# A Fast String Matching Algorithm Based on Lowlight Characters in the Pattern

Zhengjun Cao[1],      Lihua Liu[2]

[1]Department of Mathematics, Shanghai University, China. caozhj@shu.edu.cn

[2]Department of Mathematics, Shanghai Maritime University, Shanghai, China.

**Abstract**

We put forth a new string matching algorithm which matches the pattern from neither the left nor the right end, instead a special position. Comparing with the Knuth-Morris-Pratt algorithm and the Boyer-Moore algorithm, the new algorithm is more flexible to pick the position for starting comparisons. The option really brings it a saving in cost.

**Keywords.** string matching; Knuth-Morris-Pratt algorithm; Boyer-Moore algorithm; finite automata; lowlight character.

**MSC**(2010): 68Qxx, 68P10.

## 1   Introduction

There are various matching problems, such as approximate string matching [1, 5], inverse pattern matching [2], two-dimensional pattern matching [3], real scaled matching [4], scaled dictionary matching [6], property matching [7], weighted matching [7], overlap matching [8], approximate swapped matching [9], combinatorial pattern matching [10] and speculative parallel pattern matching [15]. Among these matching problems, searching for a word in a natural language text is of great importance. It is an important utility in text editors and word processers. Typically, the text is a document being edited, and the pattern searched for is a particular word supplied by the user. There are three general ways of matching technique, standard matching algorithm, Knuth-Morris-Pratt algorithm [14] and Boyer-Moore algorithm [11]. All of them are done from the perspective of character strings. These techniques could, however, be used to search for any string of bits or bytes in a binary file.

In this paper, we introduce a new string matching algorithm, more precisely, for character string matching, not for bits. It makes use of that each English alphabet has its own statistical probability. Unlike Knuth-Morris-Pratt algorithm which matches the pattern from the left end and Boyer-Moore algorithm which matches from the right end, our algorithm matches from a special position. The statistical probability of the character in the position is the smallest among that of all characters in the pattern string. We call such a character (may be not unique) a *lowlight character* of the pattern string. We shall compare the algorithm with two popular algorithms for string matching, the Knuth-Morris-Pratt algorithm and the Boyer-Moore algorithm. The flexible option to pick the position for starting comparisons really brings the new algorithm a saving in cost.

## 2 Three general algorithms for string matching

### 2.1 Standard algorithm

In the standard algorithm, we begin by comparing the first character of the text with the first character of the substring. If they match, we move to the next character of each. This process continues until the entire substring matches the text or the next characters do not match. See the following example for details.

|         |                 | comparisons |
|---------|-----------------|-------------|
| Text:   | there they are  |             |
| Pass 1: | they            | 4           |
| Text:   | there they are  |             |
| Pass 2: | they            | 1           |
| Text:   | there they are  |             |
| Pass 3: | they            | 1           |
| Text:   | there they are  |             |
| Pass 4: | they            | 1           |
| Text:   | there they are  |             |
| Pass 5: | they            | 1           |
| Text:   | there they are  |             |
| Pass 6: | they            | 1           |
| Text:   | there they are  |             |
| Pass 7: | they            | 4           |

It is easy to find that the standard algorithm wastes a lot of effort. If we have matched the beginning part of the substring, we can use that information to tell us how far to move in the text to start the next match.

## 2.2  Knuth-Morris-Pratt algorithm

The Knuth-Morris-Pratt algorithm is based on finite automata but uses a simpler method of handling the situation of when the characters don't match. In the Knuth-Morris-Pratt algorithm, we label the states with the symbol that should match at that point. We then only need two links from each state, one for a successful match and the other for a failure. The success link will take us to the next node in the chain, and the failure link will take us back to a previous node based on the word pattern. Each success link of a Knuth-Morris-Pratt automata causes the "fetch" of a new character from the text. Failure links do not get a new character but reuse the last character fetched. If we reach the final state, we know that we found the substring.

## 2.3  Boyer-Moore algorithm

The Boyer-Moore algorithm is different from the previous two algorithms in that it matches the pattern from the right instead of left end. For example, in the following example, we first compare the $y$ with the $r$ and find a mismatch. Because $r$ doesn't appear in the pattern at all, we know the pattern can be moved to the right a full four characters (the size of the pattern). We next compare the $y$ with the $h$ and find a mismatch. This time because the $h$ does appear in the pattern, we move the pattern only two characters to the right so that the $h$ characters line up. We then begin the match from the right side and find a complete match for the pattern.

|          |                  | comparisons |
|----------|------------------|-------------|
| Text:    | there they are   |             |
| Pass 1:  | they             | 1           |
| Text:    | there  they are  |             |
| Pass 2:  | they             | 1           |
| Text:    | there they are   |             |
| Pass 3:  | they             | 4           |

In the Boyer-Moore algorithm, we did 6 character comparisons verses 13 in the standard algorithm.

# 3 New algorithm

## 3.1 Description and examples

All algorithms mentioned above do not consider each English alphabet has its own statistical probability. Whereas the language property is very useful in daily life, especially in cryptanalysis [16].

Table 1: Statistical probabilities of English alphabets

| character | probability | character | probability |
|-----------|-------------|-----------|-------------|
| A | 0.082 | N | 0.067 |
| B | 0.015 | O | 0.075 |
| C | 0.028 | P | 0.019 |
| D | 0.042 | Q | 0.001 |
| E | 0.127 | R | 0.060 |
| F | 0.022 | S | 0.063 |
| G | 0.020 | T | 0.091 |
| H | 0.061 | U | 0.028 |
| I | 0.070 | V | 0.010 |
| J | 0.002 | W | 0.023 |
| K | 0.008 | X | 0.001 |
| L | 0.040 | Y | 0.020 |
| M | 0.024 | Z | 0.001 |

The basic idea behind the new algorithm is to find a character (may be not unique) which has the smallest probability among that of all characters in the pattern string. For convenience, we call such a character *lowlight character* in the pattern. It then searches the text for the lowlight character. If there is a match, then compare other characters in the pattern string with corresponding characters in the text. Usually, the method matches the pattern neither from the right nor from the left end, instead a special position.

We now describe the algorithm as follows. Suppose that the text is $T_1 \cdots T_n$, and the pattern is $P_1 \cdots P_m$, where $n \geq m$.

(1) *Find a lowlight character in the pattern.* If there are several characters in the pattern which have the same smallest probability, pick the rightmost character. For example, $P_i$ is taken as the lowlight character. Let the left segment be $P_1 \cdots P_{i-1}$, and the right segment be $P_{i+1} \cdots P_m$.

(2) *Search the text for the first mismatch.*

2-1. Compare $T_i$ with $P_i$. If $T_i \neq P_i$, go to step 3-1.

2-2. Compare the left segment. Start the comparisons from the right end of the left segment, i.e., comparing $P_{i-\ell}$ with $T_{i-\ell}$, $\ell = 1, \cdots, i-1$, one after another. Once there is a mismatch, go to step 3-2.

2-3. Compare the right segment. Start the comparisons from the right end of the right segment, i.e., comparing $P_{i+\ell'}$ with $T_{i+\ell'}$, $\ell' = m-i, \cdots, 1$, one after another. Once there is a mismatch, go to step 3-3.

(3) *Align the pattern with the text.*

3-1. If $T_i \neq P_{i-1}, \cdots, P_1$, then align $P_1$ with $T_{i+1}$. If $T_i \neq P_{i-1}, \cdots, P_{k+1}$, and $T_i = P_k, 1 \leq k < i$, then align $P_k$ with $T_i$.

3-2. Suppose that the mismatch appears at the position $s$, namely, $T_s \neq P_s$. If $T_s \neq P_{s-1}, \cdots, P_1$, then align $P_1$ with $T_{s+1}$. If $T_s \neq P_{s-1}, \cdots, P_{l+1}$, and $T_s = P_l, 1 \leq l < s$, then align $P_l$ with $T_s$.

3-3 Suppose that the mismatch appears at the position $s'$, namely, $T_{s'} \neq P_{s'}$. If $T_{s'} \neq P_{s'-1}, \cdots, P_1$, then align $P_1$ with $T_{s'+1}$. If $T_{s'} \neq P_{s'-1}, \cdots, P_{l'+1}$, and $T_{s'} = P_{l'}, 1 \leq l' < s'$, then align $P_{l'}$ with $T_{s'}$.

Now we provide some examples to explain how to use the method.

**Example 1**: Text is "there they are". Pattern is "they".

Since $P_4(y)$ has the smallest probability 0.020 in the pattern, pick it as the lowlight character. Compare $P_4(y)$ with $T_4(r)$. It is a mismatch.

Since $T_4(r)$ does not appear in the pattern, align $T_5(e)$ with $P_1(t)$. Now compare $P_4(y)$ with $T_8(h)$. It is a mismatch, too.

Since $T_8(h)=P_2(h)$, align them and compare other characters.

Like the Boyer-Moore algorithm, the new algorithm needs only 6 comparisons.

**Example 2**: Text is "attach attack attain attempt attend attention attest approve". Pattern is "attempt".

Since $P_6(p)$ has the smallest probability 0.019 in the pattern, pick it as the lowlight character. Compare $P_6(p)$ with $T_6(h)$. It is a mismatch.

Since $T_6(h)$ does not appear in the pattern, align $T_7$(blank) with $P_1(a)$. Compare $P_6(p)$ with $T_{12}(c)$.

Since $T_{12}(c)$ does not appear in the pattern, align $T_{13}$(k) with $P_1(a)$. Compare $P_6(p)$ with $T_{18}(a)$.

Since $T_{18}(a){=}P_1(a)$, align them. Compare $P_6(p)$ with $T_{23}(t)$.

Since $T_{23}(t) {=}P_3(t)$, not $P_2(t)$ (see the description of the new algorithm), align them. Compare $P_6(p)$ with $T_{26}(m)$.

Since $T_{26}(m) {=}P_5(m)$, align them. Compare $P_6(p)$ with $T_{27}(p)$. Compare other characters. Finally, we find the first appearance of the pattern in the text.

The new algorithm needs 12 comparisons. Note that the Boyer-Moore algorithm needs 10 comparisons for this example.

As mentioned earlier, the Boyer-Moore algorithm matches from the right end of the pattern in order to move right more characters once a mismatch occurs. It is more appropriate for a text of plenty of words with a same prefix. But it is insufficient for dealing with a text of plenty of words with a same suffix which is just the suffix of the pattern. Too see the shortcoming of the Boyer-Moore algorithm, we refer to the following example.

**Example 3**. The text is "bear dear fear gear hear near pear rear sear tear wear year", and the pattern is "wear".

In this example, the new algorithm picks $w$ as its lowlight character and matches from the left end of the pattern like the Knuth-Morris-Pratt algorithm. It saves much cost than the Boyer-Moore algorithm.

## 3.2   Refined algorithm

Note that there are two key factors for evaluating the new algorithm:

(1) the position for starting the comparisons in each shift is optimal;

(2) the probability of the character in the position is small enough.

It seems difficult to balance exactly the two requirements. We suggest to take the following measure.

*Refined measure*: Find a lowlight character in the right half of the pattern, instead of the whole pattern. If the chosen character is at the right end of the pattern and it is a component of a common suffix, pick the next-to-last position for starting comparisons. If there are several characters in the right half of the pattern which have the same smallest probability, pick the rightmost character in the right half.

Clearly, the refined measure has no effect on choosing the starting positions in the patterns "they" and "attempt" in example 1 and example 2, separately. But it saves much cost when we use it to deal with the example 3. In such case, the chosen character is $P_3(a)$. See the following process for details.

| Text | be**ar** dear fear gear hear near pear rear sear tear wear year |
|------|------|
| Pass 1 | wear — 3 comparisons: a&a; e&e; w&b |
| Text | bear **d**ear fear gear hear near pear rear sear tear wear year |
| Pass 2 | we**ar** — 1 comparison: a&d |
| Text | bear dea**r** fear gear hear near pear rear sear tear wear year |
| Pass 3 | we**ar** — 1 comparison: a&r |
| Text | bear dear f**e**ar gear hear near pear rear sear tear wear year |
| Pass 4 | we**ar** — 1 comparison: a&e |
| Text | bear dear fe**ar** gear hear near pear rear sear tear wear year |
| Pass 5 | we**ar** — 3 comparisons: a&a; e&e; w&f |
| $\vdots$ | $\vdots$ |

## 3.3  Complexity analysis

Generally, it is reasonable to assume that the pattern is not meaningless. Suppose that the chosen lowlight character $P$ in the pattern has the probability $\lambda$, the length of the text is $n$ and the length of the pattern is $m$. Hence, the text has about $\lambda n$ lowlight character $P$. The amount of comparisons depends essentially on the number of shifts. It is expected that the chosen lowlight character $P$ is at the $m/2$-th position in the pattern and the pattern moves right $m/4$ characters in each shift. Note that each shift is expected to has only one comparison because the first comparison happens to *the starting character $P$ in the pattern which is rarely matched*. Thus, the new algorithm needs about $4n/m$ comparisons. We conjecture that the character comparisons in the new algorithm is of order $\Theta(n/\lambda m)$ provided that $\lambda m \geq 1$. As for the matching time analysis of the Knuth-Morris-Pratt algorithm and the Boyer-Moore algorithm, we refer to [12, 13].

| Algorithm | Matching time |
|-----------|---------------|
| Standard algorithm | $\Theta((n-m+1)m)$ |
| Knuth-Morris-Pratt algorithm | $\Theta(n)$ |
| Boyer-Moore algorithm | $\Omega(n/m), O(nm)$ |
| New algorith | $\Theta(n/\lambda m), \lambda m \geq 1$ (it is conjectured) |

# 4  Conclusion

In this paper, we make use of the statistical probabilities of English alphabets in natural language texts to design a new algorithm for string matching. We hope the presentation could interest some skillful engineers to experiment on it.

# References

[1] A. Amir, Y. Aumann, O. Kapah, A. Levy, E. Porat: Approximate string matching with address bit errors. Theor. Comput. Sci. 410(51): 5334-5346 (2009)

[2] A. Amir, A. Apostolico, M. Lewenstein: Inverse Pattern Matching. J. Algorithms 24(2): 325-339 (1997)

[3] A. Amir, A. Butman, M. Crochemore, G. Landau, M. Schaps: Two-dimensional pattern matching with rotations. Theor. Comput. Sci. 314(1-2): 173-187 (2004)

[4] A. Amir, A. Butman, M. Lewenstein: Real Scaled Matching. Inf. Process. Lett. 70(4): 185-190 (1999)

[5] A. Amir, E. Eisenberg, O. Keller, A. Levy, E. Porat: Approximate string matching with stuck address bits. Theor. Comput. Sci. 412(29): 3537-3544 (2011)

[6] A. Amir, G. Calinescu: Alphabet-Independent and Scaled Dictionary Matching. J. Algorithms 36(1): 34-62 (2000)

[7] A. Amir, E. Chencinski, C. Iliopoulos, T. Kopelowitz, H. Zhang: Property matching and weighted matching. Theor. Comput. Sci. 395(2-3): 298-310 (2008)

[8] A. Amir, R. Cole, R. Hariharan, M. Lewenstein, E. Porat: Overlap matching. Inf. Comput. 181(1): 57-74 (2003)

[9] A. Amir, M. Lewenstein, E. Porat: Approximate swapped matching. Inf. Process. Lett. 83(1): 33-39 (2002)

[10] A. Amir, L. Parida: Combinatorial Pattern Matching. Inf. Comput. 213: 1 (2012)

[11] R. Boyer, J. Moore: A fast string searching algorithm, Carom. ACM 20, (10), 262-272 (1977)

[12] Z. Galil, R. Giancarlo: On the Exact Complexity of String Matching: Lower Bounds. SIAM J. Comput. 20(6): 1008-1020 (1991)

[13] Z. Galil, R. Giancarlo: On the Exact Complexity of String Matching: Upper Bounds. SIAM J. Comput. 21(3): 407-437 (1992)

[14] D. Knuth, J. Morris, V. Pratt: Fast pattern matching in strings. SIAM Journal on Computing, pp. 323-350, 1977.

[15] D. Luchaup, R. Smith, C. Estan, S. Jha: Speculative Parallel Pattern Matching. IEEE Transactions on Information Forensics and Security 6(2): 438-451 (2011)

[16] D. Stinson. Cryptography Theory and Practice, 3 edition. Chapman and Hall/CRC (2005)