# Comparing Performance of KMP-Algorithm in Serial and Parallel Environments

Joel John Thomas
USN: 1MS18CS053

## Abstract

In order to make the most out of the CPU's resources, we must find different ways to implement the algorithms. KMP (the Knuth–Morris–Pratt) algorithm is used **to find a "Pattern" in a "Text"**. This algorithm compares character by character from left to right. But whenever a mismatch occurs, it uses a preprocessed table called "Prefix Table" to skip characters comparison while matching. To activate parallel programming, we will use MPI clustering programming. MPI is a message-passing model of distributed memory. To assess parallel programming's performance, we compare it to a sequential KMP algorithm, which does not make efficient use of the CPU's resources. To determine whether an algorithm is more efficient, we evaluate the performance of sequential and parallel KMP algorithms.

## Introduction

The exponential growth of processing and network speeds means that parallel architecture isn't just a good idea; it's necessary. Big data and the IoT will soon force us to crunch trillions of data points at once. Nowadays new applications always need faster processors in today's world to operate efficiently. Parallel systems are employed in a lot of business applications nowadays. Parallel computing uses multiple computer cores to attack several operations at once. Unlike serial computing, parallel architecture can break down a job into its component parts and multi-task them. Parallel computer systems are well suited to modeling and simulating real-world phenomena. Without parallel computing, performing digital tasks would be tedious, to say the least. The advantages of parallel computing are that computers can execute code more efficiently, which can save time and money by sorting through "big data" faster than ever.

Parallel systems, divide the programs into multiple fragments and process them all at the same time in order to reduce the execution time of the program. By

constructing a parallel processing bundle or a combination of both entities, a parallel system can deal with many processors, machines, computers, or CPUs, among other things.

String searching is one of the most common problems where it tries to find whether a string is a sub-string of another particular string or not. There are several algorithms to solve this problem. These algorithms are mostly used in text edit processing, image processing, literature retrieval, natural language recognition, spell checking, and many other fields.
The most common string searching algorithm is Knuth–Morris–Pratt (KMP) algorithm. It is a prefix based searching algorithm faster than normal string searching algorithm like 'Naive String Searching Algorithm'. The algorithm uses degenerating property (pattern having same sub-patterns appearing more than once in the pattern) of the pattern and improves the worst case complexity to O(n).

We can make this algorithm faster if we make it parallelly executable. As it is linear time complexity algorithm, it will consume more time if the text data is big. If we run the KMP algorithm on multiple computers to handle larger data, it will be less time-consuming. By doing that, it can be used in large-scale applications. To achieve parallelism in KMP algorithm, we will use MPI clustering programming. MPI is a message-passing model of distributed memory. The aim of the MPI is to establish an efficient and flexible standard for message passing which is used for writing message passing programs. As such, MPI is the vendor independent and having message passing library. The benefit of developing message passing software using MPI helps to achieve design goals like portability, efficiency and flexibility. MPI performs their operations by using the methods such as MPI_send, MPI_receive, MPI_COMM_WORLD, MPI_Init and many more by Using these kind of functionality the processes communicate with themselves.

The steps of parallel implementations using MPI are given below:

1) Read the input file of Snort rules, prepare the input stream of that file, obtain the file size, allocate memory to contain whole file, copy the file into the buffer and divide it across processes using MPI_Bcast and MPI_Barrier.

2) Prepare the input stream of pattern file and get the number of lines.

3) Divide the input file among the processors.

4) Call the KMP search function and compute the longest proper suffix and at last show the matching patterns as output.

**Code:**

```
1
2 #include "mpi.h"
3 #include <stdlib.h>
4 #include <stdio.h>
5 #include <string.h>
6 #include <malloc.h>
7
8 #define master 0
9 #define MAX(a,b) ((a) > (b) ? a : b)
10 #define MILLION 1000000L
11
12 #define bug printf("trust\n");
13
14 int* kmp(char* target, char* pattern, int* table, int myrank);
15 int* kmptable(char* pattern, int len);
16 void print(char* arr, int id, int len);
17 void fillup(char** matrix, int rankid, char* msg, int len, int offset);
18 void freeDouble(char** matrix, int nproc);
19 void printMatrix(char **matrix, int nproc, int id, int length);
20 int* getRealIdx(int* answer, int len, int myrank, int x, int y, int* length);
21 void pinpoint(int* result, int* msg, int shortlen);
22
23
24 void pinpoint(int* result, int* msg, int shortlen){
25         int j = 0;
26         for(j=0; j<shortlen; j++){
27                 if(j > 0 && msg[j] == 0)
28                         break;
29                 else if(j == 0 && msg[j] == 0)
30                         continue;
31                 else{
32                         result[msg[j]] = 1;
33                 }
34         }
35 }
36
37 // get the exact idx(es) in each process
38 // return back an int pointer
39 int* getRealIdx(int* answer, int len, int myrank, int x, int y, int* length){
40         int* package = (int *) calloc(len, sizeof(int));
41         int index;
42         int real_len = 0;
43         int i;
44         for(i=0;i < len; i++){
45                 if(i > 0 && answer[i] == 0)
46                         break;
47                 else if(i == 0 && answer[i] == 0)
48                         continue;
49                 else{
50                         if(myrank == 0){
51                                 index = answer[i];
52                                 package[real_len++] = index;
53                                 // printf("This is REAL index: %d\n", index);
54                         }else{
```

```c
53              // printf("This is REAL index: %d\n", index);
54              }else{
55                  index = x * myrank - y + answer[i];
56                  package[real_len++] = index;
57                  // printf("This is REAL index: %d\n", index);
58              }
59
60          }
61      }
62      int* result = (int*)malloc((real_len) * sizeof(int));
63
64      for(i =0 ; i<real_len; i++){
65          if(i > 0 && package[i] == 0)
66              break;
67          else{
68              result[i] = package[i];
69          }
70      }
71      *length = real_len;
72      free(package);
73      return result;
74 }
75
76 void printMatrix(char **matrix, int nproc, int id, int length){
77      // printf("This matrix at %d row.\n", id);
78      int i ,j;
79      for(i=0 ;i < nproc; i++){
80          if(i == id){
81              for(j=0;j<length;j++){
82                  printf("%c", matrix[i][j]);
83                  break;
84              }
85              printf("\n");
86          }
87      }
88      printf("\n");
89 }
90
91
92 // free two-pointer matrix
93 void freeDouble(char** matrix, int nproc){
94      for (int i=0; i<nproc; ++i) {
95          free(matrix[i]);
96      }
97      free(matrix);
98 }
99
100
101 void fillup(char** matrix, int rankid, char* msg, int len, int offset){
102      int i;
103      for(i=0;i<len;i++)
104          matrix[rankid][offset + i] = msg[i];
105 }
106

107 // print out one pointer list
108 void print(char* arr, int id, int len){
109      int j;
110
111      for(j=0; j<len; j++){
112          printf("%c", arr[j]);
113      }
114
115 }
116
117 // key func kmp
118 int* kmp(char* target, char* pattern, int* table, int myrank){
119      int n = strlen(target);
120      int m = strlen(pattern);
121      int* answer = (int*) calloc(n-m+1,sizeof(int));
122      int j=0;
123      int i=0;
124      int index = 0;
125      while(i<n){
126          if(pattern[j]== target[i])
127          {
128              j++;
129              i++;
130          }
131          if(j == m){
132              printf("Start Point: %d.\n", i-j);
133              answer[index++] = i-j;
134              j = table[j-1];
135
136          }else if(i < n && pattern[j] != target[i]){
137              if(j!=0)
138                  j = table[j-1];
139              else
140                  i++;
141          }
142      }
143      return answer;
144 }
145
146
147 int main(int argc, char** argv){
148      int i,j;
149      int n;
150      int m;
151
152      struct timespec start1, end1;
153      double diff;
154
155      char target[] =
    "zztzjjsharzzfuassssssssssssssssssssssssssssssssssssssssssssssssaaaaaaaaaaaajjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjjaaaaaaaaaaaaaaaaaaaa
    k-wlew-l;;;;ppppdopsdspdspdpsdpsdspdpsdpspdpsdpspdpsdpdpspdpspdpsdpspdpsdpspdspdspdspdspdsdpspdspdpsdpsdpspdspdspap-wq=l-
    lk,ldlslldslddlsldsldlslslsldlsldsldldlsldldlsldslsldlsdslld[lsdldsll[dl[sld[sd[[sl[lds[l[lds[l[dl[sl[ls[dl[sd[sd[s[dsl[dls[d[s[dls[dl[sld[ls[dls[l[s[ewpjjjjoodgnnogndngdongodngod
156      char pattern[] = "pryyaogqzekpumxktopubehtpzeoqlkfcifewmwysgnotkbmyqhzuqujpayflinpsuvqisobabrtnsinjrofhkpzkpgvxtmpqifycubcxvwnngayxpegbvfnkm";
```

```c
156         char pattern[] = "pryyaogqzekpumxktopubehtpzeoqlkfcifewmwysgnotkbmyqhzuqujpayflinpsuvqisobabrtnsinjrofhkpzkpgvxtmpqifycubcxvwnngayxpegbvfmkm";
157
158         n = strlen(target);
159         m = strlen(pattern);
160
161         int tag = 1;
162         int tag2 = 2;
163
164         int* kmp_table = kmptable(pattern, m);
165
166
167
168         int myrank, nproc;
169         MPI_Init(&argc, &argv);                                      //initialize the MPI
170         MPI_Comm_rank(MPI_COMM_WORLD, &myrank);           //number of processes that given in command line
171         MPI_Comm_size(MPI_COMM_WORLD, &nproc);            //ranks the process rank
172
173
174
175
176         MPI_Status status;
177
178         int start_pos = n/nproc;
179         int str_per_proc = n/nproc;
180         int end = str_per_proc + n % nproc;
181
182
183         char** matrix = (char**)malloc(nproc*sizeof(char*));
184         for(i=0;i<nproc;i++){
185                 matrix[i] = (char*)malloc((end+m-1) * sizeof(char));
186         }
187         printMatrix(matrix, nproc, myrank, 10);
188
189     char send_msg[MAX(m-1, end)]; // 4-1 = 3
190         char recv_msg[MAX(m-1, end)];
191     char end_msg[MAX(m-1, end)];
192
193     double start, stop;
194
195
196         if(myrank == master){
197                 printf("Result KMP Sequential\n");
198                 start = MPI_Wtime();
199                 for(i=1;i< nproc;i++){
200                         if(i == nproc-1){
201                                 strncpy(send_msg, target + i*start_pos, end);
202                                 MPI_Send(send_msg, end, MPI_CHAR, i, tag, MPI_COMM_WORLD);
203                         }else{
204                                 strncpy(send_msg, target + i*start_pos, str_per_proc);
205                                 MPI_Send(send_msg, str_per_proc, MPI_CHAR, i, tag, MPI_COMM_WORLD);
206                         }
207
208
209                 }
```

```c
208
209                 }
210                 strncpy(send_msg, target + master * start_pos, str_per_proc);
211                 fillup(matrix, myrank, send_msg, str_per_proc, 0);
212
213         // The last process
214         }else if(myrank == nproc-1){
215                 MPI_Recv(recv_msg, end, MPI_CHAR, master, tag, MPI_COMM_WORLD, &status);
216
217                 fillup(matrix, myrank, recv_msg, end, m-1);
218         }
219         else{
220                 MPI_Recv(recv_msg, str_per_proc, MPI_CHAR, master, tag, MPI_COMM_WORLD, &status);
221                 fillup(matrix, myrank, recv_msg, str_per_proc, m-1);
222         }
223
224         MPI_Barrier(MPI_COMM_WORLD);
225
226         if(myrank == nproc-1){
227
228                 MPI_Recv(end_msg, m-1, MPI_CHAR, myrank-1, tag, MPI_COMM_WORLD, &status);
229
230                 fillup(matrix, myrank, end_msg, m-1, 0);
231                 int* answer = kmp(matrix[myrank], pattern, kmp_table, myrank);
232
233                 int len;
234                 // ---------------------
235                 int* result = getRealIdx(answer, end-m+1, myrank, str_per_proc, m-1, &len);
236                 // ---------------------
237                 free(answer);
238
239                 MPI_Send(result, len, MPI_INT, master, tag2, MPI_COMM_WORLD);
240                 free(result);
241
242         }else{
243                 strncpy(send_msg, target + str_per_proc-m + 1 + myrank * str_per_proc, m-1);
244                 MPI_Send(send_msg, m-1, MPI_CHAR, myrank+1, tag, MPI_COMM_WORLD);
245
246                 /* Processes other than master one
247                  * Re-recv more msg from the previous process whose rank is myrank-1
248                  * kmp implementation
249                  * Send the result back to master process
250                  */
251
252                 if(myrank!=master){
253                         MPI_Recv(end_msg, m-1, MPI_CHAR, myrank-1, tag, MPI_COMM_WORLD, &status);
254                         fillup(matrix, myrank, end_msg, m-1, 0);
255                         // implement kmp to get the matching result
256                         int* answer = kmp(matrix[myrank], pattern, kmp_table, myrank);
257                         int len;
258                         int* result = getRealIdx(answer, end-m+1, myrank, str_per_proc, m-1, &len);
259                         free(answer);
260                         MPI_Send(result, len, MPI_INT, master, tag2, MPI_COMM_WORLD);
261                         free(result);
```

```
271                     int* answer = kmp(matrix[myrank], pattern, kmp_table, myrank);
272                     int len;
273                     int* result = getRealIdx(answer, end-m+1, myrank, str_per_proc, m-1, &len);
274                     pinpoint(final_result,result, len);
275                     free(answer);
276                     free(result);
277                     int j;
278                     for (j = 0; j < end+m-1; j++)
279                     {
280                             recv[j] = 0;
281                     }
282
283                     for(j = 1; j <nproc; j++){
284                             MPI_Recv(recv, end+m-1, MPI_INT, j, tag2, MPI_COMM_WORLD, &status);
285                             pinpoint(final_result,recv, end+m-1);
286                     }
287                     stop = MPI_Wtime();
288                     free(recv);
289                     double timeDifference = (stop-start) * MILLION;
290                     printf("Target Length: %d\n", n);
291                     printf("Pattern Length: %d\n", m);
292                     printf("Time Elapsed: %0.05f\n",timeDifference);
293                     // for(j=0;j<n;j++){
294                     //      if(final_result[j] == 1)
295                     //          printf("Find a matching substring starting at: %d.\n", j);
296                     // }
297                     printf("\n");
298                     free(final_result);
299             }
300     }
301
302     MPI_Barrier(MPI_COMM_WORLD);
303     freeDouble(matrix, nproc);
304     MPI_Finalize();
305
306     return 0;
307 }
308
309 int* kmptable(char* pattern, int len){
310     int k = 0;
311     int i = 1;
312     int* table = (int*)malloc(len * sizeof(int));
313     table[0] = k;
314     for(i=1;i<len;i++){
315             while(k > 0 && pattern[i-1] != pattern[i]){
316                     table[i] = k-1;
317                     k = table[i];
318             }
319             if(pattern[i] == pattern[k])
320                     k++;
321             table[i] = k;
322     }
323     return table;
324 }
```

# Output:

On a Computer without using MPI

Result KMP_MPI
Start Point: 49521.
Target Length: 99044
Pattern Length: 122
Time Elapsed: 2581.42600

On both Computer using MPI Clustering

Result KMP_MPI
Start Point: 49521.
Target Length: 99044
Pattern Length: 122
Time Elapsed: 698.22000

```
data100000.txt    intrductionReport.pdf    KMPSequentical    testmake.py
data10000.txt     kmp_MPI.c                pattern.txt       text.txt
data1000.txt      kmp_MPI.out              Temporary.c
joel@joel-Inspiron-3576:~/Downloads/KMP-Parallel-Using-MPI-master$ mpicc -g -Wall -o test kmp_MPI.c
kmp_MPI.c: In function 'main':
kmp_MPI.c:153:9: warning: unused variable 'diff' [-Wunused-variable]
  153 |   double diff;
      |          ^~~~
kmp_MPI.c:152:27: warning: unused variable 'end1' [-Wunused-variable]
  152 |    struct timespec start1, end1;
      |                            ^~~~
kmp_MPI.c:152:19: warning: unused variable 'start1' [-Wunused-variable]
  152 |    struct timespec start1, end1;
      |                    ^~~~~~
kmp_MPI.c:148:8: warning: unused variable 'j' [-Wunused-variable]
  148 |   int i,j;
      |         ^
joel@joel-Inspiron-3576:~/Downloads/KMP-Parallel-Using-MPI-master$ mpiexec -n 2 ./test
Invalid MIT-MAGIC-COOKIE-1 key


Result KMP Sequential
Start Point: 49521.
Target Length: 99044
Pattern Length: 122
Time Elapsed: 2581.42600

joel@joel-Inspiron-3576:~/Downloads/KMP-Parallel-Using-MPI-master$ mpicc -0 test kmp_MPI.c
gcc: error: unrecognized command line option '-0'
joel@joel-Inspiron-3576:~/Downloads/KMP-Parallel-Using-MPI-master$ mpicc -o test kmp_MPI.c
joel@joel-Inspiron-3576:~/Downloads/KMP-Parallel-Using-MPI-master$ mpirun -np 2 ./test
Invalid MIT-MAGIC-COOKIE-1 key

Result KMP Sequential


Start Point: 49521.
Target Length: 99044
Pattern Length: 122
Time Elapsed: 698.22000

joel@joel-Inspiron-3576:~/Downloads/KMP-Parallel-Using-MPI-master$ ls
data1000000.txt  HowToRun.txt             KMPSequentical    test
data100000.txt   intrductionReport.pdf    KMPSequentical.c  testmake.py
data10000.txt    kmp_MPI.c                pattern.txt       text.txt
data1000.txt     kmp_MPI.out              Temporary.c
joel@joel-Inspiron-3576:~/Downloads/KMP-Parallel-Using-MPI-master$ mpicc -g -Wall -o test KMPSequentical.c
KMPSequentical.c: In function 'main':
KMPSequentical.c:17:8: warning: unused variable 'j' [-Wunused-variable]
   17 |   int i,j;
      |         ^
KMPSequentical.c:17:6: warning: unused variable 'i' [-Wunused-variable]
   17 |   int i,j;
      |       ^
KMPSequentical.c: In function 'kmptable':
```
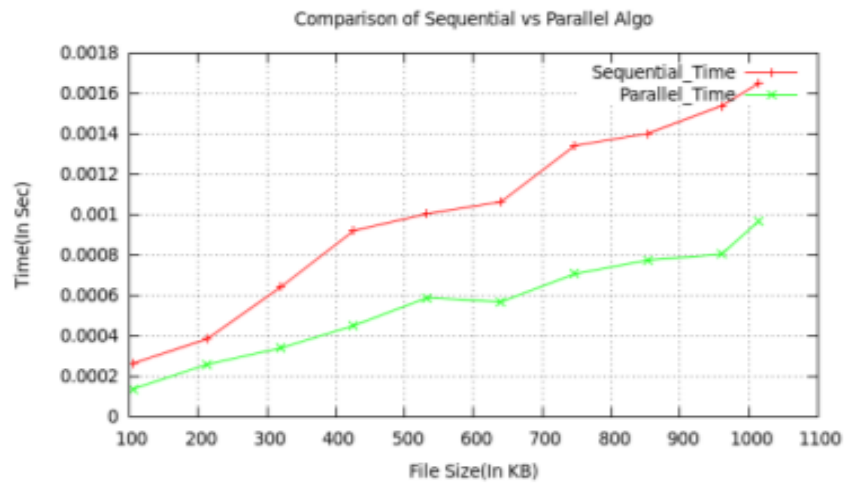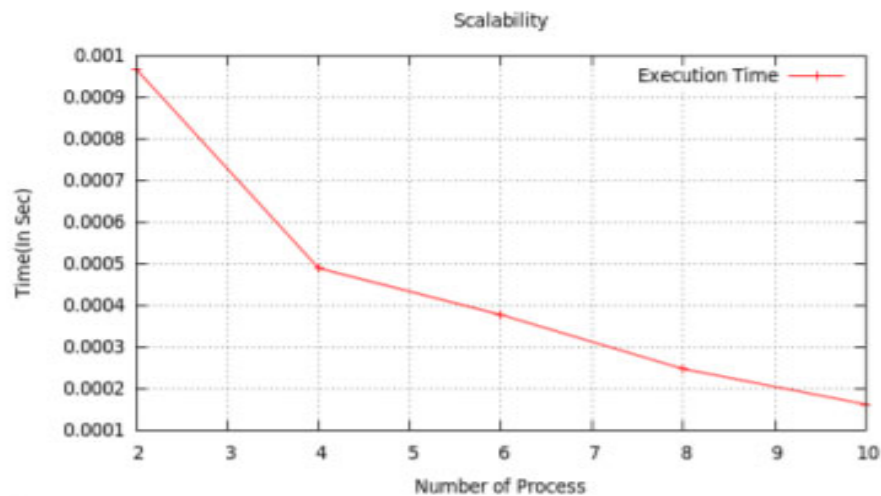
# Result

The most time consuming part of KMP search is the LPS that is Longest Prefix Suffix. The following results show the Compute time of this algorithm.

Sequential Vs Parallel Time Comparison graph:



The above graph shows us the comparison of sequential and parallel time in which we have taken the input file in KB on X-axis and Time(s) to compute in Y-axis by keeping the pattern file constant. So the above graph shows that the parallel computation time is more efficient than the sequential time and the performance of parallel is completely outplayed the sequential performance.

Scalability graph:

The above shows us the number of processes on X-axis and the time taken by that process on Y-axis. In the above graph 10 processes are taken and computed the time taken by them in seconds. It shows that if we go on increasing the number of processes, then the performance gets better. In this way the scalability is achieved.

By viewing the outputs and the graphs we can conclude by saying we implemented the Knuth-Morris-Pratt algorithm to match the Snort pattern and the performance of the parallel algorithm is far better than the sequential algorithm for multiple patterns.

# References

[1] Reetinder Sidhu and Viktor K Prasanna. Fast regular expression matching using fpgas. In Field-Programmable Custom Computing Machines, 2001. FCCM'01.The 9th Annual IEEE Symposium on, pages 227-238. IEEE, 2001.

[2] Christopher R Clark and David E Schimmel. Scalable pattern matching for high speed networks. In Field-Programmable Custom Computing Machines, 2004.FCCM 2004. 12th Annual IEEE Symposium on, pages 249-257. IEEE, 2004.

[3] Sarang Dharmapurikar and John Lockwood. Fast and scalable pattern matching for content _ltering. In Architecture for networking and communications systems, 2005. ANCS 2005. Symposium on, pages 183-192. IEEE, 2005.

[4] Alfred V Aho and Margaret J Corasick. E_cient string matching: an aid to bibliographic search. Communications of the ACM, 18(6):333-340, 1975.

[5] Nen-Fu Huang,Hsien-Wei Hung, ShengHung Lai, Yen-Ming Chu, and Wen-Yen Tsai. A gpu-based multiple-pattern matching algorithm for network intrusion detection systems. Advanced Information Networking and Applications-Workshops, 2008. AINAW 2008. 22nd International Conference on, pages 62-67. IEEE, 2008