




POLITECNICO
MILANO 1863

Benchmark porting on GPU using OpenACC

Alessandro Artoni

The objective

The scope of the project is to take a program and to parallelize it using openACC. 

Program: <https://github.com/yuhc/gpu-rodinia/tree/master/openmp/streamcluster>

Pc components:

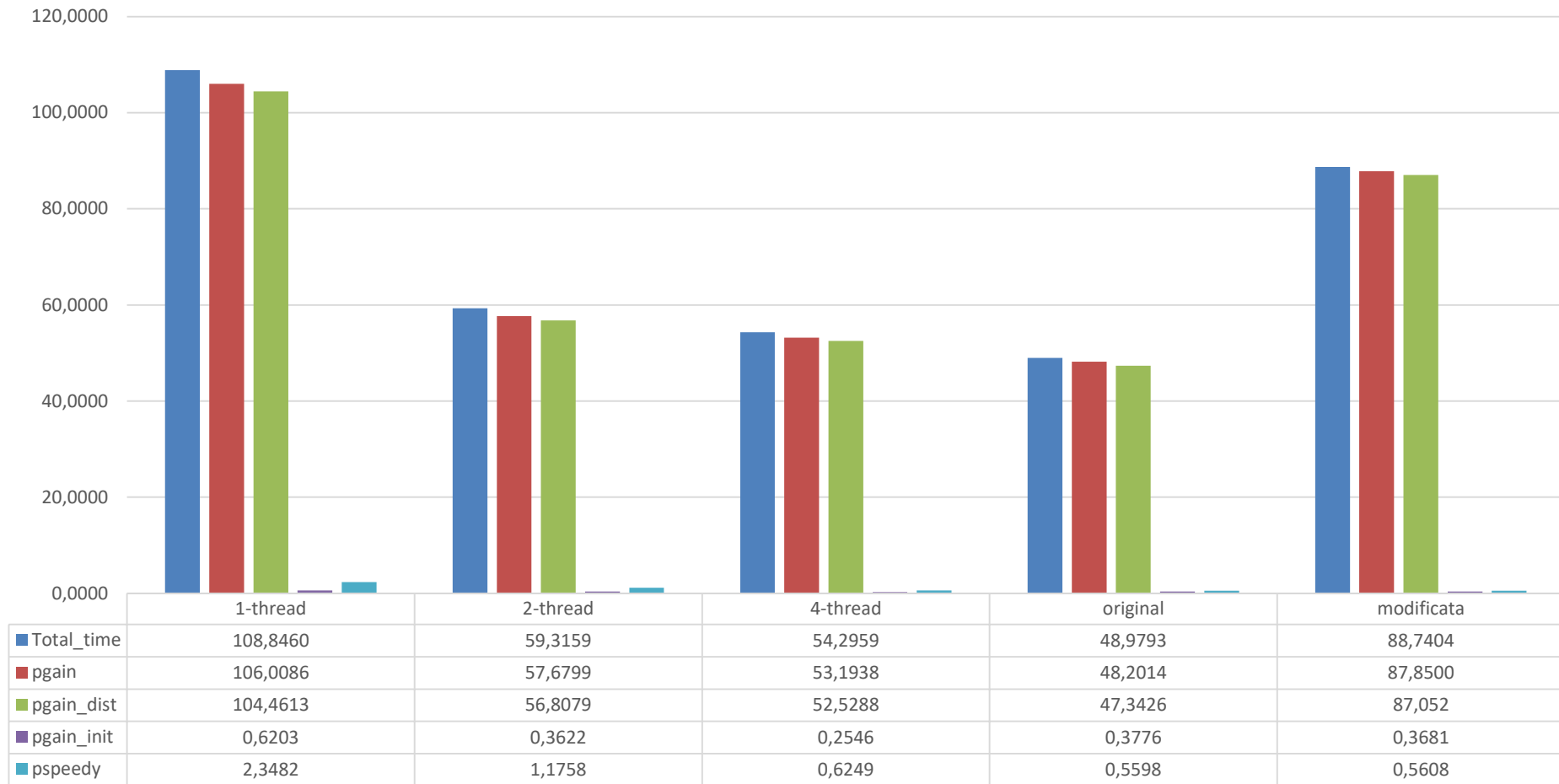
```
Detected 1 CUDA Capable device(s)

Device 0: "GeForce 920MX"
  CUDA Driver Version / Runtime Version      9.1 / 9.1
  CUDA Capability Major/Minor version number: 5.0
  Total amount of global memory:             2048 MBytes (2147483648 bytes)
  ( 2) Multiprocessors, (128) CUDA Cores/MP: 256 CUDA Cores
  GPU Max Clock rate:                        993 MHz (0.99 GHz)
  Memory Clock rate:                         900 Mhz
  Memory Bus Width:                          64-bit
  L2 Cache Size:                             1048576 bytes
```

Produttore: ASUSTek Computer Inc.
Modello: X556UV
Processore: Intel(R) Core(TM) i7-6500U CPU @ 2.50GHz 2.59 GHz
Memoria installata (RAM): 12,0 GB (11,9 GB utilizzabile)

Note: This processor has only 2 physical cores!

Analysis of the code



Time is expressed in seconds

Analysis of the code

Analysing the time the program took to execute, it has been noticed that the most expansive time part of the program was the pgain cycle. In particular this part of the loop:



```
459 for ( i = k1; i < k2; i++ ) {
460     float x_cost = dist(points->p[i], points->p[x], points->dim)
461         * points->p[i].weight;
462     float current_cost = points->p[i].cost;
463
464     if ( x_cost < current_cost ) {
465
466         // point i would save cost just by switching to x
467         // (note that i cannot be a median,
468         // or else dist(p[i], p[x]) would be 0)
469
470         switch_membership[i] = 1;
471         cost_of_opening_x += x_cost - current_cost;
472
473     } else {
474
475         // cost of assigning i to x is at least current assignment cost of i
476
477         // consider the savings that i's **current** median would realize
478         // if we reassigned that median and all its members to x;
479         // note we've already accounted for the fact that the median
480         // would save z by closing; now we have to subtract from the savings
481         // the extra cost of reassigning that median and its members
482         int assign = points->p[i].assign;
483         lower[center_table[assign]] += current_cost - x_cost;
484     }
485 }
```

Bottom-up optimization

At this point, it has been used a bottom up approach to let the GPU handle that cycle, starting from the `pgain_dist` loop using structured data region and kernels.

In order to manage better the memory from the host (cpu) to the device, some temporary arrays were introduced.

Because of this approach, the data were copied in and out of the device each time they were used inside the parallel region.

This resulted in a lot of time spent in data transfer from the CPU and the GPU.

Thus, it was decided to use unstructured data region, and «`present_or_copy`» clause. Since `present` table sometimes was giving problems, only «`update`» pragmas were used.

Then, data entering / exit were firstly made before `pgain` call (inside `pFL`), then inside `localSearch`.

Finally, compiling with `-ta=tesla:managed` gave reasonable results.

Bottom-up optimization

Structured data region

```
547  memset(temp_diff, 0, k2*sizeof(double));
548  #pragma acc data copy(temp_diff[k1:k2]) copyin(dim_points) copyin(points[0:1]) copyin(points->weight_p[0:k2])\
549    copyin(points->assign_p[0:k2]) copyin(points->cost_p[0:k2]) copyin(points->coord_p[k1:k2][0:dim_points]) \
550    copy(switch_membership[k1:k2]) copy(cost_of_opening_x) copy(result)//copyin(current_point_coords[0:dim_points]) copy(center_table[0:centerDim])
551  {
552    #pragma acc kernels// reduction(+: result)
553    {
554      for ( i = k1; i < k2; i++ ) {
555        result=0.0;
556        //#pragma acc loop reduction(+: result)
557        //{
558          for (int j=0;j<dim_points;j++)
559          {
560            result += (points->coord_p[i][j] - points->coord_p[x][j])*(points->coord_p[i][j] - points->coord_p[x][j]);
561            /*#ifdef INSERT_WASTE
562              double s = waste(result);
563              result += s;
564              result -= s;
565            #endif*/
566          }
567          float x_cost = result* points->weight_p[i];
568          float current_cost = points->cost_p[i];
569          float temp = x_cost - current_cost;
570          if ( temp < 0 ) {
571            // point i would save cost just by switching to x(note that i cannot be a median,or else dist(p[i], p[x]) would be 0)
572            switch_membership[i] = 1;
573            cost_of_opening_x += temp;
574          } else {
575            temp_diff[i] = temp*(-1);
576          }
577        }
578      }
```

Bottom-up optimization

Unstructured data region

```
#pragma acc enter data copyin(temp_diff[k1:k2]) copyin(points[0:1]) \
    copyin(points->assign_p[0:k2]) copyin(points->cost_p[0:k2]) \
    copyin(switch_membership[k1:k2]) \
    pcopyin(points->weight_p[0:k2]) pcopyin(points->coord_p[k1:k2][0:dim_points]) \
    copyin(result) copyin(dim_points) // copyin(cost_of_opening_x)

#pragma acc declare present(points[0:1], switch_membership[0:k2], temp_diff[k1:k2])

#pragma acc kernels// reduction(+: result)
{
    for ( i = k1; i < k2; i++ ) {
        result=0.0;
        // #pragma acc loop reduction(+: result)
        //{
            for (int j=0;j<dim_points;j++)
            {
                result += (points->coord_p[i][j] - points->coord_p[x][j])*(points->coord_p[i][j] - points->coord_p[x][j]);
                /*#ifndef INSERT_WASTE
                double s = waste(result);
                result += s;
                result -= s;
                #endif*/
            }
        // }
        //}
        float x_cost = result* points->weight_p[i];
        float current_cost = points->cost_p[i];
        float temp = x_cost - current_cost;
        if ( temp < 0 ) {
            // point i would save cost just by switching to x(note that i cannot be a median,or else dist(p[i], p[x]) would be 0)
            switch_membership[i] = 1;
            cost_of_opening_x += temp;
        } else {
            temp_diff[i] = temp*(-1);
        }
    }
}

#ifdef DEBUGGER
    fprintf(stdout, "PostPragma\n" );
#endif

#pragma acc exit data copyout(temp_diff[k1:k2]) copyout(switch_membership[0:k2]) \
    copyout(result) copyout(points->assign_p[0:k2]) //copyout(cost_of_opening_x)

#pragma acc exit data delete(temp_diff[k1:k2]) delete(dim_points) delete(points[0:1]) delete(points->weight_p[0:k2])\
    delete(switch_membership[k1:k2]) delete(result) \
    delete(cost_of_opening_x)
```

Bottom-up optimization

Going-up: data entering before pgain

```
#pragma acc enter data copyin(points[0:1])
#pragma acc enter data copyin(points->weight_p[0:k2])
#pragma acc enter data copyin(points->coord_p[0:k2][0:dim_points])
#pragma acc enter data create(points->assign_p[0:k2], points->cost_p[0:k2], switch_membership[0:k2])
#ifdef DEBUGGER
fprintf(stdout, "Post acc enter data main ext loop\n");
#endif
for (i=0;i<iter;i++) {
    x = i%numfeasible;
    //printf("iteration %d started*****\n", i);
    change += pgain(feasible[x], points, z, k, pid, barrier);
    c++;

    //printf("iteration %d/%d finished @@@@@"n", i,iter);
}
//sembra non sia neanche necessario copiare indietro
#pragma acc exit data copyout(points->weight_p[0:k2])
#pragma acc exit data copyout(points->coord_p[0:k2][0:dim_points])
#pragma acc exit data copyout(points->assign_p[0:k2], points->cost_p[0:k2], switch_membership[0:k2], points[0:1])
#pragma acc exit data delete(points->weight_p[0:k2])
```


Bottom-up optimization

Going-up: data entering before pFL

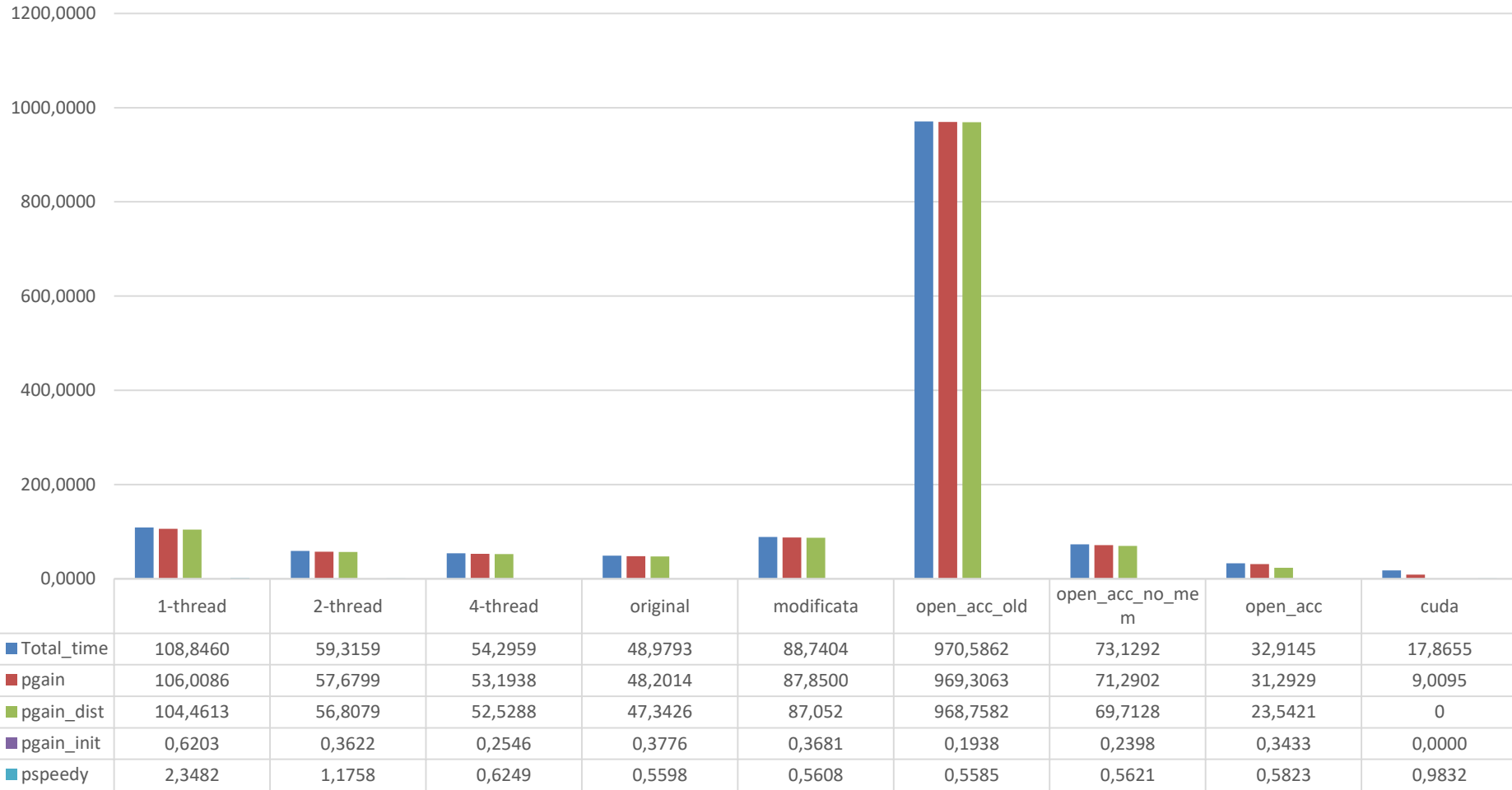
```
1148 //creating space on the device
1149 #pragma acc enter data copyin(points[0:1],points->weight_p[0:points_num],points->coord_p[0:points_num][0:dim_points])
1150 #pragma acc enter data copyin(temp_diff[0:points_num])
1151 #pragma acc enter data create(points->assign_p[0:points_num], points->cost_p[0:points_num], switch_membership[0:points_num])
1152
1153 while(1) {
1154     d++;
1155 #ifdef PRINTINFO
1156     if( pid==0 )
1157     {
1158         printf("loz = %lf, hiz = %lf\n", loz, hiz);
1159         printf("Running Local Search...\n");
1160     }
1161 #endif
1162     /* first get a rough estimate on the FL solution */
1163     // pthread_barrier_wait(barrier);
1164
1165     //lastcost = cost;
1166
1167     cost = pFL(points, feasible, numfeasible,
1168         z, &k, cost, (long)(ITER*kmax*log((double)kmax)), 0.1, pid, barrier);
1169
1170     /* if number of centers seems good, try a more accurate FL */
1171     if (((k <= (1.1)*kmax)&&(k >= (0.9)*kmin)) ||
1172         ((k <= kmax+2)&&(k >= kmin-2))) {
1173
1174 #ifdef PRINTINFO
1175         if( pid== 0)
1176         {
1177             printf("Trying a more accurate local search...\n");
1178         }
1179 #endif
```

Using GPU parallelism

At this point, it was tried to exploit the GPU thread/warp/blocks. The division in the bottom picture, resulted the best one after several tries.

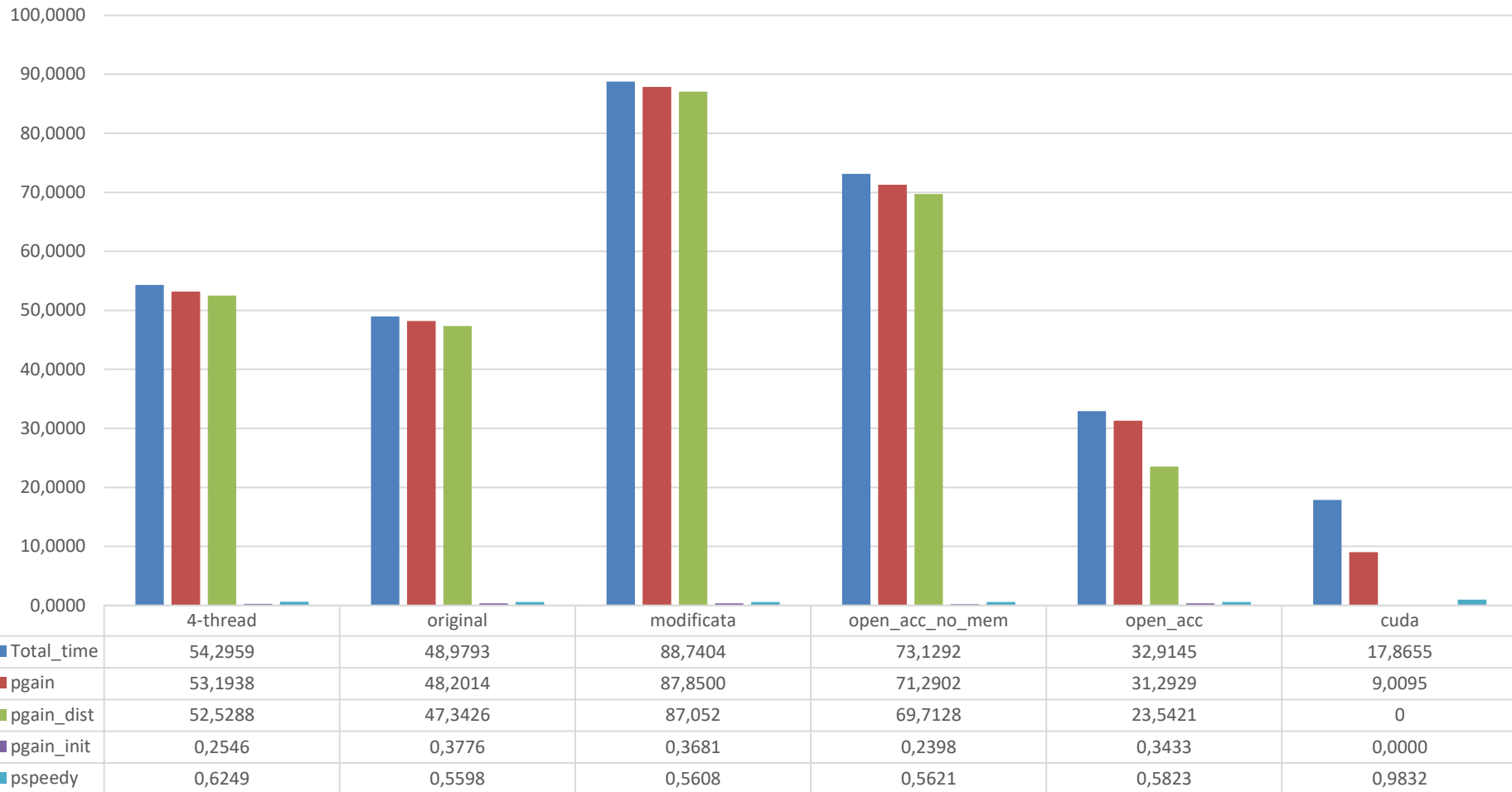
```
567 #pragma acc declare present(points[0:1], switch_membership[0:k2], temp_diff[k1:k2]) \
568 present(points->weight_p[0:k2], points->coord_p[k1:k2][0:dim_points])
569
570 #pragma acc parallel loop reduction(+: result) device_type(nvidia) gang vector_length(64)//
571 {
572     for ( i = k1; i < k2; i++ ) {
573         result=0.0;
574         #pragma acc loop gang vector //vector_length(64)//gang(16) vector_length(32) // vector
575         for (int j=0;j<dim_points;j++)
576         {
577             result += (points->coord_p[i][j] - points->coord_p[x][j])*(points->coord_p[i][j] - points->coord_p[x][j]);
```

Final results



Time is expressed in seconds.

Final results - filtered



Time is expressed in seconds.

Conclusions

Since the pc where original code using threads were tested didn't have more than two physical cores, the algorithm using thread didn't perform really well. Using other CPUs, performance significantly improves arriving to near 10 seconds – which is even better than CUDA!

The limits of OpenACC is the lack of possibility to allocate exactly the amount of block/threads needed without having to do a significant refactor of the code. This is also a reason why CUDA's code works better: the code was completely rebuilt around it.

Another limit of the solution proposed, is the limit of memory usable on the GPU.

However, OpenACC focus on portability, and the version should be able to perform well on different devices.