

6.046  
Lecture 8

Universal Hashing, Perfect Hashing

(1)

Weakness of hashing

For any choice of hash function

$\exists$  a bad set of keys that all hash to same slot.  $\hookrightarrow$

Idea: choose hash function at random,  
independently from keys

potential  
vulnerability  
against an  
adversary

Universal hashing

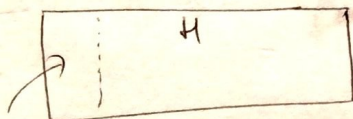
Def. Let  $U$  be a universe of keys, and  
let  $H$  be a finite collection of hash functions  
mapping  $U$  to  $\{0, 1, \dots, m-1\}$

$H$  is universal if  $\forall x, y \in U$ , where  $x \neq y$ ,

$$|\{h \in H : h(x) = h(y)\}| = \frac{|H|}{m}$$

i.e. if  $h$  is chosen randomly from  $H$ ,

the prob. of collision between  $x$  and  $y$  is  $1/m$



$\{h : h(x) = h(y)\} \leftarrow$  this subset of  $H$  would  
be different for each  
 $x, y$  distinct pair  
 $\frac{1}{m} |H|$

Thm: Choose  $h$  randomly from  $H$   
Suppose hashing  $n$  keys into  $m$  slots  
in table  $T$ . Then, for a given key  $x$ ,

$$E[\# \text{ collisions with } x] < \frac{n}{m}$$

Pf. Let  $C_x$  be r.v. denoting total  
 $\#$  collisions of keys in  $T$  with  $x$ ,

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

$$E[c_{xy}] = 1/m$$

$$C_x = \sum_{y \in T - \{x\}} c_{xy}$$

$$E[C_x] = E\left[\sum_{y \in T - \{x\}} c_{xy}\right]$$

$$= \sum_{y \in T - \{x\}} E[c_{xy}] = \sum_{y \in T - \{x\}} \frac{1}{m}$$

$$= \frac{n-1}{m} < \frac{n}{m} \leftarrow \text{exp. guarantee against adversary}$$

Constructing a universal hash function

Let  $m$  be prime. Decompose key  $k$  into  $r+1$  digits:

$$k = \langle k_0, k_1, \dots, k_r \rangle, \text{ where } 0 \leq k_i \leq m-1$$

base  $m, m^2, \dots, m^{r+1}$  then use mod

$\hookrightarrow$  Representation of  $k$   
base  $m$



- pick  $a = \langle a_0, a_1, \dots, a_r \rangle$ , each  $a_i$  is chosen randomly from  $\{0, 1, \dots, m-1\}$

- Define  $h_a(k) = \left( \sum_{i=0}^r a_i k_i \right) \bmod m$

$\nearrow$  dot product  $a$  and  $k$ , then take mod  $m$

How big is  $H$ ?

$$|H| = m^{r+1} \leftarrow \# \text{ of all } a$$

Thm:  $H$  is universal

$\nwarrow$  base  $m$  representations

A.S. let  $x = \langle x_0, x_1, \dots, x_r \rangle$

$y = \langle y_0, y_1, \dots, y_r \rangle$  be distinct keys

$\Rightarrow$  they differ in at least one digit,

wlog position 0.

For how many  $h_a \in H$  do  $x$  and  $y$  collide?

Must have  $h_a(x) = h_a(y)$

$$\Rightarrow \sum_{i=0}^r a_i x_i \equiv \sum_{i=0}^r a_i y_i \pmod{m}$$

$\uparrow$   
congruent

$$\Rightarrow \sum_{i=0}^r a_i (x_i - y_i) \equiv 0 \pmod{m}$$

$$\Rightarrow a_0(x_0 - y_0) + \sum_{i=1}^r a_i(x_i - y_i) \equiv 0 \pmod{m}$$

$$\Rightarrow a_0(x_0 - y_0) \equiv - \sum_{i=1}^r a_i(x_i - y_i) \pmod{m}$$

Number theory fact:

Let  $m$  be prime. For any  $z \in \mathbb{Z}_m$  (integers mod  $m$ )

s.t.  $z \neq 0$ ,  $\exists$  unique  $z^{-1} \in \mathbb{Z}_m$  s.t.  $z \cdot z^{-1} \equiv 1 \pmod{m}$ .

Ex  $m = 7$

$z$	1	2	3	4	5	6
$z^{-1}$	1	4	5	2	3	6

$\nearrow$  not true if  $m$  is not prime  
since any  $z \in \mathbb{Z}_m$  is relatively prime to  $m$

$a \pmod{b}$

$\nearrow$  if not relatively prime  
 $a$  does not have an inverse mod  $b$ .

$$z \equiv -5 \pmod{7}$$



6.046

Lecture 8

(2)

Since  $x_0 \neq y_0$ ,  $\exists (x_0 - y_0)^{-1}$

$$\Rightarrow a_0 \equiv \left(1 - \sum_{i=0}^r a_i (x_i - y_i)\right) \cdot (x_0 - y_0)^{-1} \pmod{m}$$

- if  $x, y$  hash to the same place (assumed)
- then  $a_0$  has a particular value as a function of other  $a_i$

- Thus, for any choice of  $a_1, a_2, \dots, a_r$
- exactly 1 of  $m$  choices of  $a_0$  that causes  $x$  and  $y$  to collide, and no collision for other  $m-1$  choices for  $a_0$ .

$$\Rightarrow \# h_a \text{'s that cause } x, y \text{ to collide} = m \cdot m \cdot \dots \cdot m \cdot 1$$

$$= m^r = \frac{141}{m} \checkmark$$

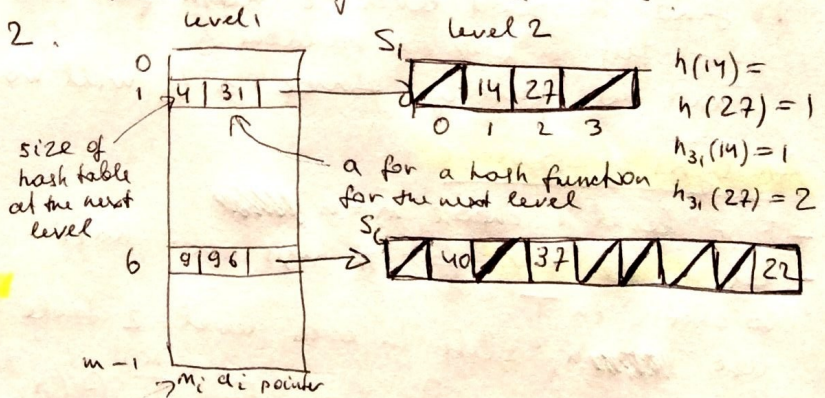
$a_1 \quad a_2 \quad \dots \quad a_r \quad a_0$

Perfect Hashing (given a fixed set of keys, build a <sup>static</sup> table with good worst case search time)  $\Rightarrow$  not in expectation

Given  $n$  keys, construct a static hash table of size  $m = O(n)$ , s.t. search takes  $O(1)$  time in the worst case.

Two-level scheme with universal hashing at both levels.  
No collisions at level 2.

If  $n_i$  items hash to level 1 slot  $i$ , then use  $m_i = n_i^2$  slots in level 2 table  $S_i$ .



$$h(14) = 1$$

$$h(27) = 1$$

$$h_{31}(14) = 1$$

$$h_{31}(27) = 2$$

use another univers. hashing w/ primes non-primes

not prime here, but could pick primes close to their values

universal hashing, pick  $h \in H$  at random

universal hashing with an  $h \in H$  picked for each slot of level 1, at random



### Level 2 analysis:

Thm: Hash  $n$  keys into  $m = n^2$  slots, using random  $h$  in universal  $H \Rightarrow E[\# \text{ collisions}] < \frac{1}{2}$

Pf: Prob. 2 given keys collide under  $h$  is  $\frac{1}{m} = \frac{1}{n^2}$   
 $\binom{n}{2}$  pairs of keys

$$E[\# \text{ collisions}] = \binom{n}{2} \frac{1}{n^2} = \frac{n(n-1)}{2} \frac{1}{n^2} = \frac{n^2}{2n^2} - \frac{n}{2n^2} =$$

Markov inequality:

For r.v.  $X \geq 0$ ,  $Pr\{X \geq t\} \leq \frac{E(X)}{t}$

$$\text{Pf. } E[X] = \sum_{x=0}^{\infty} x \cdot Pr\{X=x\} \geq \sum_{x=t}^{\infty} x \cdot Pr\{X=x\}$$

Corollary

$$Pr\{\text{no collision}\} \geq \frac{1}{2}$$

$$\geq \sum_{x=t}^{\infty} t \cdot Pr\{X=x\} = t \cdot Pr\{X \geq t\}$$

$\nearrow$  throw away lower terms

$$\text{Pf: } Pr\{\geq 1 \text{ collision}\} \leq \frac{E[\# \text{ collisions}]}{1} < \frac{1}{2} \checkmark$$

To find a good level-2 hash function, just test a few at random.

Find one quickly, since  $\geq \frac{1}{2}$  will work.

} randomized construction

### Analysis of storage

- For level 1, choose  $m = n$

- let  $n_i$  be r.v. for # keys that hash to slot  $i$  in  $T$ .

Use  $n_i^2$  slots in each level-2 table  $S_i$ .

$$E[\text{total storage}] = n + E\left[\sum_{i=0}^{n-1} \Theta(n_i^2)\right]$$

$$= \Theta(n) \text{ by bucket sort analysis}$$

check a fixed # of hash functions for each slot, s.t. the prob. to find one without collision is very high