

Case Study: Optimized Code for Neural Cell Simulations

Submitted by [Daniel Vea](https://software.intel.com/en-us/user/1234549) (<https://software.intel.com/en-us/user/1234549>) on January 4, 2016

[f Share](#) (<https://www.facebook.com/sharer/sharer.php?u=https://software.intel.com/en-us/articles/case-study-optimized-code-for-neural-cell-simulation>)

[Tweet](#) (<https://twitter.com/intent/tweet?text=Case+Study%3A+Optimized+Code+for+Neural+Cell+Simulations%3A&url=https%3A%2F%2Fsoftware.intel.com%2Fen-us%2Farticles%2Fc...>)

[g+Share](#) (<https://plus.google.com/share?url=https://software.intel.com/en-us/articles/case-study-optimized-code-for-neural-cell-simulations>)

About

Intel held the [Intel® Modern Code Developer Challenge](https://moderncodechallenge.intel.com/prizes/) (<https://moderncodechallenge.intel.com/prizes/>) that had about 2,000 students from 130 universities in 19 countries registered to participate in the Challenge. They were provided access to Intel® Xeon® processors and Intel® Xeon Phi™ coprocessors to optimize code used in a CERN openlab brain simulation research project. The goal of the research project is to find treatments and cures for neurological disorders, such as schizophrenia, epilepsy, and autism. The contestants task was to look at the code for cell clustering and 3D movement and then modify the algorithms for parallel performance by optimizing the code to reduce the runtime, all while maintaining correctness.

In this article Daniel Vea Falguera (one of the Challenge winners) shares the original code as well as the **optimized code** (Note: **Changed code line(s):**) and describes many of the optimizations he implemented. In some cases, the optimizations did not work, but he gives insight into how other changes to the code would work.

[Download \(196.73 KB\)](#) (<https://software.intel.com/sites/default/files/managed/fc/05/Case%20Study%20-Optimization%20Code.pdf>)

[Jetzt herunterladen](#) (<https://software.intel.com/sites/default/files/managed/fc/05/Case%20Study%20-Optimization%20Code.pdf>)

INCLUDES

Original Code

```
1 #include <cstring>
2 #include <cstdlib>
3 #include <ctime>
4 #include <cmath>
5 #include <getopt.h>
6 #include "util.hpp"
```

Optimized Code

```
01 #include <cstring>
02 #include <cstdlib>
03 #include <ctime>
04 #include <cmath>
05 #include <omp.h>
06 #include <getopt.h>
07 #include "util.hpp"
08 #include <malloc.h>      //Useless
09 #include <mkl.h>         //Useless
10 #include <cilk/cilk.h> //Useless
```

Changed code line: 5

Optimization Notes

The addition of #include <omp.h> is required to run OpenMP*. The OpenMP functions and clauses are used in the optimized code.

In the completed code, malloc.h, mkl.h, and cilk.h were included during the development process to optimize memory blocks (among other things), but they didn't show any improvements. They are included to show that there was an attempt to use them. More information on the Intel® Math Kernel Library (Intel® MKL) and Intel® Cilk™ Plus is given in the Other Optimizations section.

RandomFloatPos

The function RandomFloatPos() is used to generate a random number three times to generate a random 3D position inside the cellMovementAndDuplication function.

NOTE: Only the computational related functions are referenced in this paper, so functions such as “menu printings” and others aren't discussed.

Original Code

```
1 static float RandomFloatPos() {
2     // returns a random number between a given minimum and maximum
3     float random = ((float) rand()) / (float) RAND_MAX;
4     float a = 0;
5     float r = random;
6     return a + r;
```

7 }

Changed code line(s): 3, 4, 5, 6

Optimized Code

```

1 static void RandomFloatPos(float input[3],unsigned sedin) {
2     // returns a random number between a given minimum and maximum
3     __assume_aligned((float*)input, 64);
4     unsigned int seed=sedin,i;
5     for(i=0;i<3;++i){
6         input[i]=((float)rand_r(&seed))/(float)(RAND_MAX))-0.5;
7     }
8 }
```

Changed code line: 1

Optimization Notes

The original function returns a random number doing some needless steps (for instance, adding zero to the random number). The original code can't be parallelized directly due to the function rand(), which can only be executed one thread at a time. The parallelization problem was fixed using the rand_r() function. The rand_r() function allows each thread to execute at the same time.

The function RandomFloatPos() is called to generate a random 3D position minus a constant offset of 0.5, so it can be simplified as a for loop with three iterations. The optimized code does exactly this. The offset is used to make the position values range from -0.5 to 0.5, placing the (0,0,0) point in the middle of the space of cell movement.

The optimized function has two parameters: input[3] and sedin. Input[3] returns the values of the 3D random generated position and is memory aligned(64). The sedin is passed as a parameter and contains the value of the seed generated for each call to this function.

The for loop generates the 3D coordinate and stores it into input.

getNorm

Original Code

```

01 static float getNorm(float* currArray) {
02     // computes L2 norm of input array
03     int c;
04     float arraySum=0;
05     for (c=0; c<3; c++) {
06         arraySum += currArray[c]*currArray[c];
07     }
08     float res = sqrt(arraySum);
09     return res;
10 }
```

Changed code line(s): 3, 6, 8, 9

Optimized Code

```

1 static float getNorm(float* currArray) {
2     // computes L2 norm of input array
3     float arraySum=0;
4     for (int c=0; c<3; ++c) {
5         arraySum += pow(currArray[c],2);
6     }
7     return sqrt(arraySum);
8 }
```

Changed code line(s): 5, 7

Optimization Notes

This computes the norm from a given array of numbers by the input float currArray. There are two ways that the code was optimized for this function:

- The res variable was removed.
- The pow() function was added.
Extensive research done afterwards has shown that the pow() is not a significant improvement. Further, given the Intel® Xeon Phi™ coprocessor's ability to execute multiple floating point operations per cycle, the original currArray[c]*currArray[c] is provably faster in those cases. The Intel compiler vectorizes the pow() function, but adds more code to do so, resulting in slower runtime speed.

getL2Distance

This function is used to determine the linear distance between two points in 3D space.

Original Code

```

01 static float getL2Distance(float pos1x, float pos1y, float pos1z, float
02 pos2x, float pos2y, float pos2z) {
```

```

03 // returns distance (L2 norm) between two positions in 3D
04 float distArray[3];
05 distArray[0] = pos2x-pos1x;
06 distArray[1] = pos2y-pos1y;
07 distArray[2] = pos2z-pos1z;
08 float l2Norm = getNorm(distArray);
09 return l2Norm;
10 }
```

Changed code line(s): 1, 2, 5, 6, 7, 8, 9

Optimized Code

```

1 static float getL2Distance(float* pos1, float* pos2) {
2     // returns distance (L2 norm) between two positions in 3D
3     float distArray[3] __attribute__((aligned(64)));
4     distArray[0] = pos2[0]-pos1[0];
5     distArray[1] = pos2[1]-pos1[1];
6     distArray[2] = pos2[2]-pos1[2];
7     return getNorm(distArray);
8 }
```

Changed code line(s): 1-8

Optimization Notes

The original function has six inputs that represent two 3D points to calculate the distance between them. The optimized function has only two inputs; each input is an array of three elements representing the 3D point. The variable used inside is aligned.

The optimized function seems to be SIMD executable, but when using the vector notation (for example, $P[0:2]=a[0:2]*b[0:2]$) the execution time was slower. It may be that defining and using an elemental function will further optimize the function.

A white paper on elemental functions is available at:

[\(http://software.intel.com/sites/default/files/article/181418/whitepaperonelementalfunctions.pdf\)](http://software.intel.com/sites/default/files/article/181418/whitepaperonelementalfunctions.pdf)

produceSubstances

This function increases the concentration of substances for each cell position to a maximum limit of 1 unit per cell position.

Original Code

```

01 static void produceSubstances(float**** Conc, float** posAll, int* typesAll, int L, int n){
02
03     produceSubstances_sw.reset();
04     // increases the concentration of substances at the location of the cells
05     float sideLength = 1/(float)L; // length of a side of a diffusion voxel
06     int c, i1, i2, i3;
07     for (c=0; c< n; c++) {
08         i1 = std::min((int)floor(posAll[c][0]/sideLength),(L-1));
09         i2 = std::min((int)floor(posAll[c][1]/sideLength),(L-1));
10         i3 = std::min((int)floor(posAll[c][2]/sideLength),(L-1));
11         if (typesAll[c]==1) {
12             Conc[0][i1][i2][i3]+=0.1;
13             if (Conc[0][i1][i2][i3]>1) {
14                 Conc[0][i1][i2][i3]=1;
15             }
16         } else {
17             Conc[1][i1][i2][i3]+=0.1;
18             if (Conc[1][i1][i2][i3]>1) {
19                 Conc[1][i1][i2][i3]=1;
20             }
21         }
22     }
23     produceSubstances_sw.mark();
24 }
```

Changed code line(s): 1, 5, 6, 8, 9, 10, 11, 13, 14, 17, 18, 19

Optimized Code

```

01 static void produceSubstances(int L, float Conc[2][L][L][L], float posAll[][3], int* typesAll, int n) {
02
03     produceSubstances_sw.reset();
04     // increases the concentration of substances at the location of the cells
05
06     const int auxL=L;
07     --L;
08     int c,i[3] __attribute__((aligned(32))); //i array aligned
09     omp_set_num_threads(240);
10
11 #pragma omp parallel for schedule(static) private(i,c)
12     for (c=0; c< n; ++c) {
13         __assume_aligned((int*)i, 32);
14         __assume_aligned((float*)posAll, 64);
15         i[0] = std::min((int)floor(posAll[c][0]*auxL),L);
16         i[1] = std::min((int)floor(posAll[c][1]*auxL),L);
17         i[2] = std::min((int)floor(posAll[c][2]*auxL),L);
18
19         if (typesAll[c]==1) {
```

```

20         (Conc[0][i[0]][i[1]][i[2]]>0.9)?
21         Conc[0][i[0]][i[1]][i[2]]=1 :
22         Conc[0][i[0]][i[1]][i[2]]+=0.1;
23     } else {
24         (Conc[1][i[0]][i[1]][i[2]]>0.9)?
25         Conc[1][i[0]][i[1]][i[2]]=1 :
26         Conc[1][i[0]][i[1]][i[2]]+=0.1;
27     }
28 }
29 produceSubstances_sw.mark();
30 }
```

Changed code line(s): 1,6-11, 13-17, 20-22, 24-26

Optimization Notes

The optimized function inputs order changed to define in the header of the function the size of the arrays. This way we can avoid the use of pointers and work directly with the arrays so the compiler knows beforehand the size of the elements we pass into the function.

The original code used pointers to initialize the arrays, but since those arrays are static (the length is static and doesn't change) it is easier and faster to declare it directly as arrays with a defined size without using pointers.

The optimized code includes the use of one OpenMP parallel for function with static scheduling to distribute the load to the other cores equally. This is because this function can be executed in parallel without affecting the result, that is, every cycle of the for loop can be executed individually without affecting the next iterations.

During the execution of the main code (described in a later section), this function is called multiple times, each time increasing the value of n. So while it is possible to make the function parallel, this would only be useful while the values of n are low (less than 10000). Otherwise the added overhead may slow the code.

The operation posAll[c][0]/sideLength is the same as (posAll[c][0]/1/L) or (posAll[c][0]*L). Since posAll[c][0]/sideLength generates the same result with fewer calculations, it is a benefit.

The use of the L variable was changed in optimization. The original code uses the operation L-1 three times, using three extra operations. In the optimized version we simply decrement previously the L variable. The auxL constant (used to optimize as a read-only variable is faster than a read/write variable) stores the original value of L and will not change during the execution of this function.

Similarly the method for using the Conc variable was changed to optimize the code. In the original code the if clause will increment the Conc variable by 0.1, then check if the value is greater than 1 and if true, limit its value to 1, making the previous addition useless. The solution was to check if the value of Conc is greater than 0.9, and if true set the value of Conc to 1; if false then increment Conc by 0.1.

runDiffusionStep

This function has two parts. The first part of the function copies the Conc variable to tempConc. The second part iterates through the Conc array checking upper and 3D boundaries.

Original Code

```

01 static void runDiffusionStep(float**** Conc, int L, float D) {
02     runDiffusionStep_sw.reset();
03     // computes the changes in substance concentrations due to diffusion
04
05     int i1,i2,i3, subInd;
06     float tempConc[2][L][L][L];
07     for (i1 = 0; i1 < L; i1++) {
08         for (i2 = 0; i2 < L; i2++) {
09             for (i3 = 0; i3 < L; i3++) {
10                 tempConc[0][i1][i2][i3] = Conc[0][i1][i2][i3];
11                 tempConc[1][i1][i2][i3] = Conc[1][i1][i2][i3];
12             }
13         }
14     }
15
16     int xUp, xDown, yUp, yDown, zUp, zDown;
17
18     for (i1 = 0; i1 < L; i1++) {
19         for (i2 = 0; i2 < L; i2++) {
20             for (i3 = 0; i3 < L; i3++) {
21                 xUp = (i1+1);
22                 xDown = (i1-1);
23                 yUp = (i2+1);
24                 yDown = (i2-1);
25                 zUp = (i3+1);
26                 zDown = (i3-1);
27                 for (subInd = 0; subInd < 2; subInd++) {
28                     if (xUp<L) {
29                         Conc[subInd][i1][i2][i3] += (tempConc[subInd][xUp][i2][i3]-tempConc[subInd][i1][i2][i3]);
30                     }
31                     if (xDown>=0) {
32                         Conc[subInd][i1][i2][i3] += (tempConc[subInd][xDown][i2][i3]-tempConc[subInd][i1][i2][i3]);
33                     }
34                     if (yUp<L) {
35                         Conc[subInd][i1][i2][i3] += (tempConc[subInd][i1][yUp][i3]-tempConc[subInd][i1][i2][i3]);
36                     }
37                     if (yDown>=0) {
38                         Conc[subInd][i1][i2][i3] += (tempConc[subInd][i1][yDown][i3]-tempConc[subInd][i1][i2][i3]);
39                     }
40                     if (zUp<L) {
41                         Conc[subInd][i1][i2][i3] += (tempConc[subInd][i1][i2][zUp]-tempConc[subInd][i1][i2][i3]);
42                     }
43                     if (zDown>=0) {
```

```

44     }
45   }
46 }
47 }
48 }
49 runDiffusionStep_sw.mark();
50 }
51 }
```

Changed code line(s): 1, 5, 9-12, 16, 21-45

Optimized Code

```

01 static void runDiffusionStep(int L, float Conc[2][L][L][L], float D) {
02   runDiffusionStep_sw.reset();
03   // computes the changes in substance concentrations due to diffusion
04   int i1,i2,i3,auxx;
05   const float auxD=D/6;
06   const int auxL=L-1;
07   float tempConc[2][L][L][L] __attribute__((aligned(64)));
08   omp_set_num_threads(240);
09   #pragma omp parallel
10  {
11    #pragma omp for schedule(static) private(i1,i2) collapse(2)
12    for (i1 = 0; i1 < L; ++i1) {
13      for (i2 = 0; i2 < L; ++i2) {
14        memcpy(tempConc[0][i1][i2],
15               Conc[0][i1][i2],sizeof(float)*L);
16        memcpy(tempConc[1][i1][i2],
17               Conc[1][i1][i2],sizeof(float)*L);
18      }
19    }
20
21    #pragma omp for schedule(static) private(i1,i2,i3) collapse(2)
22    for (i1 = 0; i1 < L; ++i1) {
23      for (i2 = 0; i2 < L; ++i2) {
24        Conc[0][i1][i2][0] += (tempConc[0][i1][i2][1]-
25          tempConc[0][i1][i2][0])*auxD;
26        Conc[1][i1][i2][0] += (tempConc[1][i1][i2][1]-
27          tempConc[1][i1][i2][0])*auxD;
28        for (i3 = 1; i3 < auxL; ++i3) {
29          const float aux=tempConc[0][i1][i2][i3];
30          const float aux1=tempConc[1][i1][i2][i3];
31          __assume_aligned((float*)tempConc[0], 64);
32          __assume_aligned((float*)tempConc[1], 64);
33          __assume_aligned((float*)Conc[0], 64);
34          __assume_aligned((float*)Conc[1], 64);
35          if (i1<auxL) {
36            Conc[0][i1][i2][i3] +=
37              (tempConc[0][(i1+1)][i2][i3]-aux)*auxD;
38            Conc[1][i1][i2][i3] +=
39              (tempConc[1][(i1+1)][i2][i3]-aux1)*auxD;
40          }
41          if (i1>0) {
42            Conc[0][i1][i2][i3] +=
43              (tempConc[0][(i1-1)][i2][i3]-aux)*auxD;
44            Conc[1][i1][i2][i3] +=
45              (tempConc[1][(i1-1)][i2][i3]-aux1)*auxD;
46        }
47        if (i2<auxL) {
48          Conc[0][i1][i2][i3] +=
49              (tempConc[0][i1][(i2+1)][i3]-aux)*auxD;
50          Conc[1][i1][i2][i3] +=
51              (tempConc[1][i1][(i2+1)][i3]-aux1)*auxD;
52        }
53        if (i2>0) {
54          Conc[0][i1][i2][i3] +=
55              (tempConc[0][i1][(i2-1)][i3]-aux)*auxD;
56          Conc[1][i1][i2][i3] +=
57              (tempConc[1][i1][(i2-1)][i3]-aux1)*auxD;
58        }
59        Conc[0][i1][i2][i3] +=
60          (tempConc[0][i1][(i3+1)][i3]-aux)*auxD;
61        Conc[1][i1][i2][i3] +=
62          (tempConc[1][i1][(i3+1)][i3]-aux1)*auxD;
63        Conc[0][i1][i2][auxL-1] +=
64          (tempConc[0][i1][i2][auxL]-tempConc[0][i1][i2][auxL-1])*auxD;
65        Conc[1][i1][i2][auxL-1] +=
66          (tempConc[1][i1][i2][auxL]-tempConc[0][i1][i2][auxL-1])*auxD;
67    }
68  }
```

Changed code line(s): 1, 4-11, 14-17, 21, 24-64

Optimization Notes

This function has two main parts: the first copies the Conc variable (using a for loop) to tempConc, and the second one iterates through the Conc array while constantly checking the upper and lower dimensional bounds. Those two parts can be parallelized without any difficulty: all the threads copy their divisions of Conc, and then all the threads execute divisions of the next big loop. The most expensive part of this function will be in the second big loop, where during each iteration it is necessary to increase and decrease the coordinates and check the bounds. If the calculated coordinate is inside the bounds, then execute the calculations.

The optimized function creates all the threads. Within each thread their portion of the memcpy loop executes, and when all have finished then they start the execution of the dimensional loop. The use of memcpy increases the copies done by iteration.

The variables are optimized in the same way previous functions were: using constants when possible and aligning the arrays.

Both loops are 'collapsed' under the OpenMP pragma header and are considered in this way:

- The first loop is under
#pragma omp for schedule(static) private(i1,i2) collapse(2)
- The second loop is under
#pragma omp for schedule(static) private(i1,i2,i3) collapse(2)

The first loop, the tempConc copy loop, is optimized with multithreading. The loop is multithreaded with a 'parallel for' and optimized by moving larger portions of data per iteration using memcpy (as commented previously). I don't know if all the available memory bandwidth is occupied during the process of copying the Conc to tempConc, so one way to improve this is by using all the available memory bandwidth (it will be approx. $L \cdot \text{float}(64\text{bit})^2 \cdot 240$).

The second loop, handling the dimensions, has different optimizations.

The initial optimization has to do with the i3 boundary check. This check takes place in the four lines of code that start with:

```
1 | Conc[0][i1][i2][0]
2 | Conc[1][i1][i2][0]
3 |
4 | Conc[0][i1][i2][auxL-1]
5 | Conc[1][i1][i2][auxL-1]
```

The i3 bound check can be avoided if the loop is L-2 iterations (i3 goes from 1 to auxL), so we have to manually add the first and last operations. The i3 bounds are never surpassed, contained by these four lines:

```
1 | Conc[0][i1][i2][i3] += (tempConc[0][i1][i2][(i3+1)]-aux)*auxD;
2 | Conc[1][i1][i2][i3] += (tempConc[1][i1][i2][(i3+1)]-aux1)*auxD;
3 | Conc[0][i1][i2][i3] += (tempConc[0][i1][i2][(i3-1)]-aux)*auxD;
4 | Conc[1][i1][i2][i3] += (tempConc[1][i1][i2][(i3-1)]-aux1)*auxD;
```

If the bounds are unreached, then the code inside the if code is executed.

That said, this loop can be improved if all the if code inside is deleted. One possible solution would be to manually unroll the loop and take care of the bounds, then further distribute tasks among the threads with OpenMP task functions.

runDecayStep

This function iterates through all Conc variables to complete multiplication on each.

Original Code

```
01 | static void runDecayStep(float**** Conc, int L, float mu) {
02 |     runDecayStep_sw.reset();
03 |     // computes the changes in substance concentrations due to decay
04 |     int i1,i2,i3;
05 |     for (i1 = 0; i1 < L; i1++) {
06 |         for (i2 = 0; i2 < L; i2++) {
07 |             for (i3 = 0; i3 < L; i3++) {
08 |                 Conc[0][i1][i2][i3] = Conc[0][i1][i2][i3]*(1-mu);
09 |                 Conc[1][i1][i2][i3] = Conc[1][i1][i2][i3]*(1-mu);
10 |             }
11 |         }
12 |     }
13 |     runDecayStep_sw.mark();
14 | }
```

Changed code line(s): 1, 8, 9

Optimized Code

```
01 | static void runDecayStep(int L, float Conc[2][L][L][L], float mu) {
02 |     runDecayStep_sw.reset();
03 |     // computes the changes in substance concentrations due to decay
04 |     const float muu=1-mu;
05 |     int i1,i2,i3;
06 |     omp_set_num_threads(240);
07 |     #pragma omp parallel for schedule(static) private(i1,i2,i3) collapse(3)
08 |     #pragma simd
09 |     for (i1 = 0; i1 < L; ++i1) {
10 |         for (i2 = 0; i2 < L; ++i2) {
11 |             for (i3 = 0; i3 < L; ++i3) {
12 |                 __assume_aligned((float*)Conc[0][i1][i2], 64);
13 |                 __assume_aligned((float*)Conc[1][i1][i2], 64);
14 |                 Conc[0][i1][i2][i3] = Conc[0][i1][i2][i3]*muu;
15 |                 Conc[1][i1][i2][i3] = Conc[1][i1][i2][i3]*muu;
16 |             }
17 |         }
18 |     }
19 |     runDecayStep_sw.mark();
20 | }
```

Changed code line(s): 1, 4, 6-8, 12-15

Optimization Notes

There are a few simple changes to optimize this function.

- The input arrays of the function are defined on the header.
- The variables are optimized as in previous functions (using constants when possible and aligning the arrays).
- Loops are 'collapsed' under the OpenMP pragma header.
- The use of #pragma simd loop. This loop can be vectorized with SIMD instructions, which improve the execution time on the optimized code. The use of SIMD instructions gives the possibility of executing multiple operations at the same time.

Additional optimizations that could be made are to use the Intel MKL functions to their fullest possibilities. The Intel MKL functions are optimized for large arrays of data, so if the Conc array is flattened to one dimension, the Intel MKL will perform faster.

cellMovementAndDuplication

This function is used to generate the random movement of each cell and each cell's duplicate. However, not every cell duplicates.

Original Code

```

01 static int cellMovementAndDuplication(float** posAll, float* pathTraveled,
02 int* typesAll, int* numberDivisions, float pathThreshold, int
03 divThreshold, int n) {
04     cellMovementAndDuplication_sw.reset();
05     int c;
06     currentNumberCells = n;
07     float currentNorm;
08     float currentCellMovement[3];
09     float duplicatedCellOffset[3];
10     for (c=0; c<n; c++) {
11         // random cell movement
12         currentCellMovement[0]=RandomFloatPos()-0.5;
13         currentCellMovement[1]=RandomFloatPos()-0.5;
14         currentCellMovement[2]=RandomFloatPos()-0.5;
15         currentNorm = getNorm(currentCellMovement);
16         posAll[c][0]+=0.1*currentCellMovement[0]/currentNorm;
17         posAll[c][1]+=0.1*currentCellMovement[1]/currentNorm;
18         posAll[c][2]+=0.1*currentCellMovement[2]/currentNorm;
19         pathTraveled[c]+=0.1;
20         // cell duplication if conditions fulfilled
21         if (numberDivisions[c]<divThreshold) {
22             if (pathTraveled[c]>pathThreshold) {
23                 pathTraveled[c]=pathThreshold;
24                 numberDivisions[c]++;
25                 currentNumberCells++; // update number of cells in
26                 the simulation
27                 numberDivisions[currentNumberCells-1]=numberDivisions[c]; // update number of divisions the di
28                 typesAll[currentNumberCells-1]=-typesAll[c];
29                 // assign type of duplicated cell (opposite to current cell)
30
31                 // assign location of duplicated cell
32                 duplicatedCellOffset[0]=RandomFloatPos()-0.5;
33                 duplicatedCellOffset[1]=RandomFloatPos()-0.5;
34                 duplicatedCellOffset[2]=RandomFloatPos()-0.5;
35                 currentNorm = getNorm(duplicatedCellOffset);
36                 posAll[currentNumberCells-1][0]=posAll[c][0]+0.05*duplicatedCellOffset[0]/currentNorm;
37                 posAll[currentNumberCells-1][1]=posAll[c][1]+0.05*duplicatedCellOffset[1]/currentNorm;
38                 posAll[currentNumberCells-1][2]=posAll[c][2]+0.05*duplicatedCellOffset[2]/currentNorm;
39             }
40         }
41     }
42     cellMovementAndDuplication_sw.mark();
43     return currentNumberCells;
44 }
```

Changed code line(s): 1, 12-14, 16-18, 21, 22, 24-34, 36-38

Optimized Code

```

01 static int cellMovementAndDuplication(float posAll[][3], float* pathTraveled,
02 int* typesAll, int* numberDivisions, float pathThreshold, int
03 divThreshold, int n) {
04     cellMovementAndDuplication_sw.reset();
05     int c, currentNumberCells = n;
06     float currentNorm;
07     unsigned int seed=rand();
08     float currentCellMovement[3] __attribute__((aligned(64)));
09     float duplicatedCellOffset[3] __attribute__((aligned(64)));
10     omp_set_num_threads(240);
11
12     #pragma omp parallel for simd schedule (static) shared(posAll) private(c,currentNorm,currentCellMovement)
13
14     for (c=0; c<n; ++c) {
15         // random cell movement
16         RandomFloatPos(currentCellMovement,seed+c);
17         currentNorm = getNorm(currentCellMovement)*10;
18         __assume_aligned((float*)posAll, 64);
19         __assume_aligned((float*)currentCellMovement, 64);
20         posAll[c][0:3]+=currentCellMovement[0:3]/currentNorm;
21         __assume_aligned((float*)pathTraveled, 64);
22         pathTraveled[c]+=0.1;
23     }
24     seed=rand();
25     for (c=0; c<n; ++c) {
26         if ((numberDivisions[c]<divThreshold) && (pathTraveled[c]>pathThreshold)) {
27             pathTraveled[c]=pathThreshold;
28             ++numberDivisions[c]; // update number of divisions this cell has undergone
29         }
30     }
31 }
```

```

29     numberDivisions[currentNumberCells]=numberDivisions[c]; // update number of divisions the duplicate
30     typesAll[currentNumberCells]=-typesAll[c]; // assign type of
31     duplicated cell (opposite to current cell)
32
33     // assign location of duplicated cell
34     RandomFloatPos(duplicatedCell0ffset,seed+c); //The seed+c value will be different for each thread :
35     currentNorm = getNorm(duplicatedCell0ffset)*20;
36     __assume_aligned((float*)posAll, 64);
37     __assume_aligned((float*)duplicatedCell0ffset, 64);
38     posAll[currentNumberCells][0:3]=posAll[c][0:3]+duplicatedCell0ffset[0:3]/currentNorm;
39     ++currentNumberCells; // update number of cells in the simulation
40   }
41 }
cellMovementAndDuplication_sw.mark();
return currentNumberCells;
}

```

Changed code line(s): 1, 8-12, 16-21, 26, 28, 29, 34, 36-39

Optimization Notes

The optimized code for cellMovementAndDuplication has two portions that were optimized in the same way as described for previous functions:

- RandomFloatPos() functions. The optimized function can be multithreaded and returns the whole position array (the 3D position array).
- The code lines incorporating currentNorm can be vectorized with SIMD and also mathematically simplified.

Additionally, the optimized code updates the header and aligns the arrays.

The new function RandomFloatPos() requires a random generated seed before calling seed=rand(), and each thread needs a different value of this seed. To achieve this, every thread increases the random generated seed value with their corresponding for loop iteration. This way the number is both random and different for each thread.

To multithread this function I divided it into two big loops. The first (random cell movement) can be parallelized without problems. The second loop (starting with seed=rand();) must be executed in a single thread, because every iteration depends on the previous iteration's values.

runDiffusionClusterStep

This function is used to determine cell movement based on the substance gradients.

Original Code

```

01 static void runDiffusionClusterStep(float**** Conc, float** movVec, float** posAll, int* typesAll, int n, int l
02   runDiffusionClusterStep_sw.reset();
03   // computes movements of all cells based on gradients of the two substances
04   float sideLength = 1/(float)L; // length of a side of a diffusion voxel
05
06   float gradSub1[3];
07   float gradSub2[3];
08   float normGrad1, normGrad2;
09   int c, i1, i2, i3, xUp, xDown, yUp, yDown, zUp, zDown;
10
11   for (c = 0; c < n; c++) {
12     i1 = std::min((int)floor(posAll[c][0]/sideLength),(L-1));
13     i2 = std::min((int)floor(posAll[c][1]/sideLength),(L-1));
14     i3 = std::min((int)floor(posAll[c][2]/sideLength),(L-1));
15
16     xUp = std::min((i1+1),L-1);
17     xDown = std::max((i1-1),0);
18     yUp = std::min((i2+1),L-1);
19     yDown = std::max((i2-1),0);
20     zUp = std::min((i3+1),L-1);
21     zDown = std::max((i3-1),0);
22
23     gradSub1[0] = (Conc[0][xUp][i2][i3]-Conc[0][xDown][i2][i3])/(sideLength*(xUp-xDown));
24     gradSub1[1] = (Conc[0][i1][yUp][i3]-Conc[0][i1][yDown][i3])/(sideLength*(yUp-yDown));
25     gradSub1[2] = (Conc[0][i1][zUp]-Conc[0][i1][zDown])/(sideLength*(zUp-zDown));
26     gradSub2[0] = (Conc[1][xUp][i2][i3]-Conc[1][xDown][i2][i3])/(sideLength*(xUp-xDown));
27     gradSub2[1] = (Conc[1][i1][yUp][i3]-Conc[1][i1][yDown][i3])/(sideLength*(yUp-yDown));
28     gradSub2[2] = (Conc[1][i1][zUp]-Conc[1][i1][zDown])/(sideLength*(zUp-zDown));
29     normGrad1 = getNorm(gradSub1);
30     normGrad2 = getNorm(gradSub2);
31     if ((normGrad1>0)&&(normGrad2>0)) {
32       movVec[c][0]=typesAll[c]*(gradSub1[0]/normGrad1-gradSub2[0]/normGrad2)*speed;
33       movVec[c][1]=typesAll[c]*(gradSub1[1]/normGrad1-gradSub2[1]/normGrad2)*speed;
34       movVec[c][2]=typesAll[c]*(gradSub1[2]/normGrad1-gradSub2[2]/normGrad2)*speed;
35     } else {
36       movVec[c][0]=0;
37       movVec[c][1]=0;
38       movVec[c][2]=0;
39     }
40   }
41   runDiffusionClusterStep_sw.mark();
42 }

```

Changed code line(s): 1-2, 5-10, 13-39

Optimized Code

```

01 static void runDiffusionClusterStep(int L, float Conc[2][L][L][L], float movVec[][3], float posAll[][3], int* typesAll
02   runDiffusionClusterStep_sw.reset();
03   // computes movements of all cells based on gradients of the two substances

```

```

04     const float auxL=L;
05     --L;
06     float gradSub[6] __attribute__((aligned(64)));
07     float aux[3] __attribute__((aligned(64)));
08     int i[3] __attribute__((aligned(32)));
09     float normGrad[2] __attribute__((aligned(64)));
10     int c, xUp, xDown, yUp, yDown, zUp, zDown;
11     omp_set_num_threads(240);
12
13 #pragma omp parallel for schedule(static) private(i,xUp,xDown,yUp,yDown,zUp,zDown,gradSub,normGrad,aux) if
14 for (c = 0; c < n; ++c) {
15     __assume_aligned((int*)i, 32);
16     __assume_aligned((float*)posAll, 64);
17     __assume_aligned((float*)gradSub, 64);
18     __assume_aligned((float*)normGrad, 64);
19     __assume_aligned((float*)movVec, 64);
20     __assume_aligned((float*)typesAll, 64);
21     i[0:3] = std::min((int)floor(posAll[c][0:3]*auxL),L)-1;
22     xDown = std::max(i[0],0);
23     yDown = std::max(i[1],0);
24     zDown = std::max(i[2],0);
25     xUp = std::min((i[0]+2),L);
26     yUp = std::min((i[1]+2),L);
27     zUp = std::min((i[2]+2),L);
28     aux[0]=auxL/((xUp-xDown));
29     aux[1]=auxL/((yUp-yDown));
30     aux[2]=auxL/((zUp-zDown));
31     gradSub[0] = (Conc[0][xUp][i[1]][i[2]]-Conc[0][xDown][i[1]][i[2]])*aux[0];
32     gradSub[1] = (Conc[0][i[0]][yUp][i[2]]-Conc[0][i[0]][yDown][i[2]])*aux[1];
33     gradSub[2] = (Conc[0][i[0]][i[1]][zUp]-Conc[0][i[0]][i[1]][zDown])*aux[2];
34     normGrad[0] = getNorm(gradSub);
35     if (normGrad[0]>0){
36         gradSub[3] = (Conc[1][i[0]][yUp][i[2]]-Conc[1][i[0]][yDown][i[2]])*aux[1];
37         gradSub[4] = (Conc[1][xUp][i[1]][i[2]]-Conc[1][xDown][i[1]][i[2]])*aux[0];
38         gradSub[5] = (Conc[1][i[0]][i[1]][zUp]-Conc[1][i[0]][i[1]][zDown])*aux[2];
39         normGrad[1] = getNorm(gradSub+3);
40         if ( normGrad[1]>0 {
41             movVec[c][0:3]=typesAll[c]*(gradSub[0:3]/normGrad[0]-gradSub[3:5]/normGrad[1])*speed;
42         } else movVec[c][0:3]=0;
43     }
44 }
45 runDiffusionClusterStep_sw.mark();
46 }
```

Changed code line(s): 1, 5-10, 12-14, 12-43

Optimization Notes

This function is optimized using some small simplifications, vectorization and multithreading, much in the same way other functions were optimized.

- In the header the dimensions of the input arrays are specified.
- It is array aligned.
- Replacing variables with constants whenever possible.
- OpenMP parallelization.
- Simplification of operations.
- Simplification of vectorization.

The if condition can filter some of the operations out. An example is the calculation of normGrad[1] is unnecessary if normGrad[0] is not > 0. Using the if condition results in fewer comparisons and fewer useless operations are calculated.

This function can be further improved using more SIMD optimizations. However, to use the SIMD operations, the max and min functions must be implemented differently.

getEnergy

This function computes the energy measure of a subvolume of cells by assuming uniform distribution within the entire volume and determining a volume of a target number of cells.

Original Code

```

01 Original Code
02 static float getEnergy(float** posAll, int* typesAll, int n, float spatialRange, int targetN) {
03     getEnergy_sw.reset();
04     // Computes an energy measure of clusteredness within a subvolume. The
05     // size of the subvolume is computed by assuming roughly uniform
06     // distribution within the whole volume, and selecting a volume
07     // comprising approximately targetN cells.
08     int i1, i2;
09     float currDist;
10     float** posSubvol=0; // array of all 3 dimensional cell positions
11     posSubvol = new float*[n];
12     int typesSubvol[n];
13     float subVolMax = pow(float(targetN)/float(n),1.0/3.0)/2;
14     if(quiet < 1)
15         printf("subVolMax: %f\n", subVolMax);
16     int nrCellsSubVol = 0;
17     float intraClusterEnergy = 0.0;
18     float extraClusterEnergy = 0.0;
19     float nrSmallDist=0.0;
20 }
```

```

21   for (i1 = 0; i1 < n; i1++) {
22     posSubvol[i1] = new float[3];
23     if ((fabs(posAll[i1][0]-0.5)<subVolMax) && (fabs(posAll[i1][1]-0.5)<subVolMax) && (fabs(posAll[i1][2]-0.5)<subVolMax)) {
24       posSubvol[nrCellsSubVol][0] = posAll[i1][0];
25       posSubvol[nrCellsSubVol][1] = posAll[i1][1];
26       posSubvol[nrCellsSubVol][2] = posAll[i1][2];
27       typesSubvol[nrCellsSubVol] = typesAll[i1];
28       nrCellsSubVol++;
29     }
30   }
31
32   for (i1 = 0; i1 < nrCellsSubVol; i1++) {
33     for (i2 = i1+1; i2 < nrCellsSubVol; i2++) {
34       currDist = getL2distance(posSubvol[i1][0],posSubvol[i1][1],posSubvol[i1][2],posSubvol[i2][0],posSubvol[i2][1],posSubvol[i2][2]);
35       if (currDist<spatialRange) {
36         nrSmallDist = nrSmallDist+1;//currDist/spatialRange;
37         if (typesSubvol[i1]*typesSubvol[i2]>0) {
38           intraClusterEnergy = intraClusterEnergy+fmin(100.0,spatialRange/currDist);
39         } else {
40           extraClusterEnergy = extraClusterEnergy+fmin(100.0,spatialRange/currDist);
41         }
42       }
43     }
44   }
45   float totalEnergy = (extraClusterEnergy-intraClusterEnergy)/(1.0+100.0*nrSmallDist);
46   getEnergy_sw.mark();
47   return totalEnergy;
48 }

```

Changed code line(s): 1, 22, 34, 35, 37-41, 46, 48

Optimized Code

```

01 static float getEnergy(float posAll[][], int* typesAll, int n, float spatialRange, int targetN) {
02   getEnergy_sw.reset();
03   // Computes an energy measure of clusteredness within a subvolume. The
04   // size of the subvolume is computed by assuming roughly uniform
05   // distribution within the whole volume, and selecting a volume
06   // comprising approximately targetN cells.
07   int i1, i2;
08   float currDist;
09   float posSubvol[n][3] __attribute__((aligned(64)));
10   int typesSubvol[n] __attribute__((aligned(64)));
11   const float subVolMax = pow(float(targetN)/float(n),1.0/3.0)/2;
12   if(quiet < 1)printf("subVolMax: %f\n", subVolMax);
13
14   int nrCellsSubVol = 0;
15   float intraClusterEnergy = 0.0;
16   float extraClusterEnergy = 0.0;
17   float nrSmallDist=0.0;
18   for (i1 = 0; i1 < n; ++i1) {
19     __assume_aligned((float*)posAll, 64);
20     if ((fabs(posAll[i1][0]-0.5)<subVolMax) && (fabs(posAll[i1][1]-0.5)<subVolMax) && (fabs(posAll[i1][2]-0.5)<subVolMax)) {
21       __assume_aligned((float*)posSubvol[nrCellsSubVol], 64);
22       __assume_aligned((float*)typesAll, 64);
23       __assume_aligned((int*)typesSubvol, 64);
24       posSubvol[nrCellsSubVol][0] = posAll[i1][0];
25       posSubvol[nrCellsSubVol][1] = posAll[i1][1];
26       posSubvol[nrCellsSubVol][2] = posAll[i1][2];
27       typesSubvol[nrCellsSubVol] = typesAll[i1];
28       ++nrCellsSubVol;
29     }
30   }
31   omp_set_num_threads(240);
32   #pragma omp parallel for schedule(static) reduction(+:nrSmallDist,intraClusterEnergy,extraClusterEnergy) pi
33   for (i1 = 0; i1 < nrCellsSubVol; ++i1) {
34     for (i2 = i1+1; i2 < nrCellsSubVol; ++i2) {
35       currDist = getL2distance(posSubvol[i1],posSubvol[i2]);
36       if (currDist<spatialRange) {
37         ++nrSmallDist; //currDist/spatialRange;
38         (typesSubvol[i1]*typesSubvol[i2]>0)? intraClusterEnergy += fmin(100.0,spatialRange/currDist) :
39         extraClusterEnergy += fmin(100.0,spatialRange/currDist);
40       }
41     }
42   }
43   getEnergy_sw.mark();
44   return (extraClusterEnergy-intraClusterEnergy)/(1.0+100.0*nrSmallDist);
45 }

```

Changed code line(s): 1, 9-11, 19, 22-24, 32, 33, 36-40

Optimization Notes

Some of the things in this code that were easy targets for optimization are the ‘new’ constructors, some variables that are constants, and the fact that the second loop depends on the first. The first loop isn’t parallelizable since it depends on the value of the nrCellsSubVol and the if with three conditions inside this loop makes it even more difficult to parallelize. The second loop, however, can be parallelized without problems.

The specific optimizations done to this function are much the same as in the other functions:

- The dimensions of the input arrays are specified in the header.
- Use of constants when possible.
- Array aligned.

- OpenMP parallelization on the second 'for loop' using a reduction.

The next step optimization would be to parallelize the first loop if condition. It is difficult but possible and would result in provable gains.

getCriterion

This function is used to determine if the cells in the subvolume are in clusters.

NOTE: Both the original code and the optimized code are truncated versions, showing only the code portions relevant to this paper.

Original Code

```

01 static bool getCriterion(float** posAll, int* typesAll, int n, float spatialRange, int targetN) {
02     getCriterion_sw.reset();
03     // Returns 0 if the cell locations within a subvolume of the total
04     // system, comprising approximately targetN cells, are arranged as clusters, and 1 otherwise.
05     int i1, i2;
06     int nrClose=0; // number of cells that are close (i.e. within a distance of spatialRange)
07     float currDist;
08     int sameTypeClose=0; // number of cells of the same type, and that are close (i.e. within a distance of spatialRange)
09     int diffTypeClose=0; // number of cells of opposite types, and that are close (i.e. within a distance of spatialRange)
10     float** posSubvol=0; // array of all 3 dimensional cell positions in the subcube
11     posSubvol = new float*[n];
12     int typesSubvol[n];
13     float subVolMax = pow(float(targetN)/float(n),1.0/3.0)/2;
14     int nrCellsSubVol = 0;
15
16     // the locations of all cells within the subvolume are copied to array PosSubvol
17     for (i1 = 0; i1 < n; i1++) {
18         posSubvol[i1] = new float[3];
19         if ((fabs(posAll[i1][0]-0.5)<subVolMax) && (fabs(posAll[i1][1]-0.5)<subVolMax) && (fabs(posAll[i1][2]
20             -0.5)<subVolMax)) {
21             posSubvol[nrCellsSubVol][0] = posAll[i1][0];
22             posSubvol[nrCellsSubVol][1] = posAll[i1][1];
23             posSubvol[nrCellsSubVol][2] = posAll[i1][2];
24             typesSubvol[nrCellsSubVol] = typesAll[i1];
25             nrCellsSubVol++;
26         }
27     }
28
29     [section of truncated code]
30
31     for (i1 = 0; i1 < nrCellsSubVol; i1++) {
32         for (i2 = i1+1; i2 < nrCellsSubVol; i2++) {
33             currDist = getL2Distance(posSubvol[i1][0],posSubvol[i1][1],posSubvol[i1][2],posSubvol[i2][0], posS
34             if (currDist<spatialRange) {
35                 nrClose++;
36                 if (typesSubvol[i1]*typesSubvol[i2]<0) {
37                     diffTypeClose++;
38                 } else {
39                     sameTypeClose++;
40                 }
41             }
42         }
43     }
44     [section of truncated code]
45
46 }
```

Changed code line(s): 9, 10, 32-39

Optimized Code

```

01 static bool getCriterion(float posAll[][3], int* typesAll, int n, float spatialRange, int targetN) {
02     getCriterion_sw.reset();
03     // Returns 0 if the cell locations within a subvolume of the total
04     // system, comprising approximately targetN cells, are arranged as clusters, and 1 otherwise.
05     int i1, i2;
06     int nrClose=0; // number of cells that are close (i.e. within a distance of spatialRange)
07     int sameTypeClose=0; // number of cells of the same type, and that are close (i.e. within a distance of spatialRange)
08     int diffTypeClose=0; // number of cells of opposite types, and that are close (i.e. within a distance of spatialRange)
09     float posSubvol[n][3] __attribute__((aligned(64)));
10     int typesSubvol[n] __attribute__((aligned(64)));
11     const float subVolMax = pow(float(targetN)/float(n),1.0/3.0)/2;
12     int nrCellsSubVol = 0;
13
14     // the locations of all cells within the subvolume are copied to array posSubvol
15
16     for (i1 = 0; i1 < n; ++i1) {
17         __assume_aligned((float*)posAll, 64);
18         if ((fabs(posAll[i1][0]-0.5)<subVolMax) && (fabs(posAll[i1][1]-0.5)<subVolMax) && (fabs(posAll[i1][2]
19             -0.5)<subVolMax)) {
20             __assume_aligned((float*)posSubvol[nrCellsSubVol], 64);
21             __assume_aligned((float*)typesAll, 64);
22             __assume_aligned((int*)typesSubvol, 64);
23             posSubvol[nrCellsSubVol][0] = posAll[i1][0];
24             posSubvol[nrCellsSubVol][1] = posAll[i1][1];
25             posSubvol[nrCellsSubVol][2] = posAll[i1][2];
26             typesSubvol[nrCellsSubVol] = typesAll[i1];
27             ++nrCellsSubVol;
28         }
29     }
30
31     [section of truncated code]
32
33     omp_set_num_threads(240);
34     #pragma omp parallel for schedule(static)
```

```

35     reduction(+:nrClose,diffTypeClose,sameTypeClose) private(i1,i2)
36     for (i1 = 0; i1 < nrCellsSubVol; ++i1) {
37         for (i2 = i1+1; i2 < nrCellsSubVol; ++i2) {
38             if (getL2Distance(posSubvol[i1],posSubvol[i2])<spatialRange) {
39                 ++nrClose;
40                 (typesSubvol[i1]*typesSubvol[i2]<0) ? ++diffTypeClose: ++sameTypeClose;
41             }
42         }
43     }
44     [section of truncated code]
45 }
46
47 }
```

Changed code line(s): 1, 9-11, 17, 20-22, 33-35, 38-40

Optimization Notes

The optimizable parts of this function include the variable declaration and initialization. The first loop is less optimizable, but optimizing the second loop can make up the difference.

The specific optimizations done to this function are much the same as in the other functions:

- The dimensions of the input arrays are specified in the header.
- Use of constants when possible.
- OpenMP parallelization on the second ‘for loop’ using a reduction.

For further optimization, it may be possible to parallelize the first loop if condition. Additionally, flattening the posAll and posSubVol arrays can lead to a faster execution.

Main – Variable Declaration

The main code is large so for this paper it is divided into meaningful sections. This section deals with the code and optimized code for the variable declarations.

Original Code

```

01 int i,c,d;
02 int i1, i2, i3, i4;
03 float energy; // value that quantifies the quality of the cell clustering output. The smaller this value, the
04 float** posAll=0; // array of all 3 dimensional cell positions
05 posAll = new float*[finalNumberCells];
06 float** currMov=0; // array of all 3 dimensional cell movements at the last time point
07 currMov = new float*[finalNumberCells]; // array of all cell movements in the last time step
08 float zeroFloat = 0.0;
09 float pathTraveled[finalNumberCells]; // array keeping track of length of path traveled until cell divides
10 int numberDivisions[finalNumberCells]; //array keeping track of number of division a cell has undergone
11 int typesAll[finalNumberCells]; // array specifying cell type (+1 or -1)
```

Changed code line(s): 2, 4-7,

Optimized Code

```

01 int i,c,d,i1;
02 float energy; // value that quantifies the quality of the cell clustering output. The smaller this value, the l
03 float posAll[finalNumberCells][3] __attribute__((aligned(64)));
04 float currMov[finalNumberCells][3] __attribute__((aligned(64))); //array of
05 // all cell movements in the last time step
06 float pathTraveled[finalNumberCells] __attribute__((aligned(64))); // array
07 // keeping track of length of path traveled until cell divides
08 int numberDivisions[finalNumberCells] __attribute__((aligned(64))); //array
09 // keeping track of number of division a cell has undergone
10 int typesAll[finalNumberCells] __attribute__((aligned(64))); // array
11 // specifying cell type (+1 or -1)
12 float Conc[2][L][L][L] __attribute__((aligned(64))));
```

Changed code line(s): 3-12

Optimization Notes

There were many ways to optimize the code in this small section.

- Avoid the use of ‘new’ constructors.
- Delete useless variables (for example, zeroFloat).
- Avoid using pointers when we can declare an array of a known size and static.
- Memory alignment; memory access is faster though it can use more memory.
- Define the variable value at initialization time; this is faster than defining the value later.

Main – Variable Initialization

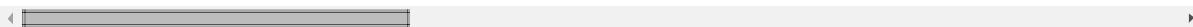
The main code is large so for this paper it is divided into meaningful sections. This section deals with the code and optimized code for the

variable initialization. The size and values of the variables are declared in this code, including the largest arrays currMov, posAll, pathTraveled, and Conc.

Original Code

```

01 // Initialization of the various arrays
02 for (i1 = 0; i1 < finalNumberCells; i1++) {
03     currMov[i1] = new float[3];
04     posAll[i1] = new float[3];
05     pathTraveled[i1] = zeroFloat;
06     pathTraveled[i1] = 0;
07     for (i2 = 0; i2 < 3; i2++) {
08         currMov[i1][i2] = zeroFloat;
09         posAll[i1][i2] = 0.5;
10     }
11 }
12 // create 3D concentration matrix
13 float*** Conc;
14 Conc = new float***[L];
15 for (i1 = 0; i1 < 2; i1++) {
16     Conc[i1] = new float**[L];
17     for (i2 = 0; i2 < L; i2++) {
18         Conc[i1][i2] = new float*[L];
19         for (i3 = 0; i3 < L; i3++) {
20             Conc[i1][i2][i3] = new float[L];
21             for (i4 = 0; i4 < L; i4++) {
22                 Conc[i1][i2][i3][i4] = zeroFloat;
23             }
24         }
25     }
26 }
```

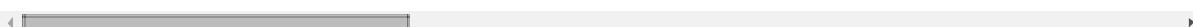


Changed code line(s): 3-5, 7-10, 13, 14, 16, 18, 20, 22

Optimized Code

```

01 // Initialization of the various arrays
02 omp_set_num_threads(240);
03 #pragma omp parallel
04 {
05     #pragma omp for simd schedule (static) private (i1) nowait
06     for (i1 = 0; i1 < finalNumberCells; ++i1) {
07         __assume_aligned((float*)posAll, 64);
08         __assume_aligned((float*)currMov[i1], 64);
09         __assume_aligned((float*)pathTraveled, 64);
10         currMov[i1][0:3]=0;
11         posAll[i1][0]=0.5;
12         posAll[i1][1]=0.5;
13         posAll[i1][2]=0.5;
14         pathTraveled[i1] = 0;
15     }
16     #pragma omp for schedule (static) private (i1,c,d) collapse(3)
17     for (i1 = 0; i1 < 2; ++i1) {
18         for (c = 0; c < L; ++c) {
19             for (d = 0; d < L; ++d) {
20                 Conc[i1][c][d][0:L]=0;
21             }
22         }
23     }
24 }
```



Changed code line(s): 2-5, 7-13, 16, 20

Optimization Notes

There are three ways the code in this portion was optimized. The most basic is the reduction of the for loops from four to three. The other optimizations are memory alignment and the OpenMP thread parallelization with SIMD.

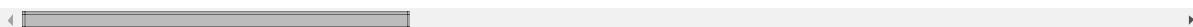
Main – Phase 1

The main code is large, so for this paper it is divided into meaningful sections. This section deals with the first steps of the simulation.

Original Code

```

01 // Phase 1: Cells move randomly and divide until final number of cells is reached
02 while (n<finalNumberCells) {
03     produceSubstances(Conc, posAll, typesAll, L, n); // Cells produce
04     // substances. Depending on the cell type, one of the two substances is produced.
05     runDiffusionStep(Conc, L, D); // Simulation of substance diffusion
06     runDecayStep(Conc, L, mu);
07     n = cellMovementAndDuplication(posAll, pathTraveled, typesAll, numberDivisions, pathThreshold, divThreshold,
08
09     for (c=0; c<n; c++) {
10         // boundary conditions
11         for (d=0; d<3; d++) {
12             if (posAll[c][d]<0) {posAll[c][d]=0;}
13             if (posAll[c][d]>1) {posAll[c][d]=1;}
14         }
15     }
16 }
```



Changed code line(s): 4, 6, 7, 14**Optimized Code**

```

01 // Phase 1: Cells move randomly and divide until final number of cells is reached
02 while (n<finalNumberCells) {
03     produceSubstances(L,Conc, posAll, typesAll, n); // Cells produce substances. Depending on the cell type, or
04     runDiffusionStep(L,Conc, D); // Simulation of substance diffusion
05     runDecayStep(L,Conc, mu);
06     n = cellMovementAndDuplication(posAll, pathTraveled, typesAll, numberDivisions, pathThreshold, divThreshold);
07     omp_set_num_threads(240);
08     #pragma omp parallel for simd schedule (static) private (c,d) if (n>500)
09     for (c=0; c<n; ++c) {
10         // boundary conditions
11         for (d=0; d<3; d++) {
12             if (posAll[c][d]<0) {
13                 posAll[c][d]=0;
14             }else if (posAll[c][d]>1) posAll[c][d]=1;
15         }
16     }
17 }
```

Changed code line(s): 4-6, 8, 9, 15**Optimization Notes**

There are two optimizations used in this portion of the code: OpenMP thread parallelization with SIMD and if clause reduction.

The OpenMP thread parallelization with SIMD is a good change because the schedule is static. Further, we know the size of the loop and the cost of each iteration is the same so the load cost will be equal for all threads. The if n>500 condition helps to reduce the added overhead for the first 500 n values.

The original code always does two checks. The optimized code reduces these if checks by doing one, and when necessary doing the second if check. This reduces the total number of comparisons done.

Main – Phase 2 Ending**Original code**

```

01 for (c=0; c<n; c++) {
02     posAll[c][0] = posAll[c][0]+currMov[c][0];
03     posAll[c][1] = posAll[c][1]+currMov[c][1];
04     posAll[c][2] = posAll[c][2]+currMov[c][2];
05
06     // boundary conditions: cells can not move out of the cube [0,1]^3
07     for (d=0; d<3; d++) {
08         if (posAll[c][d]<0) {posAll[c][d]=0;}
09         if (posAll[c][d]>1) {posAll[c][d]=1;}
10     }
11 }
```

Changed code line(s): 2-4, 9**Optimized code**

```

01 #pragma omp parallel for simd schedule (static) private (c,d) if (n>500)
02 for (c=0; c<n; ++c) {
03     __assume_aligned((float*)posAll, 64);
04     __assume_aligned((float*)currMov[c], 64);
05     posAll[c][0:3] += currMov[c][0:3];
06     // boundary conditions: cells can not move out of the cube [0,1]^3
07     for (d=0; d<3; d++) {
08         if (posAll[c][d]<0) {
09             posAll[c][d]=0;
10         }else if (posAll[c][d]>1) posAll[c][d]=1;
11     }
12 }
```

Changed code line(s): 1, 3-5, 10**Optimization Notes**

This final part of the code is optimized in the same way as the Phase 1 optimization: OpenMP thread parallelization with SIMD and if clause reduction.

The OpenMP thread parallelization with SIMD is a good change because the schedule is static. Further, we know the size of the loop and the cost of each iteration is the same so the load cost will be equal for all threads. The if n>500 condition helps to reduce the added overhead for the first 500 n values.

The original code always does two checks. The optimized code reduces these if checks by doing one, and when necessary doing the second if check. This reduces the total number of comparisons done.

Other Optimizations

Besides the optimizations listed in each section of the code previously, there were other optimizations made that were not explicitly remarked

on:

- Defining all post-increments as pre-increments. The post-increment makes a copy of the actual variable and then increments it. A pre-increment does not make this copy so it's faster. This is only a gain when the point of incrementation does not affect the result.
- Memory alignment on all arrays. When the arrays are aligned, the memory operations are faster and the CPU doesn't need to mask or convert the data.

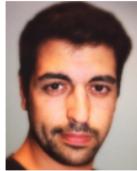
Rejected Optimization Techniques

During the development of this final code I experimented with other parallelization and vectorization techniques, such as Intel Cilk Plus and Intel® Threading Building Blocks (Intel® TBB), as well as some mathematic thread-optimized functions from the Intel MKL. Using Intel Cilk functions made the code slow, though this might be because I didn't fully understand how to implement Intel Cilk functions.

Since the OpenMP functions performed well without problems, I did not do additional testing with the Intel TBB functions, so it is possible that Intel TBB could improve performance as well.

Another possibility was using Intel MKL functions, which perform well when they work with large data. While the code moves large parts of data, the instantaneous data volume is low. Because of the low instantaneous data volume (and possibly other reasons) the Intel MKL functions were not a good option for optimization.

About the Author



[\(https://software.intel.com/sites/default/files/manage/d/4f/32/Untitled.png\)](https://software.intel.com/sites/default/files/manage/d/4f/32/Untitled.png)

Daniel Vea Falguera is an Electronic Systems Engineering student and entrepreneur with interests and knowledge about electronics and computer programming.

Links

Most of the optimizations done in this code are based on this Intel Xeon Phi coprocessor article: <https://software.intel.com/en-us/articles/optimization-and-performance-tuning-for-intel-xeon-phi-coprocessors-part-1-optimization> (<https://software.intel.com/en-us/articles/optimization-and-performance-tuning-for-intel-xeon-phi-coprocessors-part-1-optimization%20>)

Other similar commented options and solutions were provided at the Intel Modern Code for Parallel Architectures forums:

<https://software.intel.com/en-us/forums/intel-moderncode-for-parallel-architectures> (<https://software.intel.com/en-us/forums/intel-moderncode-for-parallel-architectures>)

Intel® MKL: <https://software.intel.com/en-us/mkl-reference-manual-for-c> (<https://software.intel.com/en-us/mkl-reference-manual-for-c>)

Intel® TBB: <https://www.threadingbuildingblocks.org/> (<https://www.threadingbuildingblocks.org/>)

Intel® Cilk Plus: <https://www.cilkplus.org/> (<https://www.cilkplus.org/>)

OpenMP*: <http://openmp.org/wp/> (<http://openmp.org/wp/>)

For more complete information about compiler optimizations, see our [Optimization Notice](https://software.intel.com/en-us/articles/optimization-notice#opt-en) (<https://software.intel.com/en-us/articles/optimization-notice#opt-en>).

Categories: [Academic](https://software.intel.com/en-us/search/site/field_topic/academic-78162/language/en) (https://software.intel.com/en-us/search/site/field_topic/academic-78162/language/en), [Code Modernization](https://software.intel.com/en-us/search/site/field_topic/code_modernization-80685/language/en) (https://software.intel.com/en-us/search/site/field_topic/code_modernization-80685/language/en), [Intel® Many Integrated Core Architecture](https://software.intel.com/en-us/search/site/field_topic/intel_many_integrated_core_architecture-36781/language/en) (https://software.intel.com/en-us/search/site/field_topic/intel_many_integrated_core_architecture-36781/language/en), [Threading](https://software.intel.com/en-us/search/site/field_topic/threading-20869/language/en) (https://software.intel.com/en-us/search/site/field_topic/threading-20869/language/en), [Vectorization](https://software.intel.com/en-us/search/site/field_topic/vectorization-35851/language/en) (https://software.intel.com/en-us/search/site/field_topic/vectorization-35851/language/en), [Intel® Advanced Vector Extensions](https://software.intel.com/en-us/search/site/field_technology/intel_advanced_vector_extensions-20857/language/en) (https://software.intel.com/en-us/search/site/field_technology/intel_advanced_vector_extensions-20857/language/en), [OpenMP*](https://software.intel.com/en-us/search/site/field_technology/openmp-20860/language/en) (https://software.intel.com/en-us/search/site/field_technology/openmp-20860/language/en), [C/C++](https://software.intel.com/en-us/search/site/field_programming_language/cc-20802/language/en) (https://software.intel.com/en-us/search/site/field_programming_language/cc-20802/language/en), [Server](https://software.intel.com/en-us/search/site/field_platform/server-20798/language/en) (https://software.intel.com/en-us/search/site/field_platform/server-20798/language/en), [Developers](https://software.intel.com/en-us/search/site/field_audience/developers-17152/language/en) (https://software.intel.com/en-us/search/site/field_audience/developers-17152/language/en), [Professors](https://software.intel.com/en-us/search/site/field_audience/professors-17154/language/en) (https://software.intel.com/en-us/search/site/field_audience/professors-17154/language/en), [Students](https://software.intel.com/en-us/search/site/field_audience/students-17155/language/en) (https://software.intel.com/en-us/search/site/field_audience/students-17155/language/en), [Linux*](https://software.intel.com/en-us/search/site/field_operating_system/linux-20787/language/en) (https://software.intel.com/en-us/search/site/field_operating_system/linux-20787/language/en), [Intermediate](https://software.intel.com/en-us/search/site/field_skill_level/intermediate-20808/language/en) (https://software.intel.com/en-us/search/site/field_skill_level/intermediate-20808/language/en)

Tags: [Intel® Modern Code Developer Challenge](https://software.intel.com/en-us/search/site/field_tags/intel_modern_code_developer_challenge-81768/language/en) (https://software.intel.com/en-us/search/site/field_tags/intel_modern_code_developer_challenge-81768/language/en)

[^Top](#)

Add a Comment

(For technical discussions visit our [developer forums](#). For site or software product issues [contact support](#).)

Please [sign in](#) to add a comment. Not a member?

[Join today >](#)

[Support](#) [Terms of Use](#) [*Trademarks](#) [Privacy](#) [Cookies](#)

Look for us on:

[English >](#)