

## Text searching

$T[1..n]$  - text

$P[1..m]$  - pattern

Is  $P$  a substring of  $T$ ? If so, where?

---

ALMOSTBRUTEFORCE( $T[1..n], P[1..m]$ ):

for  $s \leftarrow 1$  to  $n - m + 1$   
    *shift*  $\rightarrow$   $equal \leftarrow \text{TRUE}$   
     $i \leftarrow 1$   
    while  $equal$  and  $i \leq m$   
        if  $T[s + i - 1] \neq P[i]$   
             $equal \leftarrow \text{FALSE}$   
        else  
             $i \leftarrow i + 1$   
    if  $equal$   
        return  $s$   
return NONE

$O(mn)$  time

$s$   
~~562074~~  
AAAA

$T = \text{AAAA} \dots \text{AAAA}$

$P = \text{AAA} \dots \text{AB}$

Treat strings as integers

$$P = \sum_{i=0}^{m-1} 10^i P[i]$$

$$t_s = \text{num value of } T[s \dots s+m-1] = \sum_{i=0}^{m-1} 10^i \cdot T[s+i]$$

$$t_{s+1} = 10(t_s - 10^{m-1} T[s]) + T[s+m]$$

NUMBERSEARCH( $T[1..n], P[1..m]$ ):

$\sigma \leftarrow 10^{m-1} \bmod q$

$p \leftarrow 0$

$t_1 \leftarrow 0$

for  $i \leftarrow 1$  to  $m$

$p \leftarrow (10 \cdot p + P[i]) \bmod q$

$t_1 \leftarrow (10 \cdot t_1 + T[i]) \bmod q$

for  $s \leftarrow 1$  to  $n - m + 1$

if  $p = t_s$  IF  $P = T[s..s+m-1]$

return  $s$

$t_{s+1} \leftarrow (10 \cdot (t_s - \sigma \cdot T[s]) + T[s+m]) \bmod q$

return NONE

Time:  $O(n + Fm) = O(n + nm/q) = O(n + m)$

$\uparrow$  #False matches

If the  $\bmod q$  values are "random" — They aren't.

$E[F] = n \cdot \frac{1}{q} \Rightarrow$  choose  $q > n$

$hp \leftarrow h(P)$  //  $O(m)$  time

for  $s \leftarrow 1$  to  $n - m$

$hts \leftarrow h(T[s..s+m-1])$  //  $O(m)$  time

if  $hp = hts$

compare brute force

Rolling hash/Sliding hash

Robin/Karp

Zobrist hashing

hash chess positions

$\text{hashpiece}_1 \oplus \text{hashpiece}_2 \oplus \dots$

$\oplus \text{hashpiece}_1$

$\oplus \text{hashpiece}_1$

KARPRABIN( $T[1..n], P[1..m]$ ):

$q \leftarrow$  a random prime number between 2 and  $\lceil m^2 \lg m \rceil$

$\sigma \leftarrow 10^{m-1} \bmod q$

$\tilde{p} \leftarrow 0$

$\tilde{t}_1 \leftarrow 0$

for  $i \leftarrow 1$  to  $m$

$\tilde{p} \leftarrow (10 \cdot \tilde{p} \bmod q) + P[i] \bmod q$

$\tilde{t}_1 \leftarrow (10 \cdot \tilde{t}_1 \bmod q) + T[i] \bmod q$

for  $s \leftarrow 1$  to  $n - m + 1$

if  $\tilde{p} = \tilde{t}_s$

if  $P = T_s$   $\llbracket$ brute-force  $O(m)$ -time comparison $\rrbracket$

return  $s$

$\tilde{t}_{s+1} \leftarrow (10 \cdot (\tilde{t}_s - (\sigma \cdot T[s] \bmod q) \bmod q) \bmod q) + T[s + m] \bmod q$

return NONE

$\Theta(m^2)$  choices

Prime # Theorem: there are  $\Theta(N/\log N)$  primes less than  $N$ .

Lemma: Any integer  $x$  has  $O(\log x)$  prime factors.

$\Pr[\tilde{p} = \tilde{t}_s] \ll O(1/m)$

$p < 10^m$  so  $p$  has  $O(m)$  prime factors

$t_s < 10^m$   $\text{---} O(m)$  prime factors

$|p - t_s| < 10^m$   $\text{---}$  " " "

$O(m)$  of the  $\Theta(m^2)$  primes  $q$  cause a collision

$$E[F] \leq \frac{n}{m} \rightarrow E[\text{Time}] = O(n + \underbrace{E[F]}_{+m} \cdot m) = O(n + m)$$

CARTERWEGMANKARPRABIN( $T[1..n], P[1..m]$ ):

$q \leftarrow$  an arbitrary prime number larger than  $m^2$

$b \leftarrow \text{RANDOM}(q) - 1$   $\llcorner$ uniform between 0 and  $q - 1$  $\llcorner$

$\sigma \leftarrow b^{m-1} \bmod q$

$\tilde{p} \leftarrow 0$

$\tilde{t}_1 \leftarrow 0$

for  $i \leftarrow 1$  to  $m$

$\tilde{p} \leftarrow (b \cdot \tilde{p} \bmod q) + P[i] \bmod q$

$\tilde{t}_1 \leftarrow (b \cdot \tilde{t}_1 \bmod q) + T[i] \bmod q$

for  $s \leftarrow 1$  to  $n - m + 1$

if  $\tilde{p} = \tilde{t}_s$

if  $P = T_s$

$\llcorner$ brute-force  $O(m)$ -time comparison $\llcorner$

return  $s$

$\tilde{t}_{s+1} \leftarrow (b \cdot (\tilde{t}_s - (\sigma \cdot T[s] \bmod q) \bmod q) \bmod q) + T[s + m] \bmod q$

return NONE

Treat  $P[1..m]$  as coeffs of poly of degree  $m-1$   
evaluate at random  $b$ .

$$\Pr[\tilde{p} = \tilde{t}_s] < \frac{m}{m^2} = 1/m$$

Any poly of degree  $m$  has  $\leq m$  roots.

$\rightarrow$   $O(m+n)$  exp. time