

# String Search

- Brute force
- Rabin-Karp
- Knuth-Morris-Pratt
- Right-Left scan

## Rabin-Karp


## Right-Left scan

[illegible]

- word processors
- virus scanning
- text information retrieval (ex: Lexis/Nexis)
- digital libraries
- computational biology
- web search engines

- suffix-trie index costs  $\sim N \lg N$

## String Searching

- **match existence:** any occurrence of pattern in text?
- **enumerate:** how many occurrences?
- **match:** return index of any occurrence  focus of this lecture
- **all matches:** return indices of all occurrences
- 

•  
k v j l i x a p e j r b x e e n p p h k h t h b k w y r w a m n u g z h p p f x i y j y a n h a p f w b g h x  
m s h r l y u j f j h r s o v k v y e l n b x n a w a v g i z y v m f o h i g e a b g k s f n b k m f f x j  
f f q b u a l e y t q r p h y r b j d q j a v c t g x j i f g g y d h o i w h r v w b x g r i x y d z  
b p a j n h o p l v l a m h f [REDACTED] g g i k n g k w z i g j t l x k o z j l e f l b r b o i  
g n b z s u d s s v q y m n a p b p q v l u b d o y x k k w h c o u d v t k m i k a n g s u t d j y t h z l  
a p a w l v l i y j k m x o r z e o a f e o f b f x u h k z u k e f t n r f m o c y l c u l k s e d g r d  
i v a y j p g k r t e d e h w h r v v b b l t d k c t q

## Modelling String Searching

kvjljaxeprjbxenpphkhthbkwyrmwnugzhppfxijyanhapfwbgxmshrluyfjhrsovkvevlnbxnawagvlyzmfhoibagbksfnbkmfxjffqbualeytrphrybrjgdjvavetcxjfgggfvgdhoihrvwbqgixydzbpajnhopvlamhhfavocldftvvggkikngkwxjzjlxkoxjleflibrboi gnbzsdvsqymnnapbqdvlyubdoykxkwhcoudvtkmiansgutsutdztel

all the world's a stage and all the men and women merely players. They have their exits and their entrances, and one man in his time plays many parts. At first, the infant, mewling and puking in the nurse's arms. Then the whining schoolboy, with his satchel and shining morning face, creeping like snail unwillingly to school. And then the lover, sighing

- probability of match is less than  $N/(\text{alphabet size})^M$

.000000000000000084  
for sample problems

## Brute-force string searching

Check for pattern at every text position

using array to simplify alg descriptions  
online apps need getchar()-based implementation

```
int bruteforcesearch(char p[], char a[]) {
    int i, j;
    for (i = 0; i < strlen(a); i++) {
        for (j = 0; j < strlen(p); j++) {
            if (a[i+j] != p[j]) break;
            if (j == strlen(p)) return i;
        }
        return strlen(a);
    }
}
```

Annotations:   
- Blue arrow points to `char a[]`   
- Red arrow points to `strlen(a)` with label "mismatch"   
- Red arrow points to `strlen(p)` with label "match"   
- Red text "text loop" points to the inner `for` loop   
- Red text "pattern loop" points to the inner `if` statements

- returns `i` if leftmost pattern occurrence starts at `a[i]`

## Brute-force string searching (bug fixed)

Check for pattern at every text position

```
int bruteforcesearch(char p[], char a[]) {
    int M = strlen(p), N = strlen(a);
    int i, j;
    for (i = 0; i < N; i++) {
        for (j = 0; j < M; j++) {
            if (a[i+j] != p[j]) break;
            if (j == M) return i;
        }
        return N;
    }
}
```

Annotations:   
- Blue arrow points to `M = strlen(p)` and `N = strlen(a)` with label "lengths won't change; precompute them"   
- Red text "text loop" points to the inner `for` loop   
- Red text "pattern loop" points to the inner `if` statements

- returns `i` if leftmost pattern occurrence starts at `a[i]`

## Problem with brute-force implementation

```
for (i = 0; i < strlen(a); i++)
```

In C, `strlen` takes time proportional to string length

- evaluated every time through loop
- running time is at least  $N^2$
- same problem for simpler programs (ex: count the blanks)

### PERFORMANCE BUG

Textbook example: Performance matters in ADT design

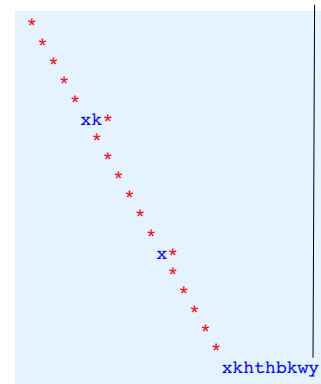
Exercise: implement string ADT with fast `strlen`

- need space to store length

## Brute-force typical case

pattern: xkthbkw

text: kvjlixkpejrbxeenppxkthbkw



```
text: 000000000000000000000000000000001
```

Backs up in text:  
need to retain M-char buffer

$$\begin{aligned} 31415 &= 84 \pmod{97} \\ 14159 &= 94 \pmod{97} \\ 41592 &= 76 \pmod{97} \\ 15926 &= 18 \pmod{97} \\ 59265 &= 95 \pmod{97} \end{aligned}$$
$$\begin{aligned} 31415 \bmod 97 &= 84 \\ 14159 \bmod 97 &= (84 - 3 \cdot 9) \cdot 10 + 9 \pmod{97} = 94 \\ 41592 \bmod 97 &= (94 - 1 \cdot 9) \cdot 10 + 2 \pmod{97} = 76 \\ 15926 \bmod 97 &= (76 - 4 \cdot 9) \cdot 10 + 6 \pmod{97} = 18 \\ 59265 \bmod 97 &= (18 - 1 \cdot 9) \cdot 10 + 5 \pmod{97} = 95 \end{aligned}$$
$$31415 \pmod{97} = 84$$
$$95 \neq 84$$
$$84 - 3 \cdot 9 \pmod{97} = 57$$

## Randomized algorithms

A randomized algorithm uses random numbers to gain efficiency

- quicksort with random partitioning element
- randomized BSTs
- Rabin-Karp

### Las Vegas algorithm

- expected to be fast
- guaranteed to be correct

Examples: quicksort, randomized BSTs, Rabin-Karp with match check

### Monte Carlo algorithm

- guaranteed to be fast

## String search implementations cost summary

Search for an M-character pattern in an N-character text

	typical	worst
brute-force	$N^M$	$N^*M$
Rabin-Karp	$N^M$	$N^M$

= assumes appropriate model

= assumed system can produce "random" numbers

## Knuth-Morris-Pratt algorithm

**Observation:** On mismatch at pattern char j we know the previous j-1 chars in the text (they are also in the pattern)

**Idea:** precompute what to do on mismatch

Example 1: mismatch 00000\* when searching for 000001 in binary text

- text had 000000
- compare next text char with last pattern char

Example 2: mismatch 000\* when searching for 000001 in binary text

- text had 0001
- compare next text char with first pattern char

char to check is completely deduced from pattern

### KMP algorithm

- precompute table of pattern char to check on mismatch, indexed by pattern position

## KMP examples

pattern: 10100110  
text: 1001110100101010100110000111

```

10*
 1*
 1*
 101001*
      1010*
        1010*
          1010*
            10100110
0120111234562343434345678
    
```

mismatch table

	0	1	2	3	4	5	6	7
0	1	0	1	3	0	2	1	

pattern: 00000001  
text: 001000001000000000000001

```

00*
 00000*
    0000000*
      0000000*
        0000000*
          0000000*
            0000000*
              0000000*
                0000000*
    
```

mismatch table

	0	1	2	3	4	5	6	7
0	0	0	0	0	0	0	7	

## KMP implementation

### Check for pattern at every text position

- char match: increment both i and j
- char mismatch: **set j to mismatch[j]**

```
int kmpsearch(char p[], char a[], int mismatch[])
{
    int M = strlen(p), N = strlen(a);
    int i, j = 0;
    for (i = 0; i < N; i++) {
        if (a[i] == p[j]) j++;
        else j = mismatch[j];
        if (j == M) return i-M+1;
    }
    return N;
}
```

← no pointer backup  
in this version

Differs from brute-force in two very significant ways

## KMP mismatch table construction

Table builds itself (!!)

mismatch table							
0	1	2	3	4	5	6	7
0	1	0	1	3	0	6	7

**Idea 1:** Simulate restart ala brute-force

Ex: mismatch[6] for 10100110

- if mismatch at 101001\* then text was 1010010
- for 010010x, x compares to p[6]
- note also: for 010011x, x compares to p[1]

100111010010  
101001\*  
1010010110

pattern: 10100110  
text: 010010x  
\*  
10\*  
10  
0012012

**Idea 2:** Remember simulation for previous entry

Ex: mismatch[7] for 10100110

- if mismatch at 1010011\* then text was 10100111
- just noted: for 010011x, x compares to p[1]

pattern: 10100110  
text: 010011x  
\*  
10\*  
1\*  
0012011

## KMP mismatch table construction implementation

**mis[j]:** index of pattern char to compare against next text char  
on mismatch on jth pattern character

**t:** index of pattern char that brute-force algorithm would compare  
against next text char on iteration after mismatch

To compute mis[j], compare p[j] with p[t]

**match:**

- mismatch[j] = mismatch[t] since mismatch action same as for t
- t = t+1 since we know that brute-force algorithm will find match

**mismatch:** opposite assignment

```
t = 0; mismatch[0] = 0;
for (int j = 1; j < M; j++)
    if (p[j] == p[t])
        { mismatch[j] = mismatch[t]; t = t+1; }
    else
```

	0	1	2	3	4	5	6	7
j	0	1	1	0	0	1	6	7
t	0	0	0	0	1	0	1	0
mismatch	0	1	0	1	3	0	6	7

## Optimized KMP implementation

Easy to create specialized program for given pattern  
(build in mismatch table)

```
int kmpsearch(char a[])
{
    int i = 0;
    s0: if (a[i] != '1') { i++; goto s0; }
    s1: if (a[i] != '0') { i++; goto s1; }
    s2: if (a[i] != '1') { i++; goto s0; }
    s3: if (a[i] != '0') { i++; goto s1; }
    s4: if (a[i] != '0') { i++; goto s3; }
    s5: if (a[i] != '1') { i++; goto s0; }
    s6: if (a[i] != '1') { i++; goto s2; }
    s7: if (a[i] != '0') { i++; goto s1; }
    return i-8;
}
```

assumes pattern in text  
(can use sentinel)

pattern 10100110

mismatch table

## String search implementations cost summary

Search for an  $M$ -character pattern in an  $N$ -character text

	typical	worst
brute-force	$N^2$	$N \cdot M$
Rabin-Karp	$N^2$	$N^2$ ← inner loop with several arithmetic instructions
Knuth-Morris-Pratt	$N$	$N$ ← tiny inner loop

= assumes appropriate model

= assumed pattern can produce "random" numbers

## Right-left pattern scan

Sublinear algorithms

- move right to left in pattern
- move left to right in text

Q: Does binary string have 9 consecutive 0s?

pattern: 000000000

text: 100111010010100010100111000111



A: No. (Needed to look at only 6 of 30 chars.)

Idea effective for general patterns, larger alphabet

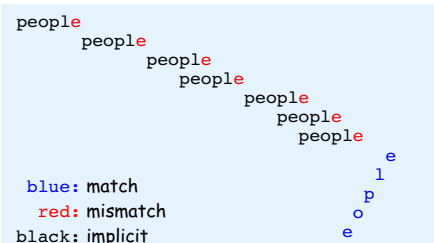
## Right-left scan example

Text char **not** in pattern: skip forward  $M$  chars

Text char **in** pattern: skip to end of pattern

pattern: people

text: now is the time for all good people to



char	skip
a	6
b	6
...	6
e	0
...	6
l	1
...	6
o	3
p	2
...	6
z	6

## Implementation of right-left pattern scan

```
void initskip(char *p)
{
    int j, M = strlen(p);
    for (j = 0; j < 256; j++) skip[j] = M;
    for (j = 0; j < M; j++) skip[p[j]] = M-j-1;
}

int misearch(char *p, char *a)
{
    int M = strlen(p), N = strlen(a);
    int i, j;
    initskip(p);
    for (i = M-1, j = M-1; j >= 0; i--, j--)
        while (a[i] != p[j])
        {
            i += max(M-j, skip[a[i]]);
            if (i >= N) return N;
            i = M-1;
        }
}
```

## String search implementations cost summary

Search for an  $M$ -character pattern in an  $N$ -character text

	typical	worst
brute-force	$N^2$	$N \cdot M$
Rabin-Karp	$N^2$	$N^2$
Knuth-Morris-Pratt	$N$	$N$
Boyer-Moore	$N/M$	$N$



beats optimal by a factor of 100 for  $M = 100$

$\equiv$  assumes appropriate model

$\equiv$  assumes system can produce "random" numbers

## String search summary

Ingenious algorithms for a fundamental problem

Rabin-Karp

- easy to implement
- extends to more general settings (ex: 2D search)

Knuth-Morris-Pratt

- quintessential solution to theoretical problem
- works well in practice, too (no backup, tight inner loop)

Right-left scan

- simple idea leads to dramatic speedup for long patterns

Tin of the iceberg (stay tuned)