

# String Matching

$\Sigma$  alphabet

$T, P$  strings over  $\Sigma$

$$n = |T|, m = |P|$$

$P$  occurs in  $T$  with shift  $s$

$\Downarrow$

$$T[s+1..s+m] = P[1..m]$$

$s$  is a valid shift when  $P$  occurs in  $T$  with shifts

String matching problem for pattern  $P$

input text  $T$

output all valid shifts / occurrences of  $P$

## Notation

$\Sigma^*$  = set of all finite length strings over  $\Sigma$

$\epsilon$  empty strings

$w \sqsubset x$   $w$  is a prefix of  $x$   $\Leftrightarrow x = wy$  for some  $y \in \Sigma^*$

$w \sqsupset x$   $w$  is a suffix of  $x$   $\Leftrightarrow x = yw$  - - -

T a|b|a|c|a|b|b|a|b|a

P b|b|a|b

T 1 2 3 4 5 | 6 7 8 9 |  
a|b|a|c|a|b|b|a|b|a

b|b|a|b

Match

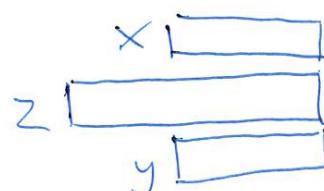
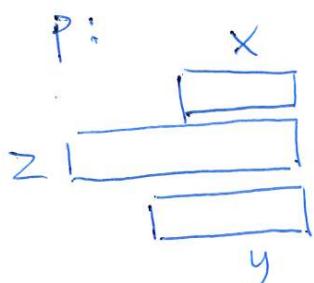
$s=5$  valid shift as  $T[6..9] = P$

### Lemma 32.1 (overlapping-suffix lemma)

Let  $x, y, z \in \Sigma^*$  such that

$x \sqsupseteq z$  and  $y \sqsupseteq z$

- If  $|x| \leq |y|$  then  $x \sqsupseteq y$
- if  $|y| \leq |x|$  then  $y \sqsupseteq x$
- if  $|x| = |y|$  then  $x = y$



### Notation

$$T[i..j] = T[i]T[i+1] \dots T[j]$$

$T[1..k]$  is denoted  $P_k$   $k \geq 0$

do for  $T[1..k]$

So  $\boxed{P \text{ occurs in } T \text{ with shift } s}$

$$\boxed{P \sqsupseteq T_{s+m} = T[1..s+m]}$$

For a given shift  $s$  we can test whether  $s$  is a valid shift in time  $\Theta(t+1)$  when  $t$  is the largest value s.t  $P_t \sqsupseteq T_{s+t}$   
 $t \leq |P|$

Naive algorithm: try all possible shifts

$$s = 0, 1, 2, \dots, n-m$$

Preprocessing: None

Running time  $O((n-m+1)m)$  worst case

Right for  $T = a^n$   $P = a^m$

problem: the algorithm completely ignores info collected while trying shifts

## The Rabin-Karp algorithm

Idea: think of  $P$  and  $T$  as numbers

and compare  $P \times T_{s+m}$  by comparing  
and  $T[s+1..s+m]$   
then numbers

$d = |\Sigma|$  so numbers base  $d$

For simplicity assume  $d = 10$   $\Sigma = \{0, 1, 2, \dots, 9\}$   
but works for all finite alphabets

$$\text{Ex } T[s+1..s+5] = 37133$$

$$P[1..5] = 37233$$

$P_5 \neq T_{s+5}$  since  $37133 \neq 37233$

### Notation

$p$  = decimal value of  $P[1..m]$

$t_s$  = decimal value of  $T[s+1..s+m]$

so  $\boxed{\text{match} \Leftrightarrow p = t_s}$

Assumption: size of the numbers  $p$  and  $t_s$  can be handled (in constant time)  
 so  $p \stackrel{?}{=} t_s$  takes  $O(1)$

Calculation of  $P$ ; use Horner's rule

see example  
 →

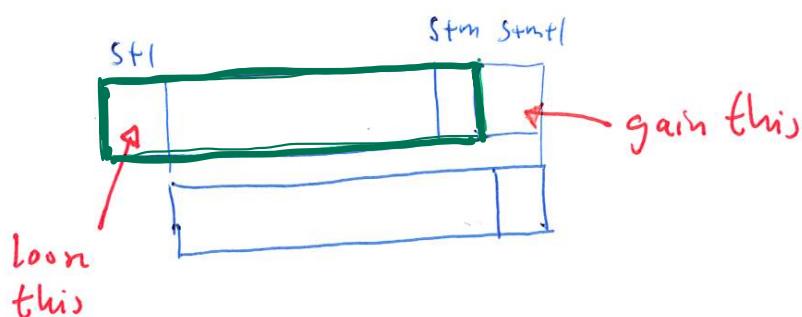
$$p = p[m] + 10(p[m-1] + 10(p[m-2] + \dots + 10p[1]) \dots)$$

similarly  $t_s$  calculated from  $T[1..m]$

in  $\Theta(m)$

Calculating  $t_{s+1}$  from  $t_s$ :

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



so  $O(1)$  to calculate next  $t_{s+1}$

when  $10^{m-1}$  is precalculated  
 (takes  $O(m)$ , can be done in  $O(\log m)$ )

eg. 33

Example of calculation using Horner's rule

$$47632 = 4 \cdot 10^4 + 7 \cdot 10^3 + 6 \cdot 10^2 + 3 \cdot 10 + 2 \cdot 10^0$$

$$= 2 \cdot 10^0 + 3 \cdot 10^1 + 6 \cdot 10^2 + 7 \cdot 10^3 + 4 \cdot 10^4$$

$$= 2 + 10(3 + 10(6 + 10(7 + 10 \cdot 4)))$$

## R-K-algorithm:

- Compute  $p$  in time  $\Theta(m)$
- Compute  $t_0$  —||— report match if  $p = t_0$
- Compute  $t_1$  —||—  $\Theta(1)$  report match if  $p = t_1$
- |
- |
- |
- |
- |
- Compute  $t_{n-m}$  in time  $\Theta(n)$  report match if  $p = t_{n-m}$

Total time  $\Theta(m)$  prep +  $\Theta(n-m)$  match

$\Theta(n-m+1)$  match time

Too good to be true?

Yes!

Unrealistic that we can compare  $p$  and  $t_s$  in constant time when  $m$  is not small

Solution Compute  $p$  and  $t_s$  modulo  $q$   
for suitable value of  $q$

Can calculate  $p \bmod q$  in  $\Theta(n)$  (by Horner's rule)

$t_0 \bmod q$  in  $\Theta(n)$  — — —

$t_1, \dots, t_{n-m}$  in total  $\Theta(n-m+1)$

good Value of  $q$ ?

as large as possible s.t. log fib in one  
computer word.

Then we can do all calculations with  
single precision arithmetic

Can of general alphabet  $\Sigma = \{0, 1, 2, \dots, d-1\}$

choose  $q$  s.t.  $d^q$  fits in one computer word

$$t_{s+1} = (d(t_s - T[s+1] \cdot h) + T[s+m+1]) \bmod q$$

when  $h = d^{m-1} \bmod q$

new problem :

Spurious hits  $t_s \equiv p \pmod{q}$  but  $t_s \neq p$

Such hits must be eliminated

↓  
need to test all shifts with  
 $t_s \equiv p \pmod{q}$        $\Theta(m)$  pr test

Hence worst case matching time is

$\Theta((n-m+1)m)$  as the naive algorithm.

happens for  $T = a^n$      $P = a^m$

Expected # hits

assumption  $\pmod{q} \underset{\text{acts like}}{\sim} \text{random map } \Sigma^* \rightarrow \mathbb{Z}_q$

(Compare with hashing )

If all the  $q$  values mod  $q$  are equally likely

then  $\Pr(t_s \equiv p \pmod{q}) = \frac{1}{q}$

so expected # spurious hits is  $\frac{\mathcal{O}(n)}{q} = \mathcal{O}\left(\frac{n}{q}\right)$

Expected matching time is thus

$$O(n) + O(m(v + n/q))$$

where  $v$  is the number of valid (correct) shifts

when  $v \in O(1)$  and  $q \geq m$

this is  $O(n)$

so total time (preprocessing + matching)

$$\text{is } O(n+m) = O(n) \text{ as } n \geq m$$

## String matching with DFA's

DFA 5-tuple  $M = (Q, q_0, A, \Sigma, \delta)$  Recall from DM565

- $Q$  finite set of states
- $q_0$  initial (start) state
- $A \subseteq Q$  the accept states
- $\Sigma$  finite input alphabet
- $\delta : Q \times \Sigma \rightarrow Q$  transition function of  $M$

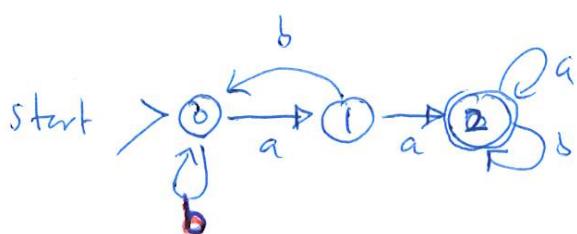
DFA  $\sim$  labelled digraph (labels on arcs)

each state is a vertex

each vertex has precisely

$|\Sigma|$  arcs leaving it

(may be loops)



$M \text{ accepts } w \in \Sigma^*$



following arcs labelled by  $w$  end  
in a state (vertex)  $q$  with  $q \in A$

$M \text{ rejects } w$  if it does not accept it

$M$  induces a function  $\phi$  the  
final state function which is  
the state in which  $M$  resides after  
reading  $w$

Hence  $M$  accepts  $w \Leftrightarrow \phi(w) \in A$

$$\phi(\epsilon) = q_0$$

$$\phi(wa) = \mathcal{J}(\phi(w), a) \quad w \in \Sigma^* \quad a \in \Sigma$$

## String matching automata

$P$  pattern  $\rightarrow M = M(P)$

Goal  $M$  will accept each prefix of  $T_{\text{str}}$  of a text  $T$  for which the last  $m=|P|$  characters of  $T_{\text{str}}$  equal  $P$ .

Show that we can build  $M$  independently of  $T$ , such that same  $M = M(P)$  works for all input texts.

### Notation

suffix function for pattern  $P$ :

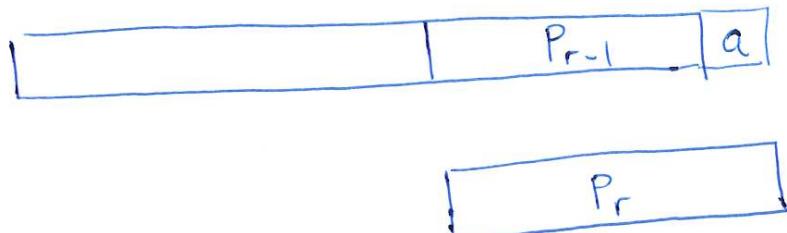
$$\sigma : \Sigma^* \rightarrow \{0, 1, 2, \dots, m\}$$

$$\boxed{\sigma(x) = \max \{k \mid P_k \sqsupseteq x\}}$$

$\sigma$  is well defined as  $\epsilon \sqsupseteq x \quad \forall x \in \Sigma^*$

Lemma 32.2

$$\forall x \in \Sigma^*, a \in \Sigma \quad \sigma(xa) \leq \sigma(x) + 1$$



If  $r = \sigma(xa)$  then the last  $r-1$  characters of  $x$  equal  $P_{r-1}$  so  $\sigma(x) \geq r-1$

Define string matching automaton

$M = M(P)$      $P = P[1..m]$  as follow

$$Q = \{q_0, q_1, \dots, q_m\} \quad q_0 = \text{start} \quad A = \{m\}$$

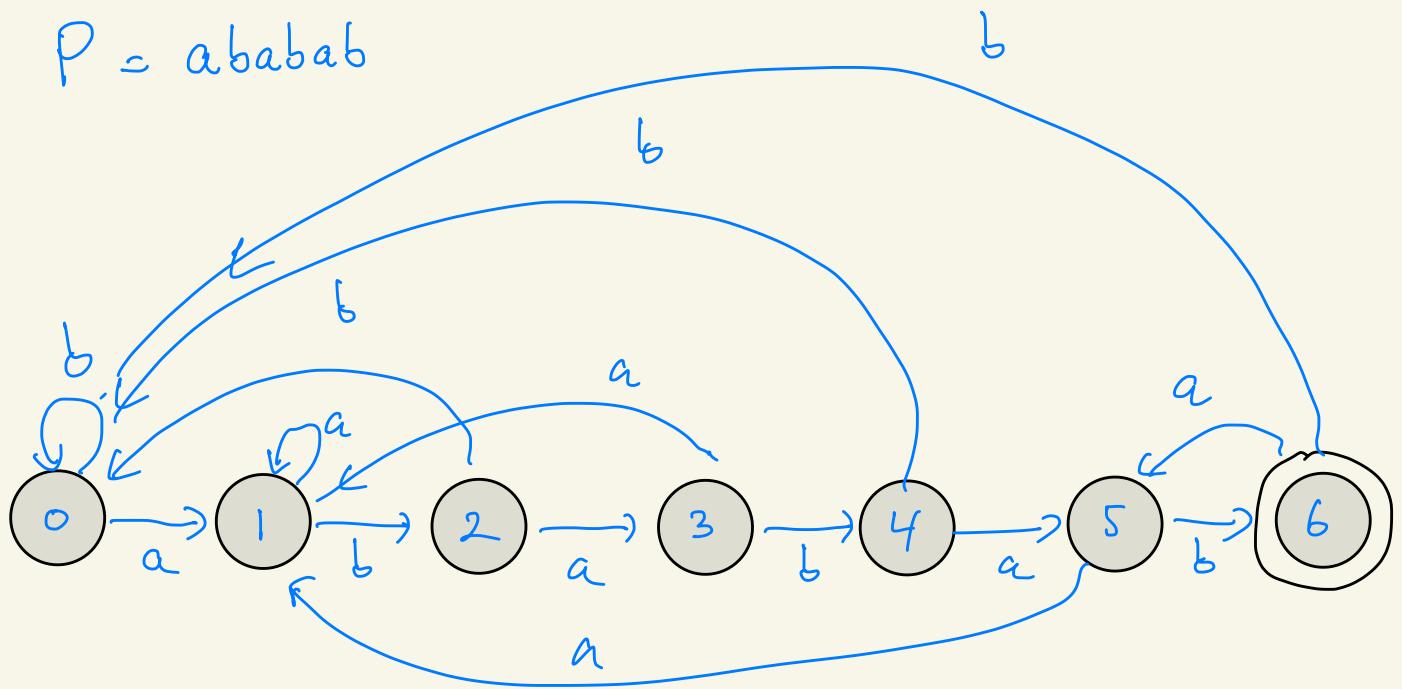
$$\delta(q, a) = \sigma(P_q a)$$

if the automata is currently

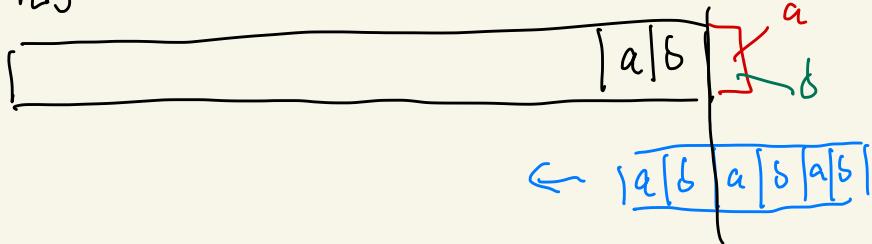
$\delta$  is chosen so that instead of the current status  $T_i$  satisfies that

$$P_q \supseteq T_i \text{ and } q = \sigma(T_i)$$

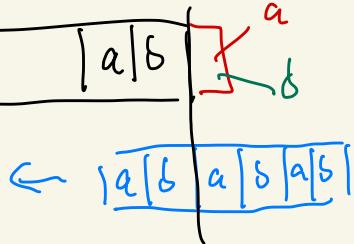
$P = ababab$



$T(I)$



$T(\omega)$



Let  $q = \phi(T_i)$  be the state we are in after reading  $T_i$

assume  $P_q$  is the longest prefix of  $P$  which is a suffix of  $T_i$

Let  $a = T[i+1]$

We want  $M$  to make a transition to the state  $q' = \sigma(T_i a)$ .

### Lemma 32.B

$\forall x \in \Sigma^*, a \in \Sigma : q = \sigma(x) \Rightarrow \sigma(xa) = \sigma(P_q a)$

P: By definition of  $\sigma$   $P_q \sqsupseteq x \Rightarrow P_q a \sqsupseteq xa$

Let  $r = \sigma(xa)$  so  $P_r \sqsupseteq xa$

Lemma 32.2  $\Rightarrow r \leq q+1$  as  $q = \sigma(x)$

$$\text{so } |P_r| = r \leq q+1 = |P_q a|$$

and thus  $P_r \sqsupseteq P_q a$  so  $r \leq \sigma(P_q a)$

but  $\sigma(P_q a) \leq \sigma(xa) \underset{\text{def}}{=} P_q a \sqsupseteq xa$

$$\text{so } \sigma(xa) = \sigma(P_q a)$$

### Theorem 32.4

If  $\phi$  is the final state function of  $M = M(P)$  and  $T[1..n]$  is the input to  $M$ , then

$$\phi(T_i) = \sigma(T_i) \text{ for } i=0, 1, 2, \dots, n$$

P: induction on  $i$

$$i=0 \quad T_0 = \varepsilon \quad \text{and} \quad \phi(T_0) = 0 = \sigma(T_0)$$

Induktion stop assume  $\phi(T_i) = \sigma(T_i)$

$$\text{let } q = \phi(T_{i+1}) \quad a = T[i+1]$$

$$\begin{aligned} \phi(T_{i+1}) &= \bar{\phi}(T_i, a) \text{ as } a = T[i+1] \\ &= \sigma(\phi(T_i), a) \text{ by def of } \phi \\ &= \sigma(q, a) \quad \text{as } q = \phi(T_i) \\ &= \sigma(p_q a) \text{ by def of } \sigma \\ &= \sigma(T_i a) \quad \text{by Lemma 32.3 + induction} \\ &= \sigma(T_{i+1}) \end{aligned}$$

□

Invariant  $\Rightarrow$  correct function of  $M$ :

when  $M$  enters state  $q$  after reading  $T_i$

$q$  is the largest value s.t.  $p_q \sqsupseteq T_i$  thus

$q = m \Leftrightarrow P \sqsupseteq T_i$  so match

## Computing the function $\delta$

Naive way: when in state  $q$

let  $k = \min\{m+1, q+2\} - 1$

try all values  $k, k-1, \dots, 0$

until  $P_k \sqsupseteq P_q a$

$O(m^2)$  work per state and  $a \in \Sigma$

total  $O(m^3 |\Sigma|)$

This can be improved to  $O(m |\Sigma|)$

so we can report all occurrences of  $P$

in  $T$  in  $O(n)$  matching time

and  $O(m |\Sigma|)$  preprocessing time.

More clever algorithm: KMP

not pseudocode