

# **Parboil:** an optimization effort

Lucas Henrique Morais  
Luís Felipe Mattos

October 11, 2013

# Contents

<b>1</b>	<b>General Methods</b>	<b>2</b>
1.1	Lock-free vector accumulation . . . . .	2
1.2	Prescient memory access . . . . .	4
<b>2</b>	<b>Optimizations</b>	<b>5</b>
2.1	sgemm . . . . .	6
2.2	stencil . . . . .	7
2.3	mri-gridding . . . . .	8
2.4	cutcp . . . . .	9
2.5	tpacf . . . . .	11

# Chapter 1

## General Methods

In this study we identified two optimization techniques that are very broad in scope and that would be recurrently employed in our solutions for each benchmark. These are as follows:

### 1.1 Lock-free vector accumulation

A very general optimization technique that was consistently employed in the studied set of benchmarks was lock-free vector accumulation. Let us introduce the topic with the following example: Usually, one could trivially avoid declaring a critical section on a piece of code like this:

```
1 int foo( ... ){  
  #pragma omp parallel  
3  for(i = 0; i < size; i++){  
    //many lines of code  
5    #pragma omp critical section  
      bin += aux;  
7  }  
}
```

Listing 1.1: 0-dimensional accumulation

Usually, one could trivially avoid declaring the critical section above by:

1. Declaring a private variable for each thread.
2. Requiring each thread to accumulate on its own version of the variable.
3. Consolidating the partial sums on the original variable.

The present report then suggests a generalization of that approach for not only lock-free accumulation on unidimensional arrays, but for n-dimensional vectors, with negligible resulting overhead.

```

1 int foo( ... ){
2   #pragma omp parallel
3   for(i = 0; i < size; i++){
4     //many lines of code
5     #pragma omp critical section
6     bin[p1][p2][p3][p4] += aux;
7   }
8 }

```

Listing 1.2: 4-dimensional accumulation

The overall strategy is comprised of three main elements, as follows:

1. Creation of an  $(n + 1)$ -dimensional auxiliar vector.
2. Private accumulation of each thread's work on its respective n-dimensional data strip.
3. Parallel consolidation of the results on the target (original) vector.

(\* Here we assume a  $n$ -dimensional target vector.) After the first one, each step is triggered by the end of the preceeding one. Each is now described in detail:

**Replication** At the end of this step, we want each thread to have its own copy of the original structure. This can be easily accomplished by allocating a vector wrapping  $k$  copies of the original D.S. We then have a  $(n + 1)$  dimensional vector on memory.

**Private accumulation** Each thread then works on its own data share and privately accumulate on its own data strip. Generally, a  $j$ -numbered thread shall access memory locations of the format

$$V_{aux}[j][*][*] \dots [*],$$

where  $*$  can assume any available value.

**Consolidation** By this time, all working data has been consumed, and we are left with the task of summing up all the partial results on the original structure, which we can fortunately do in parallel for any  $n$  greater than zero, as follows:

**Partitioning** Each thread is associated with a subset of the target structure memory positions. This can be automatically achieved with OpenMP *parallel for* pragma.

**Actual consolidation** Then, each thread computes the following expression:

$$V_{final}\{\pi_j\} = V_{aux}[0]\{\pi_j\} + V_{aux}[1]\{\pi_j\} + \dots + V_{aux}[n-1]\{\pi_j\}$$

where  $n$  is the total number of working threads and  $\pi_j$  indicates the partition of the memory positions that  $j$ -numbered thread is responsible for. This expression is iteratively computed. As partition sets are always disjoint, it follows that this step can be trivially parallelized with, for example, a *parallel for* pragma as suggested above.

For our example, we would now have:

```

1 int foo( ... ){
2   #pragma omp parallel
3   for(i = 0; i < size; i++){
4     //many lines of code
5     big_bin[thread_num][p1][p2][p3][p4] += aux;
6   }
7
8   #pragma omp parallel
9   for (i = 0; i < num_threads; i++)
10     for (j = 0; j < j_max; j++)
11       for (k = 0; k < k_max; k++)
12         for (l = 0; l < l_max; l++)
13           for (m = 0; m < m_max; m++){
14             bin = big_bin[i][j][k][l][m];
15           }
16 }

```

Listing 1.3: 4-dimensional accumulation

## 1.2 Prescient memory access

One of the greatest optimizations we applied in several algorithms was the better use of the cache memory and the main memory. Some data intensive reading algorithms can benefit from the cache memory by accessing sequential addresses from the main memory, with lesser cache misses, the algorithm has no need to search the data in the main memory, which is much slower than low level memories.

By changing the way the memory is accessed, the performance was improved in the sequential version of some algorithms, improving even the performance over some simple parallelized versions.

Several algorithms in this benchmark set have nested loops and access some structure, like a array or a matrix, inside the inner loop. Also, some algorithms have inefficient index choice for the inner loop, turning the memory access into a slow and repetitious task. Some changes that improved the performance include the remapping of the structure or/and the better choice of the order for the indexes in the nested loops.

## Chapter 2

# Optimizations

## 2.1 sgemm

```

2   for (int mm = 0; mm < m; ++mm) {
4       for (int nn = 0; nn < n; ++nn) {
6           float c = 0.0f;
7           for (int i = 0; i < k; ++i) {
8               float a = A[mm + i * lda];
9               float b = B[nn + i * ldb];
10              c += a * b;
11          }
12          C[mm+nn*ldc] = C[mm+nn*ldc] * beta +
              alpha * c;
      }
  }

```

Listing 2.1: Parboil sequential source code for sgemm

```

2   #pragma omp parallel for collapse (2)
3   for (int mm = 0; mm < m; ++mm) {
4       for (int nn = 0; nn < n; ++nn) {
5           float c = 0.0f;
6           for (int i = 0; i < k; ++i) {
7               float a = A[mm + i * lda];
8               float b = B[nn + i * ldb];
9               c += a * b;
10          }
11          C[mm+nn*ldc] = C[mm+nn*ldc] * beta +
              alpha * c;
12      }
  }

```

Listing 2.2: Parboil parallel source code for sgemm

For the matrices multiplication algorithm, it is used an array-represented structure, so, we improved the memory access by transposing the matrices. This way, the index 'i' in the inner loop, is only added instead of multiplied to calculate the correct mapping of the element in the matrix to the array-represented matrix. We needed then to allocate auxiliary arrays to represent the transposed matrices and then transpose while making the copy from one array to another.

To avoid the increase of the execution time of the algorithm, we parallelized this transposition. This can be done because it has no loop-carried dependencies and the data is distributed almost equally for all threads without concurrency, in other words, it is high parallelizable.

We created the variables aux1 and aux2 because in the inner loop, these values are constant. The original versions, the sequential and the stock parallel, calculate these constant every iteration of the inner loop, wasting some time in unnecessary multiplications.

The access to the array is almost sequential after this changes. Using better the memory, caused some great improvement in execution time of this algorithm. By applying these techniques we could run our parallel version around 10.7 times faster than sequential version and 3.9 times faster than the stock parallel version.

```

1   tb = (float *) malloc(ldb * n * sizeof(float));
2   #pragma omp parallel for collapse(2)
3   for (i = 0; i < ldb; i++)
4       for (j = 0; j < n; j++)
5           tb[j * n + i] = B[i * n + j];
6
7   ta = (float *) malloc(lda * m * sizeof(float));
8   #pragma omp parallel for collapse(2)
9   for (i = 0; i < lda; i++)
10      for (j = 0; j < m; j++)
11          ta[j * m + i] = A[i * m + j];

```

Listing 2.3: Array-represented matrices transposition

```

1   #pragma omp parallel for private(nn, c, i, a, b, mm) collapse(2) schedule (static)
2   for (mm = 0; mm < m; ++mm) {
3       for (nn = 0; nn < n; ++nn) {
4           c = 0.0f;
5           aux1 = mm * lda;
6           aux2 = nn * ldb;
7           for (i = 0; i < k; ++i) {
8               c += ta[i + aux1] * tb[i + aux2];
9           }
10          C[mm+nn*ldc] = C[mm+nn*ldc] * beta + alpha * c;
11      }
  }

```

Listing 2.4: Our parallel source code for sgemm

## 2.2 stencil

```

1  int i, j, k;
2  for(i=1;i<nx-1;i++) {
3      for(j=1;j<ny-1;j++) {
4          for(k=1;k<nz-1;k++) {
5              Anext[Index3D (nx, ny, i, j, k)] =
6                  (A0[Index3D (nx, ny, i, j, k + 1)] +
7                   A0[Index3D (nx, ny, i, j, k - 1)] +
8                   A0[Index3D (nx, ny, i, j + 1, k)] +
9                   A0[Index3D (nx, ny, i, j - 1, k)] +
10                  A0[Index3D (nx, ny, i + 1, j, k)] +
11                  A0[Index3D (nx, ny, i - 1, j, k)])*c1
12              - A0[Index3D (nx, ny, i, j, k)]*c0;
13          }
14      }
15  }
16  }

```

Listing 2.5: Parboil sequential source code for stencil

```

2  int i;
3  #pragma omp parallel for
4  for(i=1;i<nx-1;i++) {
5      int j, k;
6      for(j=1;j<ny-1;j++) {
7          for(k=1;k<nz-1;k++) {
8              // #pragma omp critical
9              Anext[Index3D (nx, ny, i, j, k)] =
10                  (A0[Index3D (nx, ny, i, j, k + 1)] +
11                   A0[Index3D (nx, ny, i, j, k - 1)] +
12                   A0[Index3D (nx, ny, i, j + 1, k)] +
13                   A0[Index3D (nx, ny, i, j - 1, k)] +
14                   A0[Index3D (nx, ny, i + 1, j, k)] +
15                   A0[Index3D (nx, ny, i - 1, j, k)])*c1
16              - A0[Index3D (nx, ny, i, j, k)]*c0;
17          }
18      }
19  }

```

Listing 2.6: Parboil parallel source code for stencil

```

1  #define Index3D(_nx,_ny,_i,_j,_k) ((_i)+_nx*((_j)+_ny*( _k)))

```

Listing 2.7: Converts from 3D indices to array-represented cube

The source codes above, specially the listing 2.7, show that the access to the cube elements are made by using the inner loop counter to multiply some constant and then map the cube to an array. The sequential access to memory is lost and the performance is impaired by successive cache misses.

Basically, we exchanged the index from the inner loop with the index of the outer loop. Then we simply parallelized the outer loop, which made our version run around 25.4 times faster than sequential version and around 15.2 times faster than the stock parallel version.

```

1  int i, j, k;
2  #pragma omp parallel for private(j, i)
3  for(k=1;k<nz-1;k++) {
4      for(j=1;j<ny-1;j++) {
5          for(i=1;i<nx-1;i++) {
6              Anext[Index3D (nx, ny, i, j, k)] =
7                  (A0[Index3D (nx, ny, i, j, k + 1)] +
8                   A0[Index3D (nx, ny, i, j, k - 1)] +
9                   A0[Index3D (nx, ny, i, j + 1, k)] +
10                  A0[Index3D (nx, ny, i, j - 1, k)] +
11                  A0[Index3D (nx, ny, i + 1, j, k)] +
12                  A0[Index3D (nx, ny, i - 1, j, k)])*c1
13              - A0[Index3D (nx, ny, i, j, k)]*c0;
14          }
15      }
16  }
17  }

```

Listing 2.8: Our parallel source code for stencil



## 2.3 mri-gridding

```

1  unsigned int k;
3  for(k=0; k<size; ++k){
4      // compute value to evaluate kernel at
5      // v in the range 0:(_width/2)^2
6      v = (((float)k)/((float)size))*cutoff2;
7
8      // compute kernel value and store
9      (*LUT)[k] = kernel_value_CPU(beta*sqrt
10         (1.0-(v/cutoff2)));
11  }

```

Listing 2.9: Auxiliary function source code for sequential mri-gridding

```

2  unsigned int k;
3  #pragma omp parallel for private(v)
4  for(k=0; k<size; ++k){
5      // compute value to evaluate kernel at
6      // v in the range 0:(_width/2)^2
7      v = (((float)k)/((float)size))*cutoff2;
8
9      // compute kernel value and store
10     (*LUT)[k] = kernel_value_CPU(beta*sqrt
11        (1.0-(v/cutoff2)));
12 }

```

Listing 2.10: Auxiliary function source code for parallel mri-gridding

```

1  int i;
3  for (i=0; i < n; i++){
4      ReconstructionSample pt = sample[i];
5
6      float kx = pt.kX;
7      float ky = pt.kY;
8      float kz = pt.kZ;
9  }

```

Listing 2.11: Parboil sequential main loop for mri-gridding

```

2  int i;
3
4  #pragma omp parallel for private(NxL, NxH,
5      NyL, NyH, NzL, NzH, dz2, nz, dx2, nx,
6      dy2, ny, idxZ, idxY, dy2dz2, idx0, v,
7      idx, w)
8  for (i=0; i < n; i++){
9      ReconstructionSample pt = sample[i];
10
11     float kx = pt.kX;
12     float ky = pt.kY;
13     float kz = pt.kZ;
14 }

```

Listing 2.12: Parboil parallel main loop for mri-gridding

```

1  /* grid data */
2  gridData[idx].real += (w*pt.real);
3  gridData[idx].imag += (w*pt.imag);
4
5  /* estimate sample density */
6  sampleDensity[idx] += 1.0;
7  }

```

Listing 2.13: Parboil sequential data processing code for mri-gridding

```

1  /* grid data */
2  #pragma omp critical (c1)
3  gridData[idx].real += (w*pt.real);
4  #pragma omp critical (c2)
5  gridData[idx].imag += (w*pt.imag);
6
7  /* estimate sample density */
8  #pragma omp critical (c3)
9  sampleDensity[idx] += 1.0;
10 }

```

Listing 2.14: Parboil parallel data processing code for mri-gridding

For this algorithm we basically allocated the lock-free structure and the accumulations are made privately and without concurrency for each thread. When all threads have finished the work, the final result is obtained by accumulating the partial results from each thread. Luckily, there is no loop-carried dependency, and the partial and the final accumulations are high parallelizable.

Our parallel version, with the lock-free vector accumulation, have a execution time around 2.5 times faster than the sequential and around 41.6 times faster than the stock parallel version. This happens because of the critical section pragmas, as we can see on listing 2.14. Each thread must execute this regions with a lock to obtain the correct result. This way, all other threads wait to get the lock and enter the critical region. By using this lock with 8 threads, there is a lot of concurrency and racing conditions, making the stock parallel the slowest version we tested, slower even the sequential version.

```

1  #pragma omp parallel private(NxL, NxH, NyL, NyH, NzL, NzH, dz2, nz, dx2, \
3  nx, dy2, ny, idxZ, idxY, dy2dz2, idx0, v, idx, w, t_id)
4  {
5  #pragma omp single
6  {
7      if (big_table == NULL)
8      {
9          //      printf("Alocou!\n");
10         size = params.gridSize[0]*params.gridSize[1]*params.gridSize[2];
11         numthreads = omp_get_max_threads();
12         big_table = (cmplx **) malloc(numthreads * sizeof(cmplx *));
13         big_table2 = (float **) malloc(numthreads * sizeof(float *));
14     }
15 }
16
17 #pragma omp for
18 for (i = 0; i < numthreads; i++)
19 {
20     big_table[i] = (cmplx *) malloc(size * sizeof(cmplx));
21     big_table2[i] = (float *) malloc(size * sizeof(float));
22 }
23 #pragma omp for collapse(2)
24 for (i = 0; i < numthreads; i++)
25     for (j = 0; j < size; j++)
26     {
27         big_table[i][j].real = 0;
28         big_table[i][j].imag = 0;
29         big_table2[i][j] = 0;
30     }
31
32 t_id = omp_get_thread_num();
33
34 #pragma omp for
35 for (i = 0; i < n; i++)
36 {
37     ReconstructionSample pt = sample[i];
38
39     float kx = pt.kX;
40     float ky = pt.kY;
41     float kz = pt.kZ;

```

Listing 2.15: Our parallel lock-free structure allocation for mri-gridding

```

1  /* grid data */
2  big_table[t_id][idx].real += (w*pt.real);
3  big_table[t_id][idx].imag += (w*pt.imag);
4
5  /* estimate sample density */
6  big_table2[t_id][idx] += 1.0;

```

Listing 2.16: Our parallel partial data processing code for mri-gridding

```

1  #pragma omp for private(j)
2  for (i = 0; i < size; i++)
3      for (j = 0; j < numthreads; j++)
4      {
5          gridData[i].real += big_table[j][i].real;
6          gridData[i].imag += big_table[j][i].imag;
7          sampleDensity[i] += big_table2[j][i];
8      }

```

Listing 2.17: Our parallel final data processing code for mri-gridding

## 2.4 cutcp

Basically, we implemented an auxiliary data structure that allowed for a lock-free version of the core routine of the benchmark, which eliminated the critical section pointed out in line 12. By doing that, we accomplished a speed up of around 25% in relation to the stock parallelization (omp\_base).

```

1  for (gindex = 0; gindex < ncell; gindex++) {
2      for (n = first[gindex]; n != -1; n = next[n]) {
3          //ommitted code
4          for (k = ka; k <= kb; k++, dz += gridspacing) {
5              //ommitted code
6              for (j = ja; j <= jb; j++, dy += gridspacing) {
7                  //ommitted code
8                  for (i = ia; i <= ib; i++, pg++, dx += gridspacing) {
9                      r2 = dx*dx + dy*dy;
10                     s = (1.f - r2 * inv_a2);
11                     e = q * (1/sqrtf(r2)) * s * s;
12                     *pg += e;
13                 }
14             }
15         }
16     }
17 }

```

Listing 2.18: cutcp serial version source

```

1  #pragma omp parallel for private ( ... )
2  for (gindex = 0; gindex < ncell; gindex++) {
3      for (n = first[gindex]; n != -1; n = next[n]) {
4          //ommitted code
5          for (k = ka; k <= kb; k++, dz += gridspacing) {
6              //ommitted code
7              for (j = ja; j <= jb; j++, dy += gridspacing) {
8                  //ommitted code
9                  for (i = ia; i <= ib; i++, pg++, dx += gridspacing) {
10                     s = (1.f - r2 * inv_a2);
11                     e = q * (1/sqrtf(r2)) * s * s;
12                     #pragma omp atomic
13                     *pg += e;
14                 }
15             }
16         }
17     }
18 }

```

Listing 2.19: Original parallel cutcp source code

```

1  nthreads = omp_get_max_threads();
2  big_table = (float **) malloc(nthreads * sizeof(float *));
3  for (i = 0; i < nthreads; i++)
4      big_table[i] = calloc(lattice->size, sizeof(float));
5
6  #pragma omp parallel private ( ... ) {
7      g = big_table[omp_get_thread_num()];
8      #pragma omp for
9      for (gindex = 0; gindex < ncell; gindex++) {
10         for (n = first[gindex]; n != -1; n = next[n]) {
11             //ommitted code
12             for (k = ka; k <= kb; k++, dz += gridspacing) {
13                 //ommitted code
14                 for (j = ja; j <= jb; j++, dy += gridspacing) {
15                     //ommitted code
16                     for (i = ia, m = index; i <= ib; i++, m++, dx += gridspacing) {
17                         pg[m] += (q/sqrtf(r2)) * (1.f - r2 * inv_a2) * (1.f - r2 * inv_a2);
18                     }
19                 }
20             }
21         }
22     }
23
24     #pragma omp for private(i)
25     for (m = 0; m < nthreads; m++)
26         for (i = 0; i < lattice->size; i++)
27             lattice->lattice[i] += big_table[m][i];
28 }

```

Listing 2.20: Our parallel source code for cutcp

## 2.5 tpacf

Here we carried out a strategy very akin to that indicated in the last section, which mainly consisted of rethinking the vector accumulation in such a way that critical sections could be circumvented.

```

1 int doCompute( ... ){
2     for (i = 0; i < ((doSelf) ? n1-1 : n1); i
3         ++){
4         for (j = ((doSelf) ? i+1 : 0); j < n2; j
5             ++){
6             while (max > min+1){
7                 k = (min + max) / 2;
8                 if (dot >= binb[k])
9                     max = k;
10                else
11                    min = k;
12            };
13            if (dot >= binb[min])
14                data_bins[min] += 1;
15            else if (dot < binb[max])
16                data_bins[max+1] += 1;
17            else
18                data_bins[max] += 1;
19        }
20    }
21 }

```

Listing 2.21: tpacf serial version source

```

1 int doCompute( ... ){
2     for (i = 0; i < ((doSelf) ? n1-1 : n1); i
3         ++){
4         #pragma omp parallel for
5         for (j = ((doSelf) ? i+1 : 0); j < n2; j
6             ++){
7             while (max > min+1){
8                 k = (min + max) / 2;
9                 if (dot >= binb[k])
10                    max = k;
11                else
12                    min = k;
13            };
14            #pragma omp critical
15            if (dot >= binb[min])
16                data_bins[min] += 1;
17            else if (dot < binb[max])
18                data_bins[max+1] += 1;
19            else
20                data_bins[max] += 1;
21        }
22    }
23 }

```

Listing 2.22: tpacf stock parallel solution

```

1 int doCompute( ... ){
2     #pragma omp parallel{
3         #pragma omp single{
4             if (doSelf){
5                 n2 = n1;
6                 data2 = data1;
7             }
8             size = ((doSelf) ? n1-1 : n1);
9             nthreads = omp_get_num_threads();
10            big_table = (int **) malloc(nthreads * sizeof(int *));
11        }
12        #pragma omp for
13        for (i = 0; i < nthreads; i++){
14            big_table[i] = (int *) malloc((nbins + 1) * sizeof(int));
15        }
16        #pragma omp for private(j)
17        for (i = 0; i < ((doSelf) ? n1-1 : n1); i++){
18            for (j = ((doSelf) ? i+1 : 0); j < n2; j++){
19                while (max > min+1){
20                    k = (min + max) / 2;
21                    if (dot >= binb[k]) max = k;
22                    else min = k;
23                };
24                if (dot >= binb[min]) k = min;
25                else if (dot < binb[max]) k = max+1;
26                else k = max;
27                big_table[omp_get_thread_num()][k]++;
28            }
29            #pragma omp for private(j)
30            for (i = 0; i < nbins + 1; i++){
31                for (j = 0; j < nthreads; j++){
32                    data_bins[i] += big_table[j][i];
33                }
34            }
35            free(big_table[i]);
36        }
37    }
38    free(big_table);
39 }

```

Listing 2.23: Our tpacf parallel version