORKER(taskid, worker)
25:      assigned.insert(taskid)
26:      taskQueue.erase(taskid)
27:    end if
28:  end for
    ▷ Drain all worker-available messages
29:  for worker in workers do
30:    msg = getMessageFromAnyWorker
31:    if msg.type == "worker is available" then
32:      activeWorkers.insert(worker)
33:    end if
34:  end for
35: end while
    6      ▷ Send stop signal to workers
36: stop = -1
37: for worker in workers do s(e)ndStopSignalToWorkerworker
38: end for
```

For instance, a time step taking 10ms and an initial data size of 10000 doubles gives a speedup of 60, when using 80 workers, i.e. 75% parallel efficiency.

Note that at every step the worker sends data to the manager, as well as a message to inform the manager that the step has been completed. This introduces a communication overhead that limits parallel scalability. Taking the example of 20,000 doubles sent by 80 workers to the manager at a cadence of 10ms, we get a data transfer rate of $1.3\text{GB/s} = 100\text{Gb/s}$, which is close to the bandwidth of the Infiniband communication fabric. To further improve parallel scalability, one would need to either increase the computational load per step, reduce the amount of data to be transferred or apply some form of data compression. For instance, one could send floats instead of double and this would improve the speedup from 40 to 60.

4 The cost of initialization

Scalability can be affected by load balancing. Some concurrently executed tasks may take longer than others to complete. When synchronization occurs, which is the case at every time step in Seapod, then workers whose task has finished need to wait for the slowest worker to finish their task. While the evolution of the reaction-diffusion equations takes the same time for all cohorts, the initialization of a new cohort, which occurs at every time step, can lead to severe load imbalance and can thus affect the parallel speedup. Fig. ?? shows the ratio of the compute time over the initialization time, the higher the better. Initialization involves reading data from file and other one off operations.

There are different ways to amortize the high cost of initialization. One possibility is to avoid constructing a new cohort at every time step. Instead a “vanilla” cohort is created initially, when the simulation starts. This cohort, reads the restart files, the input files and sets things up. Only the final stage of the construction, which involves reading forcing data and spawning, is executed when launching the new cohort.

5 Comparing blocking and non-blocking get operations

A significant part for the initialization time involves fetching data from the manager. Our library supports asynchronous get operations, which allow one to overlap communication with computation.

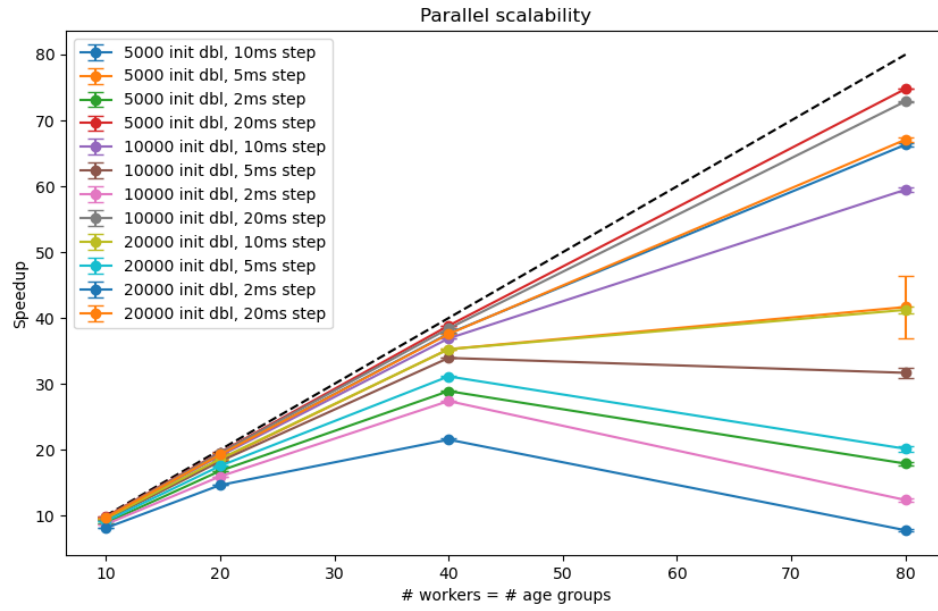
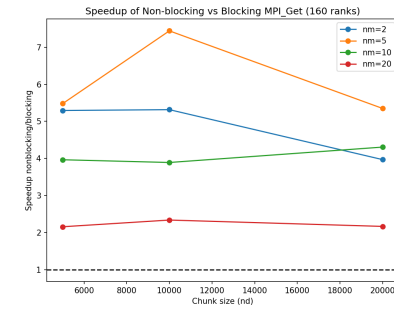
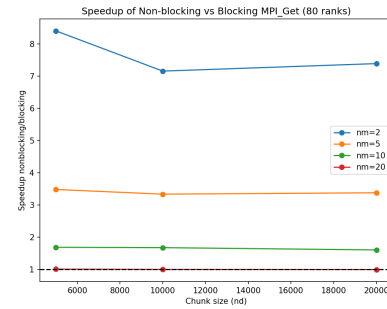
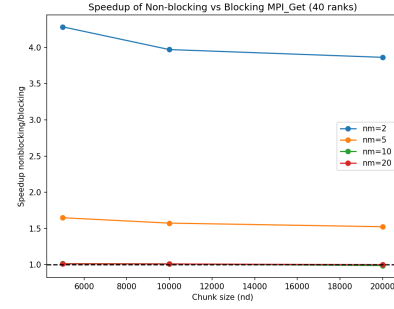
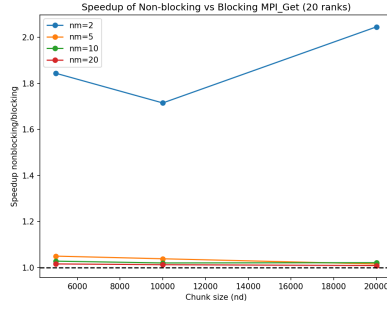


Figure 3: Ratio of compute over initialization time for Seapodym.



The default mode to `get` data is blocking, that is the retrieved data are ready when the call completes. Alternatively, the programmer may rely on `getAsync` calls. In this case the programmer indicates his/her desire to fetch the data but the data need not be immediately available. A `startEpoch` and `endEpoch` calls around the `getAsync` operations determine the time interval when remote memory access operations take place. The retrieved data are guaranteed to be available only when `endEpoch` is called. (An additional “flush” method can be invoked to synchronize operations within a period.)

In Fig. ?? we compare the execution time of MPI `get` operations using the blocking and non-blocking variants. Each test involves reading N chunks (= number of workers) of nd values from rank 0 and summing the result. In the blocking test, the chunks are read one at a time, followed by a sleep operation lasting nm milliseconds and a sum operation over the chunk. In the nonblocking test, a big array is filled by reading the chunks asynchronously, followed by the sleep operation – these two operations are within a single epoch. The final step then involves summing up the values of the big array. The benefit of using non-blocking `get` calls is particularly apparent for large chunks nd and small amount of work (low nm).

6 Summary and future work