

## Timing Study of the PageRank Algorithm Implemented with Clojure

### Introduction

The purpose of this study is to explore the efficiency of multithreading with Clojure. A PageRank algorithm, which ranks the relative importance of webpages, is used to conduct the timing study. The algorithm runs for 1000 iterations and thus is a good candidate for multithreading. The elapsed time to complete the calculations is compared with the number of requested threads. Furthermore, the study also explores how the number of hash map partitions affects execution times.

### Methods

The program starts by opening the *pages.txt* file containing information for dummy webpages, and reads in one line at a time, storing the data in a hash map. The key is the page ID, and the value is a vector of the outgoing links to other page ID's. This method produces an individual hash map for each page ID, resulting in 10,000 maps. These were combined with the custom *combine-maps* function. The next step is to find all the incoming links for each page ID. Since the data is random and not sorted, a brute force method is employed by which the entire data structure is searched for each page ID. This is by far the least efficient and most lengthy process of the entire program. In order to improve efficiency, the *some* function is used.

```
(defn search-map [targetID map] (for [[pageIndex outpages] map] (if (some # (= targetID %) (get map pageIndex))  
pageIndex)))
```

If the query is found in the vector, instead of checking the value at each remaining index of the vector, the search immediately stops and moves onto the next page entry. Another hash map is created that stores the page ID and a count of the outgoing links on that page. I chose to use another data structure to store this information rather than recalculating the count for each entry every iteration of the loop. The third hash map stores the PageRank. Before the looping begins, the PageRank of all page ID's is initialized to 1. Inside the loop, the PageRank function is called and takes the *inpagesMap*, *outpagesCountMap*, and the current *pageRankMap*. The PageRank is calculated by the following formula where  $d$  is the damping factor (0.85),  $PR(T)$  is the current PageRank of page  $T$ ,  $C(T)$  is the count of outgoing links on page  $T$ , and the summation is over all pages that have an outgoing link to page  $A$ :

$$PR(A) = (1 - d) + d \sum \frac{PR(T)}{C(T)}$$

Multithreading is incorporated into the *calc-page-rank* function. In the final code, the *modded-pmap* function, modified from the source code of the standard *pmap* function, takes the number of threads to use,  $t$ , and the sequence of hash maps, *mapPartitions*.

```
(def pageRankMap (combine-maps (modded-pmap #(calc-page-rank % outpagesCountMap pageRankMap) t  
mapPartitions)))
```

Before I determined a way to specify how many threads were called to run the PageRank loop, I approached the problem by dividing up the hash map into partitions corresponding to the number of threads that I wanted to use. The *mapPartitions* object holds a sequence of hash maps based on the number of partitions.

My idea was that *pmap*, a parallelized version *map*, which applies a function to each set of a collection, would assign each partition to a new thread. However, there is no way to guarantee that *pmap* creates the same number of threads as the partitions. Furthermore, *pmap* starts with two threads and proceeds to utilize available threads as needed. Preliminary tests suggested that the size of the map partition affected the execution time considerably. After modifying the source code for *pmap*, I did not know what number of partitions to choose in order to optimize the speed, so therefore I tested a range of map partitions from a single partition up to 10,000 partitions, or a partition for each entry in the hash map. One thing to note, if only one thread is used, the standard non-parallelized *map* function is used instead of *modded-pmap*. This is due to the way *modded-pmap* works, which requires two or more threads.

The *doseq* function allows the test to iterate through each number of threads and each number of partitions in one program execution. The first part initializes the PageRanks to 1, and partitions the hash map into *s* partitions, *mapSplits*, returning a sequence of hash maps, *mapPartitions*. After this, the loop runs for 999 more times, redefining the *pageRankMap* object during each iteration with updated values.

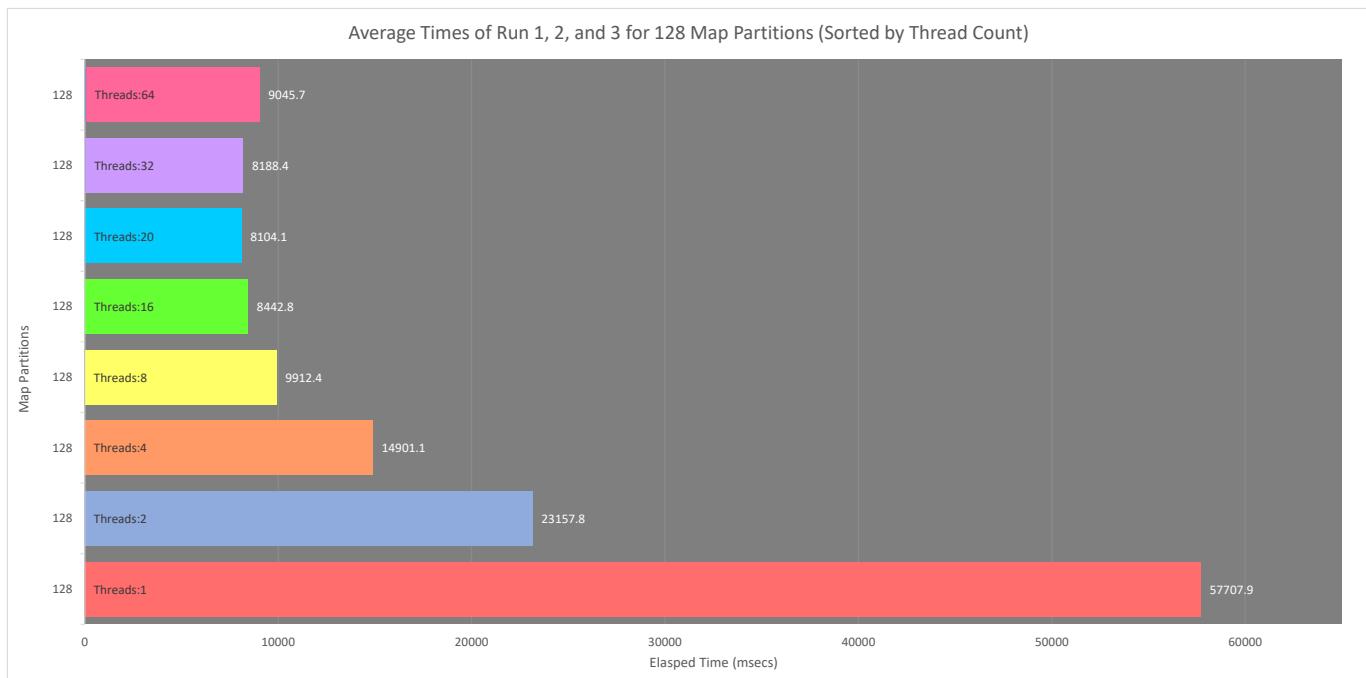
```
(defn main-body [myMap]
  (def inpagesMap (find-inpages myMap))
  (def outpagesCountMap (count-outpages myMap))
  (def threads [1 2 4 8 16 20 32 64])
  (def mapSplits [1 2 4 8 16 20 32 64 128 256 512 1024 2048 4096 8192 10000])

  (doseq [s mapSplits]
    (doseq [t threads]
      (save-times t s (with-out-str (time
          (do
            (def mapPartitions (map vec-to-map (split-map s inpagesMap))) ; returns a sequence of hash maps
            (if (= t 1)
                (def pageRankMap (combine-maps (map #(calc-page-rank % outpagesCountMap (start-rank
                  myMap 1)) mapPartitions))) ; initialize PageRanks to 1, standard map function for 1 thread
                ;; else
                (def pageRankMap (combine-maps (modded-pmap #(calc-page-rank % outpagesCountMap (start-
                  rank myMap 1)) t mapPartitions)))) ; initialize PageRanks to 1, modded pmap for > 1 thread
                (save-results pageRankMap 1 false) ; set to true if output of PageRanks is needed

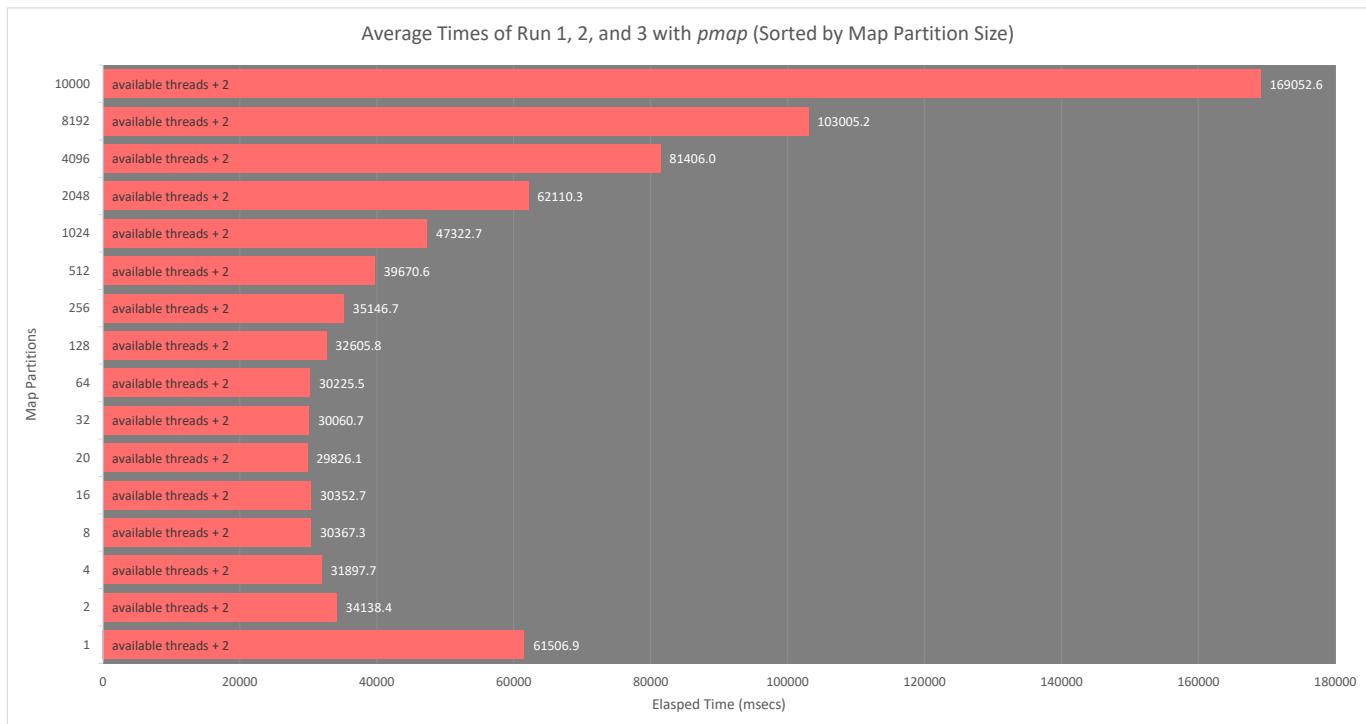
            (loop [i 1] ; already run once above
              (if (< i 1000)
                  (do
                    (if (= t 1)
                        (def pageRankMap (combine-maps (map #(calc-page-rank % outpagesCountMap pageRankMap)
                          mapPartitions))) ; standard map function for 1 thread
                        ;; else
                        (def pageRankMap (combine-maps (modded-pmap #(calc-page-rank % outpagesCountMap
                          pageRankMap) t mapPartitions)))) ; modded pmap for > 1 thread
                        (recur (+ i 1)))
                    (save-results pageRankMap i false))))))))))) ; set to true if output of PageRanks is needed
      (shutdown-agents)) ; releases threads
```

## Findings

For the timing study, the multiple thread counts and map partitions are compared. As expected, an increase in the number of threads decreased the execution time of the PageRank algorithm up until the number of physical threads is exceeded. This study was conducted on a machine with an Intel i9 7900X with 20 logical processors, the reason why a thread count of 20 is included in the tests. Requesting more threads than are physically available generally leads to a slight increase in execution time.



With 128 map partitions, 20 threads, corresponding to the number of physical threads available, performs the best at 7.12 times faster than a single thread, with an average time of just 8.10 seconds. By contrast, it takes a single thread 57.70 seconds to execute. Requesting 32 threads has a minor impact on execution time. Requesting 64 threads, however, is about 0.93 times the speed of using 16 threads, or about 600 msec slower. It is interesting that 2 threads are about 2.5 times faster than using a single thread, but 4 threads are only about 1.55 times faster than 2 threads. The execution time does not scale linearly with the thread count. The complete table of Run 1, 2, and 3 are available on page 7.



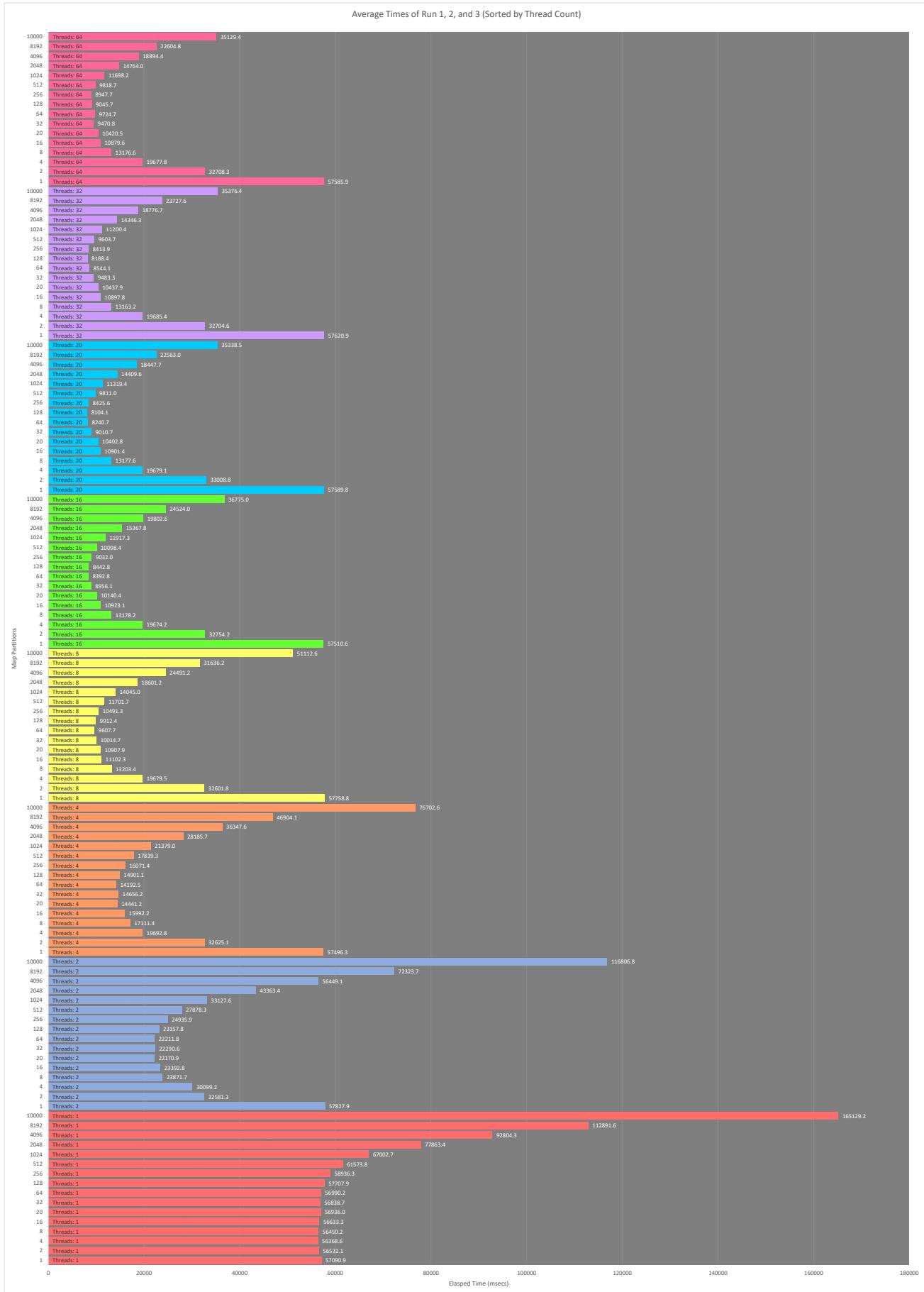
A further comparison of modded-pmap to pmap shows a substantial increase in the efficiency of modded-pmap. Although the number of threads pmap used was not recorded, with 128 map partitions, pmap is 0.77 times the speed of modded-pmap using only 2 threads.

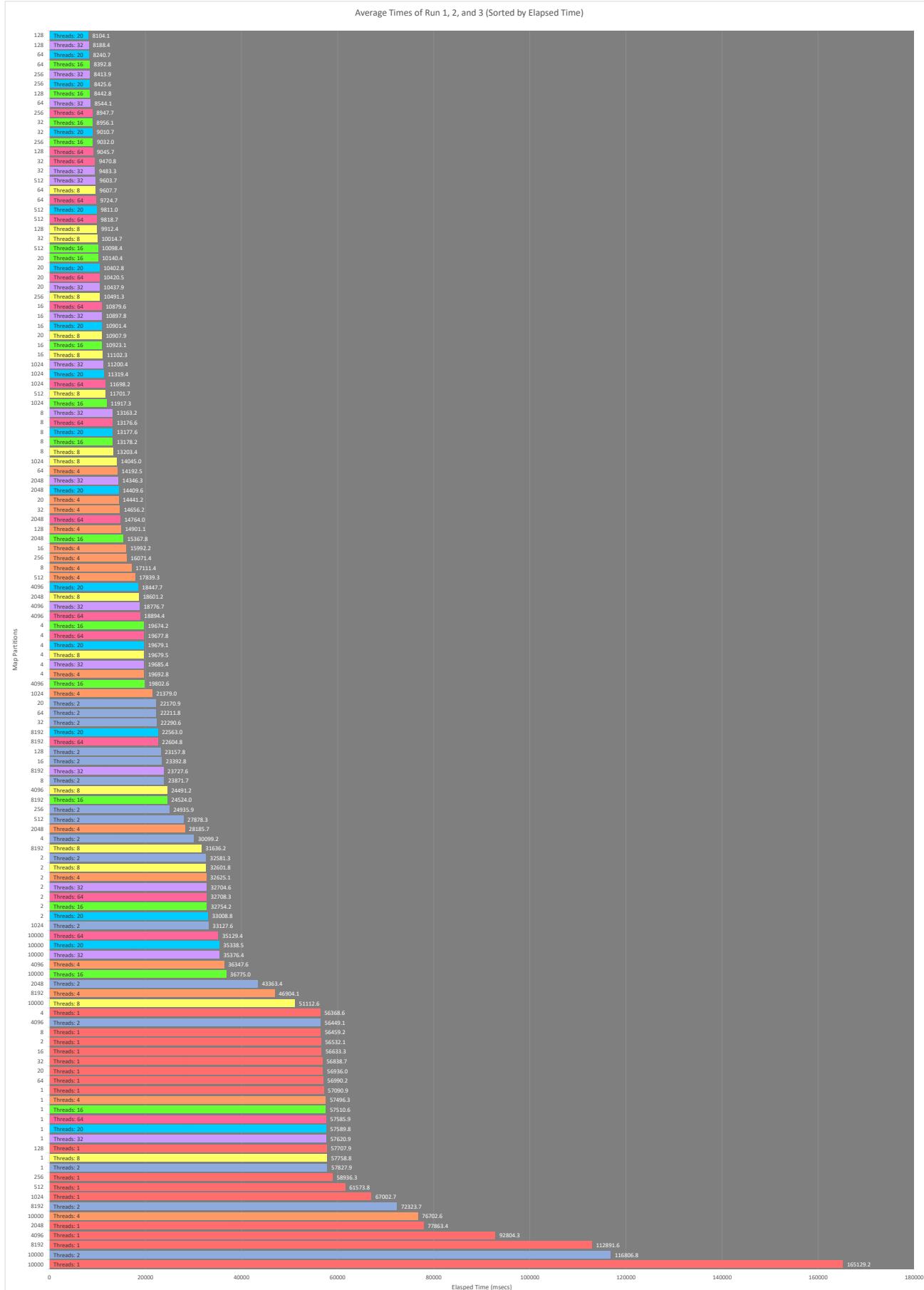
Perhaps more interesting is the effect of the number of map partitions that are used. The charts for the average times of Run 1, 2, and 3 sorted by thread count and execution time are available on pages 5 and 6, respectively.

The optimal number of map partitions appears to 64 or 128, with the fastest time reached by 20 threads and 128 map partitions at 8104.1 msecs. Using 20 threads, this is 7.11 times faster than passing the entire map in a single partition, and 4.36 times faster than passing 10,000 partitions, or each entry as an individual map. The slowest by far is using a single thread with 10,000 partitions, which is 20.38 times slower than the fastest time. The number of threads using a single map partition has no apparent effect on the execution time. This is because the map is not divided up amongst the threads. Likewise, passing 1 thread multiple map partitions has a neutral effect on execution time up to 128 map partitions, or 64 partitions for 2, 4, and 8 threads, respectively. Excluding the 1 thread count, as the number of map partitions increase, the execution time is reduced considerably up to 64 or 128 map partitions, depending on the number of threads. From 256 to 10,000 map partitions, the execution times slows considerably.

### **Conclusion**

The results show a clear relationship between the number of threads and map partitions and execution time. As the number of threads increases up to the number of physical threads, the execution time decreases. As soon as the number of physical threads is reached, a slight increase in execution time is observed. Similarly, increasing the number of map partitions past the 64 or 128 thresholds, depending on the thread count, execution time decreases. This is where the data varies per thread count, and there is no specific multiplier relating the number of available threads to the number of map partitions. However, the results demonstrate the importance of multithreading. Furthermore, they indicate that partitioning data into smaller blocks can optimize execution time considerably. However, it is important that a balance is found between threads and partitions.





Run 1 - Start: Fri Nov 22 23:54:15 CST 2019				Run 2 - Start: Sat Nov 23 06:18:55 CST 2019				Run 3 - Start: Sat Nov 23 11:39:44 CST 2019				Averages of Runs 1, 2, and 3			
Threads	Map Partitions	Time (msecs)		Threads	Map Partitions	Time (msecs)		Threads	Map Partitions	Time (msecs)		Threads	Map Partitions	Time (msecs)	
1	1	56818.5413		1	1	58820.9254		1	1	55633.2712		1	1	57090.91263	
2	1	57878.0984		2	1	59779.9896		2	1	55827.5563		2	1	57827.88477	
4	1	57761.8118		4	1	58998.49		4	1	55728.7273		4	1	57496.34303	
8	1	57906.6139		8	1	59104.2722		8	1	56265.6465		8	1	57758.8442	
16	1	57400.9371		16	1	59067.3694		16	1	56063.5713		16	1	57510.62593	
20	1	57787.4804		20	1	59017.3391		20	1	55964.8566		20	1	57589.8284	
32	1	57614.0248		32	1	59289.6294		32	1	55958.9497		32	1	57620.86797	
64	1	57766.3975		64	1	59065.5071		64	1	55925.8527		64	1	57585.9191	
1	2	56539.8151		1	2	57954.9206		1	2	55101.4196		1	2	56532.05177	
2	2	32564.4376		2	2	33492.7758		2	2	31686.6996		2	2	32581.30433	
4	2	32587.2957		4	2	33559.498		4	2	31728.4353		4	2	32625.07633	
8	2	32602.4588		8	2	33499.6007		8	2	31703.3642		8	2	32601.8079	
16	2	33038.6025		16	2	33494.5977		16	2	31729.3806		16	2	32754.1936	
20	2	33680.0089		20	2	33612.4929		20	2	31733.7892		20	2	33008.76367	
32	2	32774.9134		32	2	33514.986		32	2	31824.0299		32	2	32704.6431	
64	2	32712.332		64	2	33515.1019		64	2	31897.538		64	2	32708.32397	
1	4	56577.2214		1	4	57601.8101		1	4	54926.6555		1	4	56368.56233	
2	4	30042.3978		2	4	30932.1541		2	4	29323.0203		2	4	30099.19073	
4	4	19572.6529		4	4	20207.9263		4	4	19297.7076		4	4	19692.76227	
8	4	19549.2045		8	4	20201.7799		8	4	19287.625		8	4	19679.53647	
16	4	19583.2531		16	4	20206.6901		16	4	19232.8029		16	4	19674.2487	
20	4	19593.512		20	4	20185.0979		20	4	19258.5986		20	4	19679.0965	
32	4	19567.9733		32	4	20198.7227		32	4	19289.619		32	4	19685.43833	
64	4	19577.8179		64	4	20183.3612		64	4	19272.2981		64	4	19677.82573	
1	8	56582.4736		1	8	57655.2835		1	8	55139.8512		1	8	56459.20277	
2	8	23864.0983		2	8	24587.8128		2	8	23163.1725		2	8	23871.69453	
4	8	17035.1532		4	8	17666.0527		4	8	16632.9019		4	8	17111.36927	
8	8	13158.8581		8	8	13685.7497		8	8	12765.6436		8	8	13203.41713	
16	8	13126.7624		16	8	13661.026		16	8	12746.6819		16	8	13178.15677	
20	8	13106.3963		20	8	13665.6448		20	8	12760.82		20	8	13177.62037	
32	8	13105.336		32	8	13648.6812		32	8	12735.5101		32	8	13163.17577	
64	8	13124.8046		64	8	13657.184		64	8	12747.7493		64	8	13176.5793	
1	16	56660.0755		1	16	58133.8761		1	16	55106.0565		1	16	56633.33603	
2	16	23386.0112		2	16	24160.4878		2	16	22631.9866		2	16	23392.82853	
4	16	15943.2312		4	16	16508.3458		4	16	15524.8803		4	16	15992.15243	
8	16	11111.9694		8	16	11416.6143		8	16	10778.4164		8	16	11102.33337	
16	16	10904.1376		16	16	11157.9685		16	16	10707.1147		16	16	10923.0736	
20	16	10879.6714		20	16	11117.6464		20	16	10706.8261		20	16	10901.3813	
32	16	10863.1482		32	16	11140.3457		32	16	10690.0304		32	16	10897.84143	
64	16	10873.5953		64	16	11095.0222		64	16	10670.0897		64	16	10879.56907	
1	20	57265.5827		1	20	58300.462		1	20	55241.9837		1	20	56936.00947	
2	20	22420.5098		2	20	22933.1204		2	20	21559.1272		2	20	22170.91913	
4	20	14410.2507		4	20	14903.6479		4	20	14009.7643		4	20	14441.22097	
8	20	10936.7711		8	20	11244.7865		8	20	10542.0553		8	20	10907.87093	
16	20	10064.6874		16	20	10343.569		16	20	10012.834		16	20	10140.36347	
20	20	10369.7542		20	20	10621.1157		20	20	10217.4065		20	20	10402.7588	
32	20	10383.565		32	20	10638.3152		32	20	10291.7488		32	20	10437.87633	
64	20	10385.2663		64	20	10662.4257		64	20	10213.9121		64	20	10420.5347	
1	32	56855.4227		1	32	58219.077		1	32	55441.5059		1	32	56838.66853	
2	32	21789.2544		2	32	23147.5009		2	32	21935.0077		2	32	22290.58767	
4	32	14499.9166		4	32	15235.7498		4	32	14232.977		4	32	14656.21447	
8	32	10107.1431		8	32	10338.6542		8	32	9598.2322		8	32	10014.6765	
16	32	8968.9209		16	32	9148.8389		16	32	8750.5607		16	32	8956.106833	
20	32	8991.8654		20	32	9183.8922		20	32	8856.4142		20	32	9010.723933	
32	32	9439.6042		32	32	9687.2878		32	32	9323.0881		32	32	9483.3267	
64	32	9467.9997		64	32	9673.8981		64	32	9270.3641		64	32	9470.753067	
1	64	57440.2913		1	64	58238.7745		1	64	55291.5453		1	64	56990.2037	
2	64	22134.7644		2	64	22957.2472		2	64	21543.4456		2	64	22211.81907	
4	64	14180.5418		4	64	14629.0185		4	64	13767.8691		4	64	14192.47647	
8	64	9683.9229		8	64	9938.2156		8	64	9200.8966		8	64	9607.678367	
16	64	8404.1106		16	64	8583.0887		16	64	8191.1938		16	64	8392.7977	
20	64	8298.2011		20	64	8435.0388		20	64	7988.7658		20	64	8240.668567	
32	64	8542.985		32	64	8714.0242		32	64	8375.2633		32	64	8544.090833	
64	64	9763.1933		64	64	9913.9242		64	64	9497.0515		64	64	9724.723	
1	128	58148.682		1	128	58955.1111		1	128	56019.8911		1	128	57707.89473	
2	128	23045.9168		2	128	8376.076		2	128	7894.0946		2	128	8188.40333	
4	128	9016.7698		4	128	9249.5684		4	128	8870.8631		4	128	9045.733767	
8	128	59147.569		8	128	60338.6726		8	128	57322.5793		8	128	58936.27363	
16	128	24888.4444		16	128	6240.8723		16	128	53678.4109		16	128	24935.8792	
32	128	16000.8508		32	128	16748.6211		32	128	15464.7888		32	128	16071.42023	
64	128	10419.3111		64	128	11002.5768		64	128	10052.0302		64	128	10491.30603	
1	256	8973.059		1	256	9382.3739		1	256	8740.4694		1	256	9031.967433	
2	256	8346.6238		2	256	8714.0787		2	256	8181.1365		2	256	8425.623233	
4	256	8872.9334		4	256	9224.0834		4	256	8745.9451		4	256	8947.653967	
8	256	69147.569		8	256	10276.3793		8	256	9539.6507		8	256	9818.6966	
16	256	66653.0724		16	256	10464.2925		16	256	9271.4707		16	256	9603.686267	
32	256	1024		32	256	10464.2925		32	256	9045.9013		32	256	10739.0137	
64	256	1024		64	256	10464.2925		64	256	8730.0137		64	256	11020.42497	
1	512	61281.6498		1	512	63468.1704		1	512	59971.5328		1	512	61573.78433	
2	512	27745.4776		2	512	28934.2652		2	512	26955.1629		2	512	27878.3019	
4	512	17721.5363		4	512	18609.8317		4	512	17186.406		4	512	17839.258	
8	512	11538.5272		8	512	12396.2177		8	512	11170.3436		8	512	11701.69683	
16	51														