# String Matching

**problem:** given a text $T$ w/ $n$ char
$P$ w/ $m$ char
find the first occurance of $P$ in $T$

**solutions**
1. straight forward
2. Rabin-karp
3. Boyer-Moore

**convections:**

$P$ = pattern $\qquad$ $T$ = text



$K$
(position in P)

$j$
(position in T)

## ① Straight forward.

```
int simpleScan (char [] P, char [] T, int m) {
    int i, j, k;    ← all
            ↳ i: guess of index of P in T
              j: position in T
              k: position in P
    while (j < n)
         if (T[j] == P[k]) {
            j++;
            k++)
            if (k == m) return i; }
         else
            j = ++i)
            k = 0;
            if (j > n-m) break; }   return -1;
```

if mismatch
just shift
one to
the
right.

less than pattern left ⇒ not possible,
so leave.

n-m < j ← length check
pkm      i is index
Tjn                                            → O(mn)

① Straight forward.

```
int simpleScan ( char [] P , char [] T , int m ) {
    int i, j, k ;   ← all
        ↳ i: guess of index of P in T
        j : position in T
        k : position in P
    while ( j < n )
        if ( T[j] == P[k] ) {
            j++ ;
            k++ )                          → j-k
            if ( k == m) return i ; }
        else              → j-k+1
            j = ++i )
            k = 0 ;        ↙ less than pattern left ⇒ not possible,
            if ( j > n-m) break ; }          so leave.
                                      return -1 ;
```

if mismatch just shift one to the right. ↘

worst case
    P :   c ..... k          c for m-1
                            k for m-1 th
    T :   c ........... c        c for all

for every j , compare m times   (mth comp breaks,
                 ↓                  move on to next j)
                 n
∴  m × (n-m+1)
            ↪ last j is at n-m
    ∴ O(mn) comparisons.

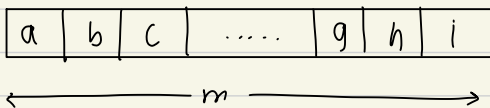best case :     T = P etc.        O(m) comparisons

$j: 1 \to m-1$   $p = p*d + P[j]$

$t_j = (t_{j-1} - T[j] \cdot d^{n-1}) \, d + T[j+m-1]$

## ② Rabin - Karp Algorithm

① convert pattern (m char) to number $p$
② convert first m char of text to number : $t$
③ if $p = t \to$ pattern found, exit
④ else    if (! end of text)                    $\checkmark^{n-m+1}$
    shift window to the right by 1 and recompute

for pattern of m characters                 $d =$ number of possible characters

| a | b | c |  .... | g | h | i |

$\longleftarrow \qquad m \qquad \longrightarrow$

$p = P[0] \cdot d^{m-1}$
$\quad + P[1] \cdot d^{m-2}$
$\quad + P[2] \cdot d^{m-3}$

ⓐ    $p = P[0] ;$
$\quad$ for $j = 1$ to $m-1$              $\vdots$
$\qquad p = p*d + P[j] ; \quad \leftarrow \theta(m)$        $+ P[m-1]$
$\qquad\qquad\qquad\qquad$ for computation

ⓑ    recomputation of $t$

$$t_j = \underbrace{(t_{j-1} \overbrace{- T[j] \cdot d^{m-1}}^{1}) \cdot d}_{2} + \overbrace{T[j+m-1]}^{3} \leftarrow O(1)$$

computation :   $\theta(m) + \theta(m) + \theta(n-m)$
$\qquad\qquad\qquad \underset{\text{for } P}{\uparrow} \qquad \underset{\text{for first } T}{\uparrow} \qquad \underset{\text{remaining } T}{\uparrow}$

① $\to \quad T[j] * d^{m-1} \leftarrow$ oldest
$\qquad d*ht - T[j]*d^m \quad \leftarrow$ oldest . removed + adv
$\qquad + T[j+m-1] \longrightarrow \quad \leftarrow$ new_added

problem w/ Robin-Karp
if the pattern is too long, the resulting number
will be huge and can cause overflow

soln: Hashing
          ↓
     take the mod w/ a prime number q
                       ‿‿‿‿‿‿‿‿‿‿‿
                        ↳
                        large to prevent multiple
① hp = pattern mod q          hashing

② ht = (text first m) mod q

③ if (ht = hp) : compare windows, if equal return

   else. if not eof, advance window

④    :C exit.

ⓐ  the hash method
   $(X+Y) \mod q = X \mod q + Y \mod q$
   $XY \mod q = [(X \mod q)(Y \mod q)] \mod q$
   $(3*10+6) \mod 13 = ((3*10) \mod 13 + 6 \mod 13) \mod 13$   mod
                     $= [(( 3 \mod 13)*(10 \mod 13) \mod 13 + 6 \mod 13)13$
                     $= (4 + 6 \mod 13) \mod 3$
                     $= 10$

$$(((a*10 + d)*10 + c) \bmod q$$
$$\underline{\phantom{(((a*10 + d)*10 + c) \bmod q}}$$
$$\bmod q \quad * \quad (c \bmod q) \bmod q$$

```
int hash( Txt , m, d )
{
    int  h = Txt[0] % q ;
    for ( int i = 1 ; i < m ; i++ )          ← θ(m)
        h = (h * d + Txt[i] ) % q ;
    return h ;                         ⌐ prevents overflow
}
```

(b)
```
int  rehash ( T, j, m, ht ){
    oldest = (T[i] * dM ) % q ;          // dM = d^{m-1} mod q
    oldest_removed = ((ht + q) - oldest ) % q ;
    return (oldest_removed * d + T[i+m] ) % q ;          ↑
}                                                         ★
                                                         ★★
                                                         ★★
```

$$t_j = \left( t_{j-1} - \overset{oldest}{T[j] \cdot d^{m-1}} \right) \cdot d + T[j+m-1]$$

$\underset{\%q}{\underbrace{\phantom{t_j}}}$

$\underset{\text{oldest\_removed} \quad \%q}{\underbrace{\phantom{oldest\_removed}}}$

$\underset{\%q}{\underbrace{\phantom{oldest\_removed oldest\_removed}}}$

```
int RKscan (P,T) {
    m = length (P)
    n = length (T)
    dM = 1 ;
    for j = 1 → m-1     dM = dM * d % q ;    } θ(m-1)

    hp = hash (P, m, d );    θ(m-1)
    ht = hash (T, m, d );    θ(m-1)

    for ( j = 0; j <= n-m ; j++)
        if ( hp == ht && equal_string (P,T,0,j,m))
            return j ;
        if ( j < n-m)   ht = rehash (T, j, m, ht );
    }                              ↘ θ(1)
    return -1;
```

n-m+1 iterations

RK
→  $\theta((n-m+1) m)$    : worst case
→  $O(n+m) +$   time for spurious hits      ⌐ real case
    ↘ $O(m)$ for 2 hash calls.
         $O(n)$ for loop

→  number of spurious hits can be kept low by using
    a large prime number q for the hash function

# ③ Boyer - Moore Algorithm

text T w/ n char       T[1] is first
pattern P w/ m char.      P[1] is first

T from left to right.
P from right to left.

preprocessing step : generates two tables based on
    ↑            which we slide the pattern
efficiency            as much as possible after a
comes from          mismatch
here
     ↘ improves with pattern size

```
int BMscan ( char[] P, char[] T, int m, int[] charJump,
                 int[] matchJump ) {

     int j ; int k ;                                    (j+m <= n)
     j = m ;  k = m ;    (right to left)
     while ( j <= n ) {                                    ↑
         if (k < 1) return j + 1 ;   match!           (j < n-m+1)
         if (T[j] == P[k])  {j--; k--; }

         else {  j + = max (charJump[T[j]], matchJump[k]);
                 k = m ; }
     }
     return -1 ;
}
```
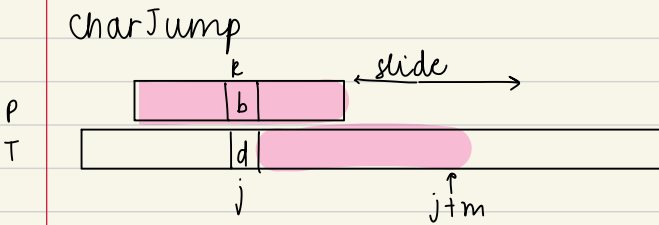
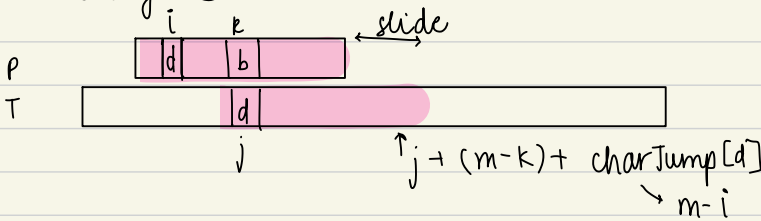starting indexis 1 →

shift window

↖ reset pattern pointer

if nothing :  j = j + m
if in array :  j = j + charjump[char]  ⟶ m-k ← starts from 1

ⓐ  charJump



if d doesn't appear in P at all , line up P after T[j]

if d does appear in P , d in P & d in T(rightmost)
are aligned .



j + (m-k) + charJump[d]
           ↘ m-i

```
                                                    charJump
void computejumps ( char [] P, int m, int alpha, int []
{
    char ch ; int k ;              → number of unique chars
    for ( ch = 0 , ch < alpha ; ch++)
        charJump [ch] = m ;              O(Σ+m)
    for (k = 1 ; k <= m ; k++ )
        charJump[P[k]] = m-k ;
}                        ↘ position from end.
```

O(alpha){  (for loop)
O(m) {  (for loop)

alpha → number of unique chars
m-k → position from end.

this can fail if the index of rightmost 'd' is greater
than k   (moves it back)
↑ check if  m-k+1 is greater and then shift
        ↘ full slide

| | | | vs | m+k-1 |
|---|---|---|---|---|
| a | m-1 | 5 | | |
| b | m-2 | 4 | | |
| c | m-3 | 3 | | |
| d | m-4 | 2 | | |
| e | m-5 | 1 | | |
| f | m-6 | 0 | | |

## ⑥ matchJump

this heuristic tries to derive the maximum shift
from the structure of the pattern.
defined for each char in P (diff from charJump)

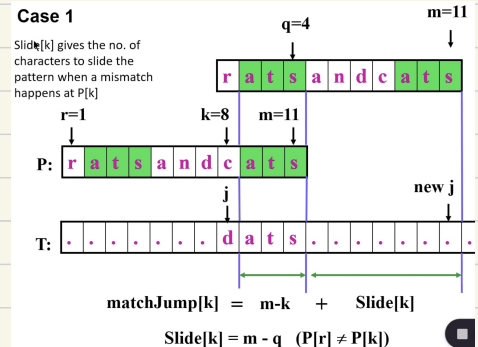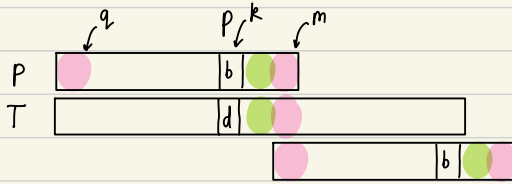case 1:



● : matching suffix
① it occurs again in P
② the preceeding character is
   different

$$j = m-k + slide[k]$$

slide by
~~~~
↓
slide by
suffix length
(j aligned w/ end of P)

m-q
~~~~
↳ slide by dist up to q

(again first
character
is at P[1])

**Case 1**

Slide[k] gives the no. of
characters to slide the
pattern when a mismatch
happens at P[k]



q=4          m=11
↓            ↓
| r | a | t | s | a | n | d | c | a | t | s |

r=1              k=8   m=11
↓                ↓     ↓
P: | r | a | t | s | a | n | d | c | a | t | s |

                 j              new j
                 ↓
T: | · | · | · | · | · | d | a | t | s | · | · | · | · | · | · | · |

matchJump[k] = m-k + Slide[k]

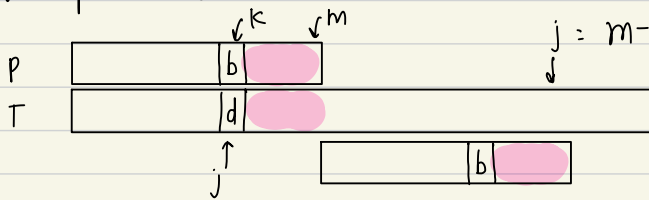Slide[k] = m - q   (P[r] ≠ P[k])

**case 2** P

T

$q$  $P$  $k$  $m$

$b$  $d$

only a part of the matching suffix is at the beginning of the pattern (prefix)

$j = m-k + slide[k]$
$\hookrightarrow m-q$

**Case 3:** there is no other occurrence of the matching suffix in the pattern



P

T

$k$  $m$

$b$  $d$

$j$  $b$

$j = m-k + slide[k]$

$\hookrightarrow m$  $(q=0)$