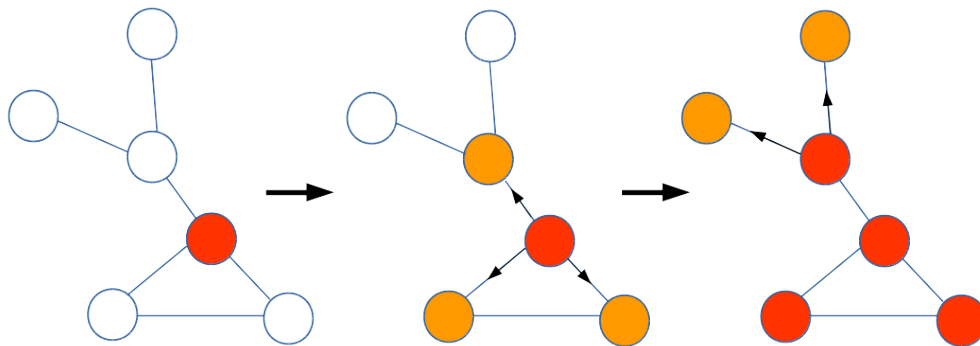


Parallel Graph Algorithms in Julia

Julian Kates-Harbeck

December 9, 2015



Abstract

We use the generic problem of Monte-Carlo simulation of stochastic graph algorithms to illustrate fine-grained multi-threading and coarse-grained multi-processing in Julia [1]. In the process we develop an automatic cluster management tool to distribute processes to cores across different machines. We develop multi-processing and multi-threading versions of our algorithms that are able to fully utilize the hardware at their disposal and achieve significant speedups over serial code. We also identify several challenging bugs and issues in the multi-processing and multi-threading components of the Julia language. The most important issues are submitted to Github as issues #14343 and #14344 in order to help with language development.

1 Introduction and Background

Our toy computational problem is concerned with the dynamics of generic “infections” on graphs. The goal of such a model is to give a realistic yet general representation of an infection process on a network. Each node of the graph is either in an *infected* (I) state or a *susceptible* (S) state. The transitions between these states happen probabilistically, where the probabilities can depend (*nonlinearly*) on the fraction of neighbors of each type. An example application might be the transmission of an infectious disease across a social network. The probability of a person (a node on the graph) becoming infected by their contacts (represented by neighbors on the graph) depends in some capacity on the number of their contacts carrying an infection. For example, here are the transition probabilities during a short time step dt used in our specific model:

$$\begin{aligned} P_{I \rightarrow S} &= (1 - y)y\beta dt \\ P_{S \rightarrow I} &= (1 - y)y^2\alpha dt \end{aligned} \tag{1}$$

Where y represents the local fraction of infected individuals, and α and β are parameters. The process we are considering might also be applied to other dynamical “infections” on graphs, such as the spread of a language pattern or a behavior across a social network, the spread of strategies among players of a game, or the spread of genes in an evolving population.

Since the process is stochastic, we are generally interested in the *statistics* of the infection process. Two important characteristics are the *fixation probability* P_{fix} of the infection, i.e. the probability that all nodes in the graph become infected, and the *distribution over infection sizes*. We define the infection size w as

$$w = \int_{T=0}^{\infty} n(t) dt$$

where $n(t)$ is the current number of infected individuals. The distribution $P(w)$ is then indicative of how frequently infections reach a certain size. We typically seed an infection with a single infected individual among all susceptible individuals. Using the micro-dynamics of the infection (i.e. the infection probabilities of a given node given its neighbors’ types, like equation 1), we can make predictions for P_{fix} and $P(w)$.

In order to test these analytical predictions, it is of course necessary to run the model many times $\gtrsim 10^5$ such that good statistics for P_{fix} and $P(w)$ may be extracted.

2 Algorithms

Equation 1 naturally suggests the following algorithm for simulating our model. On a high level, we run our graph simulation many times in order to obtain statistics over the runs, as seen in code example 1.

```
results = map(run_graph_simulation, 1:num_trials)

#analyze results...
```

Code Example 1: Monte Carlo Algorithm

On a lower level, we update every node in small time increments, applying transitions according to the transition probabilities shown in equation 1. In particular, we update all nodes in parallel, given the state of all nodes at the previous time step. While this requires a copy of the “old” state of the graph at every time step, it also makes the updates of the individual nodes independent of each other, as seen in code example 2. This is in the same spirit as the Jacobi relaxation method [3] for solving Laplace’s equation.

```

1  for t in 1:T
2      new_types *= 0
3      update_graph(g,new_types)
4  end
5
6  function update_graph{P}(g::Graph{P},new_types::Array{P,1})
7
8      #OUTER LOOP
9      for v in vertices(g)
10
11         if get_type(g,v) == INFECTED
12
13
14             #INNER LOOP
15             #possibly infect neighbors
16             for w in neighbors(g,v)
17                 if get_type(g,w) == SUSCEPTIBLE
18                     x = get_infected_neighbor_fraction(g,w)
19                     p::Float64 = p_birth(x)
20                     if rand() < p
21                         #WRITE OPERATION
22                         new_types[w] = INFECTED
23                     end
24                 end
25             end
26         end
27     end
28     set_types(g,new_types)
29 end
30

```

Code Example 2: Simplified Graph Algorithm

The two algorithms shown illustrate two very common patterns in parallel computing problems. The first is Monte Carlo simulation or sampling. Here we repeat an operation or simulation many times, where each run is independent, and aggregate the results. The second is an iterative graph algorithm, where we perform several “passes” over a complicated data structure¹, where each pass touches every component of the data structure independently and performs possible updates. It is important to note that iterations depend on previous iterations, but *within* an iteration, all operations are independent. In our case, in the loop over vertices in line 9, all loop iterations are independent. However, in the loop over times

¹In our case updating the states of the nodes on a graph, but this could also be a grid over which we are computing a stencil relaxation, or a network for which we are computing PageRank scores.

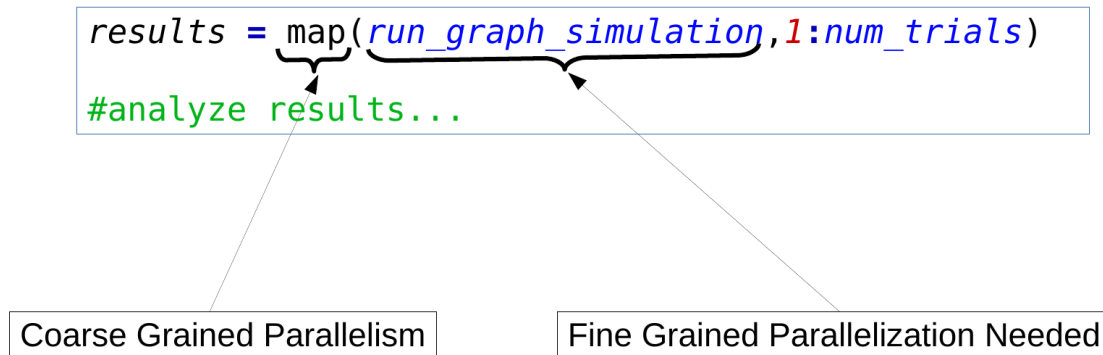


Figure 1: An illustration of the two types of parallelisms needed in our application.

steps, each time step depends on the previous one.

It is clear that both patterns lend themselves to parallelization, although the ideal way in which to parallelize these algorithms will differ, as shown in figure 1. In particular, MC sampling involves running many computationally expensive tasks independently of each other. This means that if we simply spawn a new process for every run, the overhead of spawning a process is dwarfed by the cost of the actual simulation. This allows us to use *multi-processing* to run the MC trials in parallel. On the other hand, processing each node in the graph algorithm might only involve little computational cost. Moreover, the graph itself, the type information for the nodes, and access to writable memory for updating the types must be available to all parallel workers. Thus, it will be prohibitive to duplicate memory across different processes and to spawn a process to do only a small amount of work. We thus require a very lightweight parallelization model, and chose to work with *multi-threading*². Threads are perfect for this task as they can access shared memory without duplication of the data structures and because they are very cheap to spawn.

In the following sections we will describe in more detail the steps we took to achieve high performance in our code, including serial optimization, parallelization using coarse-grained multi-processing, and parallelization using fine-grained multi-threading.

3 Serial Code

We had initially written our simulation code in Python, using the **NetworkX** graph library [2]. In the context of this project, the first step was thus to rewrite the code in serial Julia. Since we are dealing with graphs involving changing properties on the nodes, the first step was to use a simple translation of our python code into Julia using the **Graphs.jl**³ package. In particular, **Graphs.jl** is described in its documentation as being inspired by **NetworkX**. The translation was straightforward, given the similarities between Python and Julia. Using Julia’s profiling tools, we then went to optimize the Julia code in a manner

²We also tried using multi-processing for this fine grained parallelization task. However, the massive overhead of communicating and sharing data caused a significant slowdown and convinced us of the need for a more lightweight multi-threading approach.

³<https://github.com/JuliaLang/Graphs.jl>

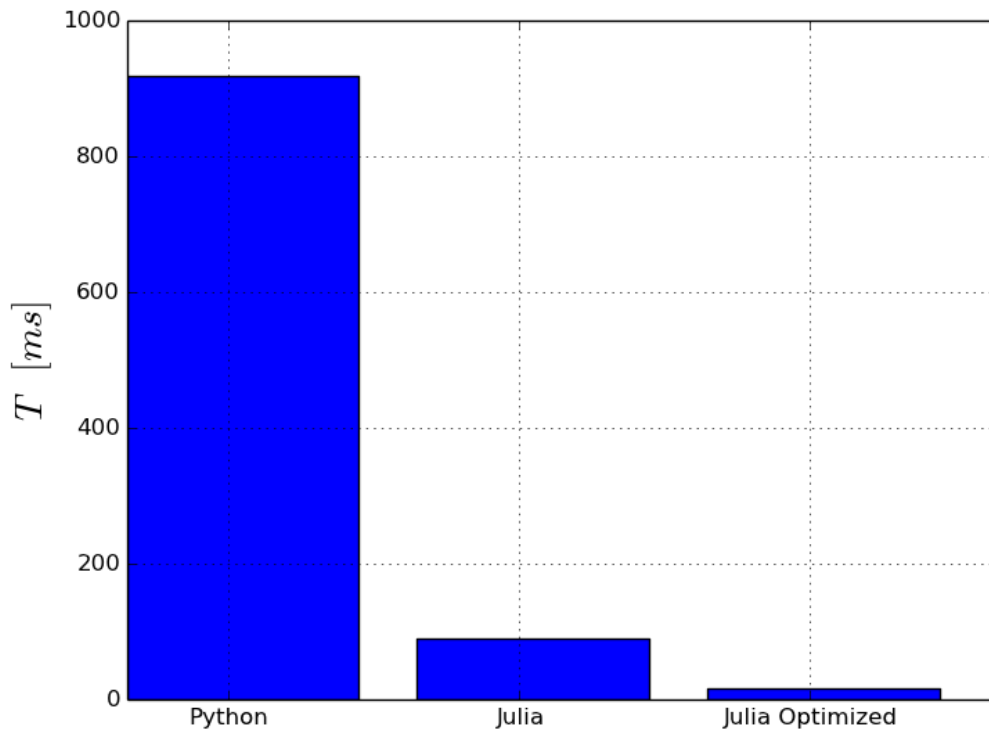


Figure 2: Serial running time on the benchmark graph problem described in the text. Moving from Python to Julia results in an order of magnitude cut in running time. However, optimizing Julia using a profile-optimize cycle resulted in another comparable speedup. The execution times are 920 ms, 89 ms, and 15 ms, respectively.

similar to Assignment 2 of this course. This entailed running the code, profiling it, finding the most time-consuming portions, and rewriting them to run faster. One important bottleneck was the iteration over neighbors as in line 16 of the simplified code sample 2 (but in the inner loop of `get_infected_neighbor_fraction()`). We found that the `Graphs.jl` package, in order to provide flexibility to the user, added significant overhead to such operations. A leaner and simpler graph data structure is provided by `LightGraphs.jl`⁴, where the graph is stored simply as an array of `Int` arrays (an array of neighbor lists). This results in much faster iteration and thus removes the main overhead in our code. We show the speedups on a benchmark problem of calculating the infected neighbor fraction (as in line 18 of the sample code) for all nodes on an Erdos-Renyi random graph with 2000 nodes and 1900 edges per node. The results are presented in figure 2. We find a speedup by a factor of about 10 simply by translating Python to Julia. Remarkably, by optimizing our Julia code in a few places, including mainly the switch to `LightGraphs.jl`, we achieve another factor of ~ 6 speedup!

⁴<https://github.com/JuliaGraphs/LightGraphs.jl>

4 Coarse-Grained Parallelization: multi-processing

It might seem initially that the problem of parallelizing across MC simulation instances is a trivial problem. Indeed the iterations are completely independent of each other and are coupled only in the final step of aggregating the statistical data, which represents a tiny fraction of the overall computational work. Thus, in an idea world, we would simply replace the `map` operation from code example 1 with a parallel `pmap`.

In Julia, `pmap` distributes the workload across all available processes in parallel and aggregates the results. Thus, once we have correctly configured our processes, this is indeed the only step we need to take to parallelize our code. However, there are two caveats:

- It is not trivial to correctly allocate one process to every available core on demand.
- The sharing of memory among different processes can be tricky.

These two issues represent the main hurdles to multi-processing in our problem, and we will describe our solutions in the following sections.

4.1 Automatic Multi-Process Management

For ideal performance, we generally want to allocate a single process⁵ per core in our computing environment. In a fixed environment, like a small personal computer or a small lab cluster, it is possible to use `addprocs` commands with explicit host names and explicit numbers of cores per host to allocate the processes. One has to essentially “hard code” these host names and the number of cores per machine into the application. While this might be feasible for a small, personal computing environment, it is unscalable, inflexible for changes in the environment, and simply unfeasible for large or shared computing environments.

We use a large, shared computing cluster⁶ as an example environment to develop a more general and scalable solution to this problem. Ideally, we want the user to be able to specify nothing but a number of processes, such that our automatic tool will handle correct allocation of these processes across different cores and machines. One challenge is that in a shared computing environment we might have access to a number of cores potentially distributed across different machines, and that this allocation can differ for every session on the shared cluster.

Our test cluster works with the SLURM resource manager. This allows a general solution of the following form:

- Use SLURM environment variables⁷ to obtain the details of the current allocation, i.e. which machines (by host name) we have access to, and how many cores on each machine.
- Read a specified number of processes from the user

⁵or in the case of hyper-threading, some small constant number of processes

⁶the Odyssey computing cluster at Harvard

⁷In this case `SLURM_NODELIST` and `SLURM_JOB_CPUS_PER_NODE`

- Find an Appropriate distribution of these processes onto the machines and cores that we have access to.
- Pass this information into an explicit call to `addprocs()` to actually add the processes in the right places.

A general allocation library would simply re-implement the first step to fit whichever resource manager is used (if other than SLURM). We implemented these steps in pure Julia. The result for the user is illustrated in figure 3. As desired, the user only has to worry about how many processes he wants, the rest is handled by the allocation manager.

4.2 Results

In this section, we use our multi-processing tool to distribute workloads of the form shown in code example 1 to processes on cores across several different machines. In particular, we are working with an allocation on the Harvard Odyssey cluster. Every machine has 64 cores, a subset of which are allocated to us. We choose an allocation of a total of 256 cores, which will then be distributed across at least 4 machines⁸. Since our actual infection algorithm is stochastic in nature, it has wildly varying run times. To give a more reliable study of timing, we choose a simpler workload of the form shown in the code example 3.

```

1 #O(N^3)
2 @everywhere myfun(N,M) = sum(randn(N,M)*randn(M,N))
3
4 map(N -> myfun(N,N),repmat([N],250)) #serial
5 pmap(N -> myfun(N,N),repmat([N],250)) #parallel

```

Code Example 3: The simple workload used to test multi-processing performance.

The scaling of this workload with N is $O(N^3)$. We run this workload for various numbers of processes and for various values of N . The results are shown in figures 4 and 5. We expect the ideal scaling to go as $\sim \frac{N^3}{n_{procs}}$. In figure 4 we compare the true scaling with n_{procs} (solid lines) to the ideal scaling (dashed lines) for various values of N . In figure 5 we summarize this data in a log-log plot, for which the ideal scaling with N^3 and n_{procs} lies on a plane (red). We find that in general the true scaling is well predicted by the ideal scaling, with the main difference being that for small problem sizes, the speedup plateaus with increasing number of processes.

Two observations are important to point out. First, we find that the scaling becomes closer to ideal ($\sim \frac{1}{n_{procs}}$) the larger N becomes. This is as expected, since we are performing more work on each of the processes as compared to the overhead of sending data around and synchronizing the processes. Second, we find that the maximum speedup obtained is $114.5\times$, which importantly is significantly greater than $64\times$ (the maximum speedup possible with one machine). This means that we are utilizing processes on more than one machine

⁸In our case, the cores per machine are 56, 16, $64(\times 2)$, and 56 on a total of 5 machines.

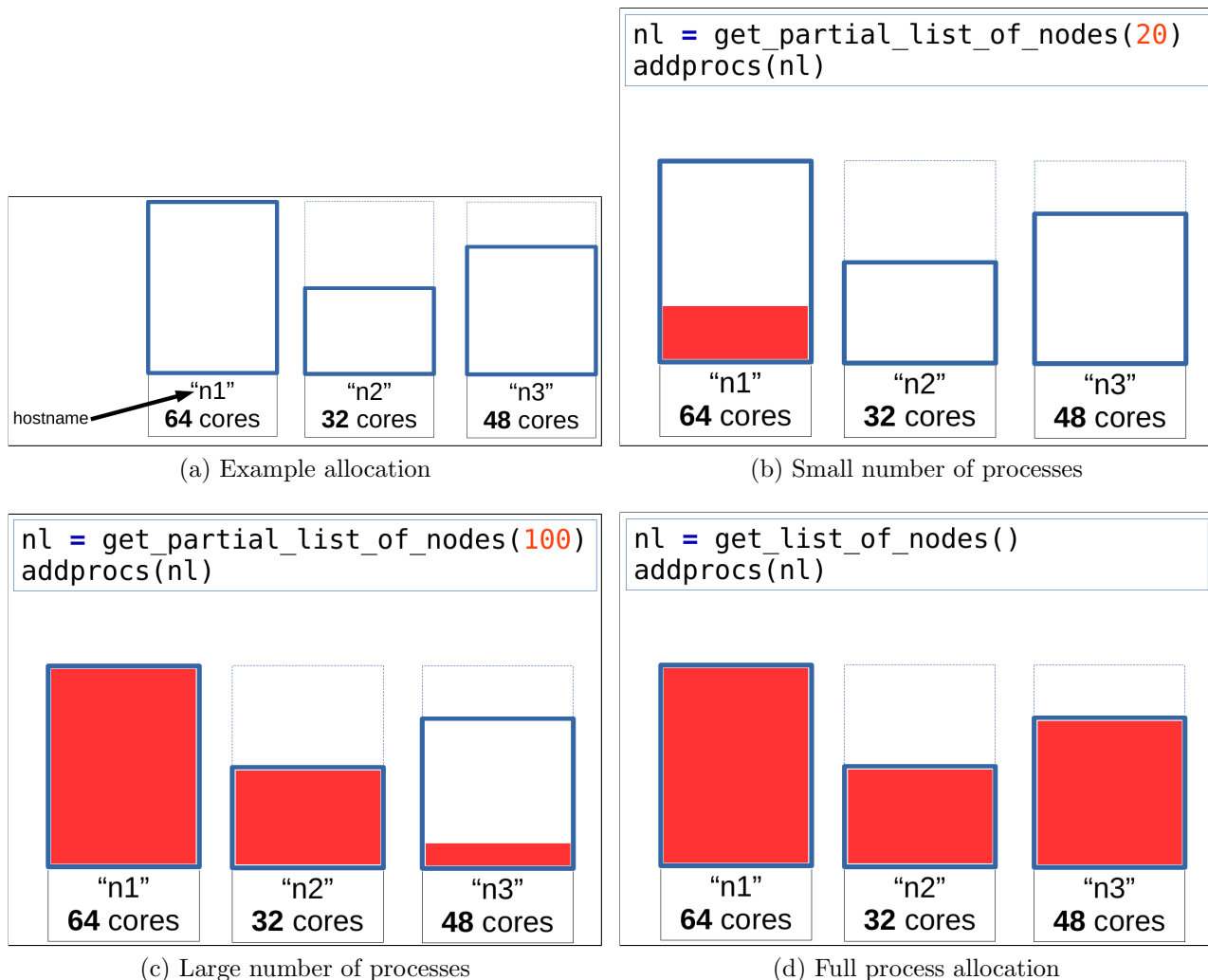


Figure 3: An illustration of the automatic multi-process management package that we developed. (a) shows an example of a possible SLURM allocation. We have received a total of 144 cores distributed as shown over nodes “n1” (64 out of 64 cores), “n2” (32 out of 64 cores) and “n3” (48 out of 64 cores). The dashed outline represents the 64 cores on each machine, while the blue box represents our example SLURM allocation. Julia processes running on a core (always one process per core) are shown in red. (b) shows the syntax for adding a small number of processes (20). These processes fit on one node and are allocated as shown in the image. (c) shows a larger allocation request for 100 processes, which don’t fit on a single machine. The machines are filled up one by one until all processes have their own core. Finally, (d) shows the shortcut syntax for allocating a process to every available core.

in parallel, which shows that our allocation procedure is *indeed distributing one process per core on many different machines*. We expect this scaling to continue even for much larger numbers of processes, given sufficiently large workloads.

4.3 Implementation Issues and Bugs

In addition to proper process allocation, the second hurdle to simply using `pmap` instead of `map` in order to parallelize our MC simulations was the proper handling of data sharing across all processes. Since the Julia documentation is often not entirely up to date on describing the details of these issues, it took many iterations of trial and error in order to obtain a functioning solution. This part took up the largest share of development time in the multi-processing component of the project. We list here some of the most challenging bugs and issues that complicated the migration to `pmap`, as well as our solutions to these issues.

It is key that all processes running the desired function have access to the data referenced by this function. This includes runtime data as well as code and modules.

Modules For modules that are used on other processes, we found the following procedure to work well. First, import the module with `using <module name>` on the master process. Then add all other processes using `addprocs()` and finally import the module on all processes with `@everywhere using <module name>`.

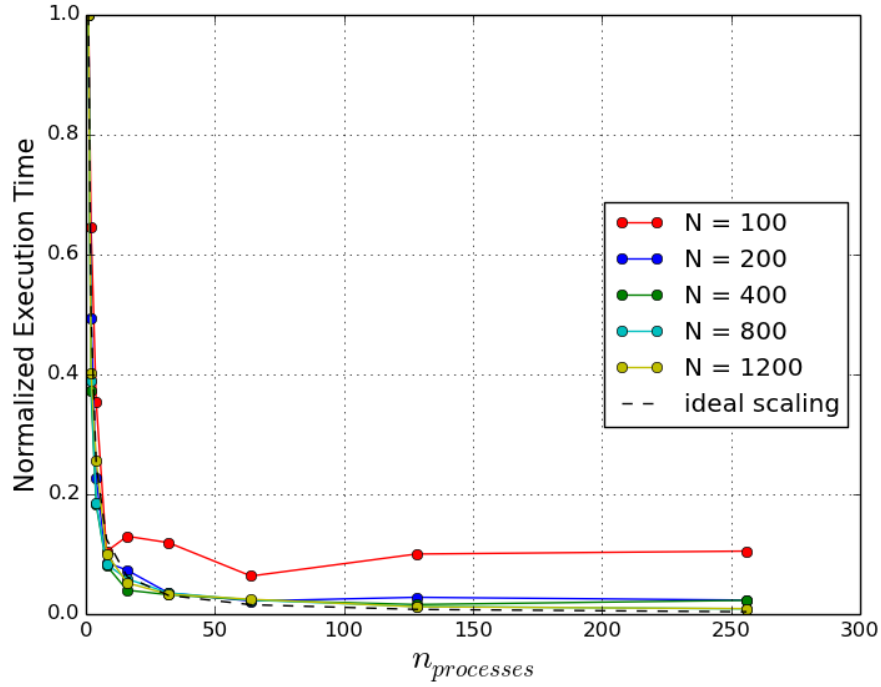
Data and Functions For data that is used on other processes, we must distinguish several different types. Functions that are used on other processes that have not been defined in one of the modules that have been imported everywhere need to be annotated with the `@everywhere` macro. Data that is passed into these functions by the `pmap` operation does *not* need to be declared `@everywhere`. Data that is captured as arguments in a curried function (which in turn is passed as the first argument to `pmap`) *does* need to be declared `@everywhere`.

In code example 4, we give a minimal example illustrating these issues. This has been submitted as Github issue #14344 to the main Julia language project.

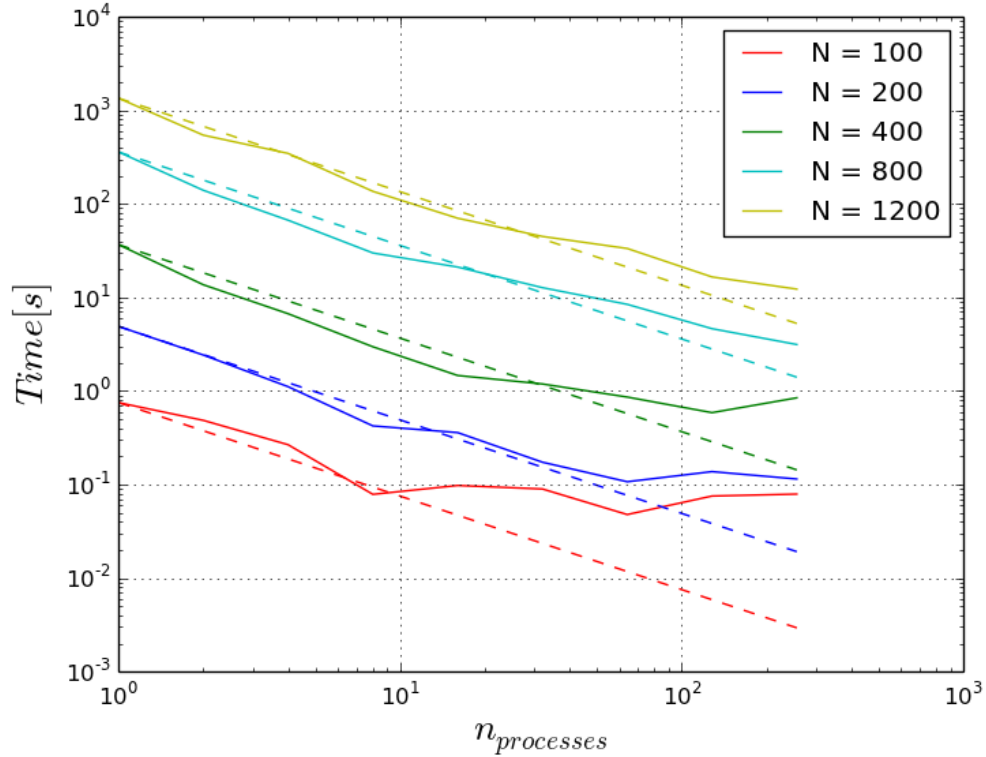
5 Fine-Grained Parallelization: multi-threading

We now turn to the issue of parallelizing *within* a graph simulation, as opposed to across several simulations. This is fundamentally a harder problem from an algorithmic point of view. However, in theory, our test problem again lends itself to simple parallelization. Inspection of the relevant algorithm (code example 2) shows that the treatments of the vertices in the for loop on line 9 are independent of each other. This means that we can use a threaded parallel for loop in this case. The only issue is that we are (potentially) writing to the common data structure `new_types` on line 22. We thus have to lock the data structure upon writing. All other operations are either reads from a common data structure or independent computations and can be left untouched⁹. Thus, in an ideal world,

⁹The function `get_infected_neighbor_fraction()` looks at all the neighbors of a vertex and computes the fraction that are infected, which is also simply a set of read operations and a computation.



(a) linear scale



(b) log-log scale

Figure 4: Execution time versus number of processes for various problem sizes. We achieve a maximum speedup of over $100\times$. See text for details.

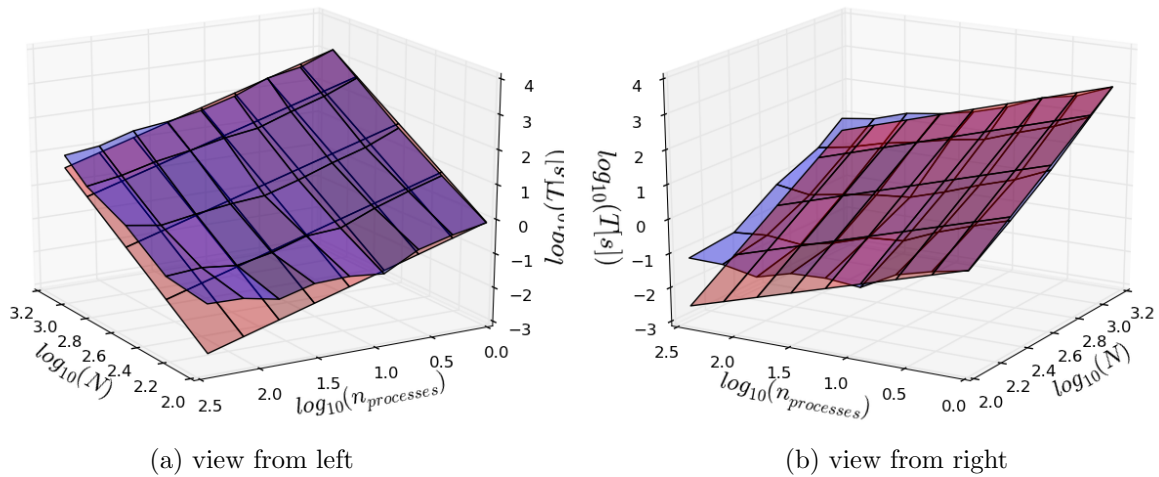


Figure 5: Execution time (blue) versus number of processes and problem size. All axes are logarithmic, such that the ideal scaling becomes a hyperplane (red). See text for details.

```

1  addprocs(2)
2
3  #no @everywhere required
4  Nlist = repmat([1000],10)
5
6  #define function to execute
7  #will not work without @everywhere
8  @everywhere myfun(N,M) = sum(randn(N,M)*randn(M,N))
9
10 #define some local variable
11 #will not work without @everywhere
12 @everywhere M = 1000
13
14 #map over curried function: make sure all captured variables are defined @everywhere!
15 @time pmap(N -> myfun(N,M),Nlist)

```

Code Example 4: Illustration of potential multi-processing issues. `Nlist` is used as the second argument in `pmap` and does not need to be declared `@everywhere`. However, the information about `myfun()` as well as the variable `M` need to be known to all other processes such that the curried function constructed as the first argument to `pmap` is well defined on all remote processors. Therefore we *do* need `@everywhere` statements when declaring these two. This has been as Github issue #14344.

we would simply add `@threads all` to the for loop in line 9, as well as a paired `lock!()` and `unlock!()` statement around the write operation on line 22. This would then parallelize

the code. While the final functioning code indeed does this, there are again several bugs and issues that came up in the development phase that took up most of the development time. We present a simplified version of the thread parallel code in code example 5. The only changes are on lines 7,9 and the lock!() - unlock!() statements around line 23.

```

1  for t in 1:T
2      new_types *= 0
3      update_graph(g,new_types)
4  end
5
6  function update_graph{P}(g::Graph{P},new_types::Array{P,1})
7      m = Mutex()
8      #OUTER LOOP, CHANGED
9      @threads all for v in vertices(g)
10
11     if get_type(g,v) == INFECTED
12
13
14         #INNER LOOP
15         #possibly infect neighbors
16         for w in neighbors(g,v)
17             if get_type(g,w) == SUSCEPTIBLE
18                 x = get_infected_neighbor_fraction(g,w)
19                 p::Float64 = p_birth(x)
20                 if rand() < p
21                     #WRITE OPERATION
22                     lock!(m)
23                     new_types[w] = INFECTED
24                     unlock!(m)
25                 end
26             end
27         end
28     end
29 end
30
31 set_types(g,new_types)
32 end

```

Code Example 5: Simplified Graph Algorithm, thread parallel version. The only changes to the non-parallel version (code example 2) are on lines 7,9 and the lock statements around line 23.

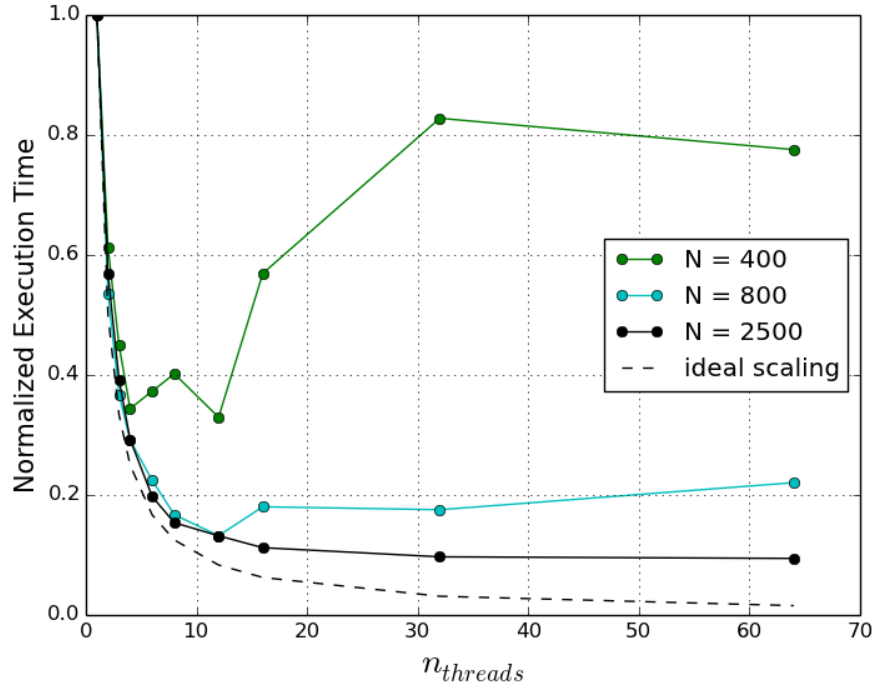
5.1 Results

In this section we test the performance of our multi-threaded code (as shown in 5) on the actual graph algorithm. We work with an 80 core machine, which is subdivided into units of 8 cores. We thus expect ideal “nice” scaling behavior up to 8 threads and some data/communication overheads beyond that. The nature of the algorithm makes this an $O(N^2)$ workload with N being the number of nodes in the graph. We run this workload for various numbers of threads and for various values of N . The results are shown in figures 6 and 7, equivalently to figures 4 and 5. We expect the ideal scaling to go as $\sim \frac{N^2}{n_{threads}}$. In figure 6 we compare the true scaling with $n_{threads}$ (solid lines) to the ideal scaling (dashed lines) for various values of N . For all problem sizes, increasing to up to 4 threads results in a gain. For small problem sizes, the overhead caused by additional threads increases running time. For large problem sizes, we have linear scaling until about 8 threads, as expected, above which the hardware limits the additional gains. We plateau at a maximum speedup of $\sim 11.5 \times > 8$. In figure 7 we summarize this data in a log-log plot, for which the ideal scaling with N^2 and n_{procs} lies on a plane (red). We find that in general the true scaling is somewhat well predicted by the ideal scaling, with the main difference being the departure due to overheads (for small problem sizes) and hardware limitations for small N and large $n_{threads}$.

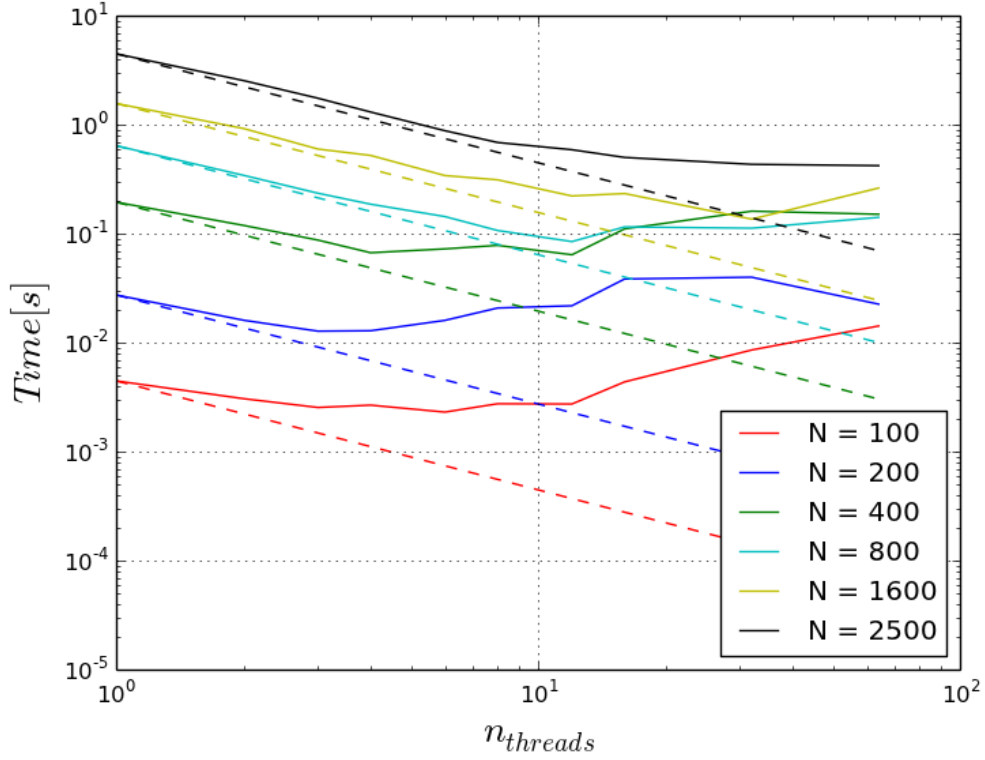
5.2 Implementation Issues and Bugs

The threading package in Julia is still in a very preliminary form. Yet, as it has recently been merged into the master branch, we were able to compile a version of Julia 0.4 that supports threading. Our development process then involved adding the simple parallelization statements as in code example 5. The bulk of the work then followed in the form of debugging the resulting code until it actually worked. We give here a summary of the threading bugs and issues that we encountered in getting the code to function.

- All errors, bugs (including Julia syntax errors!) *within* a threaded block are silently ignored. Every time there is any exception, syntax error, or other problem within a threaded section, the threaded section will simply not execute (or give undefined behavior). This can be a very confusing bug when it comes up unexpectedly.
- Any time global state is modified by threads, the behavior is buggy. It can lead to undefined behavior, `segmentation faults` or other problems. For example, each thread needs its own random number generator (because the Julia RNG requires global state). Another example is that it is not possible to `print` inside a threaded region, which complicates debugging immensely. Anytime a variable’s type is unclear, the program breaks.
- Calling functions between a `lock!()`-`unlock!()` pair can lead to problems (including `segmentation faults`).
- Calling anonymous or curried functions anywhere inside a threaded region can lead to problems (including `segmentation faults`). However, locally defined functions are fine. We illustrate this in code example 6.



(a) linear scale



(b) log-log scale

Figure 6: Execution time versus number of threads for various problem sizes. We achieve a maximum speedup of over $10\times$. See text for details.

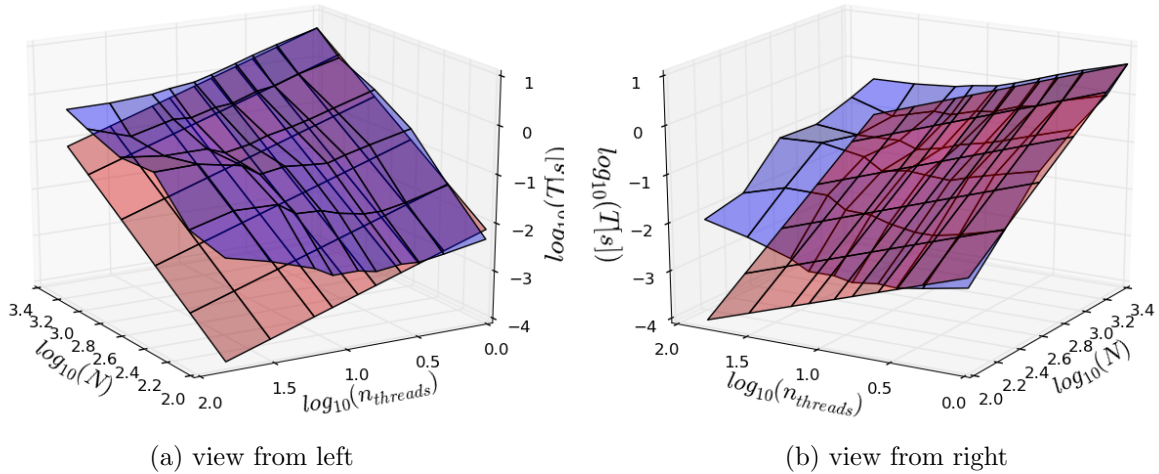


Figure 7: Execution time (blue) versus number of threads and problem size. All axes are logarithmic, such that the ideal scaling becomes a hyperplane (red). See text for details.

- Many of the above issues are *unpredictable*. This means that the issue *may* arise, but doesn't necessarily always do so. This can lead to so-called *Heisenbugs*, which sometimes appear and sometimes do not, unpredictably. In our development, a `segmentation fault` resulting from calling an anonymous function inside a threaded block only occurred about 1 in 10^5 times! One can imagine that this makes debugging extremely difficult.

We illustrate one particularly difficult issue with calling functions inside a threaded block in a minimal code example 6. This has been submitted as Github issue #14343 in the main Julia language repository.

6 Conclusion

In this project we used two archetypal problems in computing, namely Monte-Carlo simulation and iterative graph algorithms to test and illustrate Julia's parallel computing capabilities, both in a coarse-grained multi-processing setting with distributed memory, as well as in a fine-grained multi-threading setting with shared memory. We developed an automatic process management tool for allocating processes efficiently in a cluster environment, without the need to hard code these allocations. This is particularly useful on large, shared machines with changing allocations.

Along the way, we had to deal with incomplete or misleading documentation, undocumented bugs, and many challenging issues in getting the multi-threading and multi-processing to work. We are hopeful that continued development in the Julia project will alleviate these difficulties for future projects using multi-threading and multi-processing. To help in this process, we have submitted some of the most challenging issues as Github issues #14343 and #14344 to the Julia language project.

```

1  type CarryFunction
2      fn::Function
3  end
4
5  alpha = 0.1
6  fn(x) = alpha*x
7
8  function use_anonymous(N::Int, c::CarryFunction)
9      a = zeros(N)
10     @threads all for i in 1:length(a)
11         # a[i] = fn(i) #NO SEGFAULT
12         a[i] = c.fn(i) #SEGFAULT (sometimes... but not always!)
13     end
14     println(a[1],a[end])
15 end
16
17 length = 10000
18 repetitions = 100
19 for j = 1:repetitions
20     use_anonymous(length,CarryFunction(fn))
21 end

```

Code Example 6: A minimal example of how to avoid an insidious threading bug. The two statements on lines 11 and 12 perform the same work. However, in the upper case, the function is simply a locally defined function. In the lower case, the function is data as part of a type. When the function is data, it sometimes leads to **segmentation faults**. In the upper case, the code runs without problems. This has been submitted as Github issue #14343.

Our results show that both the multi-processing and multi-threading codes are truly taking advantage of their full hardware potential (several cores across several different machines in the multi-processing case, and all cores on a single machine in the multi-threading case). Putting all our speed improvements together, we achieved a $60\times$ speedup by moving from python to optimized serial Julia, a $\sim 12\times$ speedup by using multi-threading to parallelize our graph algorithm on an 80 core machine, and a $\sim O(X)$ speedup when parallelizing our Monte Carlo trials across X cores on several different machines in a cluster (in our case $X \sim 100$, but in general unlimited). Overall, if we had access to ~ 100 machines of 64 cores each, we can thus expect $60 \times 10 \times 100 \sim 4 - 5$ orders of magnitude speedup when running on all 64 threads on ~ 100 machines as compared to our original, serial python code.

7 Acknowledgments

This work was produced as an assignment for the course 6.338 at MIT, Fall Semester 2015-2016. We thank Prof. Alan Edelman and the Julia development team at MIT CSAIL for their support in this project in the form of advice, encouragement, and access to computational resources.

References

- [1] Jeff Bezanson, Stefan Karpinski, Viral B Shah, and Alan Edelman. Julia: A fast dynamic language for technical computing. *arXiv preprint arXiv:1209.5145*, 2012.
- [2] Daniel A Schult and P Swart. Exploring network structure, dynamics, and function using networkx. In *Proceedings of the 7th Python in Science Conferences (SciPy 2008)*, volume 2008, pages 11–16, 2008.
- [3] James Hardy Wilkinson, James Hardy Wilkinson, and James Hardy Wilkinson. *The algebraic eigenvalue problem*, volume 87. Clarendon Press Oxford, 1965.