

Knuth–Morris–Pratt algorithm

In computer science, the **Knuth–Morris–Pratt algorithm** (or **KMP algorithm**) is a string-searching algorithm that searches for occurrences of a "word" *W* within a main "text string" *S* by employing the observation that when a mismatch occurs, the word itself embodies sufficient information to determine where the next match could begin, thus bypassing re-examination of previously matched characters.

The algorithm was conceived by James H. Morris and independently discovered by Donald Knuth "a few weeks later" from automata theory.^{[1][2]} Morris and Vaughan Pratt published a technical report in 1970.^[3] The three also published the algorithm jointly in 1977.^[1] Independently, in 1969, Matiyasevich^{[4][5]} discovered a similar algorithm, coded by a two-dimensional Turing machine, while studying a string-pattern-matching recognition problem over a binary alphabet. This was the first linear-time algorithm for string matching.^[6]

Knuth–Morris–Pratt algorithm

Class	String search
Data structure	String
Worst-case performance	$\Theta(m)$ preprocessing + $\Theta(n)$ matching ^[note 1]
Worst-case space complexity	$\Theta(m)$

Background

A string-matching algorithm wants to find the starting index *m* in string *S* [] that matches the search word *W* [] .

The most straightforward algorithm, known as the "Brute-force" or "Naive" algorithm, is to look for a word match at each index *m*, i.e. the position in the string being searched that corresponds to the character *S* [*m*]. At each position *m* the algorithm first checks for equality of the first character in the word being searched, i.e. *S* [*m*] =? *W* [0]. If a match is found, the algorithm tests the other characters in the word being searched by checking successive values of the word position index, *i*. The algorithm retrieves the character *W* [*i*] in the word being searched and checks for equality of the expression *S* [*m*+*i*] =? *W* [*i*]. If all successive characters match in *W* at position *m*, then a match is found at that position in the search string. If the index *m* reaches the end of the string then there is no match, in which case the search is said to "fail".

Usually, the trial check will quickly reject the trial match. If the strings are uniformly distributed random letters, then the chance that characters match is 1 in 26. In most cases, the trial check will reject the match at the initial letter. The chance that the first two letters will match is 1 in 26 (1 in 26^2 chances of a match over 26 possible letters). So if the characters are random, then the expected complexity of searching string *S* [] of length *n* is on the order of *n* comparisons or *O*(*n*). The expected performance is very good. If *S* [] is 1 million characters and *W* [] is 1000 characters, then the string search should complete after about 1.04 million character comparisons.

That expected performance is not guaranteed. If the strings are not random, then checking a trial *m* may take many character comparisons. The worst case is if the two strings match in all but the last letter. Imagine that the string *S* [] consists of 1 million characters that are all *A*, and that the word *W* [] is 999 *A* characters terminating in a final *B* character. The simple string-matching algorithm will now examine 1000 characters at each trial position before rejecting the match and advancing the trial position. The simple string search example would now take about 1000 character comparisons times 1 million positions for 1 billion character comparisons. If the length of *W* [] is *k*, then the worst-case performance is *O*(*k*·*n*).

The KMP algorithm has a better worst-case performance than the straightforward algorithm. KMP spends a little time precomputing a table (on the order of the size of *W* [], *O*(*k*)), and then it uses that table to do an efficient search of the string in *O*(*n*).

The difference is that KMP makes use of previous match information that the straightforward algorithm does not. In the example above, when KMP sees a trial match fail on the 1000th character (*i* = 999) because *S* [*m*+999] ≠ *W* [999], it will increment *m* by 1, but it will know that the first 998 characters at the new position already match. KMP matched 999 *A* characters before discovering a mismatch at the 1000th character (position 999). Advancing the trial match position *m* by one throws away the first *A*, so KMP knows there are 998 *A* characters that match *W* [] and does not retest them; that is, KMP sets *i* to 998. KMP maintains its knowledge in the precomputed table and two state variables. When KMP discovers a mismatch, the table determines how much KMP will increase (variable *m*) and where it will resume testing (variable *i*).

KMP algorithm

Example of the search algorithm

To illustrate the algorithm's details, consider a (relatively artificial) run of the algorithm, where *W* = "ABCDABD" and *S* = "ABC ABCDAB ABCDABCDABDE". At any given time, the algorithm is in a state determined by two integers:

- m, denoting the position within S where the prospective match for W begins.
- i, denoting the index of the currently considered character in W.

In each step the algorithm compares $S[m+i]$ with $W[i]$ and increments i if they are equal. This is depicted, at the start of the run, like

```
      1      2
m: 01234567890123456789012
S: ABC ABCDAB ABCDABCDABDE
W: ABCDABD
i: 0123456
```

The algorithm compares successive characters of W to "parallel" characters of S, moving from one to the next by incrementing i if they match. However, in the fourth step $S[3] = ' '$ does not match $W[3] = 'D'$. Rather than beginning to search again at $S[1]$, we note that no 'A' occurs between positions 1 and 2 in S; hence, having checked all those characters previously (and knowing they matched the corresponding characters in W), there is no chance of finding the beginning of a match. Therefore, the algorithm sets $m = 3$ and $i = 0$.

```
      1      2
m: 01234567890123456789012
S: ABC ABCDAB ABCDABCDABDE
W:   ABCDABD
i:   0123456
```

This match fails at the initial character, so the algorithm sets $m = 4$ and $i = 0$

```
      1      2
m: 01234567890123456789012
S: ABC ABCDAB ABCDABCDABDE
W:   ABCDABD
i:   0123456
```

Here, i increments through a nearly complete match "ABCDAB" until $i = 6$ giving a mismatch at $W[6]$ and $S[10]$. However, just prior to the end of the current partial match, there was that substring "AB" that could be the beginning of a new match, so the algorithm must take this into consideration. As these characters match the two characters prior to the current position, those characters need not be checked again; the algorithm sets $m = 8$ (the start of the initial prefix) and $i = 2$ (signaling the first two characters match) and continues matching. Thus the algorithm not only omits previously matched characters of S (the "AB"), but also previously matched characters of W (the prefix "AB").

```
      1      2
m: 01234567890123456789012
S: ABC ABCDAB ABCDABCDABDE
W:   ABCDABD
i:       0123456
```

This search at the new position fails immediately because $W[2]$ (a 'C') does not match $S[10]$ (a ' '). As in the first trial, the mismatch causes the algorithm to return to the beginning of W and begins searching at the mismatched character position of S: $m = 10$, reset $i = 0$.

```
      1      2
m: 01234567890123456789012
S: ABC ABCDAB ABCDABCDABDE
W:   ABCDABD
i:       0123456
```

The match at $m=10$ fails immediately, so the algorithm next tries $m = 11$ and $i = 0$.

```
      1      2
m: 01234567890123456789012
S: ABC ABCDAB ABCDABCDABDE
W:   ABCDABD
i:       0123456
```

Once again, the algorithm matches "ABCDAB", but the next character, 'C', does not match the final character 'D' of the word W. Reasoning as before, the algorithm sets $m = 15$, to start at the two-character string "AB" leading up to the current position, set $i = 2$, and continue matching from the current position.

```
      1      2
m: 01234567890123456789012
S: ABC ABCDAB ABCDABCDABDE
W:   ABCDABD
i:       0123456
```

This time the match is complete, and the first character of the match is $S[15]$.

Description of pseudocode for the search algorithm

The above example contains all the elements of the algorithm. For the moment, we assume the existence of a "partial match" table T , described below, which indicates where we need to look for the start of a new match when a mismatch is found. The entries of T are constructed so that if we have a match starting at $S[m]$ that fails when comparing $S[m + i]$ to $W[i]$, then the next possible match will start at index $m + i - T[i]$ in S (that is, $T[i]$ is the amount of "backtracking" we need to do after a mismatch). This has two implications: first, $T[0] = -1$, which indicates that if $W[0]$ is a mismatch, we cannot backtrack and must simply check the next character; and second, although the next possible match will *begin* at index $m + i - T[i]$, as in the example above, we need not actually check any of the $T[i]$ characters after that, so that we continue searching from $W[T[i]]$. The following is a sample pseudocode implementation of the KMP search algorithm.

```

algorithm kmp_search:
  input:
    an array of characters, S (the text to be searched)
    an array of characters, W (the word sought)
  output:
    an array of integers, P (positions in S at which W is found)
    an integer, nP (number of positions)

  define variables:
    an integer, j  $\leftarrow 0$  (the position of the current character in S)
    an integer, k  $\leftarrow 0$  (the position of the current character in W)
    an array of integers, T (the table, computed elsewhere)

  let nP  $\leftarrow 0$ 

  while j < length(S) do
    if W[k] = S[j] then
      let j  $\leftarrow j + 1$ 
      let k  $\leftarrow k + 1$ 
      if k = length(W) then
        (occurrence found, if only first occurrence is needed,  $m \leftarrow j - k$  may be returned here)
        let P[nP]  $\leftarrow j - k$ , nP  $\leftarrow nP + 1$ 
        let k  $\leftarrow T[k]$  ( $T[\text{length}(W)]$  can't be -1)
    else
      let k  $\leftarrow T[k]$ 
      if k < 0 then
        let j  $\leftarrow j + 1$ 
        let k  $\leftarrow k + 1$ 

```

Efficiency of the search algorithm

Assuming the prior existence of the table T , the search portion of the Knuth–Morris–Pratt algorithm has complexity $O(n)$, where n is the length of S and the O is big-O notation. Except for the fixed overhead incurred in entering and exiting the function, all the computations are performed in the **while** loop. To bound the number of iterations of this loop; observe that T is constructed so that if a match which had begun at $S[m]$ fails while comparing $S[m + i]$ to $W[i]$, then the next possible match must begin at $S[m + (i - T[i])]$. In particular, the next possible match must occur at a higher index than m , so that $T[i] < i$.

This fact implies that the loop can execute at most $2n$ times, since at each iteration it executes one of the two branches in the loop. The first branch invariably increases i and does not change m , so that the index $m + i$ of the currently scrutinized character of S is increased. The second branch adds $i - T[i]$ to m , and as we have seen, this is always a positive number. Thus the location m of the beginning of the current potential match is increased. At the same time, the second branch leaves $m + i$ unchanged, for m gets $i - T[i]$ added to it, and immediately after $T[i]$ gets assigned as the new value of i , hence $\text{new_}m + \text{new_}i = \text{old_}m + \text{old_}i - T[\text{old_}i] + T[\text{old_}i] = \text{old_}m + \text{old_}i$. Now, the loop ends if $m + i = n$; therefore, each branch of the loop can be reached at most n times, since they respectively increase either $m + i$ or m , and $m \leq m + i$: if $m = n$, then certainly $m + i \geq n$, so that since it increases by unit increments at most, we must have had $m + i = n$ at some point in the past, and therefore either way we would be done.

Thus the loop executes at most $2n$ times, showing that the time complexity of the search algorithm is $O(n)$.

Here is another way to think about the runtime: Let us say we begin to match W and S at position i and p . If W exists as a substring of S at p , then $W[0..m] = S[p..p+m]$. Upon success, that is, the word and the text matched at the positions ($W[i] = S[p+i]$), we increase i by 1. Upon failure, that is, the word and the text do not match at the positions ($W[i] \neq S[p+i]$), the text pointer is kept still, while the word pointer is rolled back a certain amount ($i = T[i]$, where T is the

jump table), and we attempt to match $W[T[i]]$ with $S[p+i]$. The maximum number of roll-back of i is bounded by i , that is to say, for any failure, we can only roll back as much as we have progressed up to the failure. Then it is clear the runtime is $2n$.

"Partial match" table (also known as "failure function")

The goal of the table is to allow the algorithm not to match any character of S more than once. The key observation about the nature of a linear search that allows this to happen is that in having checked some segment of the main string against an *initial segment* of the pattern, we know exactly at which places a new potential match which could continue to the current position could begin prior to the current position. In other words, we "pre-search" the pattern itself and compile a list of all possible fallback positions that bypass a maximum of hopeless characters while not sacrificing any potential matches in doing so.

We want to be able to look up, for each position in W , the length of the longest possible initial segment of W leading up to (but not including) that position, other than the full segment starting at $W[0]$ that just failed to match; this is how far we have to backtrack in finding the next match. Hence $T[i]$ is exactly the length of the longest possible *proper* initial segment of W which is also a segment of the substring ending at $W[i - 1]$. We use the convention that the empty string has length 0. Since a mismatch at the very start of the pattern is a special case (there is no possibility of backtracking), we set $T[0] = -1$, as discussed below.

Working example of the table-building algorithm

We consider the example of $W = \text{"ABCDABD"}$ first. We will see that it follows much the same pattern as the main search, and is efficient for similar reasons. We set $T[0] = -1$. To find $T[1]$, we must discover a proper suffix of "A" which is also a prefix of pattern W . But there are no proper suffixes of "A", so we set $T[1] = 0$. To find $T[2]$, we see that the substring $W[0] - W[1]$ ("AB") has a proper suffix "B". However "B" is not a prefix of the pattern W . Therefore, we set $T[2] = 0$.

Continuing to $T[3]$, we first check the proper suffix of length 1, and as in the previous case it fails. Should we also check longer suffixes? No, we now note that there is a shortcut to checking *all* suffixes: let us say that we discovered a proper suffix which is a proper prefix (A proper prefix of a string is not equal to the string itself) and ending at $W[2]$ with length 2 (the maximum possible); then its first character is also a proper prefix of W , hence a proper prefix itself, and it ends at $W[1]$, which we already determined did not occur as $T[2] = 0$ and not $T[2] = 1$. Hence at each stage, the shortcut rule is that one needs to consider checking suffixes of a given size $m+1$ only if a valid suffix of size m was found at the previous stage (i.e. $T[x] = m$) and should not bother to check $m+2$, $m+3$, etc.

Therefore, we need not even concern ourselves with substrings having length 2, and as in the previous case the sole one with length 1 fails, so $T[3] = 0$.

We pass to the subsequent $W[4]$, 'A'. The same logic shows that the longest substring we need to consider has length 1, and as in the previous case it fails since "D" is not a prefix of W . But instead of setting $T[4] = 0$, we can do better by noting that $W[4] = W[0]$, and also that a look-up of $T[4]$ implies the corresponding S character, $S[m+4]$, was a mismatch and therefore $S[m+4] \neq \text{'A'}$. Thus there is no point in restarting the search at $S[m+4]$; we should begin at 1 position ahead. This means that we may shift pattern W by match length plus one character, so $T[4] = -1$.

Considering now the next character, $W[5]$, which is 'B': though by inspection the longest substring would appear to be 'A', we still set $T[5] = 0$. The reasoning is similar to why $T[4] = -1$. $W[5]$ itself extends the prefix match begun with $W[4]$, and we can assume that the corresponding character in S , $S[m+5] \neq \text{'B'}$. So backtracking before $W[5]$ is pointless, but $S[m+5]$ may be 'A', hence $T[5] = 0$.

Finally, we see that the next character in the ongoing segment starting at $W[4] = \text{'A'}$ would be 'B', and indeed this is also $W[5]$. Furthermore, the same argument as above shows that we need not look before $W[4]$ to find a segment for $W[6]$, so that this is it, and we take $T[6] = 2$.

Therefore, we compile the following table:

i	0	1	2	3	4	5	6	7
W[i]	A	B	C	D	A	B	D	
T[i]	-1	0	0	0	-1	0	2	0

Another example:

i	0	1	2	3	4	5	6	7	8	9
W[i]	A	B	A	C	A	B	A	B	C	
T[i]	-1	0	-1	1	-1	0	-1	3	2	0

Another example (slightly changed from the previous example):

i	0	1	2	3	4	5	6	7	8	9
W[i]	A	B	A	C	A	B	A	B	A	
T[i]	-1	0	-1	1	-1	0	-1	3	-1	3

Another more complicated example:

i	00	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24
W[i]	P	A	R	T	I	C	I	P	A	T	E		I	N		P	A	R	A	C	H	U	T	E	
T[i]	-1	0	0	0	0	0	0	-1	0	2	0	0	0	0	0	-1	0	0	3	0	0	0	0	0	0

Description of pseudocode for the table-building algorithm

The example above illustrates the general technique for assembling the table with a minimum of fuss. The principle is that of the overall search: most of the work was already done in getting to the current position, so very little needs to be done in leaving it. The only minor complication is that the logic which is correct late in the string erroneously gives non-proper substrings at the beginning. This necessitates some initialization code.

```
algorithm kmp_table:
  input:
    an array of characters, W (the word to be analyzed)
  output:
    an array of integers, T (the table to be filled)

  define variables:
    an integer, pos ← 1 (the current position we are computing in T)
    an integer, cnd ← 0 (the zero-based index in W of the next character of the current candidate substring)

  let T[0] ← -1

  while pos < length(W) do
    if W[pos] = W[cnd] then
      let T[pos] ← T[cnd]
    else
      let T[pos] ← cnd
      while cnd ≥ 0 and W[pos] ≠ W[cnd] do
        let cnd ← T[cnd]
      let pos ← pos + 1, cnd ← cnd + 1

  let T[pos] ← cnd (only needed when all word occurrences are searched)
```

Efficiency of the table-building algorithm

The time (and space) complexity of the table algorithm is $O(k)$, where k is the length of W .

- The outer loop: pos is initialized to 1, the loop condition is $pos < k$, and pos is increased by 1 in every iteration of the loop. Thus the loop will take $k - 1$ iterations.
- The inner loop: cnd is initialized to 0 and gets increased by at most 1 in each outer loop iteration. $T[cnd]$ is always less than cnd , so cnd gets decreased by at least 1 in each inner loop iteration; the inner loop condition is $cnd \geq 0$. This means that the inner loop can execute at most as many times in total, as the outer loop has executed – each decrease of cnd by 1 in the inner loop needs to have a corresponding increase by 1 in the outer loop. Since the outer loop takes $k - 1$ iterations, the inner loop can take no more than $k - 1$ iterations in total.

Combined, the outer and inner loops take at most $2k - 2$ iterations. This corresponds to $O(k)$ time complexity using the Big O notation.

Efficiency of the KMP algorithm

Since the two portions of the algorithm have, respectively, complexities of $O(k)$ and $O(n)$, the complexity of the overall algorithm is $O(n + k)$.

These complexities are the same, no matter how many repetitive patterns are in W or S .

Variants

A real-time version of KMP can be implemented using a separate failure function table for each character in the alphabet. If a mismatch occurs on character x in the text, the failure function table for character x is consulted for the index i in the pattern at which the mismatch took place. This will return the length of the longest substring ending at i matching a prefix of the pattern, with the added condition that the character after the prefix is x . With this restriction, character x in the text need not be checked again in the next phase, and so only a constant number of operations are executed between the processing of each index of the text. This satisfies the real-time computing restriction.

Booth's algorithm uses a modified version of the KMP preprocessing function to find the lexicographically minimal string rotation. The failure function is progressively calculated as the string is rotated.

Notes

1. m is the length of the pattern, which is the string we are searching for in the text which is of length n

References

1. Knuth, Donald; Morris, James H.; Pratt, Vaughan (1977). "Fast pattern matching in strings". *SIAM Journal on Computing*. **6** (2): 323–350. CiteSeerX 10.1.1.93.8147 (<https://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.93.8147>). doi:10.1137/0206024 (<https://doi.org/10.1137%2F0206024>).
2. Knuth, Donald E. (1973). "The Dangers of Computer-Science Theory". *Studies in Logic and the Foundations of Mathematics*. **74**: 189–195. doi:10.1016/S0049-237X(09)70357-X (<https://doi.org/10.1016%2FS0049-237X%2809%2970357-X>). ISBN 9780444104915.
3. Morris, J.H., Jr; Pratt, V. (1970). *A linear pattern-matching algorithm* (Technical report). University of California, Berkeley, Computation Center. TR-40.
4. Матиясевич, Юрий (1971). "О распознавании в реальное время отношения вхождения" (http://gdz.sub.uni-goettingen.de/pdfcache/PPN502141972/PPN502141972__LOG_0019.pdf) (PDF). *Записки научных семинаров Ленинградского отделения Математического института им. В.А.Стеклова* (in Russian). **20**: 104–114., translated into English as Matiyasevich, Yuri (1973). "Real-time recognition of the inclusion relation" (<http://logic.pdmi.ras.ru/~yumat/Journal/inclusion/inclusion.pdf.gz>). *Journal of Soviet Mathematics*. **1**: 64–70. doi:10.1007/BF01117471 (<https://doi.org/10.1007%2FBF01117471>). S2CID 121919479 (<https://api.semanticscholar.org/CorpusID:121919479>).
5. Knuth mentions this fact in the errata of his book *Selected Papers on Design of Algorithms* :

I learned in 2012 that Yuri Matiyasevich had anticipated the linear-time pattern matching and pattern preprocessing algorithms of this paper, in the special case of a binary alphabet, already in 1969. He presented them as constructions for a Turing machine with a two-dimensional working memory.

6. Amir, Amihood; Landau, Gad M.; Lewenstein, Moshe; Sokol, Dina (2007). "Dynamic text and static pattern matching". *ACM Trans. Algorithms*. **3** (2): 19. doi:10.1145/1240233.1240242 (<https://doi.org/10.1145%2F1240233.1240242>). S2CID 8409826 (<https://api.semanticscholar.org/CorpusID:8409826>).
- Cormen, Thomas; Leiserson, Charles E.; Rivest, Ronald L.; Stein, Clifford (2001). "Section 32.4: The Knuth-Morris-Pratt algorithm". *Introduction to Algorithms* (https://archive.org/details/introductiontoal00corm_691) (Second ed.). MIT Press and McGraw-Hill. pp. 923 (https://archive.org/details/introductiontoal00corm_691/page/n945)–931. ISBN 0-262-03293-7. Zbl 1047.68161 (<https://zbmath.org/?format=complete&q=an:1047.68161>).
 - Crochemore, Maxime; Rytter, Wojciech (2003). *Jewels of stringology. Text algorithms*. River Edge, NJ: World Scientific. pp. 20–25. ISBN 981-02-4897-0. Zbl 1078.68151 (<https://zbmath.org/?format=complete&q=an:1078.68151>).
 - Szpankowski, Wojciech (2001). *Average case analysis of algorithms on sequences*. Wiley-Interscience Series in Discrete Mathematics and Optimization. With a foreword by Philippe Flajolet. Chichester: Wiley. pp. 15–17, 136–141. ISBN 0-471-24063-X. Zbl 0968.68205 (<https://zbmath.org/?format=complete&q=an:0968.68205>).

External links

- String Searching Applet animation (<http://www.cs.pitt.edu/~kirk/cs1501/animations/String.html>)
- An explanation of the algorithm (<http://www.ics.uci.edu/~eppstein/161/960227.html>) and sample C++ code (<http://www.ics.uci.edu/~eppstein/161/kmp/>) by David Eppstein
- Knuth-Morris-Pratt algorithm (<http://www-igm.univ-mlv.fr/~lecroq/string/node8.html>) description and C code by Christian Charras and Thierry Lecroq
- Explanation of the algorithm from scratch (<https://www.inf.hs-flensburg.de/lang/algorithmen/pattern/kmpen.htm>) by H.W. Lang
- Breaking down steps of running KMP (<https://web.archive.org/web/20101227102334/http://oak.cs.ucla.edu/cs144/examples/KMPSearch.html>) by Chu-Cheng Hsieh.
- NPTELHRD YouTube lecture video (https://www.youtube.com/watch?v=Zj_er99KMb8)
- LogicFirst YouTube lecture video (<https://www.youtube.com/watch?v=4jY57Ehc14Y>)