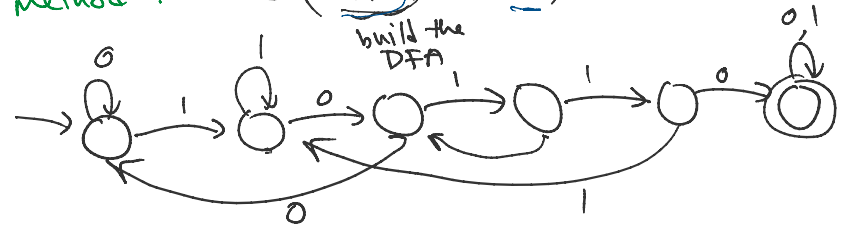


Problem Given strings  $u = a_1 a_2 \dots a_n \in \{0,1\}^*$  (text)  $\Sigma^*$   
 $v = b_1 b_2 \dots b_m \in \{0,1\}^*$  (pattern)  $\Sigma^*$   
 is  $v$  a substring of  $u$ ?  
 $(m \ll n)$

e.g.  $u = 01101011010$   
 $v = 10110$

brute force:  $O(mn)$  time

"DFA method":  $O(\#(m) + n)$  time



Knuth-Morris-Pratt '77:  $O(n)$  time  
 (by compressed version of DFA)  
 (regardless of  $|\Sigma|$ )

... a rand.  $O(n)$ -time alg'm that is simpler ...

Subproblem Alice has a string  $u \in \{0,1\}^*$   
 Bob has a string  $v \in \{0,1\}^*$  ( $|u|=|v|=n$ )

Want to test whether  $u=v$ .

Goal - minimize communication complexity  
 i.e., # bits transmitted

Obvious idea - Alice transmits  $n$  bits

another idea - "checksum"  
 check # of 1's ( $\log n$  bits)

Wrong! counterex:  $101 \neq 110$

another idea -

use a mapping  $F: \{0,1\}^* \rightarrow \{0,1,\dots,p-1\}$

e.g. think of string as base-2 number, then mod  $p$

randomized idea 1 - pick  $p$  randomly

randomized idea 2 - fix prime  $p$ .  
use rand base  $x$

let  $F_x: \{0,1\}^* \rightarrow \{0,1,\dots,p-1\}$

$$F_x(a_{n-1}a_{n-2}\dots a_0) = \left( \sum_{i=0}^{n-1} a_i x^i \right) \bmod p$$

fingerprint function  $\rightarrow$

Monte Carlo Alg'm:

$$x = \text{rand}(0, p-1)$$

$\leq \log p$  bits transmitted

Alice transmits  $F_x(u)$  to Bob

Bob says "probably equal" if  $F_x(u) = F_x(v)$   
"definitely not equal" else

Error Analysis:

if  $u = v$ , correct

if  $u \neq v$ ,

$$\text{say } u = a_{n-1}a_{n-2}\dots a_0 \\ v = b_{n-1}b_{n-2}\dots b_0,$$

$$\text{alg'm errs} \iff F_x(u) = F_x(v) \\ \iff \sum_{i=0}^{n-1} a_i x^i \equiv \sum_{i=0}^{n-1} b_i x^i \pmod{p}$$

$$\iff \sum_{i=0}^{n-1} (a_i - b_i) x^i \equiv 0 \pmod{p}$$

Counting Lemma

A non-zero polynomial of degree  $\leq n-1$   
has  $\leq n-1$  roots, mod  $p$ .

$\#$   $x$ 's satisfying (\*)

has  $\leq n-1$  roots, mod  $p$ .

# x's satisfying (\*)  
total # x's

$\Rightarrow \Pr(\text{alg errs}) \leq \frac{n-1}{p} \leq \frac{1}{n^d}$

pick  $p \approx n^{d+1} \Rightarrow \# \text{ bits transmitted } \log p \approx O(\log n)$ .



Back to string matching problem...

Rabin-Karp Rand. Alg's (Monte Carlo)

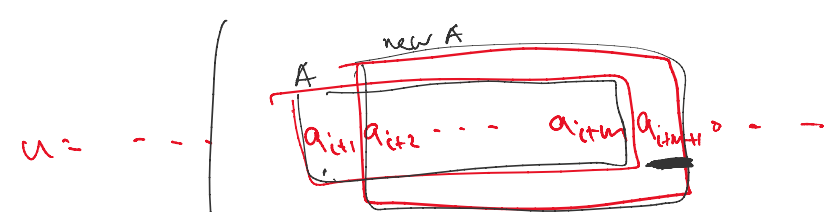
$x = \text{rand}(0, p-1)$

$B = F_x(b_1 b_2 \dots b_m)$

$A = F_x(a_1 a_2 \dots a_m)$

for  $i = 0$  to  $n-m$  {  
if  $A = B$

return "probably match at position  $i$ "



// update A from  $F_x(a_{i+1} \dots a_{i+m})$  to  $F_x(a_{i+2} \dots a_{i+m+1})$

$A = (Ax + a_{i+m+1} - a_{i+1}x^m) \text{ mod } p$

precompute

$O(1)$  time

} " ... match "

} ... ~~return~~ ~~"definitely no match"~~  $\leftarrow O(1)$  time

$\Rightarrow O(n)$  arithmetic op  
on  $(\log p)$ -bit #s  
 $\leftarrow O(\log n)$ -bit #s

### Error Analysis:

let  $E_i =$  (alg'm errs at  $i$ th iteration)

by Alice-Bob,  $\Pr(E_i) \leq \frac{m-1}{p}$

$$\Pr(\text{alg'm errs}) = \Pr\left(\bigcup_{i=0}^{n-m} E_i\right)$$

$$\leq \sum_{i=0}^{n-m} \Pr(E_i)$$

$$\leq n \frac{m}{p} \leq \frac{n^2}{p}$$

pick  $p \approx n^{d+2}$

$$\approx \frac{1}{n^d}$$

$\Rightarrow \boxed{O(n)}$  time

### Las Vegas Version:

1. run Monte Carlo version of Karp-Rabin
2. if it says "probable match at  $i$ " {
3. verify  $a_{i+1} \dots a_{i+m} = b_1 \dots b_m$  in  $O(m)$  time
4. if so, return "match at  $i$ "
5. else run brute force

Always Correct

$$\text{Prob} \geq 1 - \frac{1}{n^d}$$

Always Correct

Expected runtime analysis:

if line 3 true,

else

$$O(n+m) = O(n) \text{ time}$$

$$O(mn) \text{ time} \leftarrow \text{Prob} \leq \frac{1}{n^d}$$

$$\text{Prob} \leq 1 - \frac{1}{n^d}$$

$$\text{expected runtime} \leq O(n) \cdot \left(1 - \frac{1}{n^d}\right) + O(mn) \cdot \frac{1}{n^d}$$

$d=1$ :

$$\leq \boxed{O(n)}$$