

**Darmstadt University of Applied Sciences**  
– Faculty of Computer Science –

**Implementing a distributed lattice gas  
cellular automaton using Rust, Rayon, and  
MPI**

by  
Zebreus

## CONCEPT

---

In the fourth lab, we built a lattice gas cellular automaton (LGCA) using Rust, Rayon, and MPI. Our LGCA is based on the FHP-I model. It uses a hexagonal grid and does not support resting particles. We will also try using a lookup table to process collisions quickly. We used rayon to divide the processing of the core rows over multiple threads. In addition, we used MPI via the `rsmapi` Rust bindings for distributed memory parallelism. Our implementation uses `AVX512 SIMD` instructions to process 64 cells at once. Our solution can process 492 Gigabytes of cell data per second when running on four intel nodes on the Virgo cluster.

## REQUIREMENTS

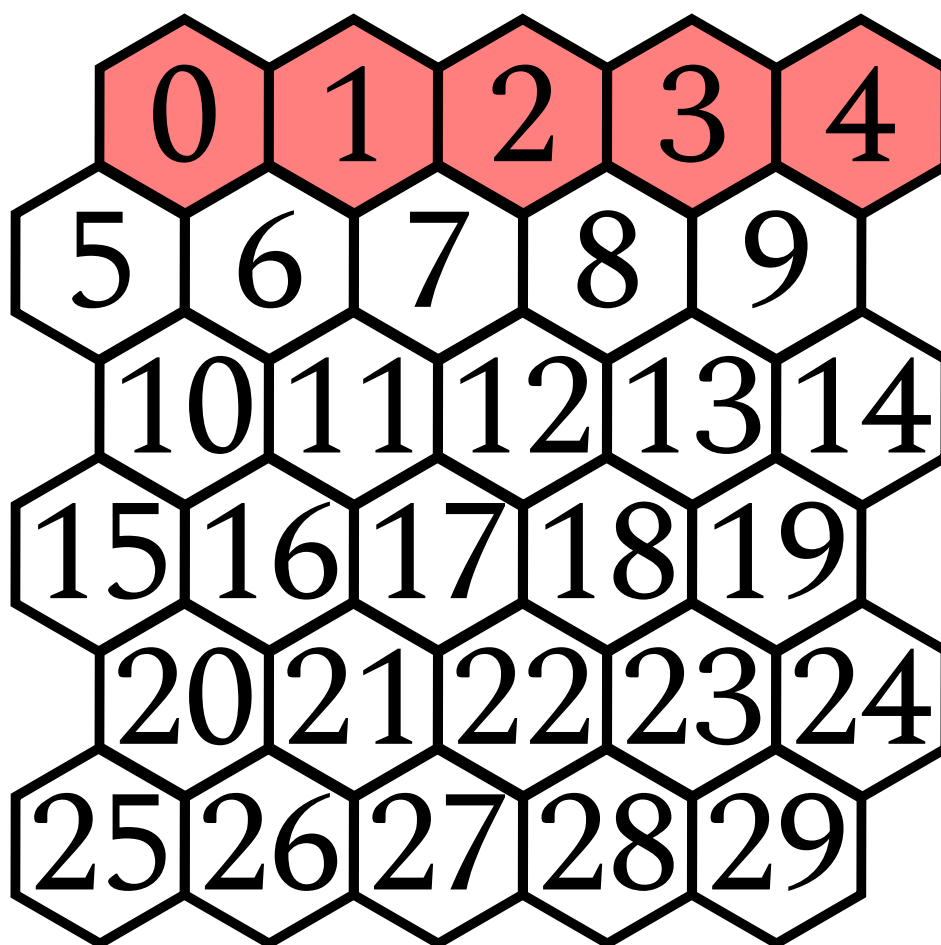
The following requirements are given by the lab instructions:

- Essentially, the lattice gas consists of a hexagonal array (the lattice), and a set of barriers.
- Each cell in the lattice is connected to 6 neighbors (see the illustrations in the report referenced above).
- Time is discrete (i.e.,  $t = 0, 1, 2, 3 \dots$ ).
- Each connection is boolean: Thus, at time  $t$ , there is either 0 or 1 particle (but never two or more particles) traveling from any given cell  $a$  to any cell  $b$  if  $a$  and  $b$  are neighbors (there might also be a particle traveling from  $b$  to  $a$ ; the two links are independent of each other).
- The state of the lattice at time  $t$  can be calculated in two phases:
- Transport: First, the particles travel linearly. For example, if cells  $a$ ,  $b$ , and  $c$  are linear neighbors (e.g.,  $c$  is to the right of  $b$ , and  $b$  is to the right of  $a$ ), then a particle traveling from  $a$  to  $b$  at time  $t = 1$  will continue to travel to  $c$  at time  $t = 2$  (in the absence of collisions, see below).
- Scattering: Second, collisions with other particles and with the barriers are taken into account. Collisions between particles are non-deterministic: More than one outcome is (sometimes) possible, and the simulation chooses between these outcomes randomly. See page 14 (188) in [Has87] for a table of the possible collisions.
- Particles reflect off barriers.

## IMPLEMENTATION CONCEPTS

The FHP-I model uses hexagonal cells. Each cell can contain up to six particles, one in each direction. The directions are called west, north-west, north-east, east, south-east, and south-west. We will store the cells row-wise in an even-r horizontal layout as shown in [Figure 1](#)

---



*Figure 1. "even-r" horizontal layout*

Cells are represented by bytes with one bit for each direction, as shown in [Figure 2](#).



*Figure 2. A single cell represented as a byte*

We will represent each row as an array of cells. The whole grid is represented as an array of rows. A single row will always be processed in one go from left to right.

First, all neighboring cells are inspected to gather the particles that will move into the current cell. At the grid's edges, the movement step reflects particles that would leave the grid. If the new direction also points outside of the grid, the particle is reflected in the direction it came from.

After the movement step, collisions inside the new cell are processed. The rules for the collisions are quite simple: \* If a cell only has two particles traveling in opposite directions, they will rotate by either  $60^\circ$  or  $-60^\circ$ . \* If a cell has three particles traveling in opposite directions, they will rotate by  $60^\circ$  (or  $-60^\circ$ , but that is the same). \* If a cell has exactly four particles with two pairs of opposite directions, they will rotate by either  $60^\circ$  or  $-60^\circ$ .

We identified four types of rows that need to be processed differently:

- Even row: The first, third, fifth, ... rows are even. The core of the row will be processed normally, but the first and last cells need to be processed differently, as they are missing some neighbors.
- Odd row: The second, fourth, ... rows are odd. The same as even rows, but the first and last cells are processed differently because other neighbors are missing.
- Top row: The first row of the grid is similar to an even row, but it has no neighbors to the north, so all particles that would move north are reflected.
- Bottom row: The last row of the grid is similar to an odd row, but it has no neighbors to the south, so all particles that would move south are reflected. Because the bottom row is an odd row, our implementation will only support grids with an even number of rows.

## IMPLEMENTATION

---

We define four functions for the four different cases. Each function has a special case for the first and last cell and delegates the core cells of the row to the `process_row_core` function. The top and bottom rows also have a special case for the core of them, as they need to reflect particles that would move out of the grid.

Nearly all particles will be processed in the grid's core for big grids, so we want to optimize that case as much as possible. We will use AVX512 SIMD instructions to process 64 cells at once. We will also try using a lookup table to process collisions quickly.

### PROCESSING ROWS

We define four functions for the four different cases. Each function has a special case for the first and last cell and delegates the core cells of the row to the `movement_core` function. The top and bottom rows also have a special case for the core of them, as they need to reflect particles that would move out of the grid.

To calculate the next state of a cell, we need to look at all of its neighboring cells and move the particles into the current cell. The movement logic for core cells is shown in [Listing 1](#). We need special logic to handle the missing neighbors for cells that are not in the core. Instead of the missing neighbor, we use a particle from the current cell as the neighbor. For example, the logic for the top-right cell is shown in [Listing 2](#).

*Listing 1. Movement calculation for core cells*

```
1      let mut new_cell = (west.raw & TO_EAST)
2          | (north_west.raw & TO_SOUTH_EAST)
3          | (north_east.raw & TO_SOUTH_WEST)
4          | (east.raw & TO_WEST)
5          | (south_east.raw & TO_NORTH_WEST)
6          | (south_west.raw & TO_NORTH_EAST);
```

RUST

*Listing 2. Movement calculation for the top-right corner cell*

```
1  // Handle border of first cell
2  result[0].raw = (below[0].raw & TO_NORTH_EAST)
3      | (below[1].raw & TO_NORTH_WEST)
4      | (current[1].raw & TO_WEST)
5      | ((current[0].raw & TO_NORTH_EAST) << 2)
6      | ((current[0].raw & TO_NORTH_WEST) << 4)
7      | ((current[0].raw & TO_WEST) << 3);
```

RUST

Our default collision handling is implemented in the `process_collision` function, as shown in [Listing 3](#). It uses a match statement to map the collision cases to the results. The default cases leave the value unchanged. In cases where there are two options, a real random number is used to select one.

*Listing 3. Default collision handling using real randomness*

RUST

```

1  pub fn process_collision(&mut self) {
2      self.raw = match self.raw {
3          // Two opposing particles
4          const { TO_WEST | TO_EAST } => {
5              if RNG.with(|f| f.borrow_mut().gen:::<bool>()) {
6                  TO_NORTH_EAST | TO_SOUTH_WEST
7              } else {
8                  TO_SOUTH_EAST | TO_NORTH_WEST
9              }
10         }
11         const { TO_SOUTH_EAST | TO_NORTH_WEST } => {
12             if RNG.with(|f| f.borrow_mut().gen:::<bool>()) {
13                 TO_NORTH_EAST | TO_SOUTH_WEST
14             } else {
15                 TO_EAST | TO_WEST
16             }
17         }
18         const { TO_SOUTH_WEST | TO_NORTH_EAST } => {
19             if RNG.with(|f| f.borrow_mut().gen:::<bool>()) {
20                 TO_SOUTH_EAST | TO_NORTH_WEST
21             } else {
22                 TO_EAST | TO_WEST
23             }
24         }
25     }
26     // Three particles
27     const { TO_SOUTH_WEST | TO_NORTH_WEST | TO_EAST } => {
28         TO_SOUTH_EAST | TO_NORTH_EAST | TO_WEST
29     }
30     const { TO_SOUTH_EAST | TO_NORTH_EAST | TO_WEST } => {
31         TO_SOUTH_WEST | TO_NORTH_WEST | TO_EAST
32     }
33
34     // Four particles with opposing holes
35     0b00110110 => {
36         if RNG.with(|f| f.borrow_mut().gen:::<bool>()) {
37             0b00011011
38         } else {
39             0b00101101
40         }
41     }
42     0b00011011 => {
43         if RNG.with(|f| f.borrow_mut().gen:::<bool>()) {
44             0b00101101
45         } else {
46             0b00110110

```

```

47         }
48     }
49     0b00101101 => {
50         if RNG.with(|f| f.borrow_mut().gen:::<bool>()) {
51             0b00011011
52         } else {
53             0b00110110
54         }
55     }
56
57     // Everything else
58     _ => self.raw,
59 }
60 }

```

While this collision implementation is correct, its performance is quite bad because we generate a random number for every collision. It also has lots of branching, and the compiler cannot auto-vectorize it.

#### OPTIMIZING

For big grids, nearly all particles will be processed in the core of the grid, so we want to optimize that case as much as possible. As the GSI has machines that support AVX512, we want to write code that utilizes these instructions to process 64 cells at once. While we could have tried to write the movement and collision logic in x86 assembly, we chose to try to write logic that can be auto-vectorized by the compiler when processing large arrays of cells. We used <https://godbolt.org> to view the generated machine code and make sure that our implementation supports AVX512. To enable auto-vectorization for AVX512 instructions, we used the rust flags shown in [Listing 4](#).

#### *Listing 4. Rust flags for enabling AVX512*

```

-C opt-level=3
-C target-
feature=+avx2,+avx,+sse2,+avx512vl,+avx512f,+avx512bw,+avx512cd,+avx512dq,+avx512vnni

```

CONSOLE

The basic function we use for processing the core of a row is shown in [Listing 5](#). It takes read-only references to the current row and the row above and below it. It then iterates through them using sliding windows so all six neighbors of the current cell are available at once. The new state of each cell is calculated and placed into the cell of the result array.

*Listing 5. Function for processing the core of a row without collisions*

```

1 pub fn process_core<const WIDTH: usize>(
2     above: &[u8; WIDTH - 1],
3     current: &[u8; WIDTH],
4     below: &[u8; WIDTH - 1],
5     result: &mut [u8; WIDTH - 2]
6 ) {
7     let context_iterator = above
8         .array_windows::<2>()
9         .zip(current.array_windows::<3>())
10        .zip(below.array_windows::<2>())
11        .zip(result.iter_mut());
12
13    context_iterator.for_each(
14        |(
15            ([north_west, north_east], [west, _current, east]), [south_west, south_e
16            ast]),
17            result,
18        )| {
19            let new_cell = (west & 0b00100000)
20                | (north_west & 0b00010000)
21                | (north_east & 0b00001000)
22                | (east & 0b00000100)
23                | (south_east & 0b00000010)
24                | (south_west & 0b00000001);
25            *result = new_cell;
26            // Collision handling omitted for now
27        },
28    )
29 }

```

**Listing 6** shows the generated assembly for the function. We can see that the compiler autovectorized the function and used the AVX512 instructions to process 64 cells at once. At first, `vpbroadcastb` is used to populate the 512-bit vector registers `zmm0 - zmm5` with the mask for each direction. The result will be placed into the `zmm6` register. The main loop first uses `vpandq` to set each cell in the result register to the value of the north-western cell masked by a mask that removes all particles except the one moving from northwest to southeast. This also resets the value of all other particles in `zmm6` to zero. Then `vpternlogq` is used with each of the five remaining directions to add the particles that move into the current cell from that direction. After that, `zmm6` is stored back into memory. This way, 64 cells are processed at once. If there are more than 64 cells left, the above is repeated. Otherwise, the program uses non-vectorized logic for the remaining cells.



Listing 6. Assembly for `process_core` without collisions

X86ASM

```

1  .LCPI0_12:
2      .byte    16
3  .LCPI0_13:
4      .byte    32
5  .LCPI0_14:
6      .byte    8
7  .LCPI0_15:
8      .byte    4
9  .LCPI0_16:
10     .byte    2
11  .LCPI0_17:
12     .byte    1
13
14  process_core:
15     vpbroadcastb    zmm0, byte ptr [rip + .LCPI0_12]
16     vpbroadcastb    zmm1, byte ptr [rip + .LCPI0_13]
17     vpbroadcastb    zmm2, byte ptr [rip + .LCPI0_14]
18     vpbroadcastb    zmm3, byte ptr [rip + .LCPI0_15]
19     vpbroadcastb    zmm4, byte ptr [rip + .LCPI0_16]
20     vpbroadcastb    zmm5, byte ptr [rip + .LCPI0_17]
21  .LBB0_1:
22     vpandq    zmm6, zmm0, zmmword ptr [rdi + rax]
23     vpternlogq    zmm6, zmm1, zmmword ptr [rsi + rax], 248
24     vpternlogq    zmm6, zmm2, zmmword ptr [rdi + rax + 1], 248
25     vpternlogq    zmm6, zmm3, zmmword ptr [rsi + rax + 2], 248
26     vpternlogq    zmm6, zmm4, zmmword ptr [rdx + rax + 1], 248
27     vpternlogq    zmm6, zmm5, zmmword ptr [rdx + rax], 248
28     vmovdqu64    zmmword ptr [rcx + rax], zmm6
29     add    rax, 64
30     cmp    rax, 9984
31     jne    .LBB0_1
32     ... ; Non-vectorized logic for the remaining cells (lt 64)
33     ret

```

This implementation is correct, but it does not handle collisions yet. The trivial way of adding the collision handling we defined above breaks auto-vectorization. We will need to find a different way to implement collisions. It is probably impossible to auto-vectorize the collision handling using true randomness; it has much conditional code that is difficult to vectorize. We tried a few different implementations, like using a simple match statement without randomness (Listing 13), a static lookup table (Listing 14) and if-else statements (Listing 15), but they all broke auto-vectorization. However, we found that if we were using multiple if statements, the compiler could transform them into vectorized assembly. Listing 7 shows the implementation we ended up with. It needs to reassign `new_cell` sometimes because autovectorization breaks otherwise. We are not entirely sure why this is the case.

**NOTE**

The lookup table solution for collisions might be possible when the CPU supports [AVX512-VBMI](#), but we did not have access to such a CPU.

*Listing 7. Non-random collision calculation using if statements*

```

1 // Three particle collisions
2 if new_cell == 0b00101010 {
3     *result = 0b00010101;
4 }
5 if new_cell == 0b00010101 {
6     *result = 0b00101010;
7 }
8 new_cell = *result;
9
10 // Two particle collisions
11 if new_cell == 0b00100100 {
12     *result = 0b00010010;
13 }
14 if new_cell == 0b00010010 {
15     *result = 0b00001001;
16 }
17 if new_cell == 0b00001001 {
18     *result = 0b00100100;
19 }
20 new_cell = *result;
21
22 // Four particle collisions
23 if new_cell == 0b00110110 {
24     *result = 0b00011011;
25 }
26 if new_cell == 0b00101101 {
27     *result = 0b00110110;
28 }
29 if new_cell == 0b00011011 {
30     *result = 0b00101101;
31 }

```

RUST

If we add that logic to the `process_core` function from above, the generated

The assembly for the `process_core` function with the collision implementation from [Listing 7](#) is shown in [Listing 8](#)

*Listing 8. Assembly for `process_core` with non-random collisions*

```

1 .LCPI0_28:
2     .byte    16
3 .LCPI0_29:
4     .byte    32
5 .LCPI0_30:
6     .byte    8

```

X86ASM

```

7  .LCPIO_31:
8      .byte 4
9  .LCPIO_32:
10     .byte 2
11 .LCPIO_33:
12     .byte 1
13 .LCPIO_34:
14     .byte 42
15 .LCPIO_35:
16     .byte 21
17 .LCPIO_36:
18     .byte 36
19 .LCPIO_37:
20     .byte 18
21 .LCPIO_38:
22     .byte 9
23 .LCPIO_39:
24     .byte 54
25 .LCPIO_40:
26     .byte 27
27 .LCPIO_41:
28     .byte 45
29
30 process_core:
31     mov     r11, rcx
32     xor     eax, eax
33     ; Masks for each direction
34     vpbroadcastb    zmm0, byte ptr [rip + .LCPIO_28]
35     vpbroadcastb    zmm1, byte ptr [rip + .LCPIO_29]
36     vpbroadcastb    zmm2, byte ptr [rip + .LCPIO_30]
37     vpbroadcastb    zmm3, byte ptr [rip + .LCPIO_31]
38     vpbroadcastb    zmm4, byte ptr [rip + .LCPIO_32]
39     vpbroadcastb    zmm5, byte ptr [rip + .LCPIO_33]
40     ; Masks for each collision case
41     vpbroadcastb    zmm6, byte ptr [rip + .LCPIO_34]
42     vpbroadcastb    zmm7, byte ptr [rip + .LCPIO_35]
43     vpbroadcastb    zmm8, byte ptr [rip + .LCPIO_36]
44     vpbroadcastb    zmm9, byte ptr [rip + .LCPIO_37]
45     vpbroadcastb    zmm10, byte ptr [rip + .LCPIO_38]
46     vpbroadcastb    zmm11, byte ptr [rip + .LCPIO_39]
47     vpbroadcastb    zmm12, byte ptr [rip + .LCPIO_40]
48     vpbroadcastb    zmm13, byte ptr [rip + .LCPIO_41]
49 .LBB0_1:
50     vpandq    zmm14, zmm0, zmmword ptr [rdi + rax]
51     vpternlogq    zmm14, zmm1, zmmword ptr [rsi + rax], 248
52     vpternlogq    zmm14, zmm2, zmmword ptr [rdi + rax + 1], 248
53     vpternlogq    zmm14, zmm3, zmmword ptr [rsi + rax + 2], 248
54     vpternlogq    zmm14, zmm4, zmmword ptr [rdx + rax + 1], 248

```

```

55     vpternlogq    zmm14, zmm5, zmmword ptr [rdx + rax], 248
56     ; Process collisions
57     vpcmpeqb     k1, zmm14, zmm6
58     vpcmpeqb     k2, zmm14, zmm7
59     vmovdqu8     zmm14 {k1}, zmm7
60     vmovdqu8     zmm14 {k2}, zmm6
61     vpcmpeqb     k1, zmm14, zmm8
62     vpblendmb    zmm15 {k1}, zmm14, zmm9
63     vpcmpeqb     k1, zmm14, zmm9
64     vmovdqu8     zmm15 {k1}, zmm10
65     vpcmpeqb     k1, zmm14, zmm10
66     vmovdqu8     zmm15 {k1}, zmm8
67     vpcmpeqb     k1, zmm15, zmm11
68     vpblendmb    zmm14 {k1}, zmm15, zmm12
69     vpcmpeqb     k1, zmm15, zmm13
70     vmovdqu8     zmm14 {k1}, zmm11
71     vpcmpeqb     k1, zmm15, zmm12
72     vmovdqu8     zmm14 {k1}, zmm13
73     vmovdqu64    zmmword ptr [r11 + rax], zmm14
74     add         rax, 64
75     cmp         rax, 9984
76     jne         .LBB0_1
77     ... ; Non-vectorized logic for the remaining cells (lt 64)
78     ret

```

This implementation always turns to the right in collisions. Implementing true randomness without sacrificing performance is not trivial. As a workaround, we adjusted our algorithm to always turn in the same direction, except when there is a particle, then we turn left. This way, two colliding particle streams will scatter, alternating between turning left and right. This behavior is not perfect, but it should be good enough for our purposes. We will revisit this when we have implemented it better and may be able to improve it then. The final implementation of `process_core` is shown in [Listing 9](#).

*Listing 9. Function for processing the core of a row with fake collisions*

```

1 pub fn process_core<const WIDTH: usize>(
2     above: &[u8; WIDTH - 1],
3     current: &[u8; WIDTH],
4     below: &[u8; WIDTH - 1],
5     result: &mut [u8; WIDTH - 2]
6 ) {
7     let context_iterator = above
8         .array_windows::<2>()
9         .zip(current.array_windows::<3>())
10        .zip(below.array_windows::<2>())
11        .zip(result.iter_mut());
12
13    context_iterator.for_each(
14        |(
15            ([north_west, north_east], [west, _current, east]), [south_west, south_e
ast]),

```

```

16         result,
17     )| {
18         let new_cell = (west & 0b00100000)
19             | (north_west & 0b00010000)
20             | (north_east & 0b00001000)
21             | (east & 0b00000100)
22             | (south_east & 0b00000010)
23             | (south_west & 0b00000001);
24         *result = new_cell;
25
26         if (new_cell == 0b00100100) && (south_east & 0b00010000 == 0) {
27             *result = 0b00010010;
28         }
29         if (new_cell == 0b00100100) && (south_east & 0b00010000 != 0) {
30             *result = 0b00001001;
31         }
32         if (new_cell == 0b00011011) && (south_east & 0b00010000 == 0) {
33             *result = 0b00101101;
34         }
35         if (new_cell == 0b00011011) && (south_east & 0b00010000 != 0) {
36             *result = 0b00110110;
37         }
38
39         if new_cell == 0b00010010 && (east & 0b00001000 == 0) {
40             *result = 0b00001001;
41         }
42         if new_cell == 0b00010010 && (east & 0b00001000 != 0) {
43             *result = 0b00100100;
44         }
45         if new_cell == 0b00101101 && (east & 0b00001000 == 0) {
46             *result = 0b00110110;
47         }
48         if new_cell == 0b00101101 && (east & 0b00001000 != 0) {
49             *result = 0b00011011;
50         }
51
52         if new_cell == 0b00001001 && (north_east & 0b00000100 == 0) {
53             *result = 0b00100100;
54         }
55         if new_cell == 0b00001001 && (north_east & 0b00000100 != 0) {
56             *result = 0b00010010;
57         }
58         if new_cell == 0b00110110 && (north_east & 0b00000100 == 0) {
59             *result = 0b00011011;
60         }
61         if new_cell == 0b00110110 && (north_east & 0b00000100 != 0) {
62             *result = 0b00101101;
63         }

```

```

64
65         if new_cell == 0b00101010 {
66             *result = 0b00010101;
67         }
68         if new_cell == 0b00010101 {
69             *result = 0b00101010;
70         }
71         new_cell = *result; // This line does nothing, but auto-
vectorization breaks without it
72     },
73 )
74 }

```

For collisions on the borders of the grid, we use a less optimized implementation which is just calling the `proccess_collision` function ([Listing 3](#)) for each cell. Only a small portion of the cells will be on the border of the grid, so this should not have a big impact on performance.

We measured the average duration for one cell for a 2000x2000 grid with 5000 rounds and got the following results:

- With real random collisions: 3.76415ns per cell
- With AVX fake random collisions: 0.22619ns per cell
- With AVX fake random collisions but with vectorization disabled: 2.50546ns per cell
- Without collisions: 0.06601ns per cell

The results show that our vectorized implementation has a 10x speedup over the non-vectorized implementation. Our test machine only supports 256-bit vector operations, so the results above only reflect that. We expect the performance to be better on machines that support AVX512.

#### MAIN LOOP

The core of the simulation is based around two grids, `grid_a` contains the current state, and `grid_b` contains the next state. We fill the top row, the core rows, and the bottom row of `grid_b` in separate functions. After `grid_b` is done, it is swapped with `grid_a`, and the process is repeated. The main loop is shown in [Listing 10](#).

*Listing 10. One step of the main loop*

```

1 movement_top_row(&grid_a[0], &grid_a[1], &mut grid_b[0]);
2
3 grid_a
4     .windows(3)
5     .zip(grid_b.iter_mut().skip(1))
6     .enumerate()
7     .for_each(|(row_index, (context, result))| {
8         let above = &context[0];

```

RUST

```

 9      let current = &context[1];
10      let below = &context[2];
11      if ((row_index + 1) % 2) == 0 {
12          movement_even_row(above, current, below, result);
13      } else {
14          movement_odd_row(above, current, below, result);
15      }
16  });
17
18 movement_bottom_row(&grid_a[WIDTH-2], &grid_a[WIDTH-1], &mut grid_b[WIDTH-1]);
19
20 std::mem::swap(&mut grid_a, &mut grid_b);

```

#### SHARED MEMORY PARALLELISM WITH RAYON

We will use rayon to divide the processing of the core rows over multiple threads. `par_windows` can be used instead of `windows`, and `par_iter_mut` can be used instead of `iter_mut` to process the rows in parallel. The main loop with rayon is shown in [Listing 11](#).

*Listing 11. One round using rayon*

```

1 movement_top_row(&grid_a[0], &grid_a[1], &mut grid_b[0]);
2
3 grid_a
4   .par_windows(3)
5   .zip(grid_b.par_iter_mut().skip(1))
6   .enumerate()
7   .for_each(|(row_index, (context, result))| {
8       let above = &context[0];
9       let current = &context[1];
10      let below = &context[2];
11      if ((row_index + 1) % 2) == 0 {
12          movement_even_row(above, current, below, result);
13      } else {
14          movement_odd_row(above, current, below, result);
15      }
16  });
17
18 movement_bottom_row(&grid_a[WIDTH-2], &grid_a[WIDTH-1], &mut grid_b[WIDTH-1]);
19
20 std::mem::swap(&mut grid_a, &mut grid_b);

```

#### *Multithreading measurements*

- Without rayon: 0.21101ns per cell
- With 1 thread: 0.21776ns per cell
- With 2 threads: 0.11265ns per cell
- With 4 threads: 0.060394ns per cell

- With 6 threads: 0.046102ns per cell
- With 7 threads: 0.042037ns per cell
- With 8 threads: 0.040239ns per cell
- With 9 threads: 0.040608ns per cell
- With 10 threads: 0.040653ns per cell
- With 10 threads: 0.038776ns per cell
- With 16 threads: 0.038706ns per cell
- With 20 threads: 0.050061ns per cell

We made measurements for a 2000x2000 grid with 5000 rounds. An intel i9-11950H was used for testing. Based on the Willow Cove microarchitecture, it has eight cpus and 16 cores. The measurements show that the performance is best with eight cores. I am not sure why this is the case. I would like to think that a single core keeps the shared vector units on its CPU busy all the time because there are nearly no branches, and most instructions in the main loop are vector instructions. However, it may as well be that there is a memory bottleneck.

Afterward, we discovered that our processor actually supports the required [AVX512](#) instructions. We redid the eight threads measurement and got a performance of 0.034543ns per cell.

#### DISTRIBUTED MEMORY PARALLELISM WITH MPI

Our project uses MPI via the `rsmpi` Rust bindings for distributed memory parallelism. Every MPI node will process a consecutive chunk of rows. For example, if we have two nodes and eight rows, node 0 will process rows 0-3, and node 1 will process rows 4-7. Our implementation is shown in [Listing 12](#). Before each round, every node sends its top row to the previous node and its bottom row to the next node. This way, all nodes that are not at the top or bottom have the context to treat their top/bottom rows like core rows. The first and last nodes still process their rows like core rows.

*Listing 12. One round using MPI and rayon*

```

1 mpi::request::scope(|scope| {
2     let mut guards = Vec::new();
3
4     if let Some(previous_rank) = previous_rank {
5         guards.push(WaitGuard::from(
6             world
7                 .process_at_rank(previous_rank)
8                 .immediate_send(scope, unsafe {
9                     std::mem::transmute:::<&mut [Cell; WIDTH], &mut [u8; WIDTH]>(&mut grid_a[0],
10                                     )
11         })),
12     },

```

RUST



```

13     ));
14     guards.push(WaitGuard::from(
15         world.process_at_rank(previous_rank).immediate_receive_into(
16             scope,
17             unsafe {
18                 std::mem::transmute:::<&mut [Cell; WIDTH], &mut [u8; WIDTH]>(
19                     receive_top,
20                 )
21             },
22         ),
23     ));
24 }
25
26 if let Some(next_rank) = next_rank {
27     guards.push(WaitGuard::from(
28         world
29             .process_at_rank(next_rank)
30             .immediate_send(scope, unsafe {
31                 std::mem::transmute:::<&mut [Cell; WIDTH], &mut [u8; WIDTH]>(
32                     &mut grid_a[height - 1],
33                 )
34             },
35     ));
36     guards.push(WaitGuard::from(
37         world
38             .process_at_rank(next_rank)
39             .immediate_receive_into(scope, unsafe {
40                 std::mem::transmute:::<&mut [Cell; WIDTH], &mut [u8; WIDTH]>(
41                     receive_bottom,
42                 )
43             },
44     ));
45 }
46 });
47
48
49 if previous_rank.is_some() {
50     movement_even_row(receive_top, &grid_a[0], &grid_a[1], &mut grid_b[0]);
51 } else {
52     movement_top_row(&grid_a[0], &grid_a[1], &mut grid_b[0]);
53 }
54
55 grid_a
56     .par_windows(3)
57     .zip(grid_b.par_iter_mut().skip(1))
58     .enumerate()
59     .for_each(|(row_index, (context, result))| {
60         let above = &context[0];

```

```

61     let current = &context[1];
62     let below = &context[2];
63     if ((row_index + 1) % 2) == 0 {
64         movement_even_row(above, current, below, result);
65     } else {
66         movement_odd_row(above, current, below, result);
67     }
68 });
69
70 if next_rank.is_some() {
71     movement_odd_row(&grid_a[height-2], &grid_a[height-
72     1], receive_bottom, &mut grid_b[height-1]);
73 } else {
74     movement_bottom_row(&grid_a[height-2], &grid_a[height-1], &mut grid_b[height-1]);
75 }
76 std::mem::swap(&mut grid_a, &mut grid_b);

```

Using this approach, we can scale our simulation to multiple nodes. We tested our implementation on the Virgo cluster.

## MEASUREMENTS

---

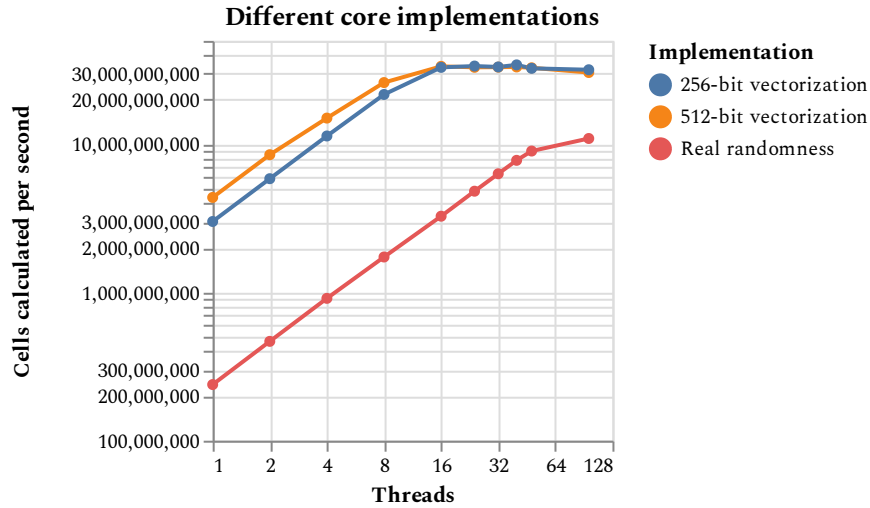
The Virgo cluster has nodes with Intel Xeon Gold 6248 CPUs, which support AVX512. Each Node has two Sockets, each with 24 CPUs and 48 cores. The measurements always reserved a whole node because the measurements were really unpredictable sometimes when sharing a node with other tasks. We think that this is caused by some other job that does not use many cores but completely saturates the memory bandwidth of the processors, but we did not investigate this further. We thought using a whole node exclusively for measurements would solve this problem. It did not, but we ran enough measurements to get some good ones.

The main metric we will focus on during our measurements is the time per cell. This metric is calculated by dividing the total runtime by the total number of cells calculated ( $\text{width} * \text{height} * \text{rounds}$ ). The runtime measurement starts after the grid is set up and ends after the last calculation is done, so it does not include memory allocation, program startup, or placing the initial cells into the grid. This metric is mostly independent of the actual test size and lets us compare different measurements with different grid sizes. It should be noted that the actual time per cell depends on the grid size because the processing of border cells is far slower than the processing of core cells.

If otherwise specified, measurements were done with a 10000x10000 cell grid over 1000 rounds. When starting a benchmark job, it tests if all nodes are working correctly and aborts if they are not.

### COMPARING DIFFERENT CORE IMPLEMENTATIONS

We wanted to know if AVX512 (512bit vector operations) brings improvement over using AVX2 (256bit vector operations). We also wanted to know if our focus on vectorization at the cost of true randomness was bringing the expected performance improvements. The results are shown in [Figure 3](#).



*Figure 3. Comparison of different implementations for calculating core cells*

The vectorized implementations process roughly ten times more cells per second than the real-random implementation. Of the vectorized implementations, the `AVX512` implementation is about 30% faster than the `AVX2` implementation. This is a bit less than we expected, but it is still a significant improvement. The non-vectorized implementation scales linearly with the number of threads, but from 48 to 96 threads, the performance only improves by 20%. This is because the node has only 48 CPUs. SMT is used when we use more threads and resources are shared between two cores on the same CPU.

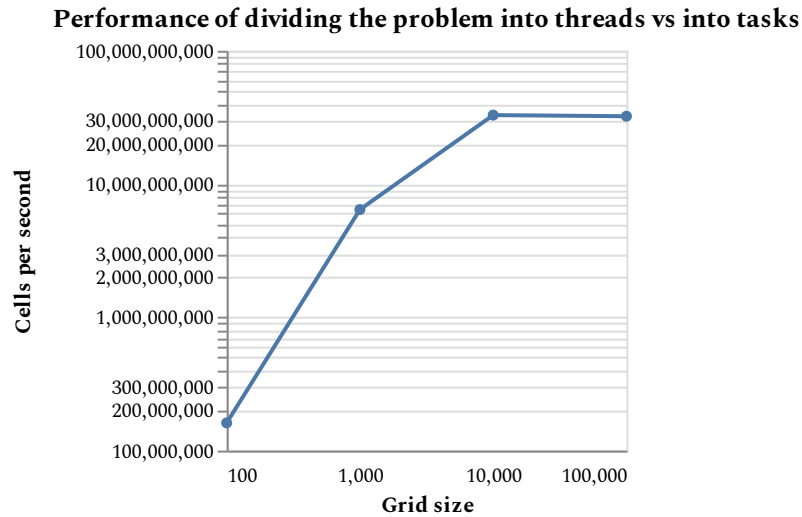
Both vectorized implementations scale about linearly with the number of cores but hit a plateau when using more than 16 cores. With more cores, the speed stays at around 35 Gigacells per second. A memory throughput bottleneck probably causes this. As one cell is 1 byte in size, the CPU is loading and storing at 35 GB/s. The theoretical maximum memory bandwidth per CPU is 143.36 GB/s (6 channels of 2933 MHz memory). The node has two sockets, so the maximum bandwidth is around 281 GB/s. Our actual maximum is about an eighth of the theoretical maximum. This is probably because we did not optimize memory access patterns to keep the code simple. We could analyze the number of page faults to improve it and ensure we are aligned with the cache lines. We could also improve caching, as our cores are split over two sockets, so they are not sharing the L3 cache. We could use Slurm to make sure every cores is on the same socket and run two tasks per node

Because the `AVX512` implementation is the fastest, the rest of the measurements will be done with that implementation.

We will also use 48 threads to avoid SMT. It should not make a big difference as we reach the memory bandwidth bottleneck at that number of threads anyway.

#### DIFFERENT GRID SIZES

We assume that the grid size significantly impacts the performance because bigger grids allow for more sequential processing of cells. We measured the performance for different grid sizes and got the results shown in [Figure 4](#).



*Figure 4. Different grid sizes compared to performance*

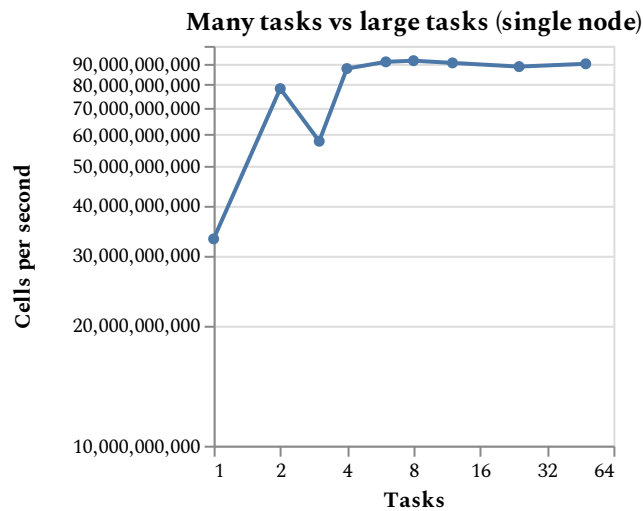
For a small 100x100 grid, we process a lot fewer cells per second than for large grids. This is probably because the core of the row is so small that the overhead is a lot bigger than the actual work.

For bigger grids, the performance increases, but it stagnates after 10000x10000 cells. We assume that, at this point, the overhead is negligible compared to the time spent doing the vectorized calculations.

For the rest of the measurements, we will use a 10000x10000 grid.

#### THREADS VS TASKS

We want to find out how running multiple openMP tasks with a low number of threads compares to running a single openMP task with a high number of threads. The results for a single node are shown in [Figure 5](#).



*Figure 5. Threads compared to tasks for a single node*

The chart does not show what we had expected initially. I would have thought that the performance would be the same for 1 task with 48 threads and 1 thread with 48 tasks. However, the performance is usually at 90 Gigacells per second, except for 3 tasks with 16 threads and 1 task with 48 threads. This could be caused by the scheduler having to split single tasks over the two available sockets for these two combinations. If a task is split among two sockets, it may have issues with caching. It seems like tasks with a large number of threads. It would also explain the split at 16 threads, as two tasks with 16 threads are not split and would run at 90, and one task is split over the two sockets and runs at 30. Taking the average  $(90 + 90 + 30) / 3$  results in 70, which is close to the actual result of 60.

#### SCALING OVER MULTIPLE NODES

We want to see how big the performance impact of distributing over multiple nodes is. We measured the performance for 1, 2, and 4 nodes. [Figure 6](#) shows the results for different numbers of tasks per node. When a node runs a low number of tasks, the tasks have more threads. In total, there are always 48 threads per node.

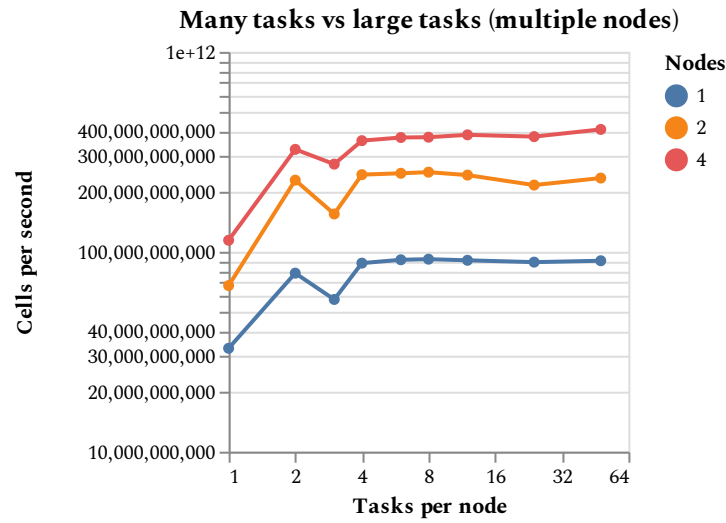


Figure 6. Performance in relation to the number of nodes

The chart shows that the performance roughly doubles when the number of nodes is doubled. The step from one to two nodes seems to be even more extreme, as it nearly triples in value. We suspect that we got a particularly good node for that measurement. Because we always use the whole node exclusively, we must wait a long time until our jobs are run. We did ten runs for the tests with one node and 5 runs with 2 and 4 nodes.

We analyzed the speedup and efficiency of distributing over multiple machines in Figure 7. The speedup is the performance gain compared to a single node. Efficiency is the performance gain in relation to the resources used.

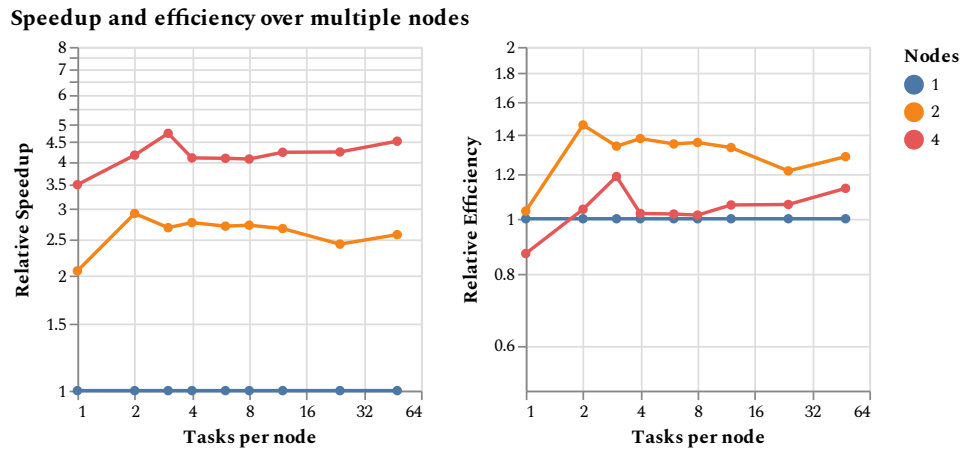


Figure 7. Speedup and efficiency

It is interesting to see that the efficiency for 2 and 4 nodes is actually higher than for 1 node. Again that is probably caused by the fact that we got particularly good nodes for the 2 measurements. For 4 nodes, the efficiency is lower than for 2 nodes, but it is still lower than for 1 node. If all nodes were

the same speed, we would expect the efficiency to be a tiny bit lower than one because of the communication overhead. However, as we are not transmitting a lot of data compared to the amount of work and only need to sync between rounds, it should not be much compared to running everything locally.



## RESULTS

---

We were able to implement a fast and scalable simulation of the LGCA. We were able to scale the simulation over multiple nodes and got a near-linear speedup for 4 nodes. We were also able to scale the simulation over multiple threads, where we achieved linear speedup as well. Our fastest run (4 nodes, 48 tasks per node) was able to **process 492 Gigabytes per second** of data, which is nearly half the theoretical maximum set by the memory bandwidth of 1126GB/s ( $140\text{GB/s} * 2 \text{ sockets} * 4 \text{ nodes}$ ).

## LIST OF ABBREVIATIONS

---

### **LGCA**

Lattice gas cellular automaton [!\[\]\(3dfb8d66e81160ad61421a3452093d1b\_img.jpg\)](#)

### **AVX512**

Advanced Vector Extensions 512 for x86 [!\[\]\(0f848bbd71cef6b345273b16f905912a\_img.jpg\)](#)

### **SMT**

Simultaneous multithreading [!\[\]\(a870788d6ed9b8fd294b7654a8c8526b\_img.jpg\)](#)

### **SIMD**

Single instruction, multiple data [!\[\]\(3211b5d1d968fc1665909b34f9f16010\_img.jpg\)](#)

## APPENDIX

---

*Listing 13. Collision calculation using match without randomness*

```
1 pub fn process_collision(cell: u8) -> u8 {
2     match cell {
3         // Two opposing particles
4         0b00001001 => 0b00100100,
5         0b00010010 => 0b00001001,
6         0b00100100 => 0b00010010,
7
8         // Three particles
9         0b00101010 => 0b00010101,
10        0b00010101 => 0b00101010,
11
12        // Four particles with opposing holes
13        0b00110110 => 0b00011011,
14        0b00011011 => 0b00101101,
15        0b00101101 => 0b00011011,
16
17        // Everything else
18        x => x,
19    }
20 }
```

RUST

*Listing 14. Collision calculation using a lookup table*

```
1 const COLLISION_TABLE: [u8; 64] = [
2     0b0000_0000,
3     0b0000_0001,
4     0b0000_0010,
5     0b0000_0011,
6     0b0000_0100,
7     0b0000_0101,
8     0b0000_0110,
9     0b0000_0111,
10    0b0000_1000,
11    0b0001_0010,
12    0b0000_1010,
13    0b0000_1011,
14    0b0000_1100,
15    0b0000_1101,
16    0b0000_1110,
17    0b0000_1111,
18    0b0001_0000,
19    0b0001_0001,
20    0b0010_0010,
21    0b0001_0011,
22    0b0001_0100,
23    0b0010_1010,
```

RUST

```

24     0b0001_0110,
25     0b0001_0111,
26     0b0001_1000,
27     0b0001_1001,
28     0b0001_1010,
29     0b0010_1101,
30     0b0001_1100,
31     0b0001_1101,
32     0b0001_1110,
33     0b0001_1111,
34     0b0010_0000,
35     0b0010_0001,
36     0b0010_0010,
37     0b0010_0011,
38     0b0000_1001,
39     0b0010_0101,
40     0b0010_0110,
41     0b0010_0111,
42     0b0010_1000,
43     0b0010_1001,
44     0b0001_0101,
45     0b0010_1011,
46     0b0010_1100,
47     0b0011_0110,
48     0b0010_1110,
49     0b0010_1111,
50     0b0011_0000,
51     0b0011_0001,
52     0b0011_0010,
53     0b0011_0011,
54     0b0011_0100,
55     0b0011_0101,
56     0b0001_1011,
57     0b0011_0111,
58     0b0011_1000,
59     0b0011_1001,
60     0b0011_1010,
61     0b0011_1011,
62     0b0011_1100,
63     0b0011_1101,
64     0b0011_1110,
65     0b0011_1111,
66 ];
67
68 pub fn process_collision(cell: u8) -> u8 {
69     COLLISION_TABLE[cell as usize]
70 }

```

*Listing 15. Collision calculation using if-else statements*

```
1  if *result == 0b00101010 {
2      *result = 0b00010101;
3  } else if *result == 0b00010101 {
4      *result = 0b00101010;
5  } else if *result == 0b00100100 {
6      *result = 0b00010010;
7  }
8
9  if *result == 0b00010010 {
10     *result = 0b00001001;
11 } else if *result == 0b00001001 {
12     *result = 0b00100100;
13 }
14
15 if *result == 0b00110110 {
16     *result = 0b00011011;
17 } else if *result == 0b00101101 {
18     *result = 0b00110110;
19 } else if *result == 0b00011011 {
20     *result = 0b00101101;
21 }
```

RUST