

# Hashing

assume  $S \subseteq \{0, 1, \dots, U-1\}$

membership queries: is  $q \in S$ ?

insert  
delete

trivial (bit vector):

$O(1)$  time  
but space is  $O(U)$



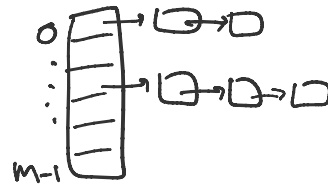
## Hashing Method 1:

pick a hash fn  $h: \{0, \dots, U-1\} \rightarrow \{0, \dots, m-1\}$   
for some  $m \ll U$ .

store array  $A[0, \dots, m-1]$

where

$A[i] = \text{list of all } x \in S \text{ with } h(x) = i$



Search(y):

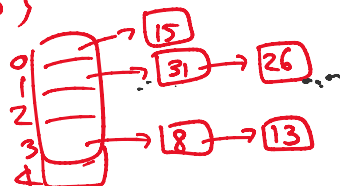
find y in  $A[h(y)]$  by linear search

insert  $O(1)$

delete  $O(1)$

e.g.  $S = \{31, 15, 8, 13, 26\}$

$h(x) = x \bmod 5$



if input is rand, unif. distrib,  
each bucket has  $\sim \frac{n}{m}$  elems on average

Set  $m \approx n \Rightarrow O(n)$  space  
 $O(1)$  "average" query time

Set  $m \approx n \Rightarrow O(n)$  space  
 $O(1)$  "average" query time

But can't assume input is random!

idea - pick a random hash fn  
from a family

Def Fix prime  $p \in [U, 2U]$ .

Pick rand  $a \in \{1, \dots, p-1\}$ ,  $b \in \{0, \dots, p-1\}$ .

Define  $h_{a,b}: \{0, \dots, U-1\} \rightarrow \{0, \dots, m-1\}$ :

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

evaluated in  $O(1)$  time

(similar to  
2-point sampling)

Prop For any fixed  $x, y \in \{0, \dots, U-1\}$ , ( $x \neq y$ ),

$$\Pr_{a,b} \left[ \underline{h_{a,b}(x) = h_{a,b}(y)} \right] \leq \underline{O\left(\frac{1}{m}\right)}.$$

called universal (Carter-Wegman '79)

More strongly, for fixed  $i, j \in \{0, \dots, m-1\}$ ,

$$\Pr_{a,b} \left[ \underline{h_{a,b}(x) = i \wedge h_{a,b}(y) = j} \right] \leq O\left(\frac{1}{m^2}\right).$$

called strong 2-universal

Pf: Fix  $i', j' \in \mathbb{Z}_p$  ( $i' \neq j'$ ).

$$\Pr_{a,b} \left[ \underline{ax+b \equiv i' \pmod{p} \wedge ay+b \equiv j' \pmod{p}} \right] = \frac{1}{p(p-1)}.$$

2x2 linear system  
in vars  $a, b$   
has unique sol'n

$$a \equiv \frac{i' - j'}{x - y}, \quad b \equiv i' - ax$$

$$\Rightarrow \Pr_{a,b} [h_{a,b}(x) = i \wedge h_{a,b}(y) = j]$$

$$\leq \left( \begin{array}{l} \# \text{ choices} \\ \text{of } i', j' \in \mathbb{Z}_p (i' \neq j') \\ \text{with } \begin{array}{l} i' \equiv i \pmod{m} \\ j' \equiv j \pmod{m} \end{array} \end{array} \right) \cdot \frac{1}{p(p-1)}$$

$$\leq \left\lceil \frac{p}{m} \right\rceil \cdot \left\lceil \frac{p}{m} \right\rceil \cdot \frac{1}{p(p-1)}$$

$$= O\left(\frac{1}{m^2}\right). \quad \square$$

**Remark:** other ex of 2-universal hash families

**Dietzfelbinger et al. '97:**

$$h_a: \{0, \dots, 2^w - 1\} \rightarrow \{0, \dots, 2^l - 1\}$$

$$h_a(x) = \left\lfloor \frac{(a \cdot x) \bmod 2^w}{2^{w-l}} \right\rfloor.$$

**Patrascu-Thorup '10:**

tabulation hashing  
(using XORs & tables of rand values)

⋮

**Analysis of query time:**

for fixed query value  $y$ ,

$\mathbb{E} [\# \text{ elems of } S \text{ that } \underline{\text{collide}} \text{ with } y]$

$$E \left( \# \text{ elems of } S \text{ that } \underline{\text{collide}} \text{ with } y \right)$$

$x \in S$        $h_{a,b}(x) = h_{a,b}(y)$

$$= E_{a,b} \left( \sum_{x \in S - \{y\}} [h_{a,b}(x) = h_{a,b}(y)] \right)$$

$$[E] = \begin{cases} 1 & \text{if } E \text{ true} \\ 0 & \text{else} \end{cases} = \sum_{x \in S - \{y\}} \Pr(h_{a,b}(x) = h_{a,b}(y))$$

$$\leq O\left(n \cdot \frac{1}{m}\right)$$

by universality

Set  $m \approx n \Rightarrow$

- $O(1)$  expected query time
- $O(1)$  insert/delete
- $O(n)$  space

Assume oblivious adversary  
(query values indep. of rand. choices made by algm)

But worst-case query time ??

Assume static ...

Alternative analysis:

$$E \left( \text{total \# of colliding pairs } (x,y) \in S \times S \right)$$

$$= E \left( \sum_{x,y \in S} [h_{a,b}(x) = h_{a,b}(y)] \right)$$

$$\leq O\left(n^2 \cdot \frac{1}{m}\right) = O\left(\frac{1}{c}\right)$$

1.

2

... related to "birthday paradox"

$$\text{Set } m = cn^2$$

← (related to "birthday paradox")

$$\Rightarrow \Pr(\text{total \# of colliding pairs} \geq 1) \leq \frac{O(\frac{1}{c})}{1} \quad \text{by Markov's in eq.}$$

$$\Rightarrow \text{with prob. } \Omega(1).$$

no collision

↔ perfect hashing

repeat till success

$$\Rightarrow \boxed{O(1)} \text{ worst-case query time}$$

$$\boxed{O(n^2)} \text{ space}$$

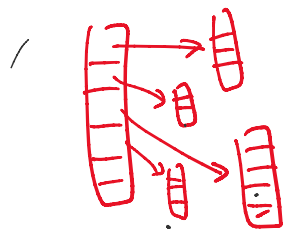
$$\boxed{O(n)} \text{ expected preproc time}$$

But can space be reduced back to  $O(n)$ ?

Final Hashing Method (Fredman, Komlos, Szemerédi '84)

idea - bootstrap

store each bucket  $A[i]$   
in the data structure with  $O(|A[i]|^2)$  space  
&  $O(1)$  worst-case query time



2-level hash table

expected space

$$O\left(m + E\left[\sum_i |A[i]|^2\right]\right)$$

$$= O(m + E[\text{total \# colliding pairs}])$$

$$= O\left(m + \frac{n^2}{m}\right)$$

$$\text{Set } m=n \Rightarrow O(n)$$

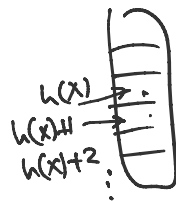
repeat till success

$$\Rightarrow \left. \begin{array}{l} \boxed{O(n)} \text{ worst-case space} \\ \boxed{O(1)} \text{ worst-case query time} \\ \boxed{O(n)} \text{ expected preproc time} \end{array} \right\}$$

**Rmk -** made dynamic by Dietzelbringer et al. '94  
with  $O(1)$  expected update time

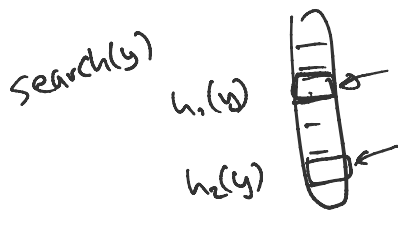
**Rmks -** alternatives with single table (open addressing)

**linear probing**



assume  $m \geq (1+\epsilon)n$

**cuckoo hashing**



use two hash fns  $h_1, h_2$