# Parallelising your OCaml code with Multicore OCaml

*Sadiq Jaffer*, Tom Kelly, Sudha Parimala, KC Sivaramakrishnan, Anil Madhavapeddy

# Overview

- Multicore OCaml
- Domains
- Domainslib
- Further optimisation

# Multicore OCaml

- Concurrency is overlapping computations
- Parallelism is simultaneous computations

Multicore OCaml ⇨ concurrency *and* shared-memory parallelism

# Multicore OCaml

- Compatible with existing OCaml code (inc ppx)
- OCaml 5 will have parallelism via Domains
- Concurrency via effects and fibers to follow

# Domains

- Unit of parallelism
- Heavyweight
- Functionality
  - Spawn/join
  - Wait/notify
  - Atomic memory operations
  - Local storage

# N-Body

- Derived from benchmarks game
- Models orbit of a number of bodies

# N-Body serial

```
1  let advance bodies n_bodies dt =
2    for i = 0 to n_bodies - 1 do
3      let b = bodies.(i) in
4      for j = i+1 to n_bodies - 1 do
5        let b_o = bodies.(j) in
6        let dx = b.x -. b_o.x  and dy = b.y -. b_o.y
7          and dz = b.z -. b_o.z in
8        let dist2 = dx *. dx +. dy *. dy +. dz *. dz in
9        let mag = dt /. (dist2 *. sqrt(dist2)) in
10
11       b.vx <- b.vx -. dx *. b_o.mass *. mag;
12       b.vy <- b.vy -. dy *. b_o.mass *. mag;
13       b.vz <- b.vz -. dz *. b_o.mass *. mag;
14
15       b_o.vx <- b_o.vx +. dx *. b.mass *. mag;
16       b_o.vy <- b_o.vy +. dy *. b.mass *. mag;
```

# N-Body serial

```
1  let advance bodies n_bodies dt =
2    for i = 0 to n_bodies - 1 do
3      let b = bodies.(i) in
4      for j = i+1 to n_bodies - 1 do
5        let b_o = bodies.(j) in
6        let dx = b.x -. b_o.x  and dy = b.y -. b_o.y
7          and dz = b.z -. b_o.z in
8        let dist2 = dx *. dx +. dy *. dy +. dz *. dz in
9        let mag = dt /. (dist2 *. sqrt(dist2)) in
10
11       b.vx <- b.vx -. dx *. b_o.mass *. mag;
12       b.vy <- b.vy -. dy *. b_o.mass *. mag;
13       b.vz <- b.vz -. dz *. b_o.mass *. mag;
14
15       b_o.vx <- b_o.vx +. dx *. b.mass *. mag;
16       b_o.vy <- b_o.vy +. dy *. b.mass *. mag;
```

# N-Body serial

```
 4      for j = i+1 to n_bodies - 1 do
 5        let b_o = bodies.(j) in
 6        let dx = b.x -. b_o.x  and dy = b.y -. b_o.y
 7          and dz = b.z -. b_o.z in
 8        let dist2 = dx *. dx +. dy *. dy +. dz *. dz in
 9        let mag = dt /. (dist2 *. sqrt(dist2)) in
10
11        b.vx <- b.vx -. dx *. b_o.mass *. mag;
12        b.vy <- b.vy -. dy *. b_o.mass *. mag;
13        b.vz <- b.vz -. dz *. b_o.mass *. mag;
14
15        b_o.vx <- b_o.vx +. dx *. b.mass *. mag;
16        b_o.vy <- b_o.vy +. dy *. b.mass *. mag;
17        b_o.vz <- b_o.vz +. dz *. b.mass *. mag;
18      done
19    done
```

# N-Body serial

All experiments on an 2x Xeon E5-2695 v4

```
real    1m23.423s
user    1m23.422s
sys      0m0.000s
```

(256 iterations, 8192 bodies)

# How fast can we go?

Amdahl's law for parallel programs:

$$\frac{1}{(1-p) + \left(\frac{p}{s}\right)}$$

*p* = proportion of parallelisable code

*s* = degree of parallelism

# Linux Perf

Sampling profile of `nbody_serial.exe` run:

| Children | Self | Command | Shared Object | Symbol |
|---|---|---|---|---|
| + 99.95% | 0.00% | nbody_serial.ex | nbody_serial.exe | [.] caml_start_program |
| + 99.95% | 0.00% | nbody_serial.ex | nbody_serial.exe | [.] caml_program |
| + 99.95% | 0.00% | nbody_serial.ex | nbody_serial.exe | [.] camlDune__exe__Nbody_serial__entry |
| + 99.55% | 99.55% | nbody_serial.ex | nbody_serial.exe | [.] camlDune__exe__Nbody_serial__advance_90 |
| 0.34% | 0.34% | nbody_serial.ex | nbody_serial.exe | [.] camlDune__exe__Nbody_serial__energy_159 |
| 0.05% | 0.05% | nbody_serial.ex | nbody_serial.exe | [.] camlDune__exe__Nbody_serial__update_152 |
| 0.03% | 0.00% | nbody_serial.ex | nbody_serial.exe | [.] _start |

Perfect scalability ⇨ ~220x speedup

If p < 0.89, max speedup single digits!

# Iteration 1: Domain per body

*Don't do this*

```
1  let advance bodies n_bodies dt =
2    let ds =
3      Array.init n_bodies (fun i -> Domain.spawn (fun _ ->
4        let b = bodies.(i) in
5        for j = 0 to n_bodies - 1 do
6          let b_o = bodies.(j) in
7          if (i!=j) then begin
8            let dx = b.x -. b_o.x  and dy = b.y -. b_o.y
9              and dz = b.z -. b_o.z in
10           let dist2 = dx *. dx +. dy *. dy +. dz *. dz in
11           let mag = dt /. (dist2 *. sqrt(dist2)) in
12           b.vx <- b.vx -. dx *. b_o.mass *. mag;
13           b.vy <- b.vy -. dy *. b_o.mass *. mag;
14           b.vz <- b.vz -. dz *. b_o.mass *. mag;
15         end
16       done
```

# Iteration 1: Domain per body

## *Don't do this*

```
2    let ds =
3      Array.init n_bodies (fun i -> Domain.spawn (fun _ ->
4        let b = bodies.(i) in
5        for j = 0 to n_bodies - 1 do
6          let b_o = bodies.(j) in
7          if (i!=j) then begin
8            let dx = b.x -. b_o.x  and dy = b.y -. b_o.y
9              and dz = b.z -. b_o.z in
10           let dist2 = dx *. dx +. dy *. dy +. dz *. dz in
11           let mag = dt /. (dist2 *. sqrt(dist2)) in
12           b.vx <- b.vx -. dx *. b_o.mass *. mag;
13           b.vy <- b.vy -. dy *. b_o.mass *. mag;
14           b.vz <- b.vz -. dz *. b_o.mass *. mag;
15         end
16       done
17     )) in
```

# Iteration 1: Domain per body

*Don't do this*

```
3     Array.init n_bodies (fun i -> Domain.spawn (fun _ ->
4         let b = bodies.(i) in
5         for j = 0 to n_bodies - 1 do
6             let b_o = bodies.(j) in
7             if (i!=j) then begin
8                 let dx = b.x -. b_o.x  and dy = b.y -. b_o.y
9                     and dz = b.z -. b_o.z in
10                let dist2 = dx *. dx +. dy *. dy +. dz *. dz in
11                let mag = dt /. (dist2 *. sqrt(dist2)) in
12                b.vx <- b.vx -. dx *. b_o.mass *. mag;
13                b.vy <- b.vy -. dy *. b_o.mass *. mag;
14                b.vz <- b.vz -. dz *. b_o.mass *. mag;
15            end
16        done
17    )) in
18  Array.iter (Domain.join) ds
```

# Iteration 1: Domain per body

Oops.

```
real    8m10.965s
user    25m24.372s
sys     11m30.816s
```

Domains are heavyweight
Aim for same number as cores
Spawn/join infrequently

# Domainslib

- Task pool

  - Parallel
    - for / reduce / scan
  - Async/await

- Channels

  https://github.com/ocaml-multicore/domainslib/

# Iteration 2: Domainslib

```
1  let advance pool n_bodies n_domains bodies dt =
2    T.parallel_for pool
3      ~chunk_size:(n_bodies/n_domains)
4      ~start:0
5      ~finish:(n_bodies - 1)
6      ~body:(fun i ->
7        let b = bodies.(i) in
8        for j = 0 to n_bodies - 1 do
9          let b_o = bodies.(j) in
10         if (i!=j) then begin
11           let dx = b.x -. b_o.x  and dy = b.y -. b_o.y
12             and dz = b.z -. b_o.z in
13           let dist2 = dx *. dx +. dy *. dy +. dz *. dz in
14           let mag = dt /. (dist2 *. sqrt(dist2)) in
15           b.vx <- b.vx -. dx *. b_o.mass *. mag;
16           b.vy <- b.vy -. dy *. b_o.mass *. mag;
```

# Iteration 2: Domainslib

```
 5      ~finish:(n_bodies - 1)
 6      ~body:(fun i ->
 7        let b = bodies.(i) in
 8        for j = 0 to n_bodies - 1 do
 9          let b_o = bodies.(j) in
10          if (i!=j) then begin
11            let dx = b.x -. b_o.x  and dy = b.y -. b_o.y
12              and dz = b.z -. b_o.z in
13            let dist2 = dx *. dx +. dy *. dy +. dz *. dz in
14            let mag = dt /. (dist2 *. sqrt(dist2)) in
15            b.vx <- b.vx -. dx *. b_o.mass *. mag;
16            b.vy <- b.vy -. dy *. b_o.mass *. mag;
17            b.vz <- b.vz -. dz *. b_o.mass *. mag;
18          end
19      done
```

# Iteration 2: Domainslib

| Cores | Time | Vs Serial | Vs Self |
|-------|------|-----------|---------|
| 1 | 128.076s | 0.65x | 1x |
| 2 | 54.987s | 1.51x | 2.32x |
| 3 | 42.577s | 1.94x | 3.00x |
| 4 | 32.753s | 2.53x | 3.90x |
| 8 | 18.868s | 4.39x | 6.78x |
| 16 | 11.438s | 7.28x | 11.19x |
| 24 | 8.465s | 9.80x | 15.12x |

# Shared state pitfalls

Minimise writes to frequently read shared state
Avoid contended writes to shared state

Use `perf stat`, `perf record` and `perf c2c`

```
perf stat -e cycles,cycle_activity.cycles_no_execute,cycle_activi

    458,432,211,938        cycles
     74,023,598,176        cycle_activity.cycles_no_execute
     67,741,246,568        cycle_activity.stalls_mem_any
        481,904,795        cache-misses
         14,350,052        mem_load_uops_l3_miss_retired.remote_hitm

       8.726805280 seconds time elapsed
```

# Iteration 3: Avoiding contention

```
1  type planet_pos = { mutable x : float;  mutable y : float;  mu
2  type planet_vec = { mutable vx: float;  mutable vy: float;  mu
3
4  let advance pool n_domains n_bodies bodies_pos bodies_vec dt =
5    T.parallel_for pool
6      ~chunk_size:(n_bodies/n_domains)
7      ~start:0
8      ~finish:(n_bodies - 1)
9      ~body:(fun i ->
10       let bp = bodies_pos.(i) in
11       let bv = bodies_vec.(i) in
12       let vx, vy, vz = ref bv.vx, ref bv.vy, ref bv.vz in
13       for j = 0 to n_bodies - 1 do
14         let bp' = bodies_pos.(j) in
15         if (i!=j) then begin
16           let dx = bp.x - bp'.x  and dy = bp.y - bp'.y  and
```

# Iteration 3: Avoiding contention

```
 5    T.parallel_for pool
 6      ~chunk_size:(n_bodies/n_domains)
 7      ~start:0
 8      ~finish:(n_bodies - 1)
 9      ~body:(fun i ->
10        let bp = bodies_pos.(i) in
11        let bv = bodies_vec.(i) in
12        let vx, vy, vz = ref bv.vx, ref bv.vy, ref bv.vz in
13        for j = 0 to n_bodies - 1 do
14          let bp' = bodies_pos.(j) in
15          if (i!=j) then begin
16            let dx = bp.x -. bp'.x  and dy = bp.y -. bp'.y  and
17            let dist2 = dx *. dx +. dy *. dy +. dz *. dz in

18            let mag = dt /. (dist2 *. sqrt(dist2)) in
19            let mass = bp'.mass in
```
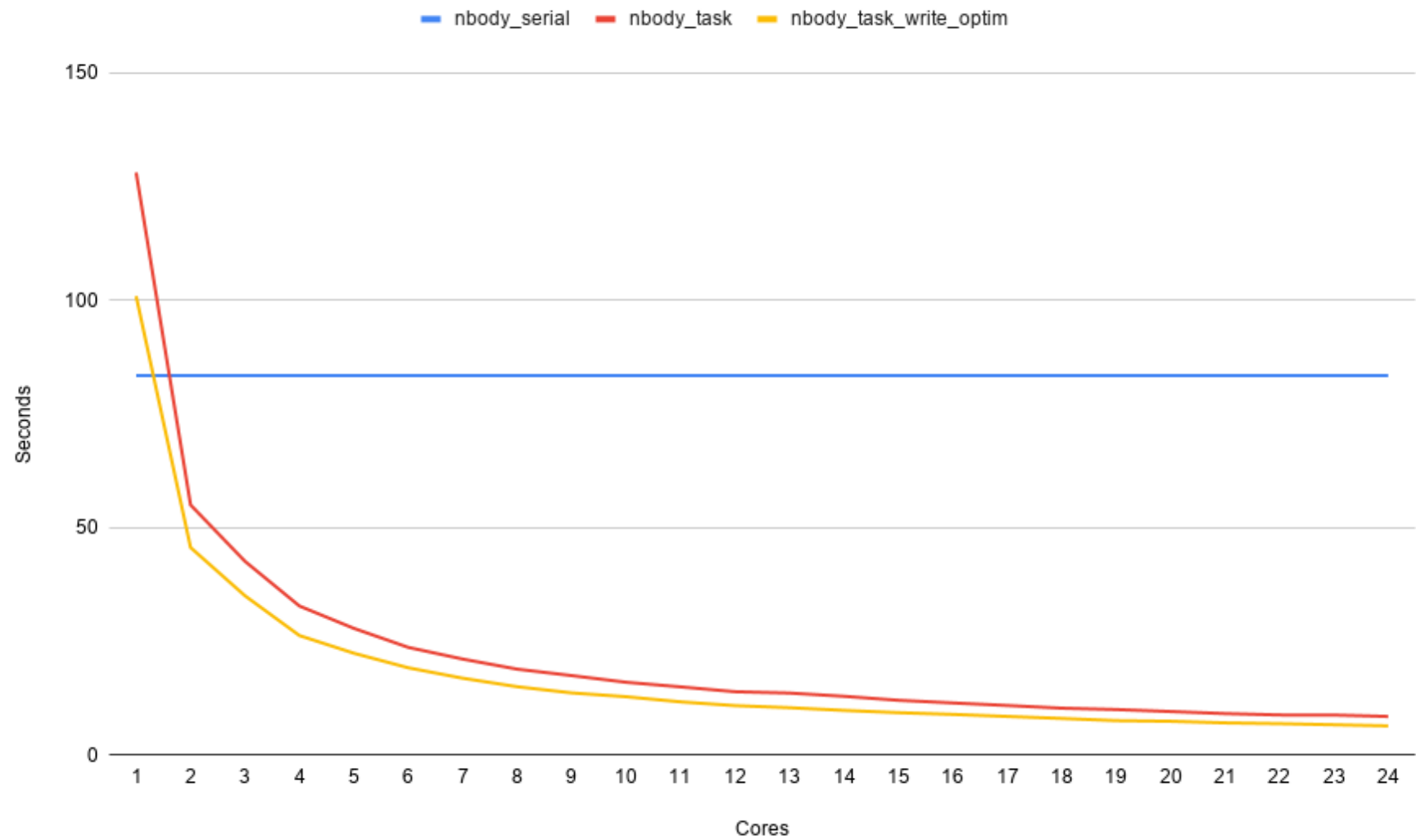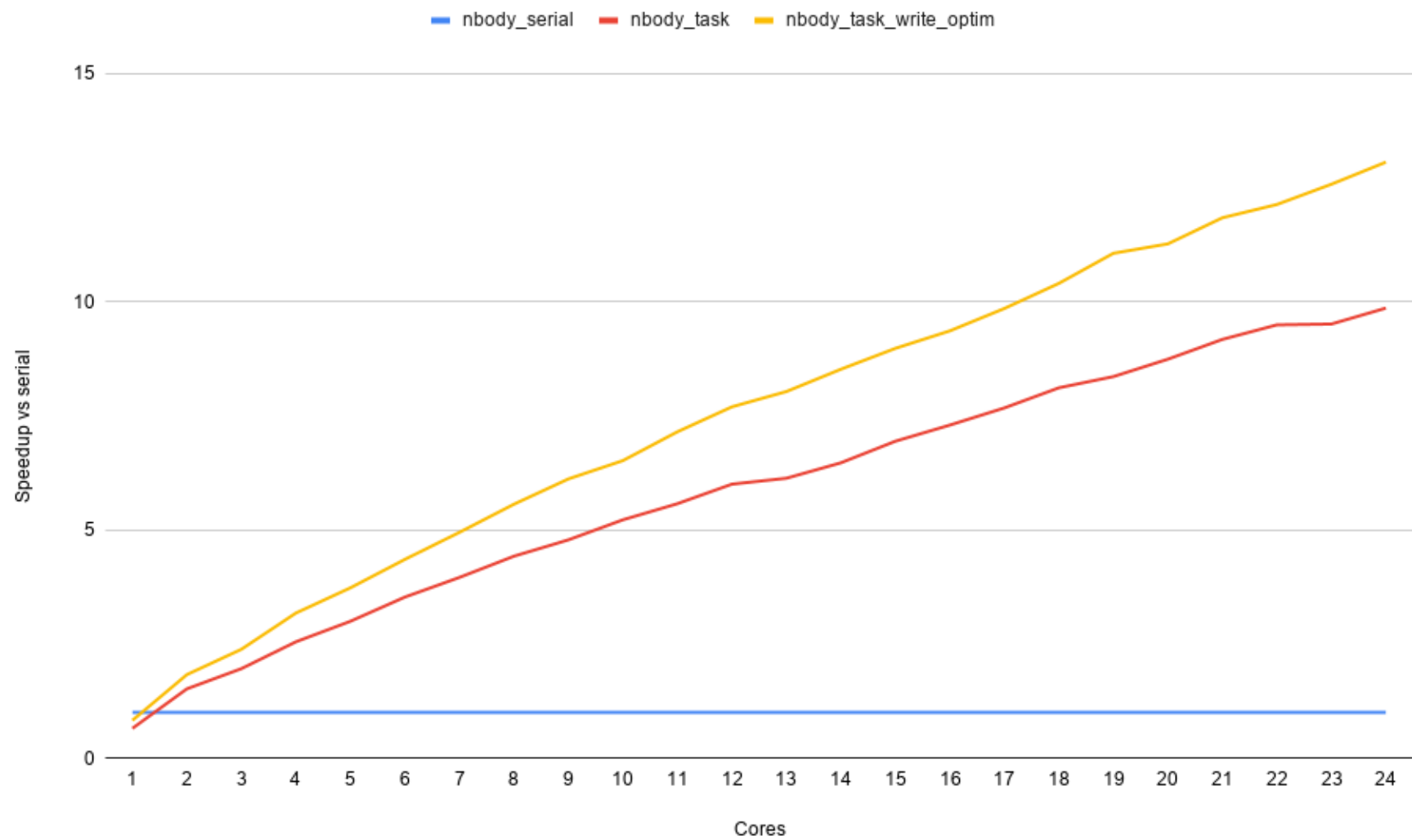
# Iteration 3: Avoiding contention

```
12      let vx, vy, vz = ref bv.vx, ref bv.vy, ref bv.vz in
13      for j = 0 to n_bodies - 1 do
14        let bp' = bodies_pos.(j) in
15        if (i!=j) then begin
16          let dx = bp.x -. bp'.x  and dy = bp.y -. bp'.y  and
17          let dist2 = dx *. dx +. dy *. dy +. dz *. dz in

18          let mag = dt /. (dist2 *. sqrt(dist2)) in
19          let mass = bp'.mass in
20          vx := !vx -. dx *. mass *. mag;
21          vy := !vy -. dy *. mass *. mag;
22          vz := !vz -. dz *. mass *. mag;
23        end
24      done;
25      bv.vx <- !vx;
26      bv.vy <- !vy;
27      bv.vz <- !vz;
```

# Iteration 3: Avoiding contention

```
perf stat -e cycles,cycle_activity.cycles_no_execute,cycle_activi

    341,067,775,833        cycles
     31,523,253,083        cycle_activity.cycles_no_execute
     29,569,589,458        cycle_activity.stalls_mem_any
         7,009,375        cache-misses
           956,792        mem_load_uops_l3_miss_retired.remote_hitm

      6.639116502 seconds time elapsed
```

~35% speedup at 24 cores

# Take aways

- Multicore is ready to use, upstreaming in progress
- Profile to understand serial performance
- Use Domainslib abstractions to add parallelism
- Share work as coarsely as you can
- Avoid writing to shared state as much as possible

Use the multicore github issue tracker if you find unexpected behaviour