

# CS7800: Advanced Algorithms

## Lecture 18: Randomized III

- String Matching
- Start Hashing

Jonathan Ullman

11-15 - 2022

# String Matching

Input: text  $T[1:n] \in \Sigma^*$

pattern  $P[1:m] \in \Sigma^*$

e.g.  $\{A, B, C, \dots, z\}$   
or  $\{0, 1\}$

Output: either  $s$  such that  $T[s:s+m-1] = P[1:m]$

or "none" if there is no match

		1	2	3	4			
T =	1	1	0	0	1	0	0	1
P =			0	1	0			

# String Matching Algorithm I

T = 11001001  
P = 010

$O(n)$  iterations

for  $s = 1, 2, \dots, n - m + 1$  :

if (  $T[s : s+m-1] = P[1:m]$  ) :  
    returns  $s$

← equality check  
takes  $O(m)$  time

return "none"

total time is  $O(nm)$

and can be  $\Omega(nm)$

T = AAAA...A  
P = AAA...AB

# String Matching Algorithm II

compute  $p$  using  $O(m)$  time  
compute  $s_i$  using

for  $s = 1, 2, \dots, n - m + 1$ :

if ( $t_s == p$ ):  
return  $s$

return "none"

$T = 11001001$

$P = 010$

$t_s = T[s:s+m-1]$  as a number

$p = P[1:m]$  as a number

Strings to Numbers

1	1	0	1
$P[1]$	$P[2]$	$P[3]$	$P[4]$

$$p = 13 = 2^3 \cdot P[1] + 2^2 \cdot P[2] + 2^1 \cdot P[3] + 2^0 \cdot P[4]$$

$$p = \left( (P[1] \times 2 + P[2]) \times 2 + P[3] \right) \times 2 + P[4]$$

# String Matching Algorithm II

compute  $p$

compute  $t_1$

//  $O(m)$  time

$T = 11001001$

$P = 010$

$t_s = T[s:s+m-1]$  as a number

$p = P[1:m]$  as a number

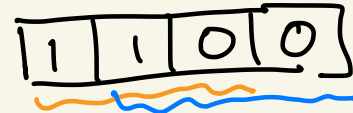
for  $s = 1, 2, \dots, n-m+1$ :

if ( $t_s == p$ ):  $O(m)$  time per iteration  
    L returns  $s$   
 $t_{s+1} = 2 \cdot (t_s - 2^{m-1} \cdot T[s]) + T[s+m]$

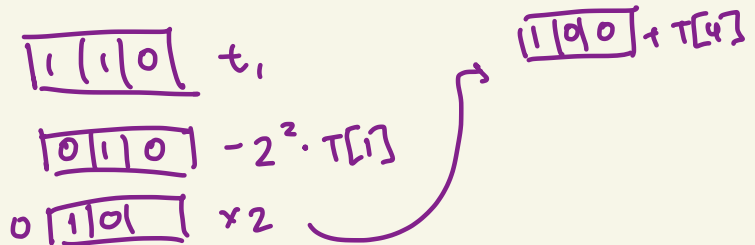
return "none"

No improvement yet

Sliding Window



$t_1 = 6$      $t_2 = 4$



# Modular Arithmetic

- $x \bmod y$  : "the remainder when you divide  $x$  by  $y$ "
  - $4 \bmod 3 = 1$
  - $7 \bmod 2 = 1$
  - $13 \bmod 5 = 3$
  - $x \bmod y = b$  means  
 $x = ay + b$  for integer  $a, b$
- If  $x = x'$  then  $(x \bmod y) = (x' \bmod y)$  for every  $y$

# String Matching Algorithm III

$T = 11001001$

$P = 010$

choose  $z$  ???

compute  $p \bmod z$   
compute  $t_1 \bmod z$   
 $\sigma = \text{compute } 2^{m-1} \bmod z$  } time  $O(m)$

for  $s = 1, 2, \dots, n-m+1$ :

if  $(t_s == p \bmod z)$  } time  $O(\log z)$

↳ check if  $t_s == p$  and if so, return  $s$  } time  $O(m)$  if we have to do it

$t_{s+1} = 2 \cdot (t_s - \sigma \cdot T[s]) + T[s+m] \bmod z$  } time  $O(\log z)$

return "none"

$O(n \log z + m \cdot (\# \text{ of matches}))$

# Random Prime Numbers

① (Prime Number Theorem) The number of primes  $\leq u$  is  $\Theta\left(\frac{u}{\log u}\right)$

Difficult to prove

② Every integer  $u$  has at most  $\log_2 u$  distinct prime factors

$$u = p_1^{a_1} \cdot p_2^{a_2} \cdot \dots \cdot p_k^{a_k} \geq 2^k$$

$$\Rightarrow k \leq \log_2 u$$

③ There is an efficient randomized algorithm to test primality

Difficult to prove



# String Matching Algorithm III

T = 11001001  
P = 010

$$u \approx m^2 \log m$$

choose  $z$  to be a random prime in  $\{2, 3, \dots, u\}$

compute  $p \bmod z \leftarrow z$  has  $O(\log m)$  digits

compute  $t_1 \bmod z$

$\sigma =$  compute  $2^{m-1} \bmod z$

for  $s = 1, 2, \dots, n - m + 1$ :

if  $(t_s == p \bmod z)$

└ check if  $t_s == p$  and if so, return  $s$

$t_{s+1} = 2 \cdot (t_s - \sigma \cdot T[s]) + T[s+m] \bmod z$

return "none"

Time:  
 $O(n \log m + m \cdot \frac{n}{m})$

think about running time of this step

•  $P(t_s == p \bmod z)$

when  $z$  is a random prime  
 $m \in \{2, 3, \dots, u\}$

•  $t_s - p \geq 0$  is an  $m$  bit number

•  $(t_s - p)$  has  $\leq m$  prime factors

•  $P(z \text{ is a factor of } (t_s - p))$

$\leq \frac{m}{\left(\frac{u}{\log u}\right)} \leq \frac{1}{m}$

• if  $u = m^2 \log m$

$\frac{n}{m}$  in expectation

# Hash Tables

Goal: Store a set of  $n$  elements  $S \subseteq U$  so that we can look up whether  $x \in U$  is in  $S$

universe



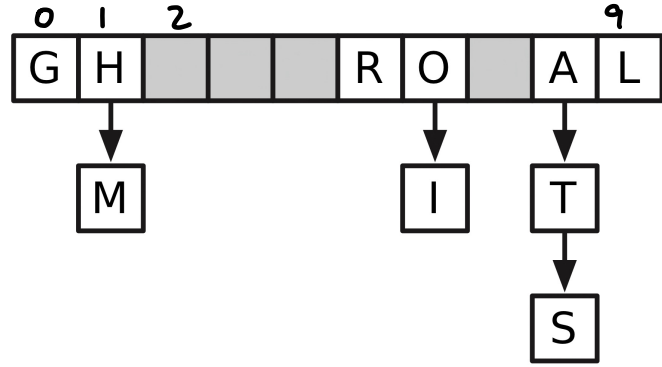
- A dictionary also lets us associate a value with key,  $x$

- A hash table  $T[1:m]$  stores the elements

- A hash function  $h: U \rightarrow \{0, 1, \dots, m-1\}$  maps elements to slots  $x \rightarrow T[h(x)]$

# Linear Chaining

A method for dealing with hash collisions



load factor  
 $n/m$

$$\begin{cases} m = 10 \\ n = 10 \end{cases}$$

$$U = \{A, B, C, \dots, Z\}$$

$$h(G) = 0$$

$$h(M) = 1$$

$$h(A) = h(T) = 8$$

- Let  $l(x)$  be the number of elements  $y \in S$  such that  $h(x) = h(y)$
- Time to lookup  $x$  is  $O(l(x))$

# Randomized Hash Functions

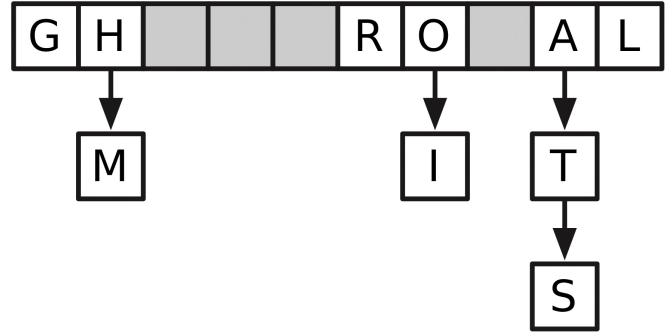
- A hash family  $\mathcal{H} = \{ h: \mathcal{U} \rightarrow \{0, 1, \dots, m-1\} \}$
- Choose a hash function  $h$  uniformly at random from  $\mathcal{H}$ .

• If  $|\mathcal{H}| = 1$  ( $h$  is deterministic) then there is always a set of size  $|\mathcal{U}|/m$  that all hash to the same bucket

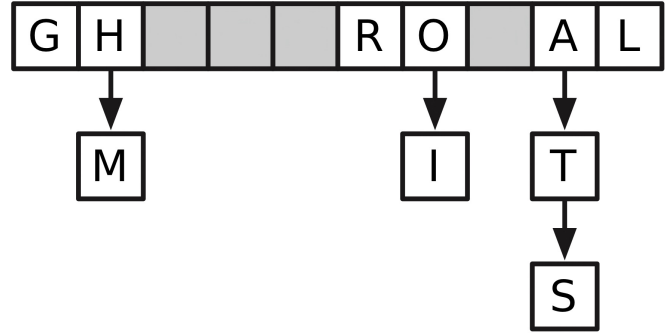
• A uniformly random function is a "good hash family"

• These simple hash families that "good enough"

Linear Chaining  
with Ideal Hashing



# Linear Chaining with Universal Hashing



# Constructing Universal Hash Families

Fix some prime  $p > |U|$  and table size  $m$

$$h_{a,b}(x) = (ax + b \bmod p) \bmod m$$

$$\mathcal{H}_{p,m} = \{ h_{a,b} \text{ for } a \in \mathbb{Z}_p^+, b \in \mathbb{Z}_p \}$$

Thm:  $\mathcal{H}_{p,m}$  is a universal hash family

# Constructing Universal Hash Families

Thm:  $H_{p,m}$  is a universal hash family

$$h_{a,b}(x) = (ax + b \bmod p) \bmod m$$

$$H_{p,m} = \{ h_{a,b} \text{ for } a \in \mathbb{Z}_p^+, b \in \mathbb{Z}_p \}$$

Lemma 1: If  $p$  is prime and  $a \neq 0$  then there is a unique value  $a^{-1} \in \{1, \dots, p-1\}$  such that  $a \cdot a^{-1} = 1 \bmod p$

(Division mod  $p$  is well defined)



# Constructing Universal Hash Families

Thm:  $H_{p,m}$  is a universal hash family

$$h_{a,b}(x) = (ax + b \bmod p) \bmod m$$

$$H_{p,m} = \{ h_{a,b} \text{ for } a \in \mathbb{Z}_p^+, b \in \mathbb{Z}_p \}$$

Lemma 2: If  $x \neq y$  and  $r \neq s$  then there is a unique solution  $(a,b)$  to the system

$$ax + b = r \bmod p$$

$$ay + b = s \bmod q$$

# Constructing Universal Hash Families

Thm:  $H_{p,m}$  is a universal hash family

$$h_{a,b}(x) = (ax + b \bmod p) \bmod m$$

$$H_{p,m} = \{ h_{a,b} \text{ for } a \in \mathbb{Z}_p^+, b \in \mathbb{Z}_p \}$$

Proof: