

CP - TP1

1st Gonalo dos Santos
Departamento de Informtica
Universidade do Minho
Braga, Portugal
pg50396

2nd Lus Faria
Departamento de Informtica
Universidade do Minho
Braga, Portugal
pg50560

I. INTRODUCTION

This project’s objective was to develop, in C, a sequential and optimized version of a simple k-means algorithm, based on Lloyd’s algorithm.

The algorithm can be resumed in the following steps:

- 1) Create samples and initiate *clusters*;
 - a) Initiate a vector with random values of N size in a (x,y) space;
 - b) Initiate the the K clusters with the coordinates of the first sample;
 - c) Attribute each sample to the nearest cluster using the euclidean distance;
- 2) Calculate each *cluster*’s centroid (also known as the geometric center);
- 3) Attribute each sample to the nearest *cluster* using the euclidean distance;
- 4) Repeat the second and third steps until no point changes *cluster*.

II. IMPLEMENTATION AND OPTIMIZATION

While planning the implementation of the project, the group decided that the first version of the program must already be as optimized as our knowledge allows it.

Which made our first version be able to vectorize all of the iterations in the algorithm and make the most use of space and temporal locality by using the same array to algorithm related variables and guarantee the continuity by not altering the same variable in consecutive operations.

A. Algorithm

After the first implementation, we tried to to optimize the algorithm and decided to calculate the distances with this formula:

$$x^2 + y^2 \quad (1)$$

That is based on the euclidean distance formula:

$$\sqrt{x^2 + y^2} \quad (2)$$

The group did that because the algorithm only compares the distances and doesn’t need the specific values and supporting our decision is the following equivalence:

$$x < y \equiv \sqrt{x} < \sqrt{y} \quad x, y \in \mathbb{R} \geq 0 \quad (3)$$

This decision was also based on the fact that the mathematical function **sqr()** is computationally heavy.

The stats after this optimization can be seen in figure 1.

B. Code

After optimizing the algorithm, the next step was to analyse the code, where we noticed that we were iterating the same list twice.

We tried to integrate those two iterations into one and noticed that the execution time had gone down and the iteration was no longer vectorized but unrolled, which was could be seen in the increase of instructions, as we expected.

In this optimization the group obtained the stats shown in figure 2.

After the previous optimization, the group didn’t find any part of the code that could be optimized with our experience at that moment, so we experimented changing some parts and found a few cases of optimization.

One of the cases we found was changing the expression

flag = (old_points[i] != ind)?1:flag

to

if(flag) flag = (old_points[i] != ind)?1:flag

and noticed that even though we used **?**, that is less computationally heavy than **if**, because we had used a memory access as the sole condition, it created a lot of cache misses and by introducing a **if(flag)**, cache misses decreased massively. We end up simplifying the expression to

flag = (flag or old_points[i] != ind)

and obtained the stats shown in figure 3.

The group’s last optimization was to delete one of the the two lists we had dedicated to store the old and the attribution of the points to the clusters, to do the comparisons. We did it because we noticed that one list was enough to serve it’s purpose and that by doing so there were marginally less cache misses and branch misses. The figure 4 shows the program final stats.

C. Compilation

To compile the program, we used the following flags:

-std=c99 -O2 -ftree-vectorize -msse4 -funroll-loops

We used the flag **-std=99** because we were required to compile the code using the c99 version.

We picked the **O2** compilation flag to apply the most generic optimizations that the compiler has available, without it increasing the size of the generated binary file, and because it gave us the best results, comparing to the **O1** and **O3** flags, as we can see from comparing the stats in the figures ??, ?? and ??.

To support the **O2** flag, the group used **ftree-vectorize** to allow the vectorization of instructions together with **msse4** that indicates to the compiler that it has available records of 256 bits size and with **funroll-loops** to unroll the loops that the compiler considers beneficial to performance.

III. CONCLUSION

We finished this first project satisfied with what we accomplished and pleased about reinforcing the previously learned knowledge about the subjects of optimization of code and it's compilation, as well as further deepening our understanding of the subjects.

We would like to further improve this project with the the knowledge that will be deepened with the next objectives of this class.

IV. ANEXOS

```
Performance counter stats for './bin/k_means' (3 runs):
      112,894,394      L1-dcache-load-misses      ( +- 0.02% ) (71.42%)
      15,950,295,765      cycles      ( +- 0.06% ) (71.42%)
      1,515,067,860      branches      ( +- 0.03% ) (71.42%)
      12,341,554      branch-misses      = 0.81% of all branches      ( +- 0.19% ) (71.43%)
      (not supported)      stalled-cycles-backend      ( +- 0.08% ) (85.72%)
      8,754,862,520      stalled-cycles-frontend      = 54.80% frontend cycles idle      ( +- 0.08% ) (85.72%)
      18,310,765,177      inst_retired.any      = 18310765177.0 Instructions      ( +- 0.04% ) (71.44%)
      15,947,691,856      cycles      = 0.9 CPI      ( +- 0.06% ) (71.44%)
      18,328,072,540      inst_retired.any      ( +- 0.03% ) (71.43%)

5.0397 +- 0.0658 seconds time elapsed ( +- 1.31% )
[pg50396@searchTedu2 TP-1]$
```

Fig. 1. Algorithm Optimization

```
Performance counter stats for './bin/k_means' (3 runs):
      115,491,399      L1-dcache-load-misses      ( +- 0.07% ) (71.41%)
      14,134,280,014      cycles      ( +- 0.06% ) (71.41%)
      1,519,097,424      branches      ( +- 0.06% ) (71.41%)
      12,535,565      branch-misses      = 0.83% of all branches      ( +- 0.14% ) (71.43%)
      (not supported)      stalled-cycles-backend      ( +- 0.17% ) (85.72%)
      3,384,405,171      stalled-cycles-frontend      = 37.53% frontend cycles idle      ( +- 0.04% ) (71.44%)
      24,321,383,374      inst_retired.any      = 24321383373.7 Instructions      ( +- 0.07% ) (71.44%)
      14,131,755,238      cycles      = 0.6 CPI      ( +- 0.05% ) (71.42%)
      24,310,329,090      inst_retired.any      ( +- 0.05% ) (71.42%)

4.432 +- 0.108 seconds time elapsed ( +- 2.43% )
[pg50396@searchTedu2 TP-1]$
```

Fig. 2. Iterations joined

```
Performance counter stats for './bin/k_means' (3 runs):
      85,469,446      L1-dcache-load-misses      ( +- 0.05% ) (71.42%)
      12,391,447,078      cycles      ( +- 0.04% ) (71.43%)
      1,136,437,343      branches      ( +- 0.11% ) (71.43%)
      749,811      branch-misses      = 0.07% of all branches      ( +- 0.57% ) (71.43%)
      (not supported)      stalled-cycles-backend      ( +- 0.13% ) (85.72%)
      4,278,358,218      stalled-cycles-frontend      = 34.53% frontend cycles idle      ( +- 0.01% ) (71.43%)
      24,382,508,407      inst_retired.any      = 24382508406.7 Instructions      ( +- 0.04% ) (71.43%)
      12,389,127,345      cycles      = 0.5 CPI      ( +- 0.02% ) (71.42%)
      24,379,305,272      inst_retired.any      ( +- 0.02% ) (71.42%)

3.8491 +- 0.0605 seconds time elapsed ( +- 1.57% )
[pg50396@searchTedu2 TP-1]$
```

Fig. 3. If statements simplified

```
Performance counter stats for './bin/k_means' (3 runs):
      84,845,064      L1-dcache-load-misses      ( +- 0.00% ) (71.41%)
      12,370,665,459      cycles      ( +- 0.04% ) (71.42%)
      1,136,157,361      branches      ( +- 0.03% ) (71.44%)
      742,429      branch-misses      = 0.07% of all branches      ( +- 0.52% ) (71.44%)
      (not supported)      stalled-cycles-backend      ( +- 0.12% ) (85.72%)
      4,461,376,328      stalled-cycles-frontend      = 36.07% frontend cycles idle      ( +- 0.01% ) (71.42%)
      24,759,232,157      inst_retired.any      = 24759232156.7 Instructions      ( +- 0.04% ) (71.42%)
      12,369,381,375      cycles      = 0.5 CPI      ( +- 0.01% ) (71.41%)
      24,756,920,115      inst_retired.any      ( +- 0.01% ) (71.41%)

3.7894 +- 0.0104 seconds time elapsed ( +- 0.28% )
[pg50396@searchTedu2 TP-1]$
```

Fig. 4. One point-cluster attribution list

```
src/k_means.c:10:5: note: loop turned into non-loop; it never loops
[pg50396@searchTedu2 TP-1]$ make test
srtn --partition=par perf stat -e L1-dcache-load-misses,cycles,branches,branch-misses,,stalled-cycles-backe
nd,stalled-cycles-frontend -M cpi,instructions ./bin/k_means > /dev/null

Performance counter stats for './bin/k_means':
      82,572,994      L1-dcache-load-misses      (71.43%)
      23,157,310,264      cycles      (71.43%)
      1,138,881,332      branches      (71.43%)
      874,321      branch-misses      = 0.08% of all branches
      (not supported)      stalled-cycles-backend      (71.43%)
      15,057,321,990      stalled-cycles-frontend      = 65.03% frontend cycles idle (85.71%)
      25,653,366,181      inst_retired.any      = 25653366181.0 Instructions
      23,151,947,411      cycles      = 0.9 CPI
      25,654,424,163      inst_retired.any

7.540661503 seconds time elapsed
7.501778000 seconds user
0.037998000 seconds sys
```

Fig. 5. flag -O1 compilation

```
src/k_means.c:10:5: note: loop turned into non-loop; it never loops
[pg50396@searchTedu2 TP-1]$ make test
srtn --partition=par perf stat -e L1-dcache-load-misses,cycles,branches,branch-misses,,stalled-cycles-backe
nd,stalled-cycles-frontend -M cpi,instructions ./bin/k_means > /dev/null

Performance counter stats for './bin/k_means':
      85,962,280      L1-dcache-load-misses      (71.41%)
      13,446,113,916      cycles      (71.41%)
      1,139,497,594      branches      (71.41%)
      786,177      branch-misses      = 0.07% of all branches      (71.44%)
      (not supported)      stalled-cycles-backend      (71.44%)
      4,540,137,697      stalled-cycles-frontend      = 33.81% frontend cycles idle (85.73%)
      25,170,249,511      inst_retired.any      = 25170249511.0 Instructions
      13,444,231,777      cycles      = 0.5 CPI
      25,164,064,307      inst_retired.any

4.332410159 seconds time elapsed
4.302426000 seconds user
0.029002000 seconds sys
```

Fig. 6. Of flag -O2 compilation

```
pg50271 pts/27 2022-10-28 15:46 (b147-156-134.ds1.telepac.pt)
pg50740 pts/37 2022-10-28 15:45 (a89-153-208-33.cpe.netcabo.pt)
[pg50396@searchTedu2 TP-1]$ make test
srtn --partition=par perf stat -e L1-dcache-load-misses,cycles,branches,branch-misses,,stalled-cycles-backe
nd,stalled-cycles-frontend -M cpi,instructions ./bin/k_means > /dev/null

Performance counter stats for './bin/k_means':
      85,328,358      L1-dcache-load-misses      (71.42%)
      16,432,483,153      cycles      (71.42%)
      1,615,597,252      branches      (71.42%)
      103,091,279      branch-misses      = 6.38% of all branches      (71.42%)
      (not supported)      stalled-cycles-backend      (71.42%)
      4,884,017,363      stalled-cycles-frontend      = 29.72% frontend cycles idle (85.71%)
      24,858,767,487      inst_retired.any      = 24858767487.0 Instructions
      16,430,260,926      cycles      = 0.7 CPI
      24,863,984,209      inst_retired.any

5.320393559 seconds time elapsed
5.297534000 seconds user
0.022002000 seconds sys
```

Fig. 7. flag -O3 compilation