

# COSC3500

## 2D Orbital Simulation Report

Maxwell Bo (43926871)

August 17, 2018

### Description

The task was to create a stock-standard, 2-dimensional gravitational  $n$ -body simulator. All bodies were to be assumed to be point masses. The simulation was to be accurate, maintaining a constant total energy, and exhibiting phenomena such as apsidal precession. The simulation was to accept arguments specifying the granularity of the simulation (number of time steps, number of instances of data export), and a file specifying masses and their initial positions and velocities. The simulation was to produce its output as fast as possible, with minimal slowdown with an increasing  $n$  number of bodies.

The simulator did not need to handle collisions between bodies. The MS1 simulator was to be free of OpenMP multiprocessing, which would be implemented in preparation for the MS2 submission.

### Implementation

At a high-level, the initial naive simulator:

1. Parsed input parameters and files
2. Constructed `Body` class instances, representing each point mass
3. Packed the `Bodys` into a `std::vector<Body>`, to maximize cache locality
4. Calculated forces between all pairwise combinations of  $n$ -bodies ( $O(n^2)$ )
5. Performed Euler's method to derive new velocities and positions
6. Output all necessary data
7. GOTO 4

By using a Quadtree, ('a tree datastructure in which each internal node has exactly four children') and the *Barnes-Hut* algorithm[1], the total cost of force calculation could be reduced to  $O(n \log n)$ , by grouping close-together bodies and approximating forces between the singular grouped pseudo-body, and distant bodies. New Quadtrees were constructed on each seperate simulation step.

Dehen and Read note that the Euler method 'performs very poorly in practice', further noting that 'errors are proportional to  $\Delta t^2$ '. They contrast it with the second-order *Leapfrog* symplectic integrator, which is 'heavily used in collisionless N-body applications'.[2]

*Leapfrog* can be expressed many in forms[3] including a synchronised form:

$$\begin{aligned}
x_i &= x_{i-1} + v_{i-1/2} \Delta t \\
a_i &= F(x_i) \\
v_{i+1/2} &= v_{i-1/2} + a_i \Delta t
\end{aligned}$$

which only requires a single acceleration calculation per every two half timesteps (the timestep  $\Delta t$  must be constant to maintain stability), and a ‘kick-drift-kick’ form

$$\begin{aligned}
v_{i+1/2} &= v_i + a_i \frac{\Delta t}{2} \\
x_{i+1} &= x_i + v_{i+1/2} \Delta t \\
v_{i+1} &= v_{i+1/2} + a_{i+1} \frac{\Delta t}{2}
\end{aligned}$$

that is stable with variable timstepping, but incurs an additional acceleration calculation per every two half timesteps.

The synchronised form was implemented, but attempts to implement the kick-drift-kick form, and variable timestepping, were left unfinished.

Thus, the final implementation:

1. Parsed input parameters and files
2. Constructed **Body** class instances, representing each point mass
3. Packed the **Bodys** into a `std::vector<Body>`, to maximize cache locality
4. Inserted all **Bodys** into a fresh **QuadTree** on full timesteps, traversing the **QuadTree** with every **Body** to calculate forces ( $O(n \log n)$ )
5. Performed the appropriate *Leapfrog* step to derive new velocities *or* positions
6. Output all necessary data
7. GOTO 4

## Correctness

I personally believe that the simulation is relatively accurate. By visualising the results with `matplotlib`, we see something that resembles an  $n$ -body simulator. The total energy is flat, observing coefficients of variation as low as 0.001%. Euler’s method consistently demonstrated coefficients of variation three times higher than that of *Leapfrog*. Due to recommendations by literature, and observed data, the use of Euler’s method was gradually phased out during my testing to speed up the process.

*Barnes-Hut* caused a significant increase in observed coefficient of variation, averaging 0.08% across multiple runs. Furthermore, total energy was observed to step up and down at varying intervals<sup>1</sup>. I suspect that this was because certain force calculations were causing groups of bodies to be approximated as a single pseudobody, after other had strayed too far from the pseudobody’s quadtree’s node’s centre of mass.

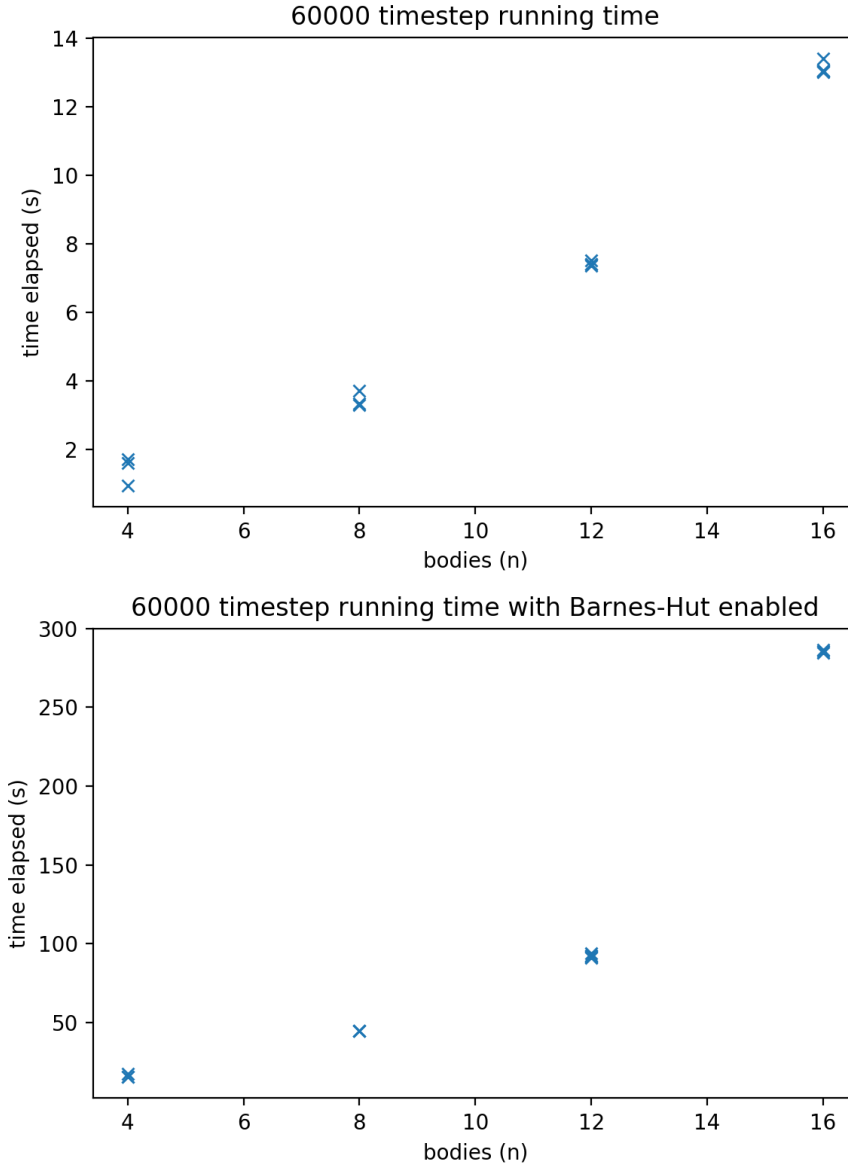
When bodies are in close proximity, anomalous total energies are observed<sup>1</sup>. Furthermore, simulations with higher numbers of bodies produce more random low energy outliers (likely due to a greater number of close encounters). There seems to be no significant difference between Euler method and *Leapfrog*, with respect to observation of anomalies.

## Performance & Scaling

Henceforth, the use of the ‘recognizable’ refers to eyeballing the output data, and making no significant effort to investigate the data more rigorously.

The *Barnes-Hut* algorithm appeared to be dominated by its constant factor.?? Generating input files with a high number of bodies that did not cause ‘spontaneous combustion’ (where bodies would fly away from each other in every direction), and unconstrained Quadtree growth (where bodies achieve escape velocity, rapidly expanding the size of the tree and causing a collapse in simulation accuracy) proved difficult. Further testing on *goliath* is needed to discover the  $n$  that causes *Barnes-Hut* to be more performant than the brute-force approach.

Strangely,  $n = 16$  with *Barnes-Hut* enabled took longer than expected, when compared to  $n \in \{4, 8, 12\}$ . Further investigation is required.



The addition of the `-march=native` compiler flag, which enables the use of all CPU specific instructions, provided no recognizable improvement in running time, but was left enabled in the instance that it improved performance on *goliath*.

The use of both the GCC and Clang Profile-Guided Optimisation features provided no recognizable improvement in running time.

Distressingly, -00, -01, -02, -03 showed no recognizable improvement in running time. -0fast led to an -4%-ish performance regression.

By profiling with `callgrind`, we saw that execution was dominated by only one incredibly costly user-defined method, with an exclusive cost of 26.33% of total running time.

Incl.	Self	Called	Function	Location
100.00	0.00	(0)	0x0000000000001030	ld-2.17.so
100.00	0.00	2	_dl_runtime_resolve_xsave	ld-2.17.so
100.00	0.00	1	0x00000000000040df6	nbody
100.00	0.00	(0)	(below main)	libc-2.17.so
100.00	9.33	1	main	nbody: main.cpp, basic_string.h, string_conversio...
89.22	26.33	672 000 056	Body::exert_force_unidirec...	nbody: Body.cpp
62.90	2.93	672 033 712	distance(double, double, ...	nbody: utils.cpp
59.97	11.70	672 033 711	hypot	libm-2.17.so
48.27	48.27	672 033 712	__hypot_finite	libm-2.17.so
0.63	0.63	48 000 000	Body::frog(double)	nbody: Body.cpp
0.52	0.52	48 000 008	Body::leap(double)	nbody: Body.cpp
0.21	0.21	48 000 000	Body::reset_force()	nbody: Body.cpp
0.07	0.00	601	dump_timestep(double, st...	nbody: main.cpp, stl_vector.h, stl_iterator.h

Performance fixes were devised.

```
void Body::exert_force_unidirectionally(const Body& there) {
    const double m1 = m;
    const double m2 = there.m;

    const double delta_x = there.x - x;
    const double delta_y = there.y - y;

-   const double r = distance(x, y, there.x, there.y);
+   const double r = hypot(delta_x, delta_y);
-   const double r2 = r * r;
+   const double r2 = pow(r, 2);

    const double F = (G * m1 * m2) / r2;

    // turn the displacement vector between our two points into a force vector
    // of the desired magnitude
    const double scale_factor = F / r;

    Fx += delta_x * scale_factor;
    Fy += delta_y * scale_factor;
}
```

Instead of recalculating  $\Delta x$  and  $\Delta y$  twice (the second time in `distance`), we calculate them only once and make a direct call to `hypot`, rather than making a call to `distance`. We also used the specialized `pow` provided by `cmath`. This yielded a recognizable 15%ish performance improvement.

```
diff --git a/Assignment_1/src/Body.cpp b/Assignment_1/src/Body.cpp
```

```
+void Body::exert_force_bidirectionally(Body& there) {
+   const double m1 = m;
+   const double m2 = there.m;
+
+   const double delta_x = there.x - x;
```

```

+   const double delta_y = there.y - y;
+
+   const double r = hypot(delta_x, delta_y);
+   const double r2 = pow(r, 2);
+
+   const double F = (G * m1 * m2) / r2;
+
+   // turn the displacement vector between our two points into a force vector
+   // of the desired magnitude
+   const double scale_factor = F / r;
+
+   Fx += delta_x * scale_factor;
+   Fy += delta_y * scale_factor;
+
+   there.Fx -= delta_x * scale_factor;
+   there.Fy -= delta_y * scale_factor;
+}

```

```

diff --git a/Assignment_1/src/main.cpp b/Assignment_1/src/main.cpp
@@ -258,8 +259,7 @@ int main(int argc, char **argv) {
    for (size_t j = i + 1; j < bodies.size(); j++) {
        auto& y = bodies[j];
-        x.exert_force_unidirectionally(y);
-        y.exert_force_unidirectionally(x);
+        x.exert_force_bidirectionally(y);
    }
}

```

This fix halved execution time, for obvious reasons.<sup>1</sup>.

```

diff --git a/Assignment_1/src/Body.cpp b/Assignment_1/src/Body.cpp
@@ -83,11 +82,11 @@ double Body::kinetic_energy() const {
    double Body::gravitational_potential_energy(const Body& there) const {
        const double R = distance(x, y, there.x, there.y); // final distance, aka, to edge

-    return (-G * m * there.m) / R;
+    return (-Gm * there.m) / R;
    }

@@ -101,7 +100,7 @@ void Body::exert_force_unidirectionally(const Body& there) {

-    const double F = (G * m1 * m2) / r2;
+    const double F = (Gm * m2) / r2;

@@ -113,7 +112,6 @@ void Body::exert_force_unidirectionally(const Body& there) {

-    const double F = (G * m1 * m2) / r2;
+    const double F = (Gm * m2) / r2;

```

---

<sup>1</sup>I had to throw out all my old profile data

```

@@ -189,6 +189,7 @@ std::vector<Body> parse_input_file(std::ifstream& input_fh) {

    for (size_t i = 0; i < bodies.size(); i++) {
        bodies[i].m = masses[i];
+       bodies[i].Gm = G * masses[i];
    }
}

```

Precomputing  $Gm$  yielded a 3%ish performance improvement.

## Parallelism

### Implementation

#### OpenMP

I slapped as many `#pragma omp parallel for shared(bodies)` on as many loops as I could find, and quickly discovered massive lock contention issues associated with `#pragma omp critical`. Pairwise force updates between bodies:

```

    for (size_t i = 0; i < bodies.size() - 1; i++) {
        auto& x = bodies[i];

        for (size_t j = i + 1; j < bodies.size(); j++) {
            auto& y = bodies[j];
            x.exert_force_bidirectionally(y);
        }
    }

```

where

```

void Body::exert_force_bidirectionally(Body& there) {
    ...
    #pragma omp critical
    {
        Fx += delta_x * scale_factor;
        Fy += delta_y * scale_factor;

        there.Fx -= delta_x * scale_factor;
        there.Fy -= delta_y * scale_factor;
    }
}

```

which had provided massive performance gains in the serial implementation over unidirectional updates was now causing massive parallel slowdowns, where increasing the `OMP_NUM_THREADS` caused a proportional increase in walltime. Instead, a **critical-less** comparison schema was devised,

```

#pragma omp parallel for shared(bodies)
for (size_t i = 0; i < bodies.size(); i++) {
    auto& x = bodies[i];

    for (size_t j = 0; j < bodies.size(); j++) {

```

```

    auto& y = bodies[j];

    if (&bodies[i] != &bodies[j]) {
        // XXX: Do NOT swap these around. You will cause
        // race conditions
        x.exert_force_unidirectionally(y);
    }
}
}

```

that provided the sought after parallel speedup. This parallel speedup was found to be present at where  $n$ -bodies was high. Testing focussed primarily on  $n > 5000$ ,<sup>2</sup> where force and integration calculation overhead dominated thread creation overhead. Parallel speedup vanished around the  $n \approx 500$  mark.

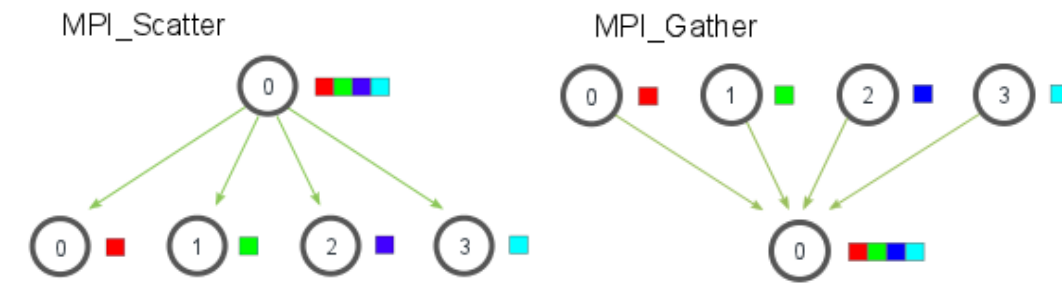
## MPI

Looking at our high-level implementation, the following two steps consume the vast majority of our computing power:

1. Inserted all **Bodys** into a fresh **QuadTree** on full timesteps, traversing the **QuadTree** with every **Body** to calculate forces ( $O(n \log n)$ )
2. Performed the appropriate *Leapfrog* step to derive new velocities *or* positions

The *Leapfrog* step is the easiest step to distribute amongst many nodes, as it does not require the comparison of a body to all other bodies - each node can compute the step on a subset of all bodies, and later combine each subset with relatively simple MPI code.

This can be achieved with `MPI.Scatterv` and `MPI.Gatherv`<sup>3</sup>. We decompose each dimension of our **Bodys** into arrays, (e.g. the  $x$  position, the  $x$  velocity, etc.), `MPI.Scatterv` these arrays, recompose them into **Bodys**, run *Leapfrog* on this subset on each node, decompose the subset into arrays, `MPI.Gatherv` these arrays back to the root node, and recompose them into **Bodys**.



```
std::vector<Body> sbodies = scatter_bodies(bodies, size, rank);
```

```

#pragma omp parallel for shared(bodies)
for (size_t i = 0; i < sbodies.size(); i++) {
    auto& body = sbodies[i];

```

<sup>2</sup>In MS1, I failed to test simulations with a large number of  $n$ , as I was focussed on testing the accuracy of the implementation, and generated input files would cause high numbers of collisions that made accuracy assessment difficult. Testing  $n = 10000$ , we find a 10x speedup with by using a serial *Barnes-Hut* implementation (Avg. 1.12s, as opposed to 9.87s without *Barnes-Hut*). *Barnes-Hut* afforded a much simpler parallelization scheme, as each body could calculate its force against the quadtree concurrently, without locks

<sup>3</sup>These `v` variants support custom chunking sizes e.g. if we have 15 bodies we can send 4, 4, 4, 3 bodies to each node respectively, which `MPI.Scatter` and `MPI.Gather` do not support dynamic chunking

```

        if (step % 2 == LEAP) {
            body.leap(timestep);
        }
        else {
            body.frog(timestep);
        }
    }
}

bodies = gather_bodies(sbodies, size, rank, bodies.size());

```

where `sbodies` is our subset of bodies distributed to each node.

Because force calculation is only performed on the root node, this messaging architecture is *master-slave*, where the root node coordinates gathers and scatters. Slave nodes, are, unfortunately, underutilized. Further work could be sent distributing force calculation work, doing away with a root node entirely, except for logging.

```

Total CPU time on rank 0 was 89.820000
Total CPU time on rank 2 was 36.750000
Total CPU time on rank 3 was 36.960000
Total CPU time on rank 1 was 37.040000

```

## Verification & Correctness

Because of my confidence in the correctness of the serial implementation, it served as a reference implementation for the verification of the parallel implementation.

Both implementations were run with identical input parameters, and their outputs were `diffd`. When run without MPI enabled, outputs were identical. When run with MPI, floating point abnormalities (initially affecting the least significant digits of each floating point value) gradually caused each diff to diverge. I hypothesize that the MPI serialization and deserialization procedure was truncating precision and causing this divergence. Perhaps some option exists to ensure that MPI guarantees floating point precision.

## Scalability

We have multiple dimensions that we can change to measure the scalability of the parallel implementation.

1. The number of bodies ( $n$ ) (512, 1024, 2048, 4096, 8192, 16384)
2. `ENABLE_BARNES_HUT` (`true`, `false`)
3. `--nodes` (Disabled, 1, 4, 16)
4. `--ntasks` (1, 4, 16)
5. `--ntasks-per-node` (1, 4)
6. `--cpus-per-task` (1, 4, 16)
7. `#pragma omp parallel schedule(static, dynamic, guided)`

Running tests with all valid combinations of these variables, with a `numTimeSteps` as low as 8, would produce useful data, across multiple dimension, with minimal load on cluster. It would be unlikely that this plan would be curtailed.



## Testing

### Changes since MS2

In MS2, we decomposed each dimension of our `Bodys` into arrays, (e.g. the  $x$  position, the  $x$  velocity, etc.), `MPI_Scattered` these arrays, recomposed them into `Bodys`, ran `Leapfrog` on this subset on each node, decomposed the subset into arrays, `MPI_Gathered` these arrays back to the root node, and recomposed them into `Bodys`. We were also allocating a fresh new `Vector` after `Body` recomposition. This was, frankly, insane, and took a program with very low memory usage and pressure in the serial implementation to one that constantly allocated and deallocated memory on every single step.

By declaring `MPI_Datatype MPI_Body` with `MPI_Type_contiguous(8, MPI_DOUBLE, &MPI_Body)`, we're now able to send and receive subsets of `Bodys` with just a single `MPI_Scatterv` and `MPI_Gatherv` call. We also `MPI_Scatterv` and `MPI_Gatherv` directly into `Vectors` allocated at the start of the program, without intermediate array allocations. This eliminated transient out-of-memory errors, and likely reduced MPI overhead.

### Testing Performed

1. The number of bodies ( $n$ ) (4, 16, 64, 256, 1024, 4096)
2. `--nodes` (1, 2, 4, 8, 12)
3. `--ntasks-per-node` (1, 2, 4, 8)
4. `--cpus-per-task` (1, 2, 4, 8)
5. `ENABLE_BARNES_HUT` (true, false)

We test with 10 `numTimeSteps`, to drown out the time dedicated to input parsing and upfront memory allocation (which is in the vicinity of 100ms).

All valid combinations<sup>4</sup> of the values in the lists above were performed. Tweaks in the number of bodies ( $n$ ) (a step from 4 to 16 bodies is double the body density; tests with 16384 bodies did not terminate), `--nodes` (`dogmatix` only has 12 nodes in the `cosc` partition; it did not make sense to disable MPI), not setting `--ntasks` parameter (it made sense to allow Slurm to dynamically allocate the required number of tasks based on our `--ntasks-per-node` and `--cpus-per-task` parameters) and adding the `--ntasks-per-node` parameter were made since our MS2 plan. Further granularity in the `--cpus-per-task` parameter was also introduced.

The OMP scheduling policy was unable to be set by program argument, and would have require three separate binaries compiled with different options - thus, this was also removed as a tested variable.

The collected data has been submitted with the report as `collected_data.json`, and is visualised at the end of the report.

### Scaling

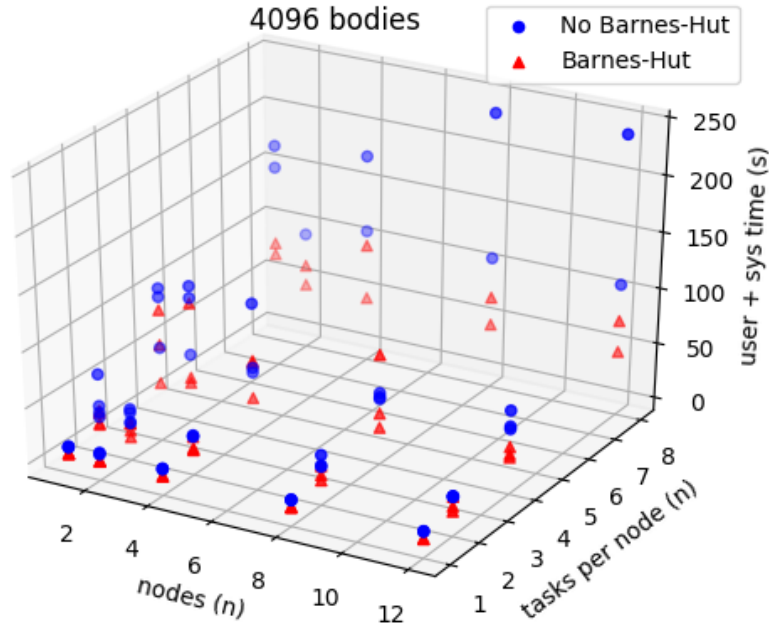
We measure time with a combination of the `sys` and `user` times reported by the `time` coreutil.

We can see a few trends emerge:

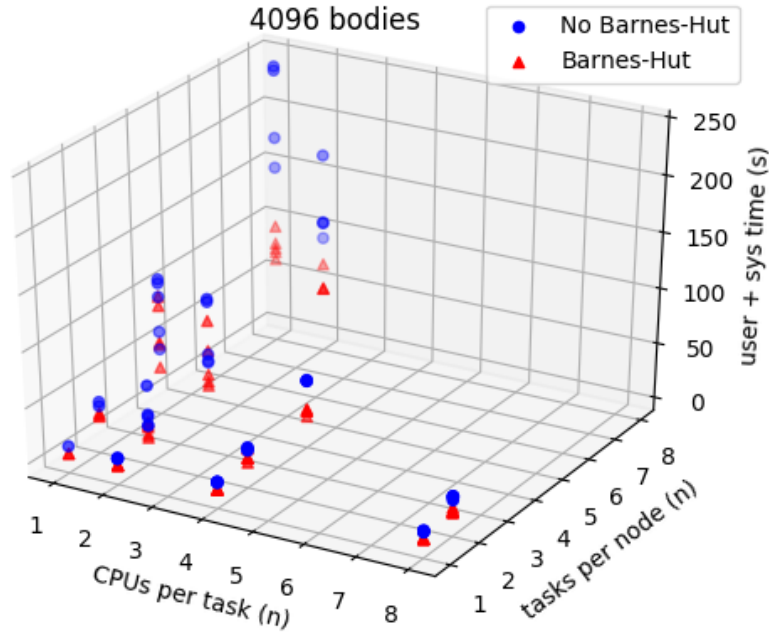
`--ntasks-per-node=8` causes time to be above an order of magnitude higher, across all results. This behaviour is most pronounced when testing with 4096 bodies, across a variable number of nodes. We can also see the pronounced effect of *Barnes-Hut* ( $O(n \log n)$ ) at 4096 bodies.

---

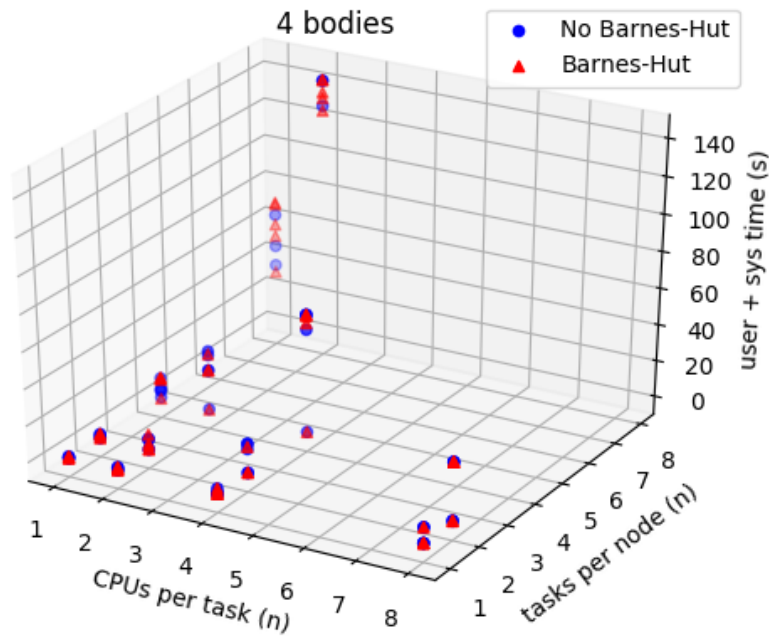
<sup>4</sup>As an example, `--ntasks-per-node=8`, `cpus-per-task=8` is not an admissible configuration, as this would require 64 CPUs on a single node. `dogmatix` only has nodes with 24 or 28 CPUs



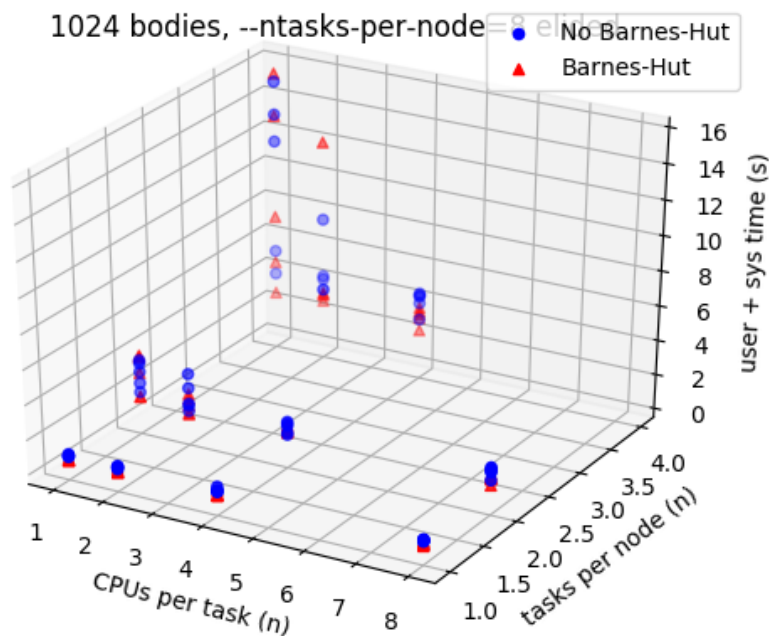
The behaviour becomes less significant when increasing the number of CPUs per task, resulting in maximum CPU utilization per node. I hypothesize that process context switch overhead is causing this behaviour, as having a greater number of tasks allocated to a node than available CPUs causes the CPU (or CPUs) to multiprocess tasks.



Furthermore, the behaviour is still pronounced at 4 bodies, and the resulting times at `--ntasks-per-node=8` do not appear to be a function of the number of bodies.



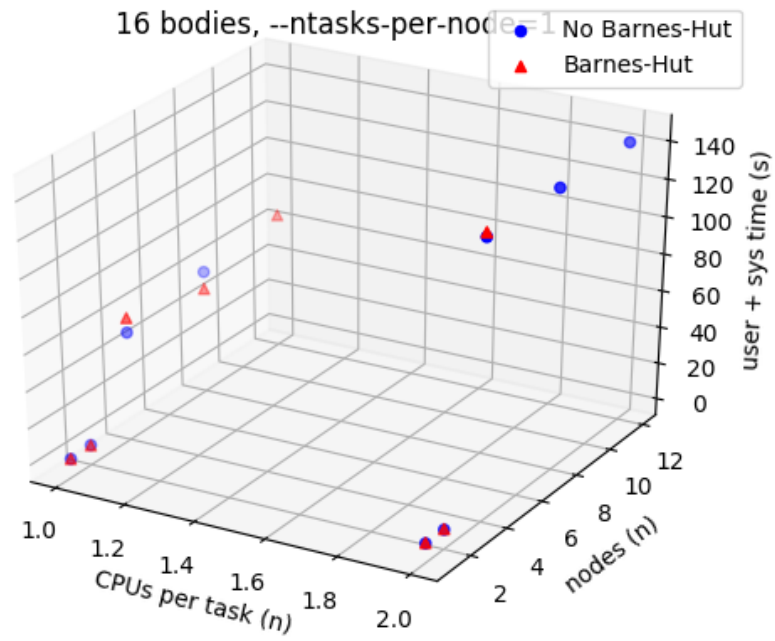
Further evidence can be seen for this multiprocessing overhead hypothesis, where decreasing the total number of processes allocated to each CPU decreases time.



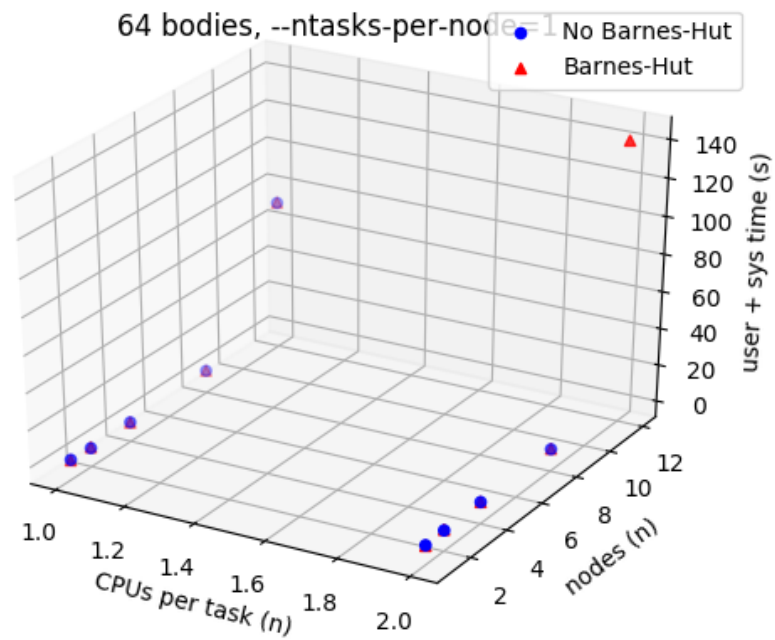
We can also see that `--ntasks-per-node=1` runs with consistently good performance. We'll drill into data where `--ntasks-per-node=1` and further analyse its scaling characteristics.

We see a large spike in running time when we move from 4 to 8 nodes when testing with 16 bodies. This is the only number of bodies that this behaviour is exhibited. I have no explanation

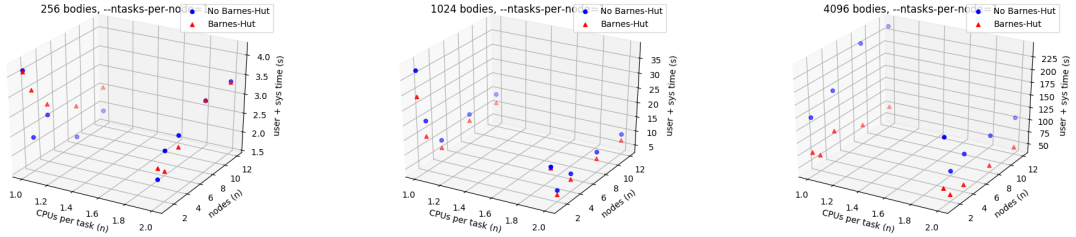
for this behaviour.



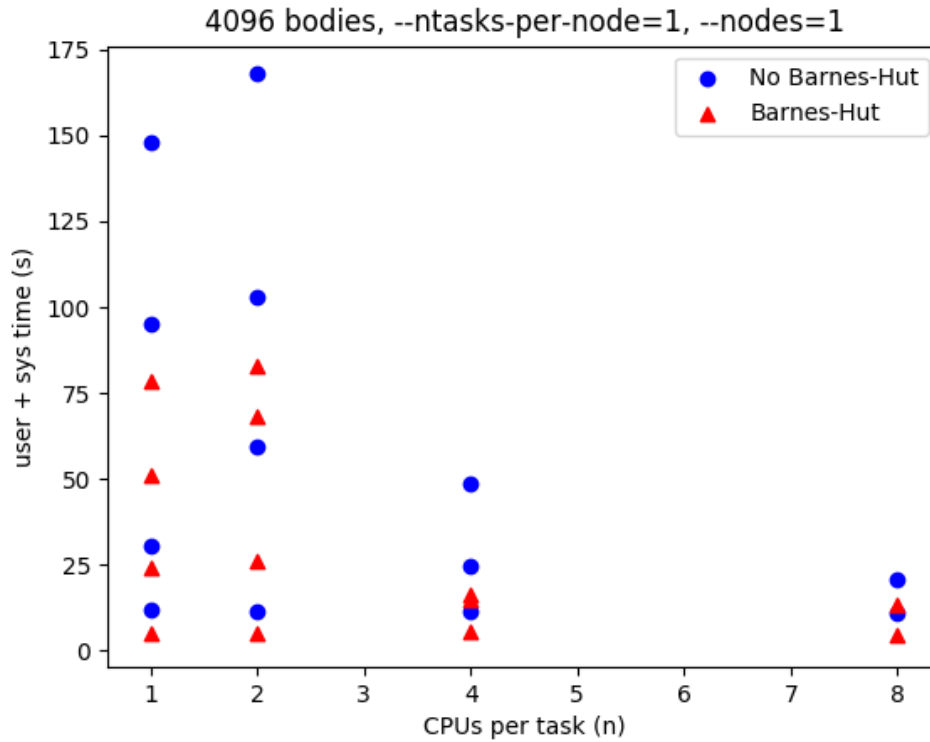
The spike occurs later for 64 bodies, instead at 12 nodes.



At 256 we see worse performance with 2 nodes and 1 CPU per task than with 4 nodes. This behaviour vanishes by 4096 bodies, and instead we see an adding more nodes increases the performance.



Perhaps the only thing that provides a parallel speedup is dedicating more CPUs to each task.



## Summary of Findings

### What Worked

*Barnes-Hut* provided an effective reduction of complexity. Small micro-optimisations made with the assistance of a profiler provided decent performance improvements. *Leapfrog* increased simulation accuracy. OpenMP provided a parallel speedup (after efforts were made to eliminate locking).

### What Didn't Work

Apart from a few anomalies, adding more nodes degraded performance. Running multiple tasks on the same CPU further degraded performance.

We failed to find the threshold where MPI overhead gave way to a parallel speedup. Future testing plans may require even higher numbers of bodies.

Testing with higher counts of bodies caused transient out-of-memory issues that persisted even after requests for more memory per node were made. This dramatically reduced the statistical significance of tests with 16384 bodies.

## What Next

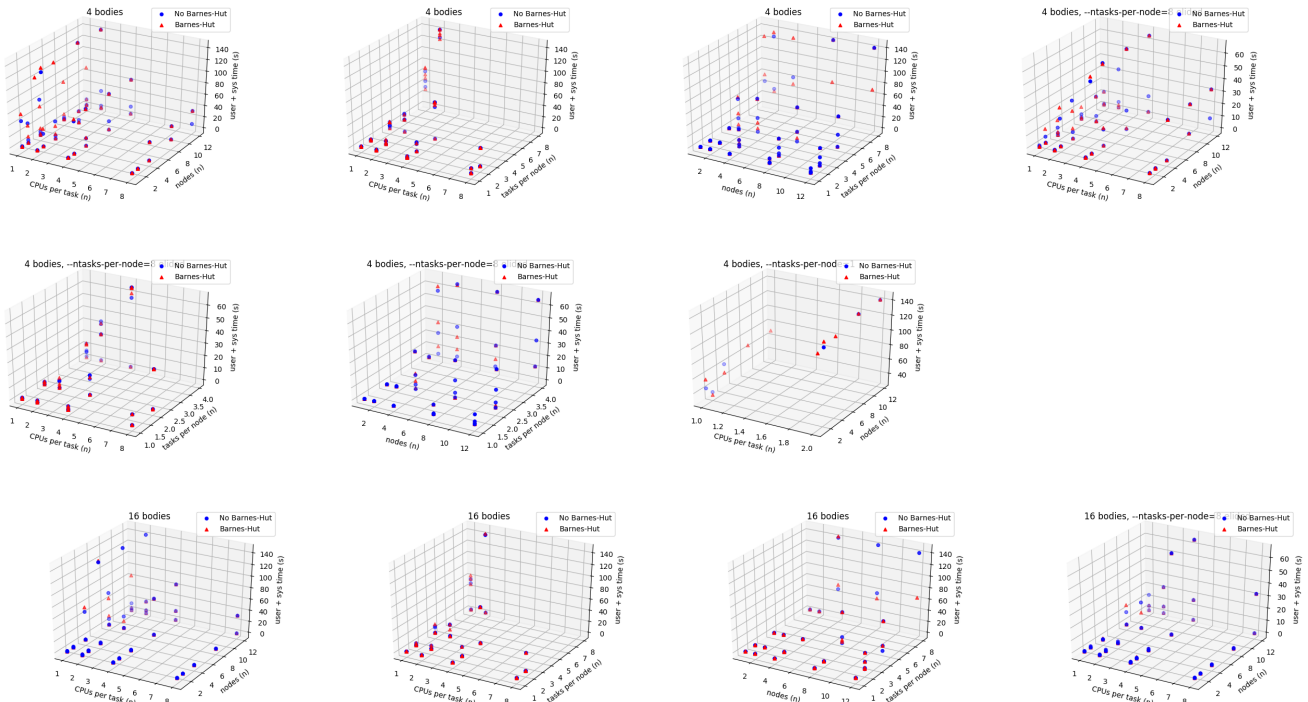
Efforts can be made to reduce the MPI overhead. Currently `Bodys` are sent in their entirety, at every single simulation step. We do not have to send `m` and the precomputed `Gm` and `m`, as each node can recover this information from the input file. The ‘leap’ step does not need `Body::Fx` and `Body::Fy`, and ‘frog’ does not need `Body::x` and `Body::y`. Thus, we do not need to `MPI_Scatterv` and `MPI_Gatherv` this information in preparation for each step. This effectively halves the amount of data that would be sent by MPI. Efforts would need to be made to determine how much of the MPI overhead is due to a hard bandwidth limit on the interconnects between nodes. I hypothesise that MPI overhead scales linearly with data transferred, and that the MPI transmission protocol is responsible for the majority of the overhead at low data transfer quantities.

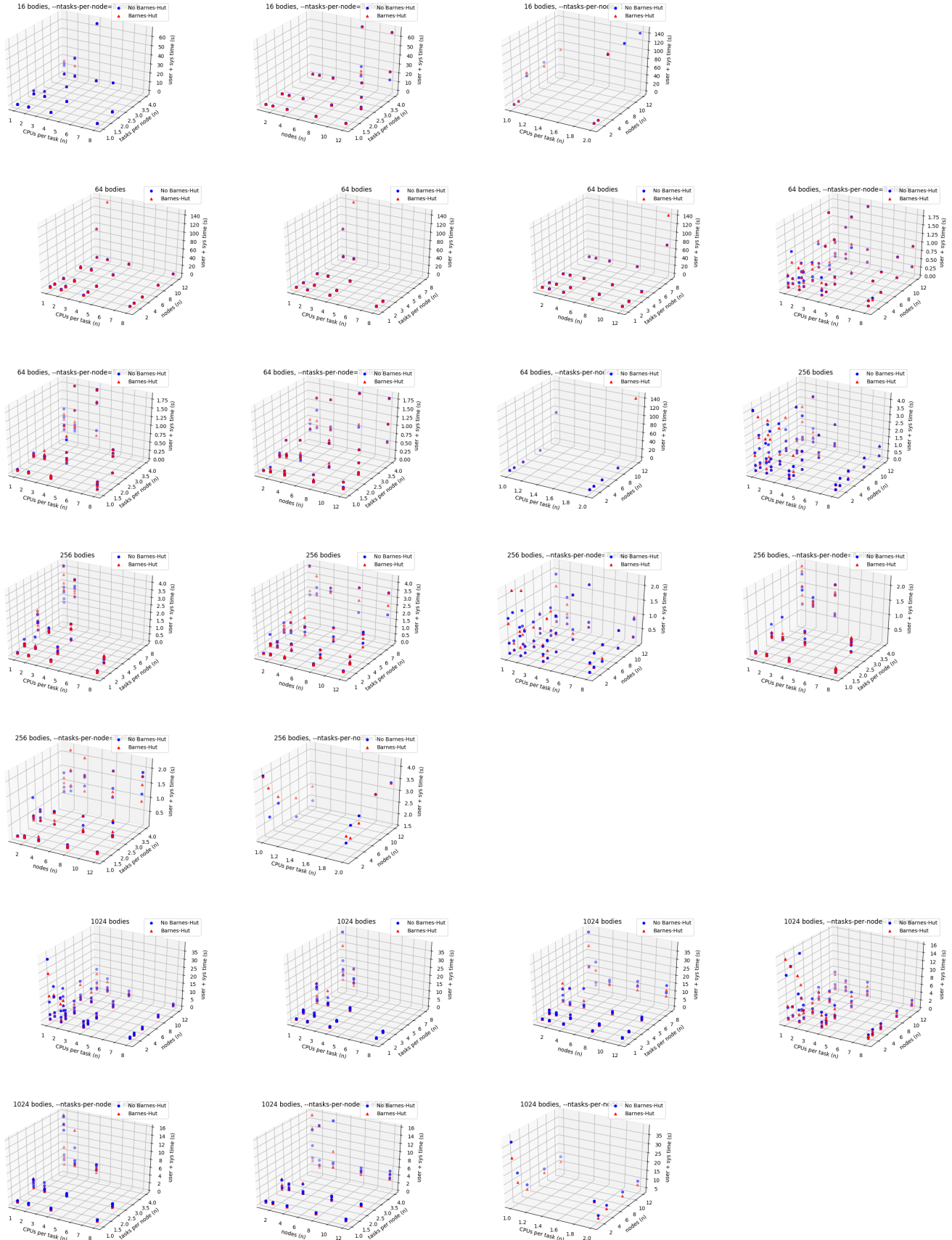
The *Barnes-Hut* algorithm provides an effective means of compartmentalizing computations, eliminating the need to `MPI_Gatherv` everything back to the root node after each step. Each Quad of the `QuadTree` could be assigned to a separate node; data would only be transferred between nodes when bodies stray over the boundary of a specific Quad. This would dramatically increase the complexity of the MPI messaging implementation.

## References

- [1] J. E. Barnes and P. Hut, “A hierarchical  $O(n\log n)$  force calculation algorithm,” *Nature*, vol. 324, p. 446, 1986.
- [2] W. Dehnen and J. I. Read, “N-body simulations of gravitational dynamics,” *European Physical Journal Plus*, vol. 126, p. 55, May 2011.
- [3] R. D. Skeel, “Variable step size destabilizes the störmer/leapfrog/verlet method,” *BIT Numerical Mathematics*, vol. 33, pp. 172–175, Mar 1993.

## Appendix







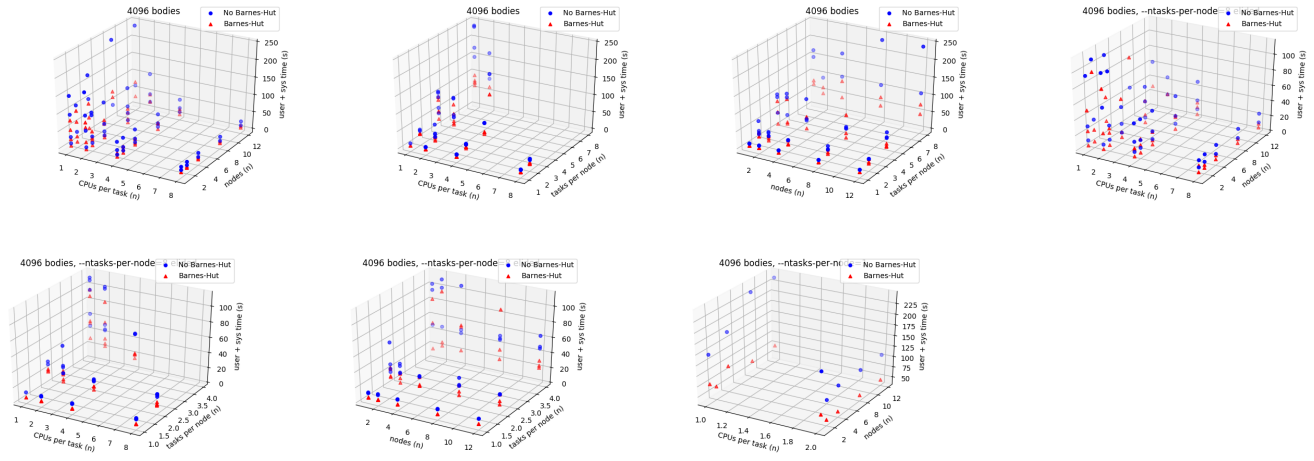




Figure 1: Observed energy anomaly while bodies in close proximity - Leapfrog

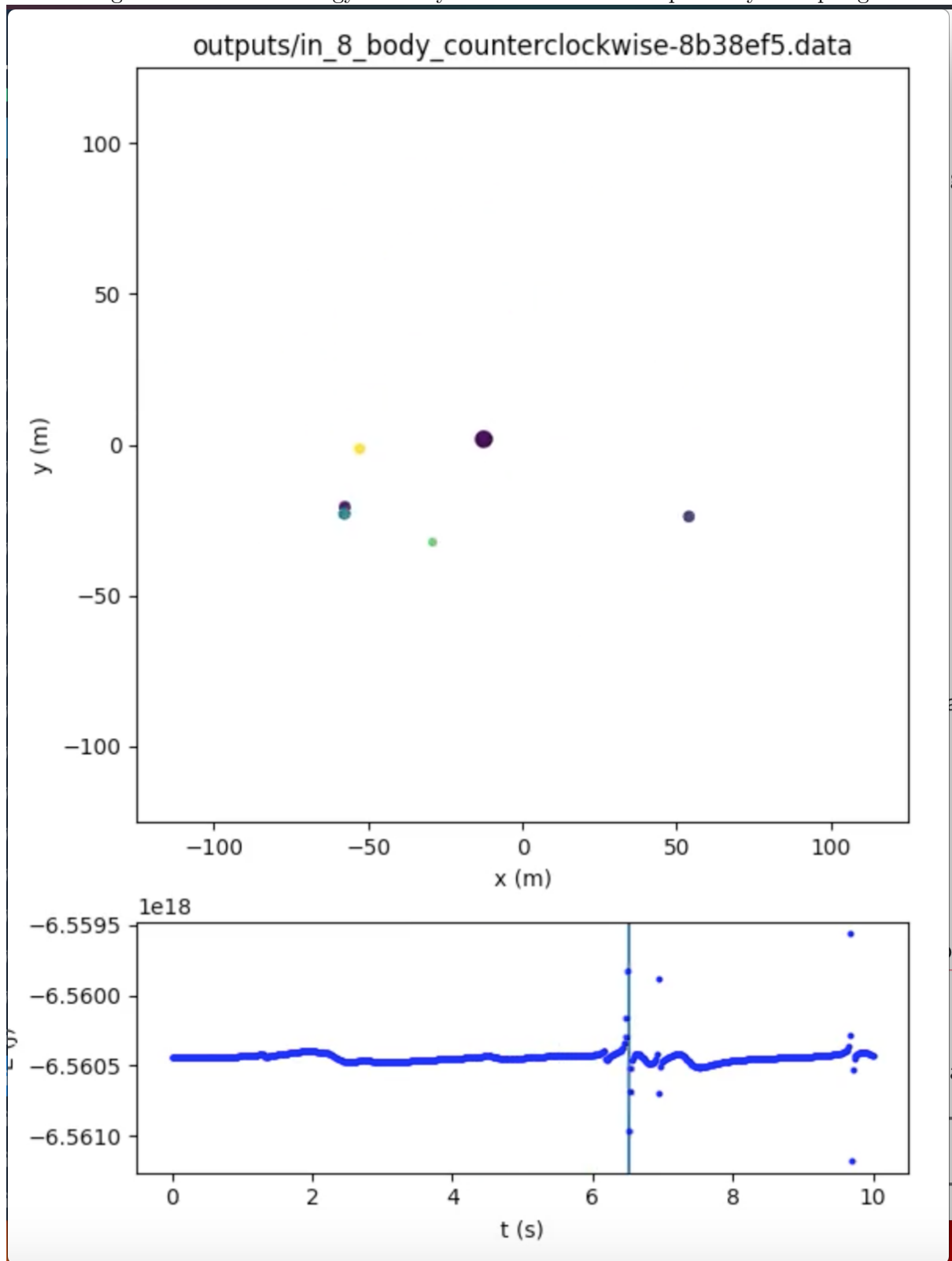


Figure 2: Barnes-Hut energy variation - Leapfrog

