

Project Report
On
Parallel String Matching Problem
(SMP)

Submitted by
PRANSHU MAHESHWARI (170001035)

PRAYAG JAIN (170001037)

Computer Science and Engineering

3rd year

Under the Guidance of

Dr. Surya Prakash



Department of Computer Science and Engineering

Indian Institute of Technology Indore

Spring 2019

Index

1. Introduction
2. Algorithm Analysis
3. Results
4. References

1. Introduction

In computer science, string-searching algorithms, sometimes called string-matching algorithms, are an important class of string algorithms that try to find a place where one or several strings (also called patterns) are found within a larger string or text. With parallel string we try to reduce the time complexity of searching a pattern in a given string. We try to reduce complexity of already existing algorithms then finally we come up with our own algorithm.

- **The Problem Statement**

The *string matching problem* (SMP) consists of finding substring (generally pattern) P in text T . In the basic form both P and T consist of characters in the same alphabet Σ .

- **Applications**

Common applications of string matching include spell checking, record linkage. With the availability of large amounts of DNA data, matching of nucleotide sequences has become an important application. String matching is also used in spam filtering.

2. Algorithm Analysis

2.1 Naïve string-matching algorithm (Brute-Force)

2.1.1 Sequential Naïve string-matching algorithm

There are many approaches to solve SMP. The simplest way to find a substring in the given text is a brute force algorithm. It compares all characters of the text with the beginning of the substring. If the characters match, it compares the next characters in the text. The brute force algorithm does not need any additional preparations or space but its time complexity is $O(n*m)$ in the worst case and $O(n)$ in the average case, where n is the size of the text and m is the size of the substring.

2.1.2 Parallel Naïve string-matching algorithm

To parallelize the Naïve string-matching algorithm we divide the text T in $NUM_THREADS$ parts and each processor is given a part of the text of size $T/NUM_THREADS + M - 1$ and the processor performs the searching using the sequential brute force algorithm. This reduces the complexity from $O(n*m)$ to $O(n*m / NUM_THREADS)$.

Code for Parallel Naïve string-matching algorithm

```
int partSize = n / NUM_THREADS;
int remainder = n % NUM_THREADS;

double start = omp_get_wtime();

#pragma omp parallel num_threads(NUM_THREADS) private(tid,start,end) shared(t,a,remainder,partSize,m)
{
    tid=omp_get_thread_num();
    long long j;

    if(tid == 0) {
        start = tid;
        end = partSize - 1;
    } else {
        start = (tid * partSize) - m;
        end = (tid * partSize) + partSize - 1;
    }
    for(long long i=start;i<=end-m;i++) {
        for(j=0;j<m;j++)
            if(t[i+j] != a[j])
                break;
        if(j == m) {
            #pragma omp critical
            v.push_back(i);
        }
    }
}

if(remainder != 0) {
    long long j;
    start = (NUM_THREADS+1) * partSize;
    end = n;
    for(long long i=start;i<=end-m;i++) {
        for(j=0;j<m;j++)
            if(t[i+j] != a[j])
                break;
        if(j == m) {
            #pragma omp critical
            v.push_back(i);
        }
    }
}
```

2.2 KMP (Knuth Morris Pratt) Pattern Searching

KMP algorithm is based on finite automaton. The states in the algorithm are marked with symbols that should match at the moment. There are two transitions from each state. One of them corresponds to the matching of the characters, while the other one corresponds to the mismatching. The automaton goes to the next state, if the comparison is successful, otherwise it goes to the previous state. Prefix-function is used for construction of the KMP automaton. This function is associated with the pattern P and gives the information about the positions of the different prefixes of the string in the pattern P . Prefix-function $\pi: \{1, 2, \dots, m\} \rightarrow \{0, 1, \dots, m-1\}$ that is associated with the pattern $P[1..m]$ is defined as $\pi[q] = \max \{k: k < q \text{ \& } P_k P_q\}$. In other words, $\pi[q]$ is the length of the longest prefix P that is the suffix of P_q . The time complexity is $O(n+m)$.

KMP can also be parallelized by dividing the text T into $NUM_THREADS$ parts of size $T/NUM_THREADS + M - 1$ and each processor performing sequential KMP on each part respectively.

Code for KMP string-matching algorithm

```
4
5 void computeLPSArray(char* pat, int M, int* lps);
6
7 void KMPSearch(char* pat, char* txt)
8 {
9     int M = strlen(pat);
10    int N = strlen(txt);
11
12    int lps[M];
13
14    computeLPSArray(pat, M, lps);
15
16    int i = 0; // index for txt[]
17    int j = 0; // index for pat[]
18    while (i < N) {
19        if (pat[j] == txt[i]) {
20            j++;
21            i++;
22        }
23
24        if (j == M) {
25            printf("Found pattern at index %d ", i - j);
26            j = lps[j - 1];
27        }
28
29        else if (i < N && pat[j] != txt[i]) {
30
31            if (j != 0)
32                j = lps[j - 1];
33            else
34                i = i + 1;
35        }
36    }
37 }
38
39 void computeLPSArray(char* pat, int M, int* lps)
40 {
41     int len = 0;
42
43     lps[0] = 0;
44
45     int i = 1;
46     while (i < M) {
47         if (pat[i] == pat[len]) {
48             len++;
49             lps[i] = len;
50             i++;
51         }
52         else
53         {
54             if (len != 0) {
55                 len = lps[len - 1];
56             }
57
58             else // if (len == 0)
59             {
60                 lps[i] = 0;
61                 i++;
62             }
63         }
64     }
65 }
66 }
67
```

2.3 Shift-AND string-matching algorithm

2.3.1 Sequential Shift-AND string-matching algorithm

In a general case we need to find all occurrences of pattern P in text T . Let us generalize the problem and assume that we should to find occurrences of all possible prefixes of pattern P : $P(1) = p_1$, $P(2) = p_1p_2$, \dots , $P(m) = p_1p_2 \dots p_m$, where p_i is the i -th symbol of pattern P . This method of string matching is known as Shift AND algorithm. The time complexity of the algorithm is $O(m*n)$.

Code for Shift-AND string-matching algorithm

```
int n = t.size(), m = a.size();
int D[m][n];
vector<int> index;

double start = omp_get_wtime();

for(int i=0; i<m; i++) D[i][0] = 0;
for(int i=0; i<n; i++) D[0][i] = (t[i] == a[0]) ? 1 : 0;

for(int i=1; i<m; i++) {
    for(int l=1; l<n; l++) {
        D[i][l] = (a[i] == t[l] && D[i-1][l-1]) ? 1 : 0;
    }
}

double end = omp_get_wtime();

for(int l=0; l<n; l++) {
    if(D[m-1][l])
        index.push_back(l);
}

for(auto it : index) {
    cout<<it<<" ";
}

cout<<"\n";
```


2.3.2 Parallel Shift-AND string-matching algorithm

We can parallelize the Shift AND algorithm by parallelizing the inner loop for each row and thus the time complexity can be reduced to $O(m)$ time with $O(n)$ processors on CREW PRAM model.

Code for Parallel Shift-AND string-matching algorithm

```
int n = t.size(), m = a.size();
int D[m][n];
vector<int> index;

double start = omp_get_wtime();

#pragma omp parallel for
for(int i=0;i<m;i++) {
    D[i][0] = 0;
}

#pragma omp parallel for
for(int i=0;i<n;i++) {
    D[0][i] = (t[i] == a[0]) ? 1 : 0;
}

for(int i=1;i<m;i++) {
    #pragma omp parallel for shared(a[i])
    for(int l=1;l<n;l++) {
        D[i][l] = (a[i] == t[l] && D[i-1][l-1]) ? 1 : 0;
    }
}

double end = omp_get_wtime();

for(int l=0;l<n;l++) {
    if(D[m-1][l])
        index.push_back(i);
}

for(auto it : index) {
    cout<<it<<" ";
}
cout<<"\n";
```

2.4 Dynamic Programming Based Solution

2.4.1 Sequential Solution

A matrix $D[m+1][n+1]$ is constructed, where $D_{i,j}$ is the minimum number of differences between a_1, a_2, \dots, a_i and any contiguous substring of the text ending at t_j . If $D_{l,m} \leq k$ then there must be an occurrence of the pattern in the text with at most k differences that ends at t_m . Time complexity of this algorithm $O(n*m)$. The last row of D matrix gives the number of mismatches for pattern ending at that index. Thus it can be used to perform fuzzy(approximate) search as well.

Code for DP based string-matching algorithm

```
15
16     int n = t.size(), m = a.size();
17     int D[m+1][n+1];
18     vector<int> v;
19
20     double start = omp_get_wtime();
21
22     for(int i=0; i<=n; i++) D[0][i] = 0;
23     for(int i=0; i<=m; i++) D[i][0] = i;
24
25     for(int i=1; i<=m; i++) {
26         for(int l=1; l<=n; l++) {
27             int temp = (a[i-1] == t[l-1]) ? D[i-1][l-1] : D[i-1][l-1] + 1;
28             D[i][l] = min({ D[i-1][l] + 1, D[i][l-1] + 1, temp });
29         }
30     }
31
32     double end = omp_get_wtime();
33
34     for(int l=0; l<=n; l++) {
35         if(!D[m][l]) {
36             v.push_back(l);
37         }
38     }
39
```

2.4.2 Parallel Solution

This solution can be parallelized by running the $n*m$ for-loop twice parallelly, once row wise and the second time column wise. The time complexity of this is $O(n+m)$ and it requires $O(\max(n, m))$ processors on a CREW PRAM model.

Code for Parallel DP based string-matching algorithm

```
#pragma omp parallel for
for(int i=0;i<=n;i++) {
    D[0][i] = 0;
}

#pragma omp parallel for
for(int i=0;i<=m;i++) {
    D[i][0] = i;
}

for(int i=1;i<=m;i++) {
    #pragma omp parallel for shared(a[i-1])
    for(int l=1;l<=n;l++) {
        int temp = (a[i-1] == t[l-1]) ? D[i-1][l-1] : D[i-1][l-1] + 1;
        D[i][l] = min({ D[i-1][l] + 1, temp });
    }
}

for(int l=1;l<=n;l++) {
    #pragma omp parallel for shared(a[l-1])
    for(int i=1;i<=m;i++) {
        int temp = (a[l-1] == t[i-1]) ? D[i-1][l-1] : D[i-1][l-1] + 1;
        D[i][l] = min({ D[i][l-1] + 1, D[i][l], temp });
    }
}

double end = omp_get_wtime();

#pragma omp parallel for
for(int i=m;i<=n;i++) {
    if(!D[m][i]) {
        #pragma omp critical
        index.push_back(i);
    }
}
```

2.5 Our CREW PRAM Solution

We allocate m processors to each character in the text T and these m processors check if the pattern starts from that character or not. It can also be used to perform fuzzy searches in the text.

The time complexity is $O(\lg(m))$ with $O(n*m)$ processors on a CREW PRAM model.

Code for our Parallel CREW string-matching algorithm

```
19     int n = t.size(), m = a.size();
20     vector<vector<int>> > check(m, vector<int>(n, 0));
21     int SUM[n];
22
23     double start = omp_get_wtime();
24
25     #pragma omp parallel for
26     for(long long i=0;i<n*m;i++) {
27         int id = i / m;
28         int j = i % m;
29         if(a[j] == t[id+j]) {
30             check[j][id] = 1;
31         }
32     }
33
34     #pragma omp parallel for
35     for(int i=0;i<n;i++) {
36         SUM[i] = 0;
37         #pragma omp parallel for reduction(+:SUM[i])
38         for(int j=0;j<m;j++) {
39             SUM[i] += check[j][i];
40         }
41     }
42
43     #pragma omp parallel for
44     for(int i=0;i<n;i++) {
45         if(SUM[i] == m) {
46             #pragma omp critical
47             v.push_back(i);
48         }
49     }
50
51     double end = omp_get_wtime();
52
53     for(auto it : v) {
54         cout<<it<<" ";
55     }
56     cout<<"\n";
```

2.6 Our CRCW PRAM Solution

We allocate m processors to each character in the text T and these m processors compares the next m characters with the pattern. If a mismatch is found the processor writes 0 to the `check[i]` variable.

The time complexity is $O(1)$ with $O(n*m)$ processors on a Common CRCW PRAM model.

Code for our Parallel CREW string-matching algorithm

```
int n = t.size(), m = a.size();
vector<int> check(n, 1);

double start = omp_get_wtime();

#pragma omp parallel for
for(long long i=0; i<n*m; i++) {
    int id = i / m;
    int j = i % m;
    if(a[j] != t[id+j]) {
        check[id] = 0;
    }
}

#pragma omp parallel for
for(int i=0; i<n; i++) {
    if(check[i]) {
        #pragma omp critical
        v.push_back(i);
    }
}

double end = omp_get_wtime();

for(auto it : v) {
    cout<<it<<" ";
}
cout<<"\n";
```

Results

Algorithm	Time Complexity	# Processors
Sequential Naïve	$O(n*m)$	$O(1)$
Parallel Naïve	$O(n*m/\text{NUM_THREADS})$	$O(\text{NUM_THREADS})$
KMP	$O(n+m)$	$O(1)$
Parallel KMP	$O((n+m)/\text{NUM_THREADS})$	$O(\text{NUM_THREADS})$
Shift And	$O(n*m)$	$O(1)$
Parallel Shift And	$O(m)$	$O(n)$
DP Solution	$O(n*m)$	$O(1)$
Parallel DP	$O(n+m)$	$O(n)$
Our CREW Solution	$O(\lg(m))$	$O(n*m)$
Our CRCW Solution	$O(1)$	$O(n*m)$

The timing analysis of the algorithms

```
Sequential Brute Force Algorithm
Text size: 10000000
Pattern size: 7
Time: 0.0573285
```

```
Parallel Brute Force Algorithm
Text size: 10000000
Pattern size: 7
NUM_THREADS: 4
Time: 0.0156179
```

```
Our CRCW Algorithm:
Text size: 10000000
Pattern size: 7
NUM_THREADS: 4
Time: 0.424392
```

```
Our CREW Algorithm
Text size: 10000000
Pattern size: 7
NUM_THREADS: 4
Time: 0.030455
```

```
DP Based Algorithm
Text size: 10000000
Pattern size: 7
Time: 0.0261375
```

```
Parallel DP Based Algorithm
Text size: 10000000
Pattern size: 7
NUM_THREADS: 4
Time: 0.11827
```

```
Sequential Shift AND Algorithm
Text size: 10000000
Pattern size: 7
Time: 0.425879
```

```
Parallel Shift AND Algorithm
Text size: 10000000
Pattern size: 7
NUM_THREADS: 4
Time: 0.133095
```

```
Sequential KMP Algorithm:
Text size: 10000000
Pattern size: 7
Time: 0.097149
```

```
Parallel KMP Algorithm:
Text size: 10000000
Pattern size: 7
NUM_THREADS: 4
0.067735
```

Project Link: [Github](#)

References

[1] Bit parallel string matching, Alina Gutnova

June 21, 2006 Available: [Bit parallel string matching - Semantic Scholar](#).

[2] General methods of sequence comparison

https://dornsife.usc.edu/assets/sites/516/docs/papers/msw_papers/msw-054.pdf

[3] <https://www.wikipedia.org/>

[4] <https://cs.stackexchange.com/>