

000
001
002
003
004
005
006
007
008
009
010
011

Abstract

In this mid-term project I have implemented K-Means clustering and then I have applied it to image in order to achieve a color quantization as visual result. I have realized three different implementations in Java, one sequential and two parallel, and then compared the execution time to obtain and consider the speed up.

Future Distribution Permission

The author(s) of this report give permission for this document to be distributed to Unifi-affiliated students taking future courses.

1. Introduction

K-Means, described in Section 2, is one of the simplest and most used algorithm for clustering points, although its performance, especially for its most basic version, is not so high. In general, the points must belong to an Euclidean space, like \mathbb{R}^n , (or a convex subset thereof) because the algorithm is based on the concepts of Euclidean distance and “average” between points; but is possible to generalize it for non-Euclidean space like string’s space, where we have a notion of distance, but there’s not an “average” between points.

In this project I have implemented K-Means considering the Euclidean case and for convenience I restrict the domain to \mathbb{R}^n and its hyper-rectangular subsets (which are the most common case, like in RGB space color introduced in following sections).

As application of this algorithm I have choosen the images: applying a clusterization method to an image (to be precise to the pixels of an im-

PC-2020/21 K-Means

Anonymous CVPR submission

Paper ID ****

age) we obtain a color quantization of the image, which can be a first result for image segmentation. I describe this aspect in Section 3.

In the following sections I present my implementation solutions, all developed in Java: the first of them is a sequential solution, the second and the third are parallel versions that use Java high level framework for concurrency, one handling explicitly the parallelism and one handling it implicitly. In the lasts sections I report the result of the performance tests, discussing the speed up and proposing my personal considerations about the efficiency of the solutions.

2. K-Means Clustering

K-Means clustering is a family of clustering methods based on the a priori knowledge of the number of clusters into which the given points must be classified. Some algorithms to deduce the optimal value of k exist, but this is outside the main topic of this project, so we assume to have k as input parameter (and we evaluate the main results also as value of k varies).

In this project we will consider the standard version of K-Means, also known as naïve K-Means [1]; the idea behind it is quite simple:

0. select k points $m_1^{(0)}, \dots, m_k^{(0)}$ from the dataset that we will consider as **centroids** of step-0 cluster;
1. assign each point of the dataset to the cluster identified with the nearest centroid;

$$S_i^{(t)} = \left\{ x_p \mid \|x_p - m_i^{(t)}\| \leq \|x_p - m_j^{(t)}\| \right\} \quad (1)$$

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

108 2. update the clusters' centroids with the aver-
 109 age of its points;

$$m_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{x_j \in S_i^{(t)}} x_j \quad (2)$$

115 3. repeat from step 1 as long as the clusters
 116 “change enough” respect to the previous it-
 117 eration.

119 This algorithm is guaranteed to converge (if we
 120 are using the Euclidean distance), though, as we
 121 can see, there are various aspects open to free
 122 choice, let's discuss them.

2.1. Clusters Initialization

127 About how to choose the initial centroids there
 128 are mainly two approaches: we can take as cen-
 129 troids k points in order to maximise their rela-
 130 tive distance or we can take as centroids the point
 131 nearest to the centroids of a raw clusterization
 132 (perhaps based on a sample of the dataset). I have
 133 adopted the first approach as in [2] and I report the
 134 algorithm:

```
136       Pick the first point at random;  

  137       while there are fewer than  $k$  points do  

  138           Add the point whose minimum distance from  

  139           the selected points is as large as possible;  

  140       end while
```

2.2. Stop Condition(s)

144 Also for the stopping rule is possible to define
 145 different conditions. Ideally we would like to stop
 146 when none of the data points change cluster from
 147 one iteration to the next one (that is we achieve the
 148 minimum sum of distances of every point from its
 149 cluster's centroid), since is guaranteed the conver-
 150 gence. As relaxation of this condition we could
 151 take as stopping rule the “convergence” of the di-
 152 ameter of the clusters (intended as the maximum
 153 distance between two points of a single cluster)
 154 or the position of centroids; obviously consider-
 155 ing as “convergence” theirs change under a cer-
 156 tain threshold.

157 In this project I assume as stop condition the
 158 latter case, taking as threshold the value 0.005.



Figure 1. An image and its color quantized version with 8 colors

3. Image Color Quantization

191 As already mentioned, with K-Means is possi-
 192 ble to get a color quantization of an image [3].
 193 This process aims to reduce the number of colors
 194 in an image, grouping the pixels in clusters and
 195 then assuming all pixels in a cluster of the same
 196 color, in particular the average color of the cluster,
 197 as it's possible to see in figure 1.

198 To do that pixels should be considered as points
 199 in the space of colors (thus loosing momentar-
 200 ily the information relative to the spatial coordi-
 201 nate inside the image) and applying clusteriza-
 202 tion algorithm to them. This process can be ap-
 203 plied to pixels from RGB color model, but also
 204 in gray scale model or others; in the first case
 205 the domain of the points is a subset of \mathbb{R}^3 (it
 206 would be $[0, 255]^3$), in the second a subset of \mathbb{R} (it
 207 would be $[0, 1]$). In this project we will consider
 208 the RGB case without loss of generality. More-
 209 over, as usual in practical use, the RGB values
 210 assume only integer values, so it isn't anymore an
 211 Euclidean space (e.g. the strings case mentioned

216 above); however during the computation I aban-
217 don the discrete representation in favor of the con-
218 tinuous one; only at the end I return to the finite
219 representation.
220

222 4. Solution Implementation

224 In my opinion the most suitable lan-
225 guage/framework for implementing K-Means
226 is CUDA, since most of calculation is indepen-
227 dent and so the algorithm doesn't require much
228 synchronization; and in addition, as generally
229 happens working with images, the typical dimen-
230 sion of datasets is really huge (*e.g.* for an image in
231 4K resolution and 16:10 ratio, the pixels amount
232 to 10.5M). Anyway I choose Java to compare the
233 sequential and the parallel versions developed
234 with the same language; and furthermore Java
235 offers, at least for the implementation of this
236 specific task, an explicit and an implicit version
237 of parallelism: it is interesting to compare them.
238

239 I have divided the algorithm in 3 subroutine:
240 the calculation of the initial centroids, then the
241 assignment of each point to a cluster and finally
242 the update of centroids. All this operation require
243 an iteration over all the points in dataset. The
244 main idea that has guided the parallelization of the
245 code is to parallelize this subroutines since they
246 are, at least this is what I assumed theoretically,
247 the hotspot in term of computation time (how-
248 ever we will discuss about this in the next Section
249 presenting the results of profiling of the sequen-
250 tial version). A fourth minor subroutine is that of
251 checking the stop condition, but this is an iteration
252 over the k centroids and typically k is a quite low
253 value; there is the possibility to parallelize also
254 this operation, but we expect a relatively small
255 improvement (if not a deterioration for too small
256 value of k due to parallelization overhead).

262 5. Sequential Solution

264 Now I present my sequential solution imple-
265 mented in Java; the code is available in [this](#)
266 [GitHub repo](#).

268 This solution is very basic and it is based
269 on the class *Point*, that represents a point in

270 \mathbb{R}^n by storing its coordinates in a float array;
271 a Point is immutable, so it can be shared be-
272 tween thread safely. The class Point implements
273 a static method *getSquaredEuclideanDistance* for
274 calculating the distance between two Points (we
275 save computation by taking the square of distance
276 because is totally indifferent in the algorithm).
277 Lastly this class can be inherited in order to rep-
278 resent particular cases of points: I have inherited
279 *RGBPoint*, that represents a tridimensional point
280 with coordinates in [0, 255], and *RGBPixel*, that
281 represents an RGBPoint with integer coordinates
282 and additional metadata for its (x, y) -coordinates
283 in the image¹.

285 To ensure that all the Points to be treated be-
286 long to the same domain, I have created the class
287 *SetOfPoints* consisting in a collection of Points
288 and a *Domain* to which the Points are constrained
289 to belong. A Domain, as initially said, it's an
290 hyper-rectangular subset of \mathbb{R}^n and so consist of
291 a dimension and two arrays of lower and upper
292 bounds for each of its dimension. A subclass of
293 *SetOfPoints* is *Cluster* that only add some meth-
294 ods typical of a cluster (like *getCenter* and *getDi-
295 ameter*).

296 The core of this implementation is the class
297 *KMeans*, genericized with a subclass of Point. Its
298 only public method *clusterize* take as parameters
299 a SetOfPoints (to this end I define a functoid class
300 *Image* to load an image and store it as SetOf-
301 Points) and the number k of desired cluster. Inter-
302 nally I abandon for convenience the SetOfPoints
303 wrapper, simply working with an array of Points.
304 This method is implemented exactly as mentioned
305 before: an initial call to the private method *initial-
306 Centroids* and then repeated calls to the private
307 methods *updateClusters* and *newCentroids*. Let's
308 see them more in depth:

- *initialCentroids*: from the dataset produces
a set of k points which are the initial cen-
troids. This method implements exactly the
algorithm shown in Section 2.1. The only

1 maybe in a general contest the composition of *RGBPixel* with *RGB-
Point* would have been better, but in this specific and simple case it turns
out to be more comfortable the inheritance: KMeans can be executed both
with points and pixels without any additional manipulation

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

simplification is that I take as first point not a random one, but exactly the first of the dataset, assuming that this is randomly sorted in the color space (in the case of images this assumption can be considered always satisfied). This trick allows to make the choice of initial centroids reproducible;

- *updateClusters*: consists simply in the iteration over the points and over the centroids, to determine which centroid is nearest to each point. It is immediately evident that this portion of code is embarrassingly parallel since every point is totally independent from the others. I represent the actual clusterization with an array of integer, where the i -nth element contains the number (from 1 to k) of the cluster to which the i -nth point is assigned.
- *newCentroids*: starting from the actual clusterization calculates the center of each cluster iterating on the points, summing its coordinates in a bidimensional array $k \times \text{dimension of the points}$ (where each row contains the sum of the coordinates of all points of that cluster) and taking into account the size of every clusters and then completing the average calculation.

The *clusterize* method then returns a list of Clusters

5.1. Profiling

Before starting the parallelization of the code, I profile the sequential version to spot the hotspot(s).

I have focused my attention on the three sub-methods, indeed together they form almost the totally of the computation. Clearly their respective percentage of the total computation change with image dimension and number of clusters. To show that and to find approximate average values I have profiled the execution in 3 different cases: a 4K image with 10 and 15 clusters and a 6K image with 10 clusters. The results can be found in the repo and they are summarized in Table 1.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

Clearly *initialCentroids* depends both on the points and k , so with the same dataset it takes more time as k grows, but since also the total execution time grows, the percentage remain almost the same. However the percentage is very low and decrease with dataset dimension.

updateClusters also depends both on points and k (it was easily predictable because at every iteration it calculates and evaluates $\text{number of points} \times \text{number of centroids}$ distances), indeed the execution time grows as dataset dimension and number of clusters grow; furthermore we can see that also the percentage over the total time (that is by far the predominant one) grows.

newCentroids depends only on points, since it was basically an iteration over them. Indeed the execution time grows with dataset's growth and remain the same with k 's growth (the percentage over total computation shrinks)

The most important aspect of these results is that, independently from percentage values, the main hotspot is by far *updateClusters*. And this is a great news, because it is embarrassingly parallel. The other two methods can be parallelized, but their small execution time (and need of synchronization) exposes us easily to a too high overhead respect to the gain.

A noteworthy fact is that analyzing the profiling results we can note that the lowest level hotspot is the method *getSquaredEuclideanDistance*: in fact it occupies more than 50% of the computation (although it was used only by *updateClusters* and *initialCentroids*).

6. Parallel Solution

As we said in Section 4, all three of this methods consist of an iteration on the dataset's points (or more iterations in the case of *initialCentroids*); so the guideline I followed to implement my solution, available on [this repo](#), is that of decompose the dataset. I have also thought to parallelize the calculation of distance (it consists of an iteration of the coordinates of the points), anyway I changed course immediately because it was

432	Image	k	Execution Time	initialCentroids	updateClusters	newCentroids	486
433	4K	10	57s	6%	85%	8.5%	487
434	4K	15	99s	6%	87%	6.5%	488
435	6K	10	322s	3.5%	88%	8.5%	489
436							490
437							491
438							492
439							493
440							494
441							495
442							496
443							497
444							498
445							499
446							500
447							501
448							502
449							503
450							504
451							505
452							506
453							507
454							508
455							509
456							510
457							511
458							512
459							513
460							514
461							515
462							516
463							517
464							518
465							519
466							520
467							521
468							522
469							523
470							524
471							525
472							526
473							527
474							528
475							529
476							530
477							531
478							532
479							533
480							534
481							535
482							536
483							537
484							538
485							539

Table 1. Profiling results

very fast in a single iteration and in all likelihood the parallelization overhead would have been too high respect to the gain. Moreover, so doing, I followed the good practice explained in [4] of parallelizing at the highest level possible.

6.1. Explicit Parallelism

The main idea was to create a single thread pool in the *clusterize* method and to submit to it dynamically three types of task for the three step of algorithm: *initialCentroidsTask*, *updateClustersTask* and *newCentroidsTask*. About the number of threads we have to consider that all this tasks are purely computational, so there should be a 1:1 ratio between threads and cores, as explained in [5]. I consider to put an extra thread to supply page fault events, but my test machines are equipped with Intel processors with Hyperthreading, so setting the number of threads equal to the number of virtual threads is sufficient to prevent page faults delays; anyway later we validate this assumption with tests.

6.1.1 initialCentroidsTask

The *initialCentroids* method executes for k times an iteration on the points. It can't be parallelized the outer iteration because every iteration depends on all the previous (that is for calculates every centroids we must know all the previous), so this task executes the iteration on a portion of the dataset. Since at every iteration of the task the main input, that is points and centroids arrays, doesn't change (at least from a Java perspective intended as the reference to the centroids array) I reduce overhead costs by instantiating the object just one time and reusing it updating/resetting only some aspects at every iteration. In particular before submitting the task I invoke the method

nextIteration passing two newly created AtomicInteger shared by all tasks representing the next identified centroid and its maximum minimum distance from the others (the reason of AtomicInteger type will be clear later).

An initial implementation expected each task to return a candidate centroid and then to choose the best sequentially. However, although this solution did not require synchronization, it was quite slow (even worse than the sequential version). The main problem is that the tasks should set *two* shared variable atomically. Anyway the solution was quite simple: I have associated the class to a Lock, every tasks find privately its candidate point to be a centroid and then in a critical section it checks if its candidate is better than the actual one and if it is, it update the shared variables. The choice of AtomicInteger is not for atomicity (they are read and written in a critical section), but because Java pass by value the primitive type (even if boxed): here I want this two value to be instantiate in main thread, written in forked threads and read again in main thread; a final note on this: the distance is a float, but can be represented with an integer relative to its bit representation.

A final consideration about the points and centroids arrays: they are shared between all tasks, but they are accessed only for reading, so they can't cause false sharing.

6.1.2 updateClustersTask

As already said this is the easier task because embarrassingly parallel. It was invoked in a while loop and so I use the previous solution of instantiating tasks just one time and reusing them, since they only need to have the centroids update at every iteration.

Also this tasks share points and centroids array

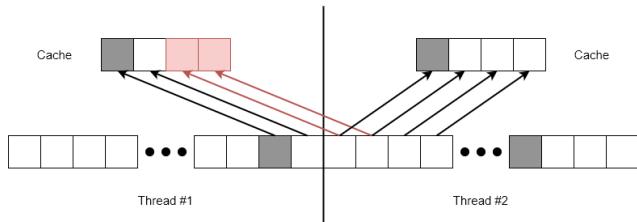


Figure 2. Example: the only case where false sharing occurs

only for reading, but, differently from *initialCentroidsTask*, they share also the clusterization Integer array. This array is accessed also for writing, so it can cause false sharing. Anyway the number of threads is infinitely smaller than dataset dimension and every task accesses (and modifies) a very big number of consecutive positions. The false sharing can occur only when it is cached a border position of a block of array: in this case it is possible that in the same line of cache there is the last positions of previous block or the first positions of the next block (as shown in figure 2). Moreover we can assume that when a thread processes the first positions of its block, also the other threads processes the first positions, since they all progress tidily on a very big number of elements, so the last positions of previous block are not written and this prevents false sharing occurs. Same reasoning for the last positions. So in light of this we can conclude that false sharing occurs rarely and its effects can be ignored.

6.1.3 newCentroidsTask

This last type of task presents any new idea: again I reduce the overhead instantiating them just one time, they share the points and the centroids but only for reading and there is a method to reset the shared variable, useful calculate the average of each clusters, before they are submitted at every iteration.

Again the best solution turned out to be that of calculating a private partial result and then in a critical section putting it in the shared variables. Differently from *initialCentroidsTask* here it was quite simple to use structures of atomic (again, the atomic float would have been replaced by an

AtomicInteger) to avoid the critical section due to absence of conditional write, but this solution was terribly worse in terms of execution time, probably for huge false sharing effects (because the shared variables were modified so frequently). So also here I have introduced a critical section.

6.2. Implicit parallelism

Since updateClustersTask consists on writing on an array in each of its position a value depending only on the position's index (this is why it is embarrassingly parallel), we can exploit a native method of Java: *Arrays.parallelSetAll*. Passing to this method the array clusterization and the lambda function that assigned an index iterate over the centroids to determine to which cluster assign the point, we obtain exactly what we want. And all done by a native function, that clearly is super-optimized and certainly is able to takes all the possible advantage from Hyperthreading.

This method is implemented with a ForkJoinPool and the number of threads (better to say degree of parallelism) can't be handled directly, but can still be changed through a system property², that I set to 4 (the number of virtual threads of my CPU) to make this solution comparable with the other.

6.3. Structure of Arrays idea

The solution I have just presented use the Array of Structure paradigm. At a certain point I have evaluated the possibility to implement the Structure of Arrays one. Anyway this idea gets bogged down shortly, indeed the results (although on a preliminary version, not optimized at all) was terribly worse, not even comparable. And this is easy to see theoretically: we access the dataset point by point, in all its coordinates, and not dimension by dimension. Just think of *getSquaredEuclideanDistance* method (that we remember to take over the 50% of total computation) where we access two, and only two, points in every coordinate. The SoA version suffers of a tremendous

²java.util.concurrent.ForkJoinPool.common.parallelism, which is by default #CPUs - 1

594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

648 number of cache miss, indeed the advantage of
649 the cache is not exploited at all.
650
651

652 7. Tests and Results

653 The execution time obviously grows as dataset
654 dimension and number of cluster grow. I have
655 considered three different dimensions for input
656 images: 4000×2000 , 5000×3000 and 6000×4000
657 for a total of, respectively, 8M, 15M and 24M of
658 pixels. I tested every dimension for different val-
659 ues of k : 2, 8, 14, 20.
660
661

662 The execution time is highly sensitive to the
663 change of input for the same size and number
664 of clusters (I have recorded a variations of up
665 to 250% with different image of the same size).
666 And sometimes I also notice relatively big casual
667 variations. So in order to present a result a bit
668 more generic I run the code on three different im-
669 ages for each size (and each k) two times. I have
670 summarized the results in Table 2, where I report
671 mean execution time, speedup and efficiency.
672
673

674 7.1. Main tests

675 First of all let us compare the two Java paral-
676 lel versions: it is clear that the version with the
677 native ForkJoinPool has generally better results.
678 The speedup increase seems to be almost constant
679 as image size increases, but it enhances with k .
680 An interesting fact is that for very small values of
681 k the version with native ForkJoinPool behaves (a
682 little) worse than the other: in all the cases the
683 speedup is worse or at least the same. This can be
684 explained, likely, with a greater instantiation over-
685 head of a ForkJoinPool respect to a classic Execu-
686 tor thread pool; clearly as k increases the algo-
687 rithm cycles more and the lightness of ForkJoin-
688 Pool's thread prevails over its overhead.
689
690

691 Excluding some cases, where probably the in-
692 put instances play an important role, the speedup
693 for the Executor thread pool version is about 1.7
694 (for an efficiency of 0.85): not a bad result but
695 not noteworthy either, consistent with my pre-
696 diction. The version with native ForkJoinPool
697 instead, again excluding the same cases, has a
698 speedup of 1.8; but not only: indeed in some cases
699
700

701 reaches a speedup of 2. This means that native
702 ForkJoinPool is able to exploit the Hyperthreading
703 the most. Clearly it is too little to talk about
704 superlinear speedup, though in some cases the
705 speedup is just above 2: indeed it could be merit
706 of random fluctuations or very lucky combina-
707 tion of image and value of k (though the speedup
708 seems to increase with k , tests with higher value
709 of k refute this possibility).
710
711

712 7.2. Optimal parallelism degree tests

713 In the main tests I assume the parameters hy-
714 potized above (4 threads and degree of paral-
715 lelism equal to 4), but I want to verify that the
716 guess is right. I run the code, for every combina-
717 tion of image size and k (just once for time rea-
718 sons), with varying number of threads/degree of
719 parallelism from 2 to 6. The figure 3 show the
720 mean speedup for every number of thread used.
721
722

723 Concerning the explicit version my suppose
724 were right: the optimal number of threads is 4,
725 that is my CPU's virtual threads; increasing the
726 number of threads we obtain a constant speed
727 up, worse than the optimal because of overhead
728 due to context switch operation. Interesting fact
729 that 3 threads do quite significantly worse than 2
730 threads: this is probably a consequence of Hyper-
731 threading.
732

733 In the implicit version I used the optimal
734 threads number for the tasks that doesn't use the
735 native function and varying the degree of paral-
736 lelism for the ForkJoinPool I have obtained an in-
737 teresting result: not only this version is always
738 better than the other, but thanks to the super op-
739 timization of the native function, setting the de-
740 gree of parallelism beyond 4 the speedup contin-
741 ues growing, reaching an average of 1.9 with 5
742 and 6 threads.
743
744

745 7.3. Extra tests

746 As I said before (from the beginning) from
747 the implicit native version I expected it to fully
748 exploit CPU advantages, i.e. cache and Hyper-
749 threading. That turned out to be true: indeed run-
750 ning the parallel versions on a CPU with more
751
752

Input	sequential	parallel						
		explicit			implicit			
	T_S	T_2	S_2	E_2	T_2	S_2	E_2	
4K	$k = 2$	5.8s	3.3s	1.76	0.88	3.4s	1.71	0.85
	$k = 8$	63.2s	38.3s	1.65	0.83	35.6s	1.78	0.89
	$k = 14$	112.7s	66.9s	1.68	0.84	61.1s	1.85	0.92
	$k = 20$	237.0s	131.0s	1.81	0.90	121.9s	1.94	0.97
5K	$k = 2$	9.0s	6.0s	1.50	0.75	6.0s	1.50	0.75
	$k = 8$	92.5s	56.7s	1.63	0.82	53.2s	1.74	0.87
	$k = 14$	202.9s	108.8s	1.86	0.93	98.7s	2.06	1.03
	$k = 20$	466.9s	253.8s	1.84	0.92	231.5s	2.02	1.01
6K	$k = 2$	17.5s	12.2s	1.43	0.71	12.2s	1.43	0.72
	$k = 8$	135.3s	77.6s	1.74	0.87	74.4s	1.82	0.91
	$k = 14$	505.8s	284.1s	1.78	0.89	267.7s	1.89	0.94
	$k = 20$	721.2s	428.6s	1.68	0.84	402.6s	1.80	0.90

Table 2. Time, speedup and efficiency for all test cases

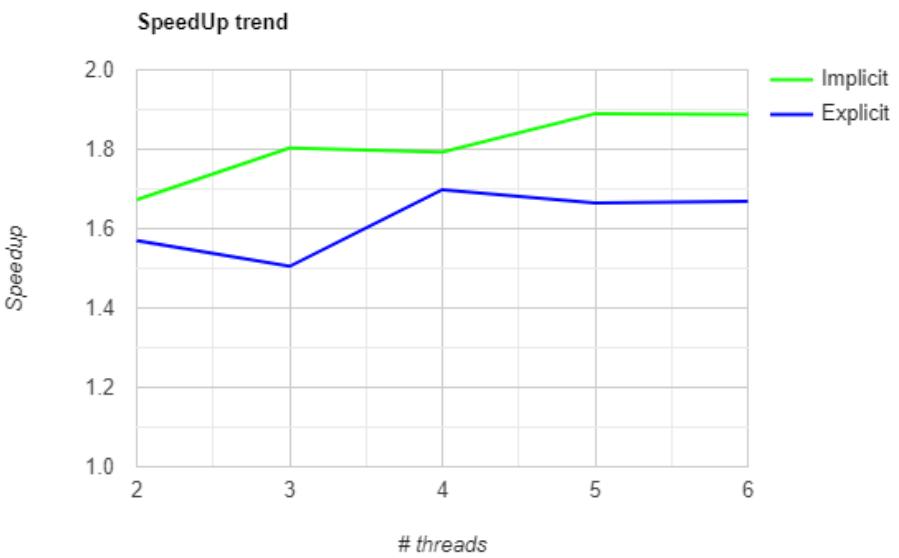


Figure 3. Speedup trend varying the number of thread/degree of parallelism

core, but no Hyperthreading, the implicit version takes longer than the explicit one. Clearly the absolute times were far better (thanks to the high level hardware), but efficiency was worse. This deterioration involve both, explicit and implicit, versions, but the implicit one was the one most affected, until to result (a little) worse than the

explicit one in all cases (as reported in table 3).

This behavior is easy to explain: not only the implicit version can't exploit the advantages of Hyperthreading, but it also presents an overhead: indeed at every iteration it has to instantiate a ForkJoinPool, instead of reusing the always same thread pool as the explicit version.

		explicit	implicit
		T_6	T_6
4K	$k = 2$	1.3s	1.3s
	$k = 8$	11.9s	12.0s
	$k = 14$	19.7s	20.2s
	$k = 20$	37.7s	39.1s
5K	$k = 2$	2.7s	3.0s
	$k = 8$	20.1s	21.7s
	$k = 14$	30.6s	32.2s
	$k = 20$	72.1s	76.2s
6K	$k = 2$	4.0s	4.2s
	$k = 8$	25.5s	26.6s
	$k = 14$	79.7s	90.2s
	$k = 20$	133.2s	146.1s

Table 3. Parallel execution time without Hyperthreading

7.4. Complete tests data

All tests result is available for consultation in table 4. From these is possible to obtain other interesting statistics. For example that the implicit parallel version has better result with the maximum degree of parallelism with a big input and worse with small input. There is also the execution time of the parallel versions with only one thread (so it can be evaluated the overhead respect to the sequential version). Anyway the reported speedup is intended relatively to the sequential version. Final note: the implicit parallel version has fixed number of threads for non-native parallel section (equal to 4 that was the best choice from previous results).

References

- [1] *K-means clustering*, Wikipedia. [Online]. Available: https://en.wikipedia.org/wiki/K-means_clustering (visited on Feb. 1, 2020).
 - [2] J. Leskovec, A. Rajaraman, and J. Ullman, *Mining Massive Datasets*. Cambridge University Press, 2011, ch. 7.3.
 - [3] *Color quantization*, Wikipedia. [Online]. Available: https://en.wikipedia.org/wiki/Color_quantization (visited on Feb. 1, 2020).
 - [4] C. Breshears, *The Art of Concurrency: A Thread Monkey's Guide to Writing Parallel Applications*. O'Reilly, 2009, ch. 4.
- [5] N. Matloff, *Programming on Parallel Machines: GPU, Multicore, Clusters and More*. University of California, 2012, ch. 3.

8. Additional Notes

8.1. Host OS

The main tests were runned on a PC equipped with Ubuntu, since some tests show that both the parallel versions get better results than with Windows. The sequential version instead gets better result on Windows. Anyway the difference were not so high. As consequence, running on Ubuntu the speedup results a little higher.

8.2. Technical Specification

Main tests were runned on/with:

- Intel Core i7-6500U (2.5 GHz up to 3.1 GHz, 2 cores, 4 threads, 4 MB of cache);
- 8 GB of RAM;
- Ubuntu 20.04/20.10;
- Java SE 11;

Extra tests were runned on/with:

- Intel Core i5-9600K (3.7 GHz up to 4.6 GHz, 6 cores, 6 threads, 9 MB of cache);
- 16 GB of RAM;
- Windows 10 Home;
- Java SE 11;

Input	sequential	explicit						parallel					
		#1	#2	#3	#4	#5	#6	#1	#2	#3	#4	#5	#6
4K	$k=2$	5.8s	6.0s	3.6s	3.4s	3.3s	3.4s	3.7s	3.6s	3.4s	3.4s	3.4s	3.3s
	$k=8$	63.2s	68.0s	40.2s	45.7s	38.3s	36.7s	37.8s	38.5s	39.7s	35.1s	35.6s	34.3s
	$k=14$	112.7s	117.9s	72.7s	82.8s	66.9s	67.8s	69.3s	65.6s	63.5s	60.6s	61.1s	58.8s
5K	$k=20$	237.0s	262.1s	141.7s	157.3s	131.0s	140.0s	139.4s	126.0s	126.7s	118.5s	121.9s	116.1s
	$k=2$	9.0s	10.2s	7.4s	6.6s	6.0s	6.0s	6.5s	6.9s	6.3s	6.5s	6.0s	5.9s
	$k=8$	92.5s	107.9s	66.0s	63.1s	56.7s	56.1s	62.8s	60.2s	60.3s	54.2s	53.2s	50.6s
6K	$k=14$	202.9s	236.3s	117.1s	122.3s	108.8s	106.3s	104.5s	108.8s	109.0s	100.4s	98.7s	102.7s
	$k=20$	466.9s	558.8s	272.4s	292.1s	253.8s	264.7s	252.7s	246.2s	255.9s	2228.8s	231.5s	240.2s
	$k=2$	17.5s	17.5s	13.3s	12.9s	12.2s	12.2s	12.5s	12.5s	14.2s	12.4s	12.3s	11.6s
7K	$k=8$	135.3s	138.4s	81.0s	83.8s	77.6s	80.6s	75.6s	89.9s	78.2s	73.2s	74.4s	65.1s
	$k=14$	505.8s	488.7s	269.4s	300.4s	284.1s	272.2s	260.3s	317.0s	270.0s	255.9s	267.7s	225.7s
	$k=20$	721.2s	673.3s	480.1s	532.3s	428.6s	501.8s	466.8s	423.6s	408.5s	384.9s	402.6s	339.9s
Speedup avg		-	-	1.57	1.51	1.70	1.67	1.67	1.62	1.67	1.80	1.79	1.89

Table 4. Execution time data of all tests cases.