



UPPSALA
UNIVERSITET

Individual Project

High Performance Programming

Author:

Moa Li Lekander

Uppsala

March 8, 2021

Contents

1	Introduction	1
2	Problem Description	1
3	Solution Method	2
3.1	Original Implementation	2
3.2	Serial Optimizations	2
3.3	Parallelization Using Pthreads	3
4	Experiments	4
4.1	Serial Optimizations	4
4.2	Parallelization Using Pthreads	6
5	Conclusions	7

1 Introduction

Code optimizations is a matter of great importance in many areas of science today. Although the performance and computing capacity of the modern computers is rapidly improving, there still exist many scenarios where it is highly auspicious to minimize the run time of programs to the greatest extent. By learning techniques of optimizing and parallelizing code, heavy simulations and computations that originally could takes days or even weeks to run, can instead be written so that its run time can be reduced significantly.

The aim of this project is to implement an efficient code using parallelization, for an application or algorithm of own choosing. In this report, a Game of Life simulation was implemented, and the code was firstly being serial optimized and thereafter parallelized using Pthreads. Also, the time complexity of the final algorithm was determined and the effect of parallellization with different numbers of threads was analysed.

2 Problem Description

A Game of Life simulation visualizes the evolution of a specific initial state in a 2D binary world. Hence, the simulation is played out on an quadratic grid of cells, where each cell is defined as either dead or alive. Depending on the amount of alive cells among the eight neighbours of each cell (see Figure 1), the next state of each cell is determined by a set of rules, that are as follows:

- If a cell is *alive* and has two or three neighbours it stays alive in the next time step, otherwise it will be *dead* in the next time step.
- If a cell is *dead* and has exactly three neighbours it will become alive in the next time step, otherwise it will stay dead in the next time step.

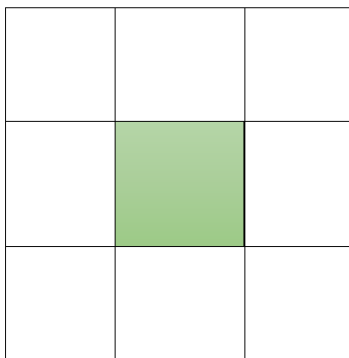


Figure 1: The current cell displayed as green and its eight surrounding neighbours

The goal of this project is to implement a Game of Life simulation, optimize the serial code and finally parallelize the code using Pthreads. Thereafter, the time complexity of the program will be analysed and the effects of the parallelization will be evaluated.

3 Solution Method

3.1 Original Implementation

The original code for the Game of Life simulation can be found in Appendix A. The script consists of a main function that asks the user for an input of the size of the grid, the number of time steps and a variable that determines if the grid state should be printed on the screen at each time step (i.e. if a visual simulation should be enabled or not). The last mentioned variable should take the value 1 if the simulation is desired or 0 if it is not. If the wrong input format is entered, the user will receive an error message describing the expected input and the program will terminate.

The size of the grid is declared and stored in the variable `grid_size`. The grid is then initiated and filled with a randomly generated initial state by the function `void initGrid(int grid_size, int** grid)`. The grid is represented by a 2D array of the datatype `int`, and required memory for this is allocated using `malloc`. The `initGrid` function uses two nested for-loops to fill the grid with a randomized initial state. The cells of the grid will randomly obtain either the number 0, indicating that the cell is *dead*, or the number 1 meaning that the cell is *alive*.

After the grid has been initialized, the goal is to update the grid at every time step, which is done by the function `void updateGrid(int grid_size, int** grid, int** temp_grid)`, that updates the grid according to the rules described in section 2. Before this function is called however, a larger grid `temp_grid` is initialized. The purpose of this grid is to make sure that the edges of the grid are handled correctly, and therefore at the beginning of every time step a deep copy of the main grid is made and stored in this larger grid. This grid is also given a padding of zeros on its edges that will ensure that the edge cells of the main grid will be handled correctly in the computations. Then, both the main grid and the padded grid is sent into the function `updateGrid`.

Inside the `updateGrid` function, two nested for loops are used to iterate over the cells, and based on the padded array `temp_grid` the number of alive neighbours are computed by the function `int countNeighbours(int grid_size, int** grid, int i, int j)`. This function takes the temporary padded grid as input together with the position of the current cell within the padded grid, and returns the sum of all the eight surrounding cells which corresponds to the number of alive neighbours (since alive is defined as 1 and dead as 0). The new state of each cell is then updated in the original (unpadded) array according to the rules described in section 2. This process for updating the grid will be repeated for the given number of time steps by using a for loop in the main function. When this process is finished, all the allocated memory is freed and the program terminates.

If the simulation option is enabled (i.e. the variable `simulation` is set to 1), the main grid will be printed at every time step by the function `void printGrid(int grid_size, int** grid)`. This function takes the grid as input and displays the cells as `'.'` if the cell is dead and `'*'` if the cell is alive. The function `nanosleep` from the library `time.h` is called after every print of the grid in order to make the simulation show the states more smoothly.

3.2 Serial Optimizations

The code for when the serial optimizations had been made is presented in Appendix B. When the code was being optimized, several different methods for optimizing the code were tested and systematically evaluated by measuring the time it took for the program to run. Those optimizations that resulted in an improved performance were kept. Even though some alterations might have made the readability of the code worse, they were kept anyways since the main focus of this project was to minimize the run time of the program. In this section a summary of all the adjustments that were made to the code is presented. A more thorough description of the optimization techniques that were tested can be found in section 4.1.

Firstly, it was found that the `-O3` compiler flag generated the fastest code, and therefore that compiler flag was used when compiling the program. Secondly, several if-statements within loops were removed and replaced by additional for-loops instead. After this change, the remaining if-statements that existed in the code were optimized by merging two else-if branches. Another optimization that was made was to remove the `countNeighbours` function that was called inside a nested for-loop. Instead, the computations that were made by the function were computed directly inside the loop instead. Also, the variable `int neighbours` was initiated in the main function instead of inside the nested for-loop in order to minimize the memory usage of the program. Loop unrolling was introduced with the `-funroll-loops` compiler flag and the datatype for the arrays used in the program was changed to `short int` instead of `int`. Also, both the `constant` keyword and the `restrict` keyword were added to multiple places in the code. Additionally, the number of `malloc` calls were reduced significantly by pointing arrays into larger buffer arrays instead. Lastly, it was confirmed that the `-O3` compiler flag successfully managed to vectorize 6 loops in the code.

3.3 Parallelization Using Pthreads

The code for the parallelized implementation is presented in Appendix C. In this implementation, the two arrays used in the program (`grid` and `temp_grid`) are made global so that all threads easily can have access to them. Also, another input argument is required by the user to enter, which states the number of threads to use when running the program. This value is stored in the variable `num_of_threads`.

In order to be able to send data to the different threads, a struct `thread_data_t` is created. This enables every thread to get access to specific values that are needed for the computations. An array of the struct `thread_data_t` is initialized for the given number of threads, and the amount of rows in the grid for every thread to work on is stored in the variable `work_size`. If the size of the grid is not evenly divisible by the number of threads entered, this work size amount will not cover the last rows in the grid. This problem is solved by computing the rest of the division and storing this in a variable called `rest`. From this, the last rows will be assigned to the last thread that is being created.

The grids are initialized and filled with values in the same manner before. Then, a global mutex and a global condition variable is initiated using the functions `pthread_mutex_init` and `pthread_cond_init`, that later are used in the synchronisation of the threads. After this step, the threads are given their needed values for the computations and their start and stop indexes that specify which rows they should do their work on. All threads are pointed to the same thread function `void* updateGrid(void* arg)`. This function updates the grid for every time step (in the same manner as in previous implementation), and uses the barrier function `void barrier(const int num_of_threads, const int grid_size, const int simulation)` to make sure that every thread has updated the current grid before they all move on to the next time step.

The `barrier` function keeps count of how many threads that have finished their computations and are ready to move on to the next time step. This value is kept in the global variable `waiting` that is increased every time the barrier function is called (which is done by every thread in the beginning of every time step). As long as all threads are not finished with their work for the current time step, the thread will be put to sleep by the function `pthread_cond_wait`. When the last thread is done with its computations and calls for the `barrier` function (i.e. when all other threads are waiting for the last thread to finish), the variable `waiting` is set to zero and a call for the function `pthread_cond_broadcast` is made in order to wake up all waiting threads that have been set to sleep. If the simulation is desired (i.e. the `simulation` variable has been set to 1), a call for the `printGrid` function is made, which displays the current grid on the screen. When this has been done, the newly awakened threads will continue to do their work for the following time steps according to the same course of action that has just been described.

When the grid has been updated for all time steps, the main function will call the function `pthread_join`

for every thread. The allocated memory is freed and the mutex and condition variables are destroyed by the functions `pthread_mutex_destroy` and `pthread_cond_destroy` respectively. Finally, the program terminates by returning 0.

4 Experiments

When the code was being serial optimized, time measurements were made in order to evaluate which optimization techniques that were effective. When these experiments have been performed, the executable file has been run with the following command: `./gameOfLife 1000 1000 0`, and the time was measured in wall seconds (by using the function `static double get_wall_seconds()`). This command indicates that the simulation was run for 1000 time steps with a grid size of 1000x1000. Also, no visual simulation was made when optimizing, i.e. the `simulation` variable was set to 0 (as can be seen in the input command). Whenever the code was changed and a time measurement was made, the code was run five times and the measurement with the best performance was considered as the final result. This approach was done for the purpose of minimizing the effects that other programs running in the background of the computer could have had on the measurements. The code was optimized on a *Macbook Pro Intel(R) Core(TM) i5-5257U CPU @ 2.70GHz* with the compiler version *clang-1001.0.46.4*.

4.1 Serial Optimizations

The first step when optimizing the serial code was to test different compiler flags and examine their effect on the time performance of the program. In Table 1 the measured times is presented for the different compiler flags. It can be noted that the flag `-O3` gave the best performance.

Table 1: Time measurements for different compiler flags, where the compiler flag with the best performance is marked in bold

Flags	Time [s]
-O0 (default)	20.085
-O1	6.290
-O2	4.790
-O3	4.403
-Ofast	4.452
-Os	7.577
-O3 -ffast-math	4.407
-O3 -march=native	9.245
-O3 -march=native -ffast-math	4.427
-O2 -march=native -ffast-math	4.712
-Ofast -march=native -ffast-math	4.485

The next step in the optimization process was to evaluate all loops used in the program and examine if they could be rewritten in any advantageous way. The first goal was to try to avoid having if-statements within loops, since that often can worsen the performance of a program. The if-statements in the nested for-loops that made the deep copy of the grid and added padding to the temporary grid was removed (see lines 54-58 and 75-79 in Appendix A). Instead, two more for-loops were added in order to make the code function correctly (see lines 55-67 and 80-92 in Appendix B). When implementing this change, the time required for the code to run was reduced to 4.185 seconds and thus this alteration was kept.

It was not possible to remove the if-statement that checked the rules of the game (lines 116-124 in Appendix A) in the `updateGrid` function. However, the two else-if statements on lines 118-120 in Appendix A could be merged together (see result in lines 125-126 in Appendix B) which resulted in an improvement of

the performance since the time reduced to 3.726 seconds. The order of the boolean expressions within the first if-statement (on line 116 in Appendix A) was changed as well, for the purpose of implementing boolean short circuits. Since there is a greater probability that the number of neighbours is not correct, compared to the 50% chance that the cell is either dead or alive, it seemed more beneficial to have the boolean expression checking the number of neighbours first. Despite this reasoning, that change made the performance slightly worse, as the time increased to 4.150 seconds, and therefore the change was not kept.

It was also tested to put the computation of the number of neighbours directly in the loop instead of calling for the function `countNeighbours` (line 114 in Appendix A). This alteration resulted in a small improvement since the time reduced to 3.650, and so the change was kept. Additionally, for the purpose of reducing the memory footprint of the program, the initiation of the variable `int neighbours` (on line 114 in Appendix A) was moved to the main function and the variable was then used as an input to the function `updateGrid` instead. From this, the same memory location for that variable could be used through the whole simulation instead of declaring a new variable every time the function `countNeighbours` was called.

Loop unrolling is also an approach that sometimes can improve the performance of a code. The compilation flag `-funroll-loops` was added in the compilation in order to test its effect. It did not give any improvement in the performance however, nor did it make the performance worse, so the change was kept.

For the purpose of minimizing the memory usage of the program, the arrays for the grids were changed to be of the datatype `short int` instead of `int`. Because the array elements only take on the values 0 or 1, it seemed unnecessary to allocate more memory for the arrays than was needed, and therefore less memory could be used by the program which often can improve the performance. When testing this change, it could be noted that the performance did improve as the time reduced to 3.590 seconds. Therefore, this alteration was kept.

The `const` keyword was added to the unchanging variables that existed in the code, i.e. the variables `double time1`, `int time_steps`, `int simulation`, `int grid_size` (on lines 20, 26, 27 and 25 in Appendix A respectively) and all `int grid_size` input parameters to the functions. Because of this keyword, the compiler would not have to assume that the variables can change during the execution of the program, which sometimes can improve the performance. In this case, a small improvement of the performance could be noted as the time reduced to 3.535 seconds. Thus, the keywords were kept.

In all functions the keyword `__restrict` was added in the function input parameters for the arrays used in the program. This keyword could be used since it was known that none of the pointers used as input in the respective functions pointed to the same memory address. However, this alteration did not have any significant improvement of the performance, nor did it make it any worse, and so the keywords were kept.

Another thing that can improve the performance of a program is to minimize the number of `malloc` and `free` calls. In order to do this, two buffer arrays for the rows and columns in the two grids used in the program were created. From this, only two `malloc` call to the buffer arrays were needed since they allocated enough memory for both arrays. Then, the two arrays for the grids were pointed into the larger buffer arrays instead (see implementation in Appendix B lines 33-46 and 99-100). This change reduced the number of `malloc` and `free` calls significantly, which could also be noticed in the performance of the program since the run time reduced to 3.249 seconds. Thus, this change was kept.

Vectorization of the code is also an approach that can improve the performance of a code. Since the `-O3` compiler flag was used, which by default includes auto-vectorization of the code, it was assumed that auto vectorization already was enabled. For the purpose of examining if the code could be rewritten in any way so that more auto vectorizations could be made, the flag `-ftree-vectorizer-verbose=2` was included when compiling. From this, it was concluded that 6 loops in the program had successfully been vectorized,

namely the loops starting on lines 38, 44, 56, 60, 81 and 85 in Appendix B. The nested for-loops inside the `updateGrid` function and inside the `printGrid` function were two of the loops that were not vectorized, which probably was due to the if-statements that existed within the loops. For the purpose of vectorizing the rest of the unvectorized loops, it was tested to add the `-march=native`, `-ffast-math` and `-mavx` flags when compiling. However, these flags made no difference, and unfortunately, since the rest of the loops did not have any distinct loop invariants, it was difficult to figure out why these loops were not being vectorized. Hence, no more loops in the code were successfully vectorized.

4.2 Parallelization Using Pthreads

In order to determine what parts of the code that consumed most of the time, the tool `gprof` was used. From this it was clear that the function `updateGrid` occupied a the majority of the time required for the code to run (more precisely 86% of the time). Thus, the main focus of the parallelization process was laid on optimizing this function.

The final code, presented in Appendix C, was run for different grid sizes on one thread for the purpose of analysing the time complexity of the algorithm. The code was run with the following command: `./gameOfLife grid_size 1000 0 1`, meaning that the variable `grid_size` was varied and that the code was run for 1000 time steps, without the simulation enabled and with only one thread. In Figure 2 the result of these measurements is displayed. It can be observed from the shape of the curve that the time complexity of the program most probably is $\mathcal{O}(N^2)$. The time for when the command `./gameOfLife 1000 1000 0 1` was run (the same command used when the serial optimized code was being developed) was measured to be 3.207 seconds which is slightly better than its previous version of the code (the one presented in Appendix B).

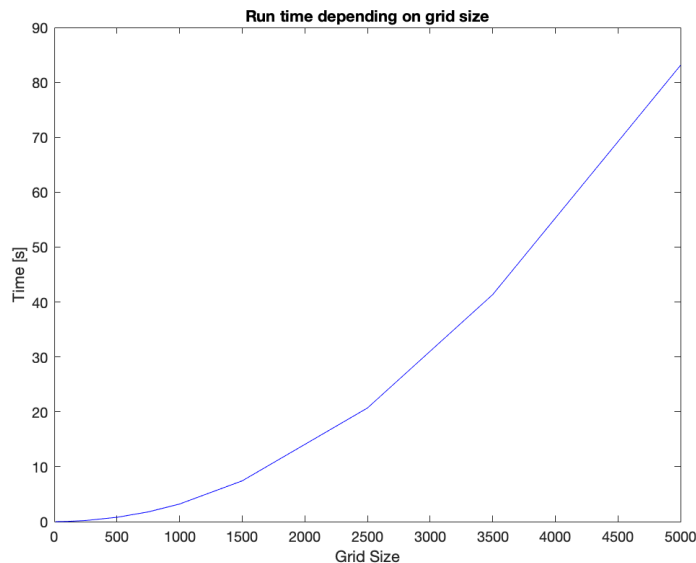


Figure 2: Run time of the program for different sizes of the grid

An analysis of the run time of the program for different number of threads was also done. The code was run with the command `./gameOfLife 1000 1000 0 num_of_threads` where the `num_of_threads` variable was varied. This experiment was executed on one of the University's Linux computers, due to the large amount of cores available on it (16 cores). More precisely, the used computer was an Intel® Xeon® E5520 @ 2.27GHz, with a total of 16 CPU:s (4 cores \times 2 sockets, 2 threads per core), and compiler *gcc version*

7.5.0 (Ubuntu 7.5.0-3ubuntu1 18.04). The result of these measurements is presented in Figure 3. It can be noted that the time reduces in an $1/x$ like fashion and that when the number of threads exceeds the number of cores available, the time does not continue to decrease anymore. The lowest time measurement was 0.670 seconds which was when 16 threads were used.

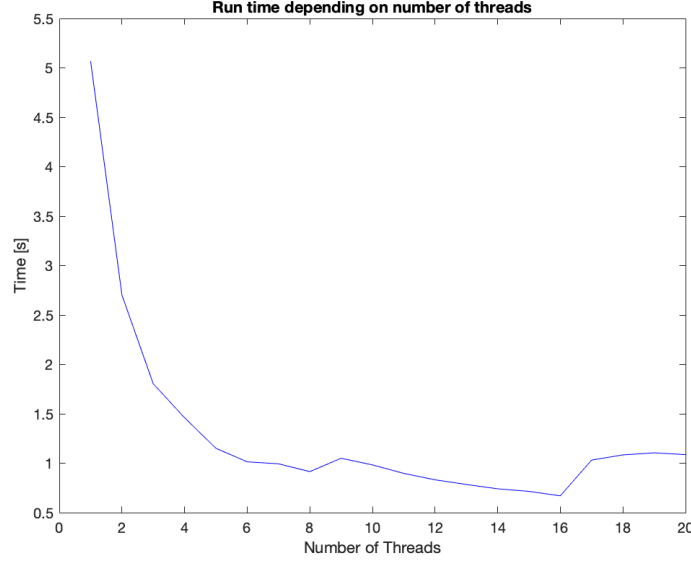


Figure 3: Run time of the program for different number of threads

5 Conclusions

In summary, the performance of the implemented Game of Life simulation has increased significantly after serial optimizations and parallelization of the code was made. Also, the time measurements when different numbers of threads were used is behaving as expected. However, the time complexity of the program was determined to be $\mathcal{O}(N^2)$ which means that the run time of the program will increase very fast as the grid size increases, which is not a preferred aspect of an algorithm. This time complexity is due to the fact that iterations over every cell in the grid is made in order to find the next state.

An approach that could be considered to improve the complexity of the algorithm could be to only count the number of neighbours of the alive cells in the grid. In order for this approach to work properly, all 3x3 neighbourhoods with exactly 3 alive cells must also be detected, so that the dead cells within that neighbourhood will be updated correctly. Since these situations only can occur in the presence of at least one alive cell however, only these potential areas would have to be searched instead of the whole grid. This algorithm might be effective for states with few alive cells in the grid (for example when specific seeds with few alive cells ought to be tried). However, for a randomized initial state as used in this implementation, this algorithm will most likely not result in any significant improvement of the time. Another thing that could be considered though, is to somehow keep track of which cells that were changed in the previous time step, so that only these areas of these cells of the grid must be examined. This approach might in some cases reduce the computations needed by the program, but depending on the initial state, it can not be ensured. In summary, it was hard to come up with an algorithm with a better time complexity than $\mathcal{O}(N^2)$, when also considering the worst case of the simulation. Therefore, the current implementation can be considered as an optimal solution of the problem.

The parallelization of the code was made using the `barrier` function, for the purpose of keeping the threads

alive through all time steps instead of creating and terminating threads on every time step. This approach is highly preferable since creating and joining threads can take up a significant amount of the run time of a program if it has to be done often. When running the program for many time steps, the barrier implementation is therefore more efficient than it would have been if the other approach would have been used. Additionally, the work amount is evenly distributed to all the threads which is favorable for parallelized work since it minimizes the time spent on waiting for the last thread to finish its work (on every time step). Thus, the parallelized implementation made in this project can be considered as an optimal solution to the problem.

Regarding the precision of the experiments in the project, it might have been more beneficial to start every simulation with the same initial seed. This would have made sure that exactly the same computations were made during all runs of the program and it would also have been easier to determine the correctness of the simulation after changes were made to the code (since the grid would always look the same after a specific number of time steps). Now, a different evolution of the states were generated every time the code was run which could have had a small impact on the time. Perhaps, in this project it would have been better to take an average value of the 5 runs that were made for every measurement, instead of considering the fastest time as the final result.

When evaluating the serial optimizations that were made, another approach could have been to allocate the arrays on stack instead of calling `malloc`. This approach can sometimes give an improved performance, however it was not chosen in this project because it would set a limit on the size of the grid. If a too large grid size would be entered by the user, the code would generate a `segmentation fault` error due to the fact that not enough memory on stack would be available for the desired grid size. Although, in similar implementations if large grid sizes are not needed for the purpose of the program, it might be beneficial to allocate the arrays on stack instead.

Another serial optimization that was not tested, was to introduce loop unrolling manually (instead of using the compiler flag `-funroll-loops`). This was not done because it seemed relatively complicated since nearly all loops in the program were dependent on the input parameter `grid_size`. In order to handle the cases where the grid size was not divisible with the unroll factor, an additional loop would be necessary to compute the remaining values. For small unroll factors, this would most likely not have made a big difference in the performance because the loop condition for the remaining values would have to be checked anyways. However, if larger unroll factors were to be used, it might have given an improved performance if it was implemented.

Appendix A

Listing 1: *Original implementation*

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <unistd.h>
4  #include <sys/time.h>
5  #include <time.h>
6
7  void initGrid(int grid_size, int** grid);
8  void updateGrid(int grid_size, int** grid, int** temp_grid);
9  int countNeighbours(int grid_size, int** grid, int i, int j);
10 void printGrid(int grid_size, int** grid);
11
12 static double get_wall_seconds(){
13     struct timeval tv;
14     gettimeofday(&tv, NULL);
15     double seconds = tv.tv_sec + (double)tv.tv_usec / 1000000;
16     return seconds;
17 }
18
19 int main(int argc, char* argv[]){
20     double time1 = get_wall_seconds(); // start time
21     if (argc != 4){
22         printf("Expected input: ./gameOfLife grid_size time_steps simulation\n");
23         return -1;
24     }
25     int grid_size = atoi(argv[1]);
26     int time_steps = atoi(argv[2]);
27     int simulation = atoi(argv[3]); // 1 = simulation on, 0 = simulation off
28     if(simulation != 0 && simulation != 1){
29         printf("Error: Last input argument must be 1 or 0. \n");
30         return -1;
31     }
32
33     // Initiate grids
34     int** grid = (int**)malloc(grid_size*sizeof(int*)); // main grid
35     for(int i = 0; i < grid_size; i++){
36         grid[i] = (int*)malloc(grid_size*sizeof(int*));
37     }
38     initGrid(grid_size, grid); // fill grid with start seed
39
40     int** temp_grid = (int**)malloc((grid_size+2)*sizeof(int*)); // "old" (padded) grid
41     for(int i = 0; i < grid_size+2; i++){
42         temp_grid[i] = (int*)malloc((grid_size+2)*sizeof(int*));
43     }
44
45     if(simulation == 1){ // Simulation ON
46         printf("Start grid:\n");
47         printGrid(grid_size, grid);
48
49         // Update grid for every timestep
50         for(int k = 0; k < time_steps; k++){
51             // Make a deep copy of current grid and add padding to temp_grid
52             for(int i = 0; i < grid_size+2; i++){
53                 for(int j = 0; j < grid_size+2; j++){
54                     if(i < 1 || i > grid_size || j < 1 || j > grid_size){
55                         temp_grid[i][j] = 0;
56                     } else {
57                         temp_grid[i][j] = grid[i-1][j-1];
58                     }
59                 }

```

```

60         }
61
62         updateGrid(grid_size, grid, temp_grid);
63
64         printf("\nUpdated grid:\n");
65         printGrid(grid_size, grid);
66         nanosleep((const struct timespec[]){0, 100000000L}, NULL);
67     }
68     } else { // Simulation OFF
69
70         // Update grid for every timestep
71         for(int k = 0; k < time_steps; k++){
72             // Make a deep copy of current grid and add padding to temp_grid
73             for(int i = 0; i < grid_size+2; i++){
74                 for(int j = 0; j < grid_size+2; j++){
75                     if(i < 1 || i > grid_size || j < 1 || j > grid_size){
76                         temp_grid[i][j] = 0;
77                     } else {
78                         temp_grid[i][j] = grid[i-1][j-1];
79                     }
80                 }
81             }
82
83             updateGrid(grid_size, grid, temp_grid);
84         }
85     }
86
87     // Free allocated memory
88     for(int i = 0; i < grid_size; i++){
89         free(grid[i]);
90     }
91     free(grid);
92     for(int i = 0; i < grid_size+2; i++){
93         free(temp_grid[i]);
94     }
95     free(temp_grid);
96
97     printf("gameOfLife took %7.3f wall seconds.\n", get_wall_seconds()-timel);
98     return 0;
99 }
100
101 void initGrid(int grid_size, int** grid){
102     srand(1);
103     for(int i = 0; i < grid_size; i++){
104         for(int j = 0; j < grid_size; j++){
105             grid[i][j] = rand()%2; // 0 = dead, 1 = alive
106         }
107     }
108 }
109
110 void updateGrid(int grid_size, int** grid, int** temp_grid){
111     // Compute new grid
112     for(int i = 1; i ≤ grid_size; i++){
113         for(int j = 1; j ≤ grid_size; j++){
114             int neighbours = countNeighbours(grid_size+2, temp_grid, i, j);
115
116             if(temp_grid[i][j] == 0 && neighbours == 3){
117                 grid[i-1][j-1] = 1;
118             } else if(temp_grid[i][j] == 1 && neighbours == 2){
119                 grid[i-1][j-1] = 1;
120             } else if(temp_grid[i][j] == 1 && neighbours == 3){
121                 grid[i-1][j-1] = 1;
122             } else {
123                 grid[i-1][j-1] = 0;
124             }

```

```

125     }
126 }
127 }
128
129 int countNeighbours(int grid_size, int** grid, int i, int j){
130     return grid[i-1][j-1]+grid[i-1][j]+grid[i-1][j+1]+grid[i][j-1]+grid[i][j+1]+
131     grid[i+1][j-1]+grid[i+1][j]+grid[i+1][j+1];
132 }
133
134 void printGrid(int grid_size, int** grid){
135     char c;
136     for(int i = 0; i < grid_size; i++){
137         for(int j = 0; j < grid_size; j++){
138             if(grid[i][j] == 1){ // alive
139                 c = '*';
140             } else { // dead
141                 c = '.';
142             }
143             printf("%c ", c);
144         }
145         printf("\n");
146     }
147 }

```

Appendix B

Listing 2: *Serial optimized implementation*

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <unistd.h>
4  #include <sys/time.h>
5  #include <time.h>
6
7  void initGrid(const int grid_size, short int** __restrict grid);
8  void updateGrid(const int grid_size, short int** __restrict grid, short int** __restrict ...
    temp_grid, int neighbours);
9  void printGrid(const int grid_size, short int** __restrict grid);
10
11 static double get_wall_seconds(){
12     struct timeval tv;
13     gettimeofday(&tv, NULL);
14     double seconds = tv.tv_sec + (double)tv.tv_usec / 1000000;
15     return seconds;
16 }
17
18 int main(int argc, char* argv[]){
19     const double time1 = get_wall_seconds(); // start time
20     if (argc != 4){
21         printf("Expected input: ./gameOfLife grid_size time_steps simulation\n");
22         return -1;
23     }
24     const int grid_size = atoi(argv[1]);
25     const int time_steps = atoi(argv[2]);
26     const int simulation = atoi(argv[3]); // 1 = simulation on, 0 = simulation off
27     if(simulation != 0 && simulation != 1){
28         printf("Error: Last input argument must be 1 or 0. \n");
29         return -1;
30     }
31     int neighbours;
32
33     // Initiate grids
34     short int** buffer_rows = (short int**)malloc((grid_size+grid_size+2)*sizeof(short int*));
35     short int* buffer_col = (short int*)
        malloc((grid_size*grid_size+(grid_size+2)*(grid_size+2))*sizeof(short int*));
36
37     short int** grid = buffer_rows;
38     for(int i = 0; i < grid_size; i++){
39         grid[i] = (buffer_col + i*grid_size);
40     }
41     initGrid(grid_size, grid); // fill grid with start seed
42
43     short int** temp_grid = (buffer_rows+grid_size); // "old" (padded) grid
44     for(int i = 0; i < grid_size+2; i++){
45         temp_grid[i] = (buffer_col + grid_size*grid_size + i*(grid_size+2));
46     }
47
48     if(simulation == 1){ // Simulation ON
49         printf("Start grid:\n");
50         printGrid(grid_size, grid);
51
52         // Update grid for every timestep
53         for(int k = 0; k < time_steps; k++){
54             // Make a deep copy of current grid and add padding
55             for(int i = 1; i ≤ grid_size; i++){ // fill middle of grid
56                 for(int j = 1; j ≤ grid_size; j++){
57                     temp_grid[i][j] = grid[i-1][j-1];
```

```

58         }
59     }
60     for(int j = 0; j < grid_size+2; j++){ // fill horizontal padding
61         temp_grid[0][j] = 0;
62         temp_grid[grid_size+1][j] = 0;
63     }
64     for(int i = 1; i ≤ grid_size; i++){ // fill vertical padding
65         temp_grid[i][0] = 0;
66         temp_grid[i][grid_size+1] = 0;
67     }
68
69     updateGrid(grid_size, grid, temp_grid, neighbours);
70
71     printf("\nUpdated grid:\n");
72     printGrid(grid_size, grid);
73     nanosleep((const struct timespec[]){0, 100000000L}, NULL);
74 }
75 } else { // Simulation OFF
76
77     // Update grid for every timestep
78     for(int k = 0; k < time_steps; k++){
79         // Make a deep copy of current grid and add padding
80         for(int i = 1; i ≤ grid_size; i++){ // fill middle of grid
81             for(int j = 1; j ≤ grid_size; j++){
82                 temp_grid[i][j] = grid[i-1][j-1];
83             }
84         }
85         for(int j = 0; j < grid_size+2; j++){ // fill horizontal padding
86             temp_grid[0][j] = 0;
87             temp_grid[grid_size+1][j] = 0;
88         }
89         for(int i = 1; i ≤ grid_size; i++){ // fill vertical padding
90             temp_grid[i][0] = 0;
91             temp_grid[i][grid_size+1] = 0;
92         }
93
94         updateGrid(grid_size, grid, temp_grid, neighbours);
95     }
96 }
97
98 // Free allocated memory
99 free(buffer_col);
100 free(buffer_rows);
101
102 printf("gameOfLife took %7.3f wall seconds.\n", get_wall_seconds()-timel);
103 return 0;
104 }
105
106 void initGrid(const int grid_size, short int** __restrict grid){
107     srand(1);
108     for(int i = 0; i < grid_size; i++){
109         for(int j = 0; j < grid_size; j++){
110             grid[i][j] = rand()%2; // 0 = dead, 1 = alive
111         }
112     }
113 }
114
115 void updateGrid(const int grid_size, short int** __restrict grid, short int** __restrict ...
temp_grid, int neighbours){
116     // Compute new grid
117     for(int i = 1; i ≤ grid_size; i++){
118         for(int j = 1; j ≤ grid_size; j++){
119             neighbours = temp_grid[i-1][j-1]+temp_grid[i-1][j]+temp_grid[i-1][j+1]+
120 temp_grid[i][j-1]+temp_grid[i][j+1]+temp_grid[i+1][j-1]+temp_grid[i+1][j]+
121 temp_grid[i+1][j+1];

```

```

122
123         if(temp_grid[i][j] == 0 && neighbours == 3){
124             grid[i-1][j-1] = 1;
125         } else if(temp_grid[i][j] == 1 && (neighbours == 2 || neighbours == 3)){
126             grid[i-1][j-1] = 1;
127         } else {
128             grid[i-1][j-1] = 0;
129         }
130     }
131 }
132 }
133
134 void printGrid(const int grid_size, short int** __restrict grid){
135     char c;
136     for(int i = 0; i < grid_size; i++){
137         for(int j = 0; j < grid_size; j++){
138             if(grid[i][j] == 1){ // alive
139                 c = '*';
140             } else { // dead
141                 c = '.';
142             }
143             printf("%c ", c);
144         }
145         printf("\n");
146     }
147 }

```


Appendix C

Listing 3: *Parallelized implementation*

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <unistd.h>
4  #include <sys/time.h>
5  #include <time.h>
6  #include <pthread.h>
7
8  // Structs
9  typedef struct thread_data {
10     int start;
11     int stop;
12
13     int grid_size;
14     int time_steps;
15     int num_of_threads;
16     int simulation;
17
18     pthread_t thread;
19 } thread_data_t;
20
21 // Global variables
22 short int** grid;
23 short int** temp_grid;
24
25 pthread_mutex_t lock;
26 pthread_cond_t cv;
27 int state = 0;
28 int waiting = 0;
29
30 // Functions
31 void initGrid(const int grid_size, short int** __restrict grid);
32 void* updateGrid(void* arg);
33 void printGrid(const int grid_size, short int** __restrict grid);
34 void barrier(const int num_of_threads, const int grid_size, const int simulation);
35
36 static double get_wall_seconds(){
37     struct timeval tv;
38     gettimeofday(&tv, NULL);
39     double seconds = tv.tv_sec + (double)tv.tv_usec / 1000000;
40     return seconds;
41 }
42
43 int main(int argc, char* argv[]){
44     const double time1 = get_wall_seconds(); // start time
45     if (argc != 5){
46         printf("Expected input: ./gameOfLife grid_size time_steps simulation ...
47             num_of_threads\n");
48         return -1;
49     }
50     const int grid_size = atoi(argv[1]);
51     const int time_steps = atoi(argv[2]);
52     const int simulation = atoi(argv[3]); // 1 = simulation on, 0 = simulation off
53     if(simulation != 0 && simulation != 1){
54         printf("Error: simulation input argument must be 1 or 0. \n");
55         return -1;
56     }
57     const int num_of_threads = atoi(argv[4]);
58
59     thread_data_t threads[num_of_threads];
```

```

59  const int work_size = grid_size/num_of_threads;
60  const int rest = grid_size%num_of_threads; // extra work amount for last thread
61
62  // Initiate grids
63  short int** buffer_rows = (short int**)malloc((grid_size+grid_size+2)*sizeof(short int*));
64  short int* buffer_col = (short ...
        int*)malloc((grid_size*grid_size+(grid_size+2)*(grid_size+2))*sizeof(short int*));
65
66  grid = buffer_rows;
67  for(int i = 0; i < grid_size; i++){
68      grid[i] = (buffer_col + i*grid_size);
69  }
70  initGrid(grid_size, grid); // fill grid with start seed
71
72  temp_grid = (buffer_rows+grid_size); // "old" (padded) grid
73  for(int i = 0; i < grid_size+2; i++){
74      temp_grid[i] = (buffer_col + grid_size*grid_size + i*(grid_size+2));
75  }
76
77  // Make deep copy of grid to temp_grid and add padding on edges
78  for(int i = 1; i ≤ grid_size; i++){ // fill middle of grid
79      for(int j = 1; j ≤ grid_size; j++){
80          temp_grid[i][j] = grid[i-1][j-1];
81      }
82  }
83  for(int j = 0; j < grid_size+2; j++){ // fill horizontal padding
84      temp_grid[0][j] = 0;
85      temp_grid[grid_size+1][j] = 0;
86  }
87  for(int i = 1; i ≤ grid_size; i++){ // fill vertical padding
88      temp_grid[i][0] = 0;
89      temp_grid[i][grid_size+1] = 0;
90  }
91
92  // Initiate mutex and condition variables
93  pthread_mutex_init(&lock, NULL);
94  pthread_cond_init(&cv, NULL);
95
96  // Create threads
97  for(int k = 0; k < num_of_threads-1; k++){
98      threads[k].start = (k+1) * work_size - work_size;
99      threads[k].stop = threads[k].start + work_size;
100
101      threads[k].grid_size = grid_size;
102      threads[k].time_steps = time_steps;
103      threads[k].num_of_threads = num_of_threads;
104      threads[k].simulation = simulation;
105
106      pthread_create(&(threads[k].thread), NULL, updateGrid, &threads[k]);
107  }
108  threads[num_of_threads-1].start = num_of_threads * work_size - work_size;
109  threads[num_of_threads-1].stop = threads[num_of_threads-1].start + work_size + rest;
110  threads[num_of_threads-1].grid_size = grid_size;
111  threads[num_of_threads-1].time_steps = time_steps;
112  threads[num_of_threads-1].num_of_threads = num_of_threads;
113  threads[num_of_threads-1].simulation = simulation;
114  pthread_create(&(threads[num_of_threads-1].thread), NULL, updateGrid, ...
        &threads[num_of_threads-1]);
115
116  // Join threads
117  for(int k = 0; k < num_of_threads; k++){
118      pthread_join(threads[k].thread, NULL);
119  }
120
121  // Free allocated memory

```

```

122     free(buffer_col);
123     free(buffer_rows);
124
125     // Destroy mutex and condition variables
126     pthread_mutex_destroy(&lock);
127     pthread_cond_destroy(&cv);
128
129     printf("gameOfLife took %7.3f wall seconds.\n", get_wall_seconds()-timel);
130     return 0;
131 }
132
133 void initGrid(const int grid_size, short int** __restrict grid){
134     srand(1);
135     for(int i = 0; i < grid_size; i++){
136         for(int j = 0; j < grid_size; j++){
137             grid[i][j] = rand()%2; // 0 = dead, 1 = alive
138         }
139     }
140 }
141
142 void* updateGrid(void* arg){
143     thread_data_t* info = (thread_data_t *) arg;
144     int neighbours;
145
146     for(int k = 0; k < info->time_steps; k++){ // For every time step
147         barrier(info->num_of_threads, info->grid_size, info->simulation);
148
149         // Update temp_grid with most recent values
150         for(int i = info->start+1; i <= info->stop; i++){ // fill middle of grid
151             for(int j = 1; j <= info->grid_size; j++){
152                 temp_grid[i][j] = grid[i-1][j-1];
153             }
154         }
155
156         // Compute new grid
157         for(int i = info->start+1; i <= info->stop; i++){ // specific rows
158             for(int j = 1; j <= info->grid_size; j++){ // all columns
159                 neighbours = temp_grid[i-1][j-1]+temp_grid[i-1][j]+temp_grid[i-1][j+1]+
160                     temp_grid[i][j-1]+temp_grid[i][j+1]+temp_grid[i+1][j-1]+temp_grid[i+1][j]+
161                     temp_grid[i+1][j+1];
162
163                 if(temp_grid[i][j] == 0 && neighbours == 3){
164                     grid[i-1][j-1] = 1;
165                 } else if(temp_grid[i][j] == 1 && (neighbours == 2 || neighbours == 3)){
166                     grid[i-1][j-1] = 1;
167                 } else {
168                     grid[i-1][j-1] = 0;
169                 }
170             }
171         }
172     }
173     return NULL;
174 }
175
176 void printGrid(const int grid_size, short int** __restrict grid){
177     char c;
178     for(int i = 0; i < grid_size; i++){
179         for(int j = 0; j < grid_size; j++){
180             if(grid[i][j] == 1){ // alive
181                 c = '*';
182             } else { // dead
183                 c = '.';
184             }
185             printf("%c ", c);
186         }

```

```

187         printf("\n");
188     }
189 }
190
191 void barrier(const int num_of_threads, const int grid_size, const int simulation){
192     int mystate;
193     pthread_mutex_lock(&lock);
194     mystate = state;
195     waiting++; // Number of threads ready for next time step
196
197     if(waiting == num_of_threads){ // All threads ready for next timestep
198         waiting = 0;
199         state = 1-mystate;
200
201         if(simulation == 1){ // simulation enabled
202             printf("Current state:\n");
203             printGrid(grid_size, grid); // Print grid
204             nanosleep((const struct timespec[]){0, 100000000L}, NULL);
205         }
206
207         pthread_cond_broadcast(&cv); // Wake up sleeping threads
208     }
209     while(mystate == state){ // Not all threads are done
210         pthread_cond_wait(&cv, &lock); // Let thread sleep
211     }
212     pthread_mutex_unlock(&lock);
213 }

```