

EVALUATING THE EFFECTIVENESS OF KEY ALGORITHMS OF PATTERN MATCHING TO ENHANCE THE EFFICACY OF BOYER-MOORE-HORSPOOL ALGORITHM AND RABIN KARP ALGORITHM IN MALICIOUS PROGRAM DETECTION

Ruchika Goel

ABSTRACT

Computer viruses are serious threats to society. There are numerous procedures that are used for resolving this downside. Our purpose here is to create a collective look on numerous pattern matching approaches that are employed in the domain of malicious program detection and mitigation. In virus detection, the program starts matching with some known or trained patterns which tell us that the attack has been done. Keywords: Pattern Matching Algorithms

INTRODUCTION

Now a day's pattern matching has numerous approach which includes solving the problem of classification. Lots of algorithms are used in pattern matching concept. Various approaches are used in detecting computer viruses like signature scanning based approach. In pattern matching process antiviruses maintain a database of commonly used file signature and start searching for the pattern in systems. An everyday number of viruses increases, it requires updating of virus databases. In this approach, it takes lots of time to search and match the files which make the process slow. There should be some optimization technique which enhances the process of matching and takes lesser time.

PATTERN MATCHING ALGORITHMS

There are a number of pattern matching algorithms that are employed in the domain of malicious program detection. These algorithms are often classified as [1]

Single Keyword pattern matching –

a) Brute force formula

- b) Karp-Rabin formula
- c) Boyer Moore formula
- d) Backward oracle matching formula
- e) Knuth Morris Pratt formula

Multiple Keyword Pattern Matching –

- a) Aho Corasick formula
- b) Commentz conductor formula
- c) Wu-Manber formula
- d) Fan-Su formula

The Aho-Corasick formula wants a tri-like DFA and a failure perform. the development of those is given below.

Sandeep Kumar in his research A Common Bug Scanner

The formula is analogous to the Aho –Corasick formula, however, it's stretched regex character [3].

? compare any nibble with the input data

%n avoid 0-n nibbles within the input stream

*n skip precisely n nibbles

** skip associate degree absolute variety, as well as zero

for(each nibble within the input stream)

if(cross(inputNibbleStream, TreeRoot))

node *cross(ifNibbleStream& i, node& n)

{

if(there area unit virus signatures related to this node) come &n;

```

ch = upcoming input data ; pos = current position of nibble i;
if(link on ch from n(approx))
if(ret = traverse(i, the node found by following the link from n on ch)) come ret;
for(help, *d, ? & ** links of node n)
sort %d)
price of d; for(uint L = 0; L finish of file; L++)
{
restore file position to posi; skip precisely l nibbles;
if(ret = traverse(i, the node obtained by following the link from n on **))return ret;
}
return data location to posi;
}
}
return 0;
}

```

Pattern matching formula employed in common scanner in c++[three]

Another approach is taken by another formula that create use of Boyer Moore Horspool formula for pattern match. The Boyer Moore Horspool calculation or Horspool calculation is utilized to discover the substring in the information content or substantial record accumulations. In 1980, this calculation was distributed by Nigel Horspool. It is an improvement of the Boyer Moore calculation which is related to the Knuth-Morris-Pratt calculation (KMP)⁶. In the Boyer-Moore-Horspool calculation, it contrasts the content character t_i and the last character p_m of the example. In the event that they coordinate, at that point, it analyzes the past characters of the content with comparing characters in the example successively appropriate to left until to recognize either a recurrence of the example or a befuddle on a content character. Assume, regardless of the match is happened, it slides the example as per the following event of the character t_i in the pattern δ . The quantity of positions to be moved is controlled by the estimation of skip (t_i). Calculation of the skirt Table in the Boyer-Moore Horspool calculation has an

unpretentious distinction with the first avoid table definition proposed in the Boyer-Moore algorithm 11. In the Boyer-Moore calculation, the estimation of skip (pm) is dependable 0. In the Horspool adaptation, skip (pm) = m if pm is special inside the example; generally skip (pm) = m - k, where pm-k is the penultimate (furthest right) appearance of the character pm in the pattern 14. Boyer Moore Horspool calculation [a couple of]

```

1. Initialize pattern length  $m \leftarrow |p|$ ;
2. Initialize the text length  $n \leftarrow |t|$ ;
3. Compute skip table GENERATE-SKIP-TABLE( $\Sigma, p$ ):
   a. Set pattern length  $m \leftarrow |p|$ ;
   b. Initialize skip table skip ( $\sigma$ ) = m for all symbols  $\sigma \in \Sigma$ ;
   c. Initialize pattern index  $j \leftarrow 1$ ;
   d. For jth character  $P_j$  in the pattern, set skip ( $P_j$ )  $\leftarrow m-j$ ;
   e. Increment pattern index  $j \leftarrow j+1$ ;
   f. If  $j < m-1$  then go to step 4;
   g. Stop.
4. Initialize text pointer  $i \leftarrow 0$ ;
5. Initialize pattern pointer  $j \leftarrow m$ ;
6. While  $j > 0$  and  $t_{i+j} = P_j$ 
   Do move pattern pointer to left  $j \leftarrow j-1$ ;
7. If  $j=0$  then
   Print "pattern occurs at text index"  $i+1$ ;
8. Shift the text pointer  $i \leftarrow i + \text{skip}(t_{i+m})$ ;
9. If  $i \leq n - m$  then
   Go to step 5 to continue matching process.
10. Stop

```

Performance per size of information [a pair of] is completed on the idea of databank amounts and a discount within the period is restrained. There is a unit numerous alternative pattern matching formula that we tend to mention that may be changed as shown by alternative authors. alternative a brand new formula is often developed for quick looking out however it needs a deep analysis of existing algorithms intimately.

Rabin- Karp Algorithm

Rabin-Karp Algorithm is the least difficult string seeking calculation. This calculation was produced by Michael O. Rabin and Richard M. Karp in 1987. This calculation utilizes the hash capacity to find the potential example in the info content. For the length of content n and example p of shared length m, its normal and best case running time is $O(n+m)$ in space $O(p)$, and furthermore, the most pessimistic scenario time is $O(nm)$ in space $O(m)$ ¹⁵. It is utilized to find the hash estimation of the specific example substring and afterward, it finds the hash estimation of all conceivable m length substring of the info content. In the event that the hash estimation of the example and content substring match than it restores the esteem generally next substring esteem is coordinated to compute the string of length m.

```

Rabin_Karp_Matcher(T,P,d,q)
1. n ← length [T]
2. m ← length [P]
3. h ← dm-1 mod K
4. p ← 0
5. t0 ← 0
6. for I ← 1 to m
7. do P ← (dp + P[i]) mod K
8. t0 ← (d + T[i]) mod K
9. for s ← 0 to n-m
10. Do if p = ts
11. then if P[1.....m] = T[s+1.....s+m]
12. then print "pattern occurs with shift s "
13. if s < n-m
14. then Ts+1 = (d(ts - T[s+1])h) + T[s+m+1]) mod k

```

Knuth-Morris-Pratt Algorithm

The Knuth–Morris–Pratt has produced a straight time string seeking calculation by examination of the beast drive calculation or innocent calculation. The calculation was created in 1974 by Donald Knuth and Vaughan Pratt, and freely by James H. Morris and they distributed it mutually in 1977. The Knuth-Morris-Pratt calculation directs the aggregate number of examinations of the example against the information string 16. A coordinating time of $O(n)$ is refined by dodging relationship with fundamentals of 'S' that have prior been engaged with the examination with a portion of the particular component of the example 'p' to be coordinated. i.e., backtracking on the string 'S' surely not happens.

```

1. Initialize pattern length j ← 1;
2. Initialize the text length k ← 1;
3. Set length of the pattern m ← |p|;
4. Set length of the text n ← |t|;
5. While j > 0 and pj ≠ tj do
    Shift pattern pointer (j ← Next(j));
6. Advance text pointer i ← i+1
7. Advance pattern pointer j ← j+1;
8. If j > m then
    Print "pattern occurs at text index" i-m
    Else shift pattern pointer j ← Next(j);
9. If i ≤ n and j ≤ m then
    Goto step 5 to continue pattern matching
    Else stop

```

Enhanced Algorithms Enhanced

Boyer-Moore-Horspool Algorithm

This calculation is a capable string coordinating calculation, and it has been the standard perspective for the string coordinating issues. This calculation checks the characters of the example from ideal to left order 6. On these terms of confound or a total match of full example, it utilizes the two capacities to move the window from left to one side and the two capacities are great postfix move an awful character move. The seeking period of the calculation in $O(nm)$ time many-sided quality and the best execution is $O(n/m)$ [18]. First, compute the state change table S from the example P , the example might be a single line, various lines or a record. At that point set the pointer esteem and state esteem. On the off chance that the pointer esteem is littler than the example and content esteem at that point read the character from appropriate to left, starting the furthest right one. For this situation if coordinate happens, the arrival the list of the character, generally move the pointer esteem again a similar procedure will be finished by each level until the point that the example found or not found. [5]

```

1. Input: m- Length of Pattern; n- Length of Text
2. Output: text index
3. Calculate state transition table S from the pattern P
4. Set pointer = 0, state = 1 and ppointer = m - 1
5. While pointer < n - m do
  a. Read character T(aligned + ppointer) into c
  b. (shift, state, ppointer, match, state) = S[state, c]
  c. If match = true then
    i. Return pointer
  d. End if
  e. If state = true AND remaining characters match then
    i. Return pointer
  f. End if
  g. pointer = pointer + shift
6. end while
7. return NIL

```

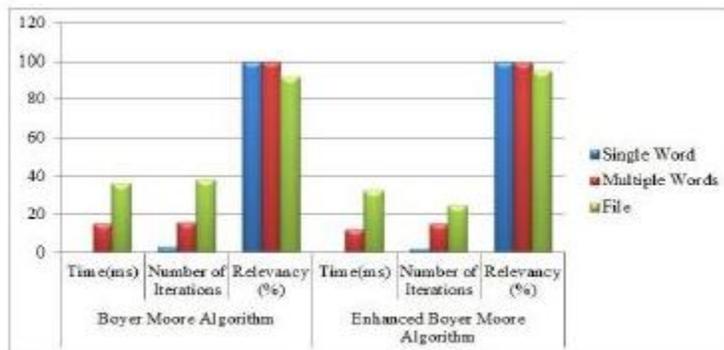
Enhanced Rabin Karp Algorithm

This searching algorithm that uses the hashing function to find any one of a set of pattern in the input text. Hashing offers a simple method to avoid a total number of character comparisons [19]. Instead of checking at each position of the text, it checks only the content of the window whether the pattern occurs or not. For the length of text N and the pattern P of combined length M , its best case running time is $O(N+M)$. And the worst case time is $O(NM)$. First, the algorithm used to find the hash value of the pattern. Then it checks the input text along with its hash value. If a mismatch occurs, shift the window to the next character then calculate the hash value and the same process will continue. Otherwise, it returns the index position of the particular character. [5]

```

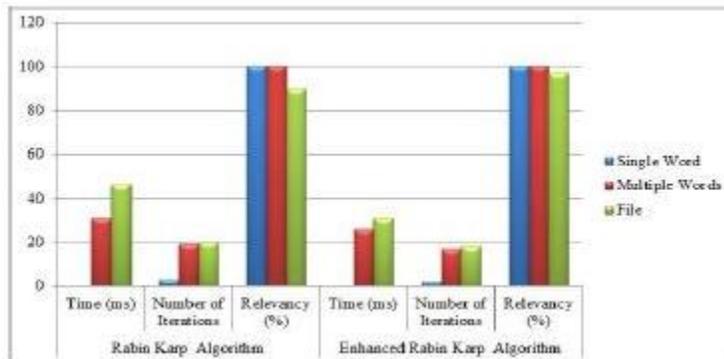
1. Function RabinKarp (S, P, n, m, K, q)
2. Begin
3.  $h \leftarrow K^{m-1} \bmod q$ ;
4.  $p \leftarrow 0$ ;
5.  $t_0 \leftarrow 0$ ;
6. for  $i = 1$  to  $m$  do
7.  $p \leftarrow (K \cdot p + P[i], \bmod q)$ ;
8.  $t_0 \leftarrow (K \cdot t_0 + S[i]) \bmod q$ ;
9. end for
10. for  $j = 0$  to  $n - m$  do
    a. if  $p = t_j$  then
        i. if  $P = S[j+1, j+m]$  then
            1. output  $j+1$ ;
        ii. end if
    b. end if
11. if  $j < n - m$  then
12.  $t_{j+1} = ((K \cdot (t_j - S[j+1]) \cdot h) + S[j+m+1]) \bmod q$ ;
13. end for
14. end
    
```

Performance Analysis of Boyer-Moore and Enhanced Boyer-Moore [5]



Input	Boyer Moore Algorithm			Enhanced Boyer Moore Algorithm		
	Time(ms)	Number of Iterations	Relevancy (%)	Time(ms)	Number of Iterations	Relevancy (%)
Single Word	0	3	100	0	2	100
Multiple Words	15	16	100	12	15	99
File	36	38	92	32	25	95

Rabin Karp Algorithm Performance Analysis[5]



Input	Rabin Karp Algorithm			Enhanced Rabin Karp Algorithm		
	Time (ms)	Number of Iterations	Relevancy (%)	Time (ms)	Number of Iterations	Relevancy (%)
Single Word	0	2	100	0	1	100
Multiple Words	22	16	100	20	12	100
File	45	24	95	30	23	96

Conclusion

In our research, we tend to create an evaluation of numerous algorithms that area unit used for malicious program recognition and display the outcomes achieved by writers to indicate the potency of their changed algorithms. There is a unit heap of labor that may be tired this field. we tend to solely mention 2 pattern matching formulas not bushed the pc virus domain as a result of our purpose here to create a review study that offers an initial concept.

References

[1] www.wikipedia.com.

[2] Sunita Kanaujiya, Dr. S.P.Tripathi, N.C.Sharma, " rising Speed of the Signature Scanner victimization BMH Algorithm" International Journal of pc Applications.

[3] Sandeep Kumar Eugene H. Spafford, " A GenericVirusScannerin C++."Department of pc Sciences Purdue University

[4] R. Janani* and S. Vijayarani IJST, Vol 9(43), DOI: 10.17485/ijst/2016/v9i43/95454, November 2016