

PC40 Hands-on: UPC

Adrien Burgun

Automne 2021

Abstract

For this hands-on assignment, I wanted to measure the performance of the different versions of the code and compare them to measure speedups. For this reason, my code diverges slightly from the template that was given to us, as I needed to insert a timing method and avoid refactoring the code.

You will find in this report listings of the different versions of the code, alongside experimental measurements and commentaries.

The source code itself, this report's source code and instructions on how to build and run the code yourself can be found on this project's git repository: <https://github.com/adri326/pc40-upc/>.

Contents

1	Simplified 1D Laplace solver	3
1.1	C implementation	3
1.2	Porting the C code to UPC	4
1.3	Optimizing the inner for loop	6
1.4	Optimizing the array blocking factor	6
1.4.1	Note on the update/copy loop	7
1.5	Detecting convergence	9
1.5.1	Using reduction operations	10
1.6	Conclusion (1D Laplace Solver)	10
2	2D Heat conduction	10
2.1	First C implementation	11
2.2	Porting the C code to UPC	13
2.3	Optimizing the array accesses	14
2.4	Array privatization	14
2.5	Dynamic problem size	17
3	Conclusion	19
3.1	Speed comparison	19
3.2	Personal remarks	21
	References	22

1 Simplified 1D Laplace solver

1.1 C implementation

The C implementation of the laplace solver has been slightly modified to time the main loop. It is single-threaded but the implementation allows for compiler SIMD optimizations. Settings common to every version (vector size, epsilon, max number of iterations) have been placed in a file called `settings.h`.

This code has been compiled with `gcc v4.9.0` and run on the `mesoshared` server, yielding the following results:

Options	Time (avg)	CI ($\sigma = 0.01$)
-00	450 μ s/iter	$\pm 10\mu$ s/iter
-03	140 μ s/iter	$\pm 10\mu$ s/iter

Table 1: Timing results for the C implementation of the 1D Laplace solver (Listing 1)

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <time.h>
4 #include "settings.h"
5
6 #define LEN DESIRED_LEN
7
8 void init();
9 void iteration();
10 void copy_array();
11
12 double x_new[LEN];
13 double x[LEN];
14 double b[LEN];
15
16 int main() {
17     init();
18
19     clock_t begin = clock();
20     for (size_t n = 0; n < ITERATIONS; n++) {
21         iteration();
22         if (n != ITERATIONS - 1) copy_array();
23     }
24     clock_t end = clock();
25
26     double diff = (double)(end - begin) / CLOCKS_PER_SEC / ITERATIONS;
27     double mean_squared = 0.0;
28     if (DISPLAY) {
29         printf("| b          | x          | x_new      | diff        |\n");
30         printf("| :----- | :----- | :----- | :----- |\n");
31         for (size_t i = 0; i < LEN; i++) {
32             printf("| %1.4lf | %1.4lf | %1.4lf | %1.4lf |\n", b[i], x[i], x_new[i],
33                 x_new[i] - x[i]);
34         }
35         for (size_t i = 0; i < LEN; i++) mean_squared += (x_new[i] - x[i]) * (x_new[i] - x[i]);
36         printf("\nIterations: %d\n", ITERATIONS);
37         printf("Mean squared difference: %1.4lf\n", mean_squared / LEN);
38         printf("Took %.3lfms/iter!\n", diff * 1000);
```

```

39 }
40
41 void init() {
42     srand(time(0));
43     for (size_t i = 0; i < LEN; i++) {
44         b[i] = (double)rand() / RAND_MAX;
45         x[i] = (double)rand() / RAND_MAX;
46         x_new[i] = 0;
47     }
48 }
49
50 void iteration() {
51     for (size_t i = 1; i < LEN - 1; i++) {
52         x_new[i] = 0.5 * (x[i-1] + x[i+1] + b[i]);
53     }
54     x_new[0] = x[0];
55     x_new[LEN - 1] = x[LEN - 1];
56 }
57
58 void copy_array() {
59     for (size_t i = 0; i < LEN; i++) {
60         x[i] = x_new[i];
61     }
62 }

```

Listing 1: C implementation of the 1D Laplace solver

1.2 Porting the C code to UPC

The base C implementation was designed to be able to quickly port it to UPC. A new `if` statement had to be inserted in the `for` loop within `iteration()`. Additionally, several lines had to only be executed by the thread 0, so additional conditionals were added when needed.

Finally, the `x`, `xnew` and `b` arrays were made shared and `upc_barrier` statements were added at the end of `init`, `iteration` and `copy_array`: this is to prevent threads from beginning processing the next iteration when the `x` and `xnew` arrays aren't ready, and to prevent the thread 0 from stopping early and printing the wrong timing.

This code was compiled with `upcc v2.22.0 + gcc v4.2.4` and run on the `mesoshared` server, yielding the following results:

Threads	Time (avg)	CI ($\sigma = 0.01$)
2	515.2 μ s/iter	$\pm 3.9\mu$ s/iter
3	536.2 μ s/iter	$\pm 2.4\mu$ s/iter
4	312.4 μ s/iter	$\pm 2.0\mu$ s/iter
8	204.0 μ s/iter	$\pm 2.0\mu$ s/iter
16	166.2 μ s/iter	$\pm 2.5\mu$ s/iter
32	150.4 μ s/iter	$\pm 3.2\mu$ s/iter

Table 2: Timing results for the first UPC implementation of the 1D Laplace solver (Listing 2)

When compiled and run with 3 threads, the code runs noticeably slower. For curiosity, I ran the code with 24 and 31 threads and obtained a similar slowdown:

Threads = 24, Time = 252.6 μ s/iter \pm 3.7 μ s/iter (expected $\approx 160 \mu$ s/iter)

Threads = 31, Time = $317.6\mu\text{s}/\text{iter} \pm 4.4\mu\text{s}/\text{iter}$ (expected $\approx 150\mu\text{s}/\text{iter}$)

```
1 #include <upc_relaxed.h>
2 #include <stdio.h>
3 #include <stdlib.h>
4 #include <time.h>
5 #include "settings.h"
6
7 #define LEN DESIRED_LEN
8
9 void init();
10 void iteration();
11 void copy_array();
12
13 shared double x_new[LEN];
14 shared double x[LEN];
15 shared double b[LEN];
16
17 int main() {
18     init();
19
20     clock_t begin = clock();
21     for (size_t n = 0; n < ITERATIONS; n++) {
22         iteration();
23         if (n != ITERATIONS - 1) copy_array();
24     }
25     clock_t end = clock();
26
27     if (MYTHREAD == 0) {
28         double diff = (double)(end - begin) / CLOCKS_PER_SEC / ITERATIONS;
29         double mean_squared = 0.0;
30         if (DISPLAY) {
31             printf("| b          | x          | x_new      | diff        |\n");
32             printf("| :----- | :----- | :----- | :----- |\n");
33             for (size_t i = 0; i < LEN; i++) {
34                 printf("| %1.4lf | %1.4lf | %1.4lf | % 1.4lf |\n", b[i], x[i],
35 x_new[i], x_new[i] - x[i]);
36             }
37             for (size_t i = 0; i < LEN; i++) mean_squared += (x_new[i] - x[i]) * (
38 x_new[i] - x[i]);
39             printf("\nIterations: %d\n", ITERATIONS);
40             printf("Mean squared difference: %1.4lf\n", mean_squared / LEN);
41             printf("Took %.3lfms/iter!\n", diff * 1000);
42         }
43     }
44
45     void init() {
46         if (MYTHREAD == 0) {
47             srand(time(0));
48             for (size_t i = 0; i < LEN; i++) {
49                 b[i] = (double)rand() / RAND_MAX;
50                 x[i] = (double)rand() / RAND_MAX;
51                 x_new[i] = 0;
52             }
53         }
54         upc_barrier;
55     }
56
57     void iteration() {
58         for (size_t i = 1; i < LEN - 1; i++) {
```

```

58     if (i % THREADS == MYTHREAD) {
59         x_new[i] = 0.5 * (x[i-1] + x[i+1] + b[i]);
60     }
61 }
62
63 if (MYTHREAD == 0) x_new[0] = x[0];
64 if (MYTHREAD == (LEN - 1) % THREADS) x_new[LEN - 1] = x[LEN - 1];
65
66 upc_barrier;
67 }
68
69 void copy_array() {
70     for (size_t i = 0; i < LEN; i++) {
71         if (i % THREADS == MYTHREAD) x[i] = x_new[i];
72     }
73     upc_barrier;
74 }

```

Listing 2: First UPC implementation of the 1D Laplace solver

1.3 Optimizing the inner for loop

The first step in optimizing the UPC implementation of the 1D Laplace equation solver is to replace the `for (...) if (...)` with a single, more efficient `for` loop.

To achieve this, only the `iteration()` function had to be changed. This transformation is shown in Figure 1

This change greatly increases the speed of the program:

Threads	Time (avg)	CI ($\sigma = 0.01$)
2	482.0 μ s/iter	$\pm 3.2\mu$ s/iter
4	265.2 μ s/iter	$\pm 2.3\mu$ s/iter
8	128.8 μ s/iter	$\pm 1.6\mu$ s/iter
16	74.6 μ s/iter	$\pm 3.1\mu$ s/iter
32	46.0 μ s/iter	$\pm 2.5\mu$ s/iter

Table 3: Timing results for the second UPC implementation of the 1D Laplace solver

1.4 Optimizing the array blocking factor

Our next optimization is to increase the blocking factor B for the x , x_{new} and b arrays. Each operation on the item of index j accesses $b[j]$, $x[j-1]$, $x[j+1]$ and $x_{new}[j]$.

With the default block factor of 1, accesses to both $x[j-1]$ and $x[j+1]$ will always be outside of the current thread's affinity. With a greater block factor, these outside accesses can be reduced to only $2/B$ accesses on average.

Two parts of the code had to be changed: the array declarations (Figure 2) and the loop in `iteration()` and `copy_array()` (Figure 3).

The effects of this change is dependent on B : similar to the results found in Table 2, the code runs fastest when $B = 2^n$. Whichever value is chosen, however, the program did not run faster

```
// ex3.upc:
void iteration() {
    for (size_t i = 1; i < LEN - 1; i++) {
        if (i % THREADS == MYTHREAD) {
            x_new[i] = 0.5 * (x[i-1] + x[i
+1] + b[i]);
        }
    }

    if (MYTHREAD == 0) x_new[0] = x[0];
    if (MYTHREAD == (LEN - 1) % THREADS)
        x_new[LEN - 1] = x[LEN - 1];

    upc_barrier;
}
```

```
// ex4.upc:
void iteration() {
    size_t i = MYTHREAD;
    if (MYTHREAD == 0) {
        i += THREADS;
        x_new[0] = x[0];
    }

    for (; i < LEN - 1; i += THREADS) {
        x_new[i] = 0.5 * (x[i-1] + x[i
+1] + b[i]);
    }

    if (MYTHREAD == THREADS - 1) x_new[
LEN - 1] = x[LEN - 1];

    upc_barrier;
}
```

(a) Previous `iteration()` function from Listing 2 (b) New `iteration()` function without the inner `if`

Figure 1: Modification in `iteration()` to remove the reduce the number of loop iterations.

than `ex4.upc` (Table 3), but rather ran much slower ([the source code of this version can be found here](#)):

Threads	Time (avg)	CI ($\sigma = 0.01$)
2	959.2 μ s/iter	$\pm 3.9\mu$ s/iter
4	863.9 μ s/iter	$\pm 16.1\mu$ s/iter
8	691.6 μ s/iter	$\pm 6.1\mu$ s/iter
16	581.5 μ s/iter	$\pm 4.1\mu$ s/iter
32	536.9 μ s/iter	$\pm 4.0\mu$ s/iter

Table 4: Timing results for the third UPC implementation of the 1D Laplace solver, $B = 32$

When ran with $B = 1$, $THREADS = 32$, the code takes 318 μ s/iter \pm 7.8 μ s/iter. I do not know what causes this slowdown, but it seems like `upc_forall` could contribute to it.

When limiting the modification of `ex4.upc` to only replace the `for` loop with a `upc_forall` loop, the example becomes 83% slower (206 μ s/iter \rightarrow 376 μ s/iter, 8 threads, AMD Ryzen 2700X, `upcc v2.28.2 + gcc v9.3.0`).

1.4.1 Note on the update/copy loop

So far, the program has been running many iterations of the "main" loop, which calls `iteration()` and `copy_array()`. A `upc_barrier`; call is necessary at the end of each of those two operations, as `copy_array()` modifies x and depends on the processed value of x_{new} , while `iteration()` modifies x_{new} and depends on the copied value in x .

```
// ex5.upc:
void iteration() {
    upc_forall (size_t i = 1; i < LEN - 1; i++; &x_new[i]) {
        x_new[i] = 0.5 * (x[i-1] + x[i+1] + b[i]);
    }
}
```

```
// ex4.upc:
shared double x_new[LEN];
shared double x[LEN];
shared double b[LEN];
```

(a) Previous array declarations

```
// settings.h:
#define BLOCKSIZE 32

// ex5.upc:
shared[BLOCKSIZE] double x_new[LEN];
shared[BLOCKSIZE] double x[LEN];
shared[BLOCKSIZE] double b[LEN];
```

(b) New array declarations with blocking factor

Figure 2: Modification of the array declarations to increase the blocking factor B .

```
if (MYTHREAD == 0) x_new[0] = x[0];
if (MYTHREAD == THREADS - 1) x_new[LEN - 1] = x[LEN - 1];

upc_barrier;
}

void copy_array() {
    upc_forall (size_t i = 0; i < LEN; i++; &x_new[i]) {
        x[i] = x_new[i];
    }
    upc_barrier;
}
```

Listing 3: New `iteration()` and `copy_array()` implementations

1.5 Detecting convergence

One of my last modifications is to measure $\delta_{max} = \max_{i \in \llbracket 0; n-1 \rrbracket} |x[i] - x_{new}[i]|$ (named `diffmax` in the source code) and to stop once $\delta_{max} \leq \varepsilon$.

Because of the shared nature of x and x_{new} , we need to compute a partial δ_{max}^t and compute $\delta_{max} = \max_{t \in [0; T]} (\delta_{max}^t)$, with:

$$\delta_{max}^t = \max_{i \in \llbracket 0; n-1 \rrbracket, \text{affinity}(x_{new}[i])=t} |x[i] - x_{new}[i]|$$

The new timings are:

Threads	Time (avg)	CI ($\sigma = 0.01$)
2	1400.4 μ s/iter	$\pm 6.0\mu$ s/iter
4	1092.8 μ s/iter	$\pm 17.0\mu$ s/iter
8	803.4 μ s/iter	$\pm 7.7\mu$ s/iter
16	647.8 μ s/iter	$\pm 18.9\mu$ s/iter
32	572.0 μ s/iter	$\pm 4.6\mu$ s/iter

Table 5: Timing results for the fourth UPC implementation of the 1D Laplace solver, $B = 32$, `MAX_ITER` = 1000

```

1 // ex6.upc (partial):
2 shared double diff[THREADS];
3 shared double diffmax;
4
5 void handle_diff(double d) {
6     if (d < 0.0) d = -d;
7     if (diff[MYTHREAD] < d) diff[MYTHREAD] = d;
8 }
9
10 int main() {
11     // ...
12     while (true) {
13         iteration();
14
15         upc_barrier;
16
17         if (MYTHREAD == 0) {
18             for (size_t n = 0; n < THREADS; n++) {
19                 if (diffmax < diff[n]) diffmax = diff[n];
20             }
21         }
22
23         upc_barrier;
24
25         if (diffmax <= EPSILON) break;
26         if (++iter > MAX_ITER) break;
27
28         copy_array(); // upc_barrier; at the end of copy_array();
29     }
30     // ...
31 }
32
33 void iteration() {
34     upc_forall (size_t i = 1; i < LEN - 1; i++; &x_new[i]) {

```

```

35     x_new[i] = 0.5 * (x[i-1] + x[i+1] + b[i]);
36     handle_diff(x_new[i] - x[i]);
37 }
38 }

```

Listing 4: Excerpt from ex6.upc: new `handle_diff` function and calculation of `diffmax`

1.5.1 Using reduction operations

We can further optimize this code by using reduction operations from `upc_collective.h`. The optimized variant is available at [laplace/optimized.upc](#).

By removing the calls to `upc_forall`, I could reduce the time taken by a further 35% ($\pm 5\%$ pt).

Threads	Time (avg)	CI ($\sigma = 0.01$)
2	1216.0 μ s/iter	$\pm 27.1\mu$ s/iter
4	778.0 μ s/iter	$\pm 15.1\mu$ s/iter
8	415.8 μ s/iter	$\pm 10\mu$ s/iter
16	225.8 μ s/iter	$\pm 4.8\mu$ s/iter
32	130.0 μ s/iter	$\pm 3.1\mu$ s/iter
64	98.2 μ s/iter	$\pm 2.7\mu$ s/iter

Table 6: Timing results for the fully optimized UPC implementation of the 1D Laplace solver, $B = 32$, `MAX_ITER` = 1000

1.6 Conclusion (1D Laplace Solver)

The first implementation in UPC of the 1D Laplace solver algorithm and its subsequent first optimization (flattening the `for (...) if (...)`) gave very promising results for the improvement of the speed of the program, bringing the timings down from 140 μ s/iteration to only 46 μ s/iteration.

Unfortunately, the next optimization attempt, which introduces a blocking factor to the primary arrays and used `upc_forall`, brought the speeds down and made the UPC implementation much slower than the C implementation.

The addition of δ_{max} , as a way to stop the algorithm once a sufficiently accurate solution is found, inevitably added another overhead to the program. With my best efforts, I could only bring its performance to somewhere between those of the first implementation and its first optimization.

Nonetheless, the resulting program runs slightly faster than the original C implementation, given enough threads).

2 2D Heat conduction

For this second algorithm, I will try to improve on the knowledge gathered from the first algorithm and implement a simple, 2D heat conduction simulation. As with Section 1, I begin with a simple C implementation of the algorithm, which will work as my performance base value. The provided

template came bundled with a performance measurement method, but I swapped it out with `clock()` in the subsequent UPC implementations for consistency with the last section's code.

2.1 First C implementation

Following are the performance results for the first C implementation, ran on `mesoshared` with `gcc v4.9.0` and on an AMD Ryzen 2700X with `clang v12.0.1`:

Options	Time (avg)	CI ($\sigma = 0.01$)
-00	89.1 μ s/iter	$\pm 0.3\mu$ s/iter
-03	15.0 μ s/iter	$\pm 0.1\mu$ s/iter

Table 7: Timing results for `heat_c.c` (mesoshared, Listing 5)

Options	Time (avg)	CI ($\sigma = 0.01$)
-00	80.9 μ s/iter	$\pm 2.6\mu$ s/iter
-03	17.9 μ s/iter	$\pm 2.2\mu$ s/iter

Table 8: Timing results for `heat_c.c` (Ryzen 2700X, Listing 5)

```

1 #include <stdio.h>
2 #include <math.h>
3 #include <sys/time.h>
4 #include <stdbool.h>
5 #include "settings.h"
6
7 double grid[N+2][N+2], new_grid[N+2][N+2];
8
9 void initialize(void) {
10     int j;
11
12     /* Heat one side of the solid */
13     for (j=1; j<N+1; j++) {
14         grid[0][j] = 1.0;
15         new_grid[0][j] = 1.0;
16     }
17 }
18
19 int main(void) {
20     struct timeval ts_start, ts_end;
21     double dTmax, dT, time;
22     int i, j, k, l;
23     bool finished = false;
24     double T;
25     int n_iter = 0;
26
27     initialize();
28
29     /* Set the precision wanted */
30     finished = 0;
31
32     /* and start the timed section */
33     gettimeofday(&ts_start, NULL);
34
35     do {
36         dTmax = 0.0;
37         for (i=1; i<N+1; i++) {

```

```

38     for (j=1; j<N+1; j++) {
39         T = 0.25 *
40             (grid[i+1][j] + grid[i-1][j] +
41              grid[i][j-1] + grid[i][j+1]); /* stencil */
42         dT = T - grid[i][j]; /* local variation */
43         new_grid[i][j] = T;
44         if (dTmax < fabs(dT)) dTmax = fabs(dT); /* max variation in this
iteration */
45     }
46 }
47 if (dTmax < EPSILON) { // is the precision reached good enough ?
48     finished = true;
49 } else { // It isn't: preparing for a new iteration
50     for (k=0; k<N+2; k++) {
51         for (l=0; l<N+2; l++) {
52             grid[k][l] = new_grid[k][l];
53         }
54     }
55 }
56 n_iter++;
57 } while (!finished);
58
59 gettimeofday(&ts_end, NULL); /* end the timed section */
60
61 /* compute the execution time */
62 time = ts_end.tv_sec + (ts_end.tv_usec / 1000000.0);
63 time -= ts_start.tv_sec + (ts_start.tv_usec / 1000000.0);
64
65 printf("%d iterations in %.3lf sec\n", n_iter, time);
66 printf("Took %.3lf ms/iter\n", time * 1000.0 / n_iter);
67
68 return 0;
69 }

```

Listing 5: C implementation of the 2D Heat simulation

2.2 Porting the C code to UPC

Porting the template C version to UPC is straightforward. The performance, however, takes a hit, as the default work sharing causes a lot of remote accesses to memory, despite the blocking factor that was put in place. We obtain the following measurements on `mesoshared`:

Threads	Time (avg)	CI ($\sigma = 0.01$)
2	259.9 μ s/iter	$\pm 1.9\mu$ s/iter
4	193.6 μ s/iter	$\pm 1.5\mu$ s/iter
8	164.1 μ s/iter	$\pm 1.2\mu$ s/iter
16	150.6 μ s/iter	$\pm 1.4\mu$ s/iter
32	151.7 μ s/iter	$\pm 1.7\mu$ s/iter

Table 9: Timing results for `heat_1.upc`

```
// heat_c.c:
for (i=1; i<N+1; i++) {
  for (j=1; j<N+1; j++) {
    T = 0.25 *
      (grid[i+1][j] + grid[i-1][j] +
       grid[i][j-1] + grid[i][j+1]); /*
      stencil */
    dT = T - grid[i][j]; /* local
    variation */
    new_grid[i][j] = T;
    if (dTmax < fabs(dT)) dTmax = fabs(
    dT); /* max variation in this
    iteration */
  }
}
```

(a) Main loop in `heat_c.c`

```
// heat_1.upc:
for (size_t i = 1; i <= N; i++) {
  upc_forall (size_t j = 1; j <= N; j++;
    &grid[i][j]) {
    double T = 0.25 * (grid[i+1][j] +
      grid[i-1][j] + grid[i][j-1] + grid[i][j+1]);
    double dT = fabs(T - grid[i][j]);
    new_grid[i][j] = T;
    if (dTmax < dT) dTmax = dT;
  }
}
```

(b) Main loop in `heat_1.upc`

Figure 3: Porting the main loop from C to UPC

2.3 Optimizing the array accesses

One simple optimization is to avoid copying the destination array (`new_grid`) into the source array (`grid`).

To do this, I store two shared pointers that will point on either array, and I swap them at the end of each iteration, simulating a copy of `grid` into `new_grid`, with a $\mathcal{O}(1)$ time complexity instead of $\mathcal{O}(n)$. Doing so effectively halves the amount of synchronization that needs to happen for each loop, which can be observed in the timing results:

Threads	Time (avg)	CI ($\sigma = 0.01$)
2	134.0 μ s/iter	$\pm 1.4\mu$ s/iter
4	103.9 μ s/iter	$\pm 4.1\mu$ s/iter
8	87.4 μ s/iter	$\pm 1.0\mu$ s/iter
16	81.4 μ s/iter	$\pm 1.2\mu$ s/iter
32	82.9 μ s/iter	$\pm 1.5\mu$ s/iter

Table 10: Timing results for `heat_3.upc`

This optimization was also done on the C version for comparison, and I obtain the following speeds:

Options	Time (avg)	CI ($\sigma = 0.01$)
-O0	64.0 μ s/iter	$\pm 0.5\mu$ s/iter
-O3	11.3 μ s/iter	$\pm 0.4\mu$ s/iter

Table 11: Timing results for `heat_c_ptr.c`

```
// heat_3.upc:
shared[BLOCKSIZE] double (*ptr)[N+2] = grid;
shared[BLOCKSIZE] double (*new_ptr)[N+2] = new_grid;

// ...

if (n_iter % 2 == 0) {
    ptr = grid;
    new_ptr = new_grid;
} else {
    ptr = new_grid;
    new_ptr = grid;
}

// ...

double T = 0.25 * (ptr[i+1][j] + ptr[i-1][j] + ptr[i][j-1] + ptr[i][j+1]);
double dT = fabs(T - ptr[i][j]);
new_ptr[i][j] = T;
```

Listing 6: Excerpt from `heat_3.upc`, where pointer-swapping was implemented

2.4 Array privatization

As observed in the later examples of the 1D Laplace equation solver (Table 4), `upc_forall` has a significant performance cost over a regular `for` loop. We also know that local shared

memory accesses have a considerable cost when the shared pointer math is done without dedicated hardware support [1].

To make use of less shared accesses, the program can copy a chunk of the *grid* array into private memory using `upc_memget` (into `ptr_priv` and `new_ptr_priv`). We specifically only retrieve the memory that has the affinity to the current thread, so that the code can write to it and broadcast it in one go. The memory is copied also duplicated to make use of pointer-swapping as in Section 2.3.

We then operate on the private array, only using remote accesses when necessary, and store the results in the new private array. Once the work is finished, the program put the new private array into `new_grid` with `upc_mempout` and synchronize δ_{max} .

Threads	Time (avg)	CI ($\sigma = 0.01$)
2	13.5 μ s/iter	$\pm 1.4\mu$ s/iter
4	12.6 μ s/iter	$\pm 1.2\mu$ s/iter
8	13.3 μ s/iter	$\pm 1.0\mu$ s/iter
16	15.0 μ s/iter	$\pm 1.2\mu$ s/iter
32	21.5 μ s/iter	$\pm 1.1\mu$ s/iter

Table 12: Timing results for `heat_4.upc` (Listing 7)

```

1 #include <upc_relaxed.h>
2 #include <bupc_collectivev.h>
3 #include <stdio.h>
4 #include <math.h>
5 #include <time.h>
6 #include <stdbool.h>
7 #include "settings.h"
8
9 #if ((N+2) % THREADS) != 0
10 #error N+2 must be divisible by THREADS
11 #endif
12
13 // Change blocksize to (N+2)^2 / THREADS?
14 #define BLOCKSIZE ((N+2) * (N+2) / THREADS)
15 // #define BLOCKSIZE ((N+2) * THREADS)
16 #define LOCALWIDTH ((N+2) / THREADS)
17 #define LOCALSIZE (LOCALWIDTH * sizeof(double) * (N+2))
18
19 shared[BLOCKSIZE] double grid[N+2][N+2];
20 shared[BLOCKSIZE] double new_grid[N+2][N+2];
21 shared double dTmax[THREADS];
22
23 void init() {
24     for (size_t j = 1; j < N+1; j++) {
25         grid[0][j] = 1.0;
26         new_grid[0][j] = 1.0;
27     }
28 }
29
30 int main() {
31     shared[BLOCKSIZE] double (*ptr)[N+2] = grid;
32     shared[BLOCKSIZE] double (*new_ptr)[N+2] = new_grid;
33     double (*ptr_priv)[N+2] = malloc(LOCALSIZE);
34     double (*new_ptr_priv)[N+2] = malloc(LOCALSIZE);
35
36     if (MYTHREAD == 0) {
37         init();

```

```

38 }
39
40 upc_barrier;
41
42 upc_memget(ptr_priv, &grid[LOCALWIDTH * MYTHREAD], LOCALSIZE);
43 upc_memget(new_ptr_priv, &new_grid[LOCALWIDTH * MYTHREAD], LOCALSIZE);
44
45 bool finished = false;
46 int n_iter = 0;
47
48 clock_t begin = clock();
49 do {
50     double dTmax = 0.0;
51
52     size_t o = LOCALWIDTH * MYTHREAD;
53     size_t i = 0;
54
55     // Local block start
56     if (i + o > 0) {
57         for (size_t j = 1; j <= N; j++) {
58             double T = 0.25 * (ptr_priv[i+1][j] + ptr[o+i-1][j] + ptr_priv[i][j-1] + ptr_priv[i][j+1]);
59             // printf("%zu+%zu %zu: %lf\n", o, i, j, T);
60             double dT = fabs(T - ptr_priv[i][j]);
61             new_ptr_priv[i][j] = T;
62             if (dTmax < dT) dTmax = dT;
63         }
64     }
65
66     // Local block middle
67     for (i += 1; i < LOCALWIDTH - 1; i++) {
68         for (size_t j = 1; j <= N; j++) {
69             double T = 0.25 * (ptr_priv[i+1][j] + ptr_priv[i-1][j] + ptr_priv[i][j-1] + ptr_priv[i][j+1]);
70             // printf("%zu+%zu %zu: %lf\n", o, i, j, T);
71             double dT = fabs(T - ptr_priv[i][j]);
72             new_ptr_priv[i][j] = T;
73             if (dTmax < dT) dTmax = dT;
74         }
75     }
76
77     // Local block end
78     if (i + o < N + 1) {
79         for (size_t j = 1; j <= N; j++) {
80             double T = 0.25 * (ptr[o+i+1][j] + ptr_priv[i-1][j] + ptr_priv[i][j-1] + ptr_priv[i][j+1]);
81             // printf("%zu+%zu %zu: %lf\n", o, i, j, T);
82             double dT = fabs(T - ptr_priv[i][j]);
83             new_ptr_priv[i][j] = T;
84             if (dTmax < dT) dTmax = dT;
85         }
86     }
87
88     // printf("%d: %lf\n", MYTHREAD, dTmax);
89
90     // upc_barrier;
91
92     // Update ptr_priv and new_ptr_priv
93     upc_mempu(&new_ptr[LOCALWIDTH * MYTHREAD], new_ptr_priv, LOCALSIZE);
94
95     // printf("%lf %lf\n", new_ptr_priv[1][1], new_ptr[o+1][1]);

```



```

96 // Implicit barrier here:
97
98 double dTmax_g = bupc_allv_reduce_all(double, dTmax, UPC_MAX);
99
100 if (dTmax_g < EPSILON) {
101     finished = true;
102 } else {
103     // Swap ptr and new_ptr
104     shared[BLOCKSIZE] double (*ptr_tmp)[N+2] = ptr;
105     ptr = new_ptr;
106     new_ptr = ptr_tmp;
107
108     double (*ptr_tmp_priv)[N+2] = ptr_priv;
109     ptr_priv = new_ptr_priv;
110     new_ptr_priv = ptr_tmp_priv;
111 }
112 n_iter++;
113 // upc_barrier;
114 } while (!finished);
115 clock_t end = clock();
116
117 if (MYTHREAD == 0) {
118     double seconds = (double)(end - begin) / CLOCKS_PER_SEC;
119     printf("%d iterations in %.3lf sec\n", n_iter, seconds);
120     printf("Took %.3lf ms/iter\n", seconds * 1000.0 / n_iter);
121 }
122 }

```

Listing 7: Optimized UPC implementation of the 2D Heat algorithm (heat_4.upc)

2.5 Dynamic problem size

This last version of the code is not an optimization: we had to adapt it so that we could specify the number of threads and the size of the grid at runtime, rather than at compile time. This requires changing a big portion of the code, as it was previously using pre-defined constants for these two parameters.

Because speed for this part is not a worry, I based the code from `heat_3.upc`, as the previous optimization step made the code harder to work with.

The declaration of `grid` and `new_grid` have been deleted and are done implicitly through the compiler-supplied function `upc_all_alloc`. The matrix notation, while useful, had to be gotten rid of, as it was dependent on N , which is now dynamic. This meant re-writing all of the matrix accesses to now multiply the y coordinate by $(n + 2)$ and to add its result to the x coordinate.

```

1 #include <upc_relaxed.h>
2 #include <upc_collective.h>
3 #include <stdio.h>
4 #include <math.h>
5 #include <time.h>
6 #include <stdbool.h>
7 #include <stdlib.h>
8 #include "settings.h"
9
10 shared double dTmax_g;
11
12 void init(shared[] double* grid, shared[] double* new_grid, size_t n) {
13     for (size_t i = 0; i < n+2; i++) {

```

```

14     for (size_t j = 0; j < n+2; j++) {
15         grid[i * (n+2) + j] = 0.0;
16         new_grid[i * (n+2) + j] = 0.0;
17     }
18 }
19
20 for (size_t j = 1; j < n+1; j++) {
21     grid[j] = 1.0;
22     new_grid[j] = 1.0;
23 }
24 }
25
26 int main(int argc, char* argv[]) {
27     if (argc != 2) {
28         if (MYTHREAD == 0) fprintf(stderr, "Error: expected two arguments, got %d\
n", argc);
29         exit(1);
30     }
31
32     size_t n = atoi(argv[1]);
33
34     shared[] double* dTmax = upc_all_alloc(1, THREADS);
35     shared[] double* ptr = upc_all_alloc((n+2) * (n+2) / THREADS, (n+2) * (n+2));
36     shared[] double* new_ptr = upc_all_alloc((n+2) * (n+2) / THREADS, (n+2) * (n
+2));
37
38     // Handy way to access ptr from now on
39     #define ptr_get(x, y) ptr[(x) * (n+2) + (y)]
40
41     if (MYTHREAD == 0) {
42         init(ptr, new_ptr, n);
43     }
44
45     upc_barrier;
46     bool finished = false;
47     int n_iter = 0;
48
49     clock_t begin = clock();
50     do {
51         dTmax[MYTHREAD] = 0.0;
52         size_t i = (n+2) * MYTHREAD / THREADS;
53         if (i == 0) i = 1;
54         size_t imax = (n+2) * (MYTHREAD + 1) / THREADS;
55         if (imax > n+1) imax = n+1;
56
57         for (; i < imax; i++) {
58             for (size_t j = 1; j <= n; j++) {
59                 double T = 0.25 * (ptr_get(i+1, j) + ptr_get(i-1, j) + ptr_get(i,
j-1) + ptr_get(i, j+1));
60                 double dT = fabs(T - ptr_get(i, j));
61                 new_ptr[i * (n+2) + j] = T;
62                 if (dTmax[MYTHREAD] < dT) dTmax[MYTHREAD] = dT;
63             }
64         }
65
66         upc_all_reduceD(&dTmax_g, dTmax, UPC_MAX, THREADS, 1, NULL, UPC_IN_ALLSYNC
| UPC_OUT_ALLSYNC);
67
68         if (dTmax_g < EPSILON) {
69             finished = true;
70         } else {

```

```

71     shared[] double* tmp = ptr;
72     ptr = new_ptr;
73     new_ptr = tmp;
74 }
75 n_iter++;
76 // upc_barrier;
77 } while (!finished);
78 clock_t end = clock();
79
80 if (MYTHREAD == 0) {
81     double seconds = (double)(end - begin) / CLOCKS_PER_SEC;
82     printf("%d iterations in %.3lf sec\n", n_iter, seconds);
83     printf("Took %.3lf ms/iter\n", seconds * 1000.0 / n_iter);
84 }
85 }

```

Listing 8: Dynamic UPC implementation of the 2D Heat algorithm (`heat_5.upc`)

3 Conclusion

For both the Simplified 1D Laplace Equation Solver (Section 1) and the 2D Heat Algorithm (Section 2), I was able to take a simple C implementation and port it to UPC with minimal effort. These two problems are easily parallelizable, as all of the work within an iteration can be done independently (or vertically).

We then went along the process of refining the UPC implementation and optimizing the various accesses to the arrays to reduce the overhead of parallelization induced by UPC. Unfortunately, the overhead encountered here is mainly caused by the shared pointer operations (as highlighted by the results in Section 1.4 and Section 2.4). Minimizing this overhead required heavy modification to the code.

With all of my efforts, the UPC implementation only came shortly before the C implementation when benchmarked and required at least 8 threads to overcome the language's overhead. This difficulty to beat the speed of C is explained by the nature of the problems at hand: each iteration consists of a small amount of work over a large set of data, that needs to be synchronized before the next iteration. Even with a perfect parallelization of the algorithm (with no synchronization overhead), the implementation would still be bottlenecked by the data transfer speeds and the slowest thread contributing to the calculation.

3.1 Speed comparison

Following are two comparative tables of all of the measurements of this report.

The Simplified 1D Laplace Equation Solver scored great until I introduced δ_{max} , which required some costly synchronization.

The 2D Heat Algorithm did worse, only achieving speeds close to those of C. The synchronization costs made it unfeasible to use a greater amount of threads, as can be seen with the last few rows.

Implementation	Options	Time (avg)	CI ($\sigma = 0.01$)
<code>laplace/ex2.c</code>	-00	450 μs /iter	$\pm 10\mu s$ /iter
	-03	140 μs /iter	$\pm 10\mu s$ /iter
<code>laplace/ex3.upc</code> (w/o δ_{max})	$T = 2$	515.2 μs /iter	$\pm 3.9\mu s$ /iter
	$T = 4$	312.4 μs /iter	$\pm 2.0\mu s$ /iter
	$T = 8$	204.0 μs /iter	$\pm 2.0\mu s$ /iter
	$T = 16$	166.2 μs /iter	$\pm 2.5\mu s$ /iter
	$T = 32$	150.4 μs /iter	$\pm 3.2\mu s$ /iter
<code>laplace/ex4.upc</code> (w/o δ_{max})	$T = 2$	482.0 μs /iter	$\pm 3.2\mu s$ /iter
	$T = 4$	265.2 μs /iter	$\pm 2.3\mu s$ /iter
	$T = 8$	128.8 μs /iter	$\pm 1.6\mu s$ /iter
	$T = 16$	74.6 μs /iter	$\pm 3.1\mu s$ /iter
	$T = 32$	46.0 μs /iter	$\pm 2.5\mu s$ /iter
<code>laplace/ex5.upc</code> (with δ_{max})	$T = 2$	959.2 μs /iter	$\pm 3.9\mu s$ /iter
	$T = 4$	863.9 μs /iter	$\pm 16.1\mu s$ /iter
	$T = 8$	691.6 μs /iter	$\pm 6.1\mu s$ /iter
	$T = 16$	581.5 μs /iter	$\pm 4.1\mu s$ /iter
	$T = 32$	536.9 μs /iter	$\pm 4.0\mu s$ /iter
<code>laplace/ex6.upc</code> (with δ_{max})	$T = 2$	1400.4 μs /iter	$\pm 6.0\mu s$ /iter
	$T = 4$	1092.8 μs /iter	$\pm 17.0\mu s$ /iter
	$T = 8$	803.4 μs /iter	$\pm 7.7\mu s$ /iter
	$T = 16$	647.8 μs /iter	$\pm 18.9\mu s$ /iter
	$T = 32$	572.0 μs /iter	$\pm 4.6\mu s$ /iter
<code>laplace/optimized.upc</code> (with δ_{max})	$T = 2$	1216.0 μs /iter	$\pm 27.1\mu s$ /iter
	$T = 4$	778.0 μs /iter	$\pm 15.1\mu s$ /iter
	$T = 8$	415.8 μs /iter	$\pm 10\mu s$ /iter
	$T = 16$	225.8 μs /iter	$\pm 4.8\mu s$ /iter
	$T = 32$	130.0 μs /iter	$\pm 3.1\mu s$ /iter
	$T = 64$	98.2 μs /iter	$\pm 2.7\mu s$ /iter

Table 13: Comparative table of the benchmarks of the Simplified Laplace 1D Solver

Implementation	Options	Time (avg)	CI ($\sigma = 0.01$)
heat_c.c (w/o ptr-swap)	-00	80.9 μs /iter	$\pm 2.6\mu s$ /iter
	-03	17.9 μs /iter	$\pm 2.2\mu s$ /iter
heat_c_ptr.c (with ptr-swap)	-00	64.0 μs /iter	$\pm 0.5\mu s$ /iter
	-03	11.3 μs /iter	$\pm 0.4\mu s$ /iter
heat_1.upc (w/o ptr-swap)	$T = 2$	259.9 μs /iter	$\pm 1.9\mu s$ /iter
	$T = 4$	193.6 μs /iter	$\pm 1.5\mu s$ /iter
	$T = 8$	164.1 μs /iter	$\pm 1.2\mu s$ /iter
	$T = 16$	150.6 μs /iter	$\pm 1.4\mu s$ /iter
	$T = 32$	151.7 μs /iter	$\pm 1.7\mu s$ /iter
heat_3.upc (with ptr-swap)	$T = 2$	134.0 μs /iter	$\pm 1.4\mu s$ /iter
	$T = 4$	103.9 μs /iter	$\pm 4.1\mu s$ /iter
	$T = 8$	87.4 μs /iter	$\pm 1.0\mu s$ /iter
	$T = 16$	81.4 μs /iter	$\pm 1.2\mu s$ /iter
	$T = 32$	82.9 μs /iter	$\pm 1.5\mu s$ /iter
heat_4.upc (with ptr-swap)	$T = 2$	13.5 μs /iter	$\pm 1.4\mu s$ /iter
	$T = 4$	12.6 μs /iter	$\pm 1.2\mu s$ /iter
	$T = 8$	13.3 μs /iter	$\pm 1.0\mu s$ /iter
	$T = 16$	15.0 μs /iter	$\pm 1.2\mu s$ /iter
	$T = 32$	21.5 μs /iter	$\pm 1.1\mu s$ /iter

Table 14: Comparative table of the benchmarks of the 2D Heat Algorithm

3.2 Personal remarks

Overall, I enjoyed working with UPC and optimizing the parallelized code. However, I feel like UPC is lacking behind in terms of what the compiler can optimize for the developer. For instance, it could unwrap a `upc_forall` loop into two `for` loops when the increment is a divisor of the blocking factor: this can easily be done with macros already and can only make the code faster.

The hardest part of working on this project was finding documentation for the language itself. The documentation is scattered across internet, with some tutorials explaining how basic functions work but missing the more complicated functions, some documentation hosted on code mirroring websites and the rest burried inside of the compiler's source code.

Searching for example `"upc_all_reduceD"` online brings up [an 11-year old mail archive](#), [a benchmark of UPC](#) and quite further down the summary of the UPC Collective specification. The only way to truly know what this function does is to [read its source code](#).

This contrasts with other, smaller languages that I have encountered until now, whose documentation is usually centered in one place. With the similarly niche language "Pony", searching `"ponylang hashmap"` online yields [a rendered version of the documentation of the language's standard library](#) as the first result.

The errors outputted by the language were also of little help, as they did not show any additional information. The `-verbose` flag only displayed information unrelated to the error.

Nonetheless, I am glad that I got to try out this language and work with parallelized code.

```
1 $ upcc 2d_heat/heat_5.upc -o build/heat_5
2 upcc: error during UPC-to-C translation (sgiupe stage):
3 2d_heat/heat_5.upc: In function 'main':
4 2d_heat/heat_5.upc:66: incompatible type 'for' argument 1 of 'bupc_all_reduceD'
5
```

Figure 4: *"I guess the compiler wants me to learn the function signature by myself."*

References

- [1] Olivier Serres. *Hardware Support for Productive Partitioned Global Address Space (PGAS) Programming*. PhD thesis, George Washington University, 2015.