

# CS7800: Advanced Algorithms

## Lecture 19: Randomized IV

- Universal hash functions
- Perfect hashing

Jonathan Ullman

11-18 - 2022

# Hash Tables

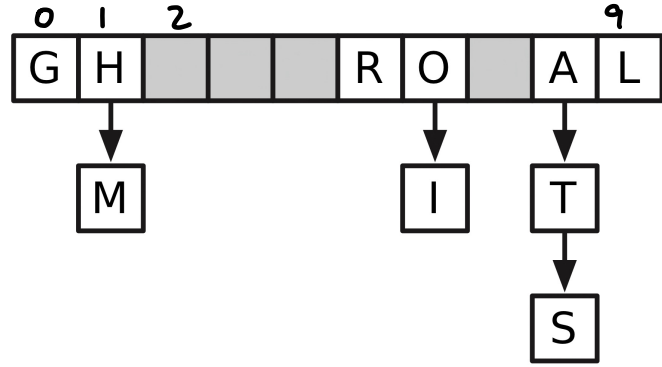
Goal: Store a set of  $n$  elements  $S \subseteq U$  so that we can look up whether  $x \in U$  is in  $S$

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

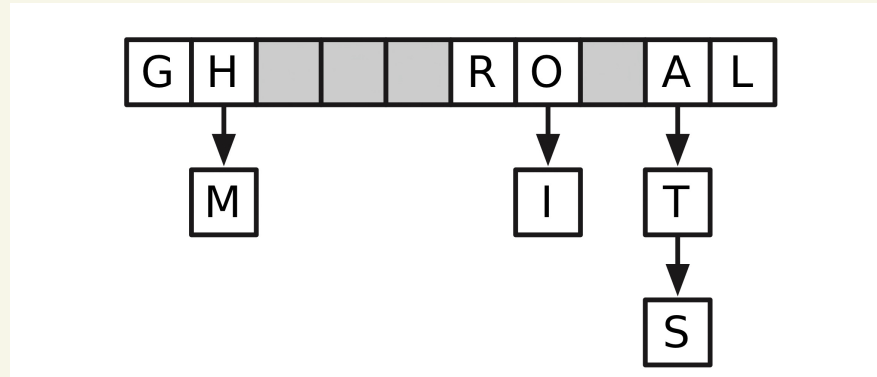
• A uniformly random function is a "good hash family"

• These simple hash families that "good enough"

# Linear Chaining with Ideal Hashing

Property of an ideal hash fn:

For any set of <sup>distinct</sup> elts  $x_1, \dots, x_k$   
the values  $h(x_1), \dots, h(x_k)$  are independent



$$l(x) = \# \text{ of elts } y \in S \text{ s.t. } h(x) = h(y)$$

Expected Lookup Time: (fixed  $S$ , fixed  $x$ , random  $h$ )

$$\mathbb{E}(l(x)) = \mathbb{E}\left(\sum_{y \in S} C_{x,y}\right) = \sum_{y \in S} \mathbb{E}(C_{x,y}) = \sum_{y \in S} \mathbb{P}(h(x) = h(y))$$

$$C_{x,y} = \begin{cases} 1 & \text{if } h(x) = h(y) \\ 0 & \text{otherwise} \end{cases}$$

Assume  $x \neq y$

$$= \sum_{y \in S} \frac{1}{m} = \frac{|S|}{m} = \frac{n}{m}$$

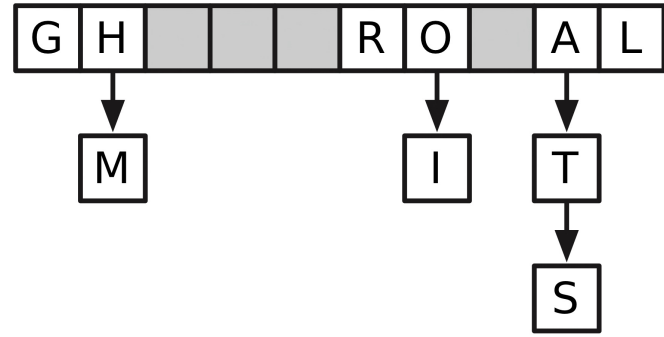
load factor

# Linear Chaining with Ideal Hashing

Property of an ideal hash fn:

For any set of <sup>distinct</sup> elts  $x_1, \dots, x_k$

the values  $h(x_1), \dots, h(x_k)$  are independent



$$l(x) = \# \text{ of elts } y \in S \text{ s.t. } h(x) = h(y)$$

Worst-Case Lookup Time: ( $S$  is fixed,  $h$  is random)

$$\mathbb{E}(\max_{x \in U} l(x)) = \mathbb{E}(\max \text{ items in any bucket}) \leftarrow \text{balls and bins}$$

$$\text{(if } \frac{m}{n} = 1 \text{) then } = \Theta\left(\frac{\log n}{\log \log n}\right)$$

# Pseudorandom Hash Families

A hash family  $\mathcal{H} = \{ h: \mathcal{U} \rightarrow \{0, 1, \dots, m-1\} \}$

- A hash family  $\mathcal{H}$  is 2-wise uniform if for every  $x \neq y \in \mathcal{U}$   
$$\mathbb{P}(h(x) = i \text{ and } h(y) = j) = \frac{1}{m^2}$$

- A hash family  $\mathcal{H}$  is 2-wise universal if for every  $x \neq y \in \mathcal{U}$

$$\mathbb{P}(h(x) = h(y)) \leq \frac{1}{m}$$

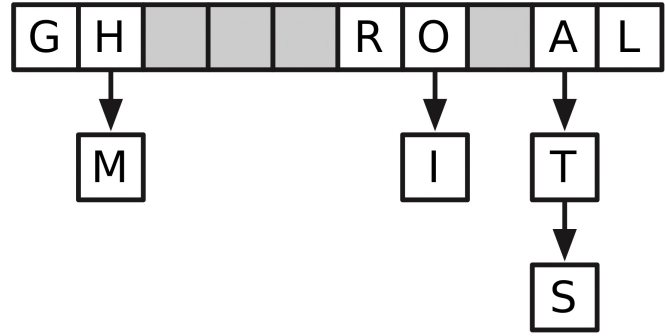
$c_1, c_2$  are com's

$$h_{c_1, c_2}(x) = \begin{cases} c_1 & x=1 \\ c_2 & x=2 \\ c_1 \oplus c_2 & x=3 \end{cases}$$

# Linear Chaining with Universal Hashing

If  $H$  is 2-universal then

$$\mathbb{E}(l(x)) = \frac{m}{n}$$





# 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.)

Pf: Suppose that  $az = az' \bmod p \Rightarrow a(z-z') = 0 \bmod p$   
for  $z, z' \in \mathbb{Z}_p^+$

$\Rightarrow z-z'$  is divisible by  $p$

$\Downarrow$

$\Rightarrow z-z' = 0$

$z-p \leq z-z' \leq p-2$

must be  $z$  s.t.  
 $a \cdot z = 1 \bmod p$

$az = h \bmod p$  has  
at most one solution

$az = 0 \bmod p$  has  
no solutions

# 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

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

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

# Constructing Universal Hash Families

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

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

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

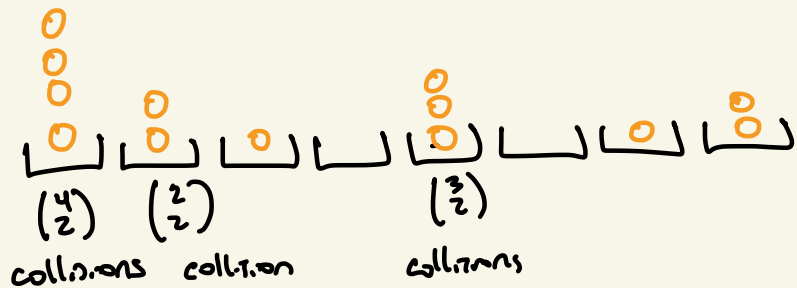
Proof:

By Lemma 2  $\mathbb{P}(ax + b = r \bmod p \text{ and } ay + b = s \bmod p) = \frac{1}{p(p-1)}$

$$\mathbb{P}(h_{a,b}(x) = h_{a,b}(y)) = \frac{N}{p(p-1)} \quad \text{where } N \text{ is number of } r \neq s \in \mathbb{Z}_p \text{ such that } r = s \bmod m$$
$$\leq \frac{p(p-1)}{p(p-1)} \cdot \frac{1}{m}$$

$$N \leq \underbrace{p}_{\text{choices of } r} \cdot \underbrace{\frac{p-1}{m}}_{\text{choices of } s \text{ for given } r}$$

# Maximum Load for Universal Hashing



if max load  $\geq k \approx \sqrt{c}$  then # of collisions  $\geq \binom{k}{2} = c$

$$\Rightarrow \text{max load} \approx \sqrt{\# \text{ collisions}}$$

$$\mathbb{E}(\text{max load}) \approx \sqrt{\mathbb{E}(\# \text{ collisions})}$$

$$\approx \sqrt{n + \frac{n^2}{m}}$$

$$\mathbb{E}(\# \text{ collisions}) = \sum_{x, y} \mathbb{P}(h(x) = h(y)) = \sum_x 1 + \sum_{x \neq y} \frac{1}{m} = n + \frac{n(n-1)}{m} \approx n + \frac{n^2}{m}$$