

6) Lecture 8-12 Dictionary Problem

12 August 2024 15:29

Motivating example:

○ Bank account scenario:

- Each account is identified by an account number
- Each account stores additional information such as
 - Balance,
 - Name & address of account holder
 - History of transactions

15 digit key - account no.

Val1
Val2
Val3

○ Any application that wish to operate on account need to provide the account no as key

data := (key, value)

- Search(key)

- Insert(key, value)

- delete(key)

The dictionary:

- Dictionaries is a data structure that store elements so that they can be located quickly using keys

- It provides and abstract model of database such that

○ Stores (key, element) pair

○ The main operations supported by dictionary are

- FindElement(key),
- insertItem(key, Element),
- deleteItem(Key),

Following data structure can be used to implement dictionary:

- Array & linked list (inefficient)
- Hash tables
- Binary tree
- Red black tree
- AVL tree
- B tree

Searching problem:

Input: Sequence of numbers (database)

A single number (query)

Output: Index of the number, if found or NIL

Standard Techniques:

Data = { (key, value)}

- [] Insert (key, value)
- [] Search (key)
- [] Delete (key)

```
struct node
{
    int key;
    int value;
}
```

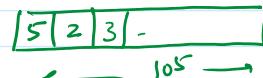
Certain ordering can be defined on key.

(i) Create an Unordered Sequence

5, 2, 3, 10 -

Use array - append new element at the end of array

 └ insertion $\Theta(1)$ — Advantage



 └ search $\Theta(n)$ (linear search - list is unordered)

 └ deletion $\Theta(n)$ (due to search + shift)

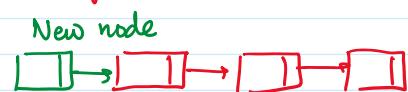
 └ Space complexity is high

disadvantage

-
-
-

Use linked list - append new element at beginning of linked list

 └ insertion $\Theta(1)$ Advantage



 └ search / deletion $\Theta(n)$ (linear search)

 └ Lesser

 └ Space complexity

Useful to applications where frequent insertion & rare search / deletion.

(ii) Create an ordered sequence

- create an array of struct - insertion sort

- search $\Theta(\log n)$ (Binary Search)

- insertion $\Theta(n)$

 └ Mainly due to shift

- deletion $\Theta(n)$

4



- Deletion $O(n)$ / Mainly due to shift



Useful in application where frequent search and rare insertion/deletion

Input $\{(key, value)\} = \text{Data}$

Search (Key)

Insert (Key, value)

Delete (Key)

Direct address table

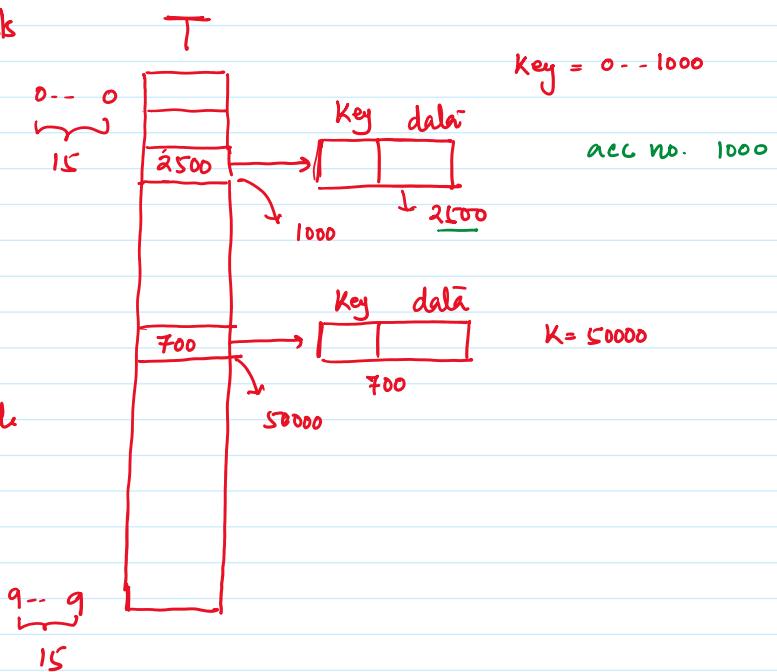
Key: Account No. 15 digits

0..0

:

9..9

One location for each possible acc. no.



- Search $O(1)$

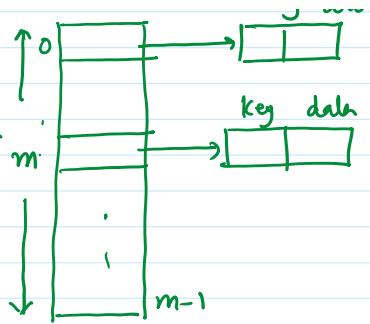
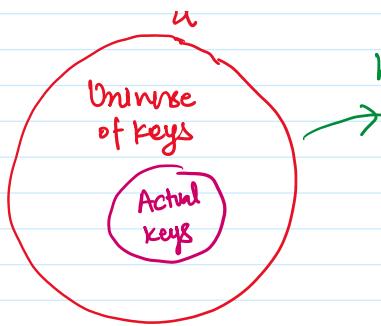
- Insertion $O(1)$

- Deletion $O(1)$

Disadvantage: Space complexity is high

u
Universe





Universe of Keys
:= set of all possible
16 digit no.

$|U| := \# \text{ of elements in universe}$

$N := \text{set of actual key}$

$$\approx \frac{n}{m}$$

$m := \text{size of array } T$

We used $m = |U| \rightarrow \text{higher space complexity}$ Hash fn.

we want $m \ll |U|$

$$[h: U \rightarrow \{0 \dots m-1\}]$$

$m = 100$

Naive approach:

$$h(x) = x \bmod m$$

$x := \text{acc. no.}$

$$h: U \rightarrow \{0 \dots m-1\} := \text{hash function.}$$

Key $\in \mathbb{Z}_{+}^+$

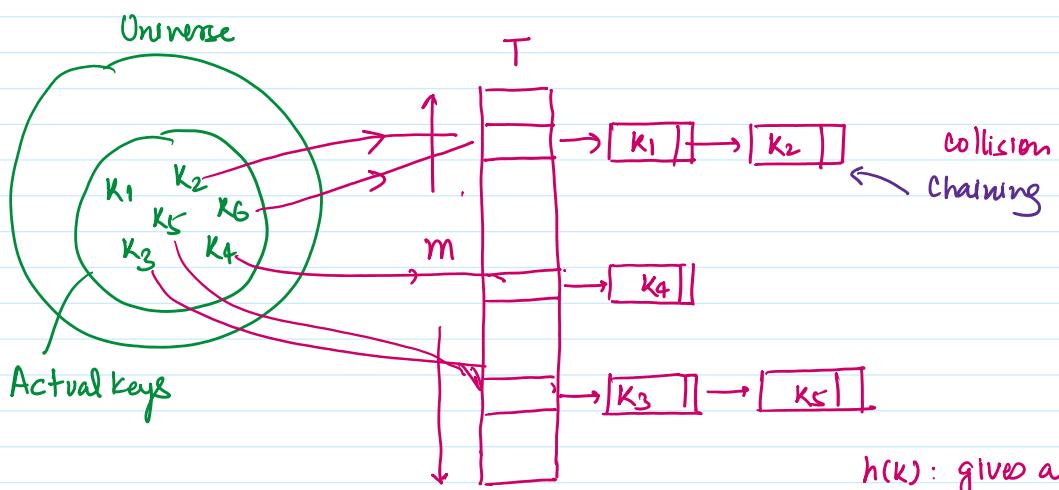
$m = 100$

$$h: \text{acc. no.} \rightarrow \{0 \dots 99\}$$

15 digit

acc. no. mod 100

$m \ll |U|$



$h(K) : \text{gives a slot}$
 $\{0 \dots m-1\} \text{ in } T$

Insert (Key, Val) $O(1)$

Search (Key) $> O(n)$

Delete (Key)

Hash Table

: n Keys

$\approx \underline{m}$ in each slot

Delete (key) /

: n Keys $\approx \frac{n}{m}$ in each slot
: m slots in T

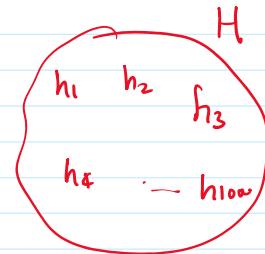
if $m = \sqrt{n}$

Universal Hashing

Design a family of hash function H

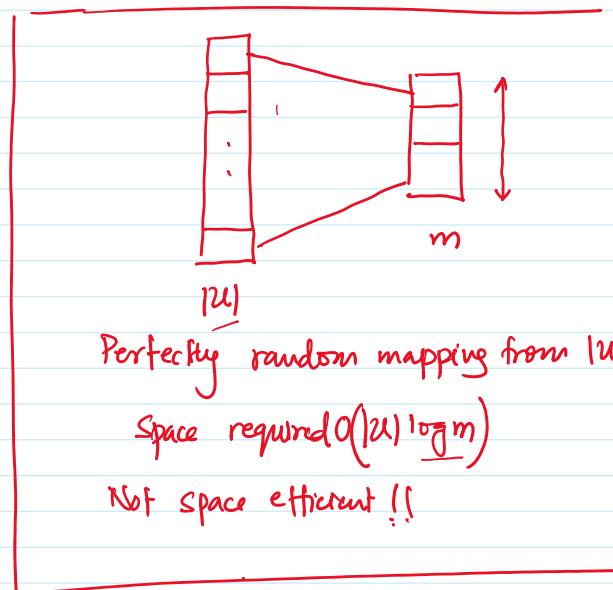
for all pair of keys $k \neq k' \in \mathcal{U}$,

$$\text{if } \Pr_{h \in H} [h(k) = h(k')] \leq \frac{1}{m}. - (1)$$



Then H is called as Universal Hash family.

$h : \mathcal{U} \rightarrow \text{Universal hash fn.}$



Th: Let n arbitrary distinct key $\in \mathcal{U}$, Let H be a Universal Hash family

$$E [\# \text{keys colliding on particular slot } (0 \dots m-1)] \leq \frac{n}{m} + 1.$$

Proof: Let $k_0 \dots k_{n-1}$ (n - distinct keys)

Consider some arbitrary key k_i

$$E [\# \text{of key colliding with } k_i]$$

Let $I_{ij} = \begin{cases} 1 & \text{if } \underbrace{h(k_i) = h(k_j)}_{h \in H} \\ 0 & \text{o.w.} \end{cases}$

$h \in H$

h is sampled u.a.r. from H

$$E[\# \text{ of key colliding with } k_i] = E\left[\sum_{i \neq j} I_{ij} + I_{ii}\right]$$

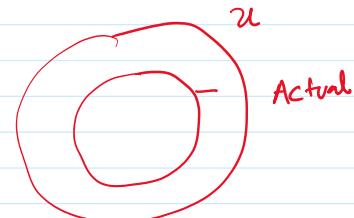
$$= E\left[\sum_{i \neq j} I_{ij}\right] + 1$$

$$= \sum_{i \neq j} E[I_{ij}] + 1$$

$$= \sum_{i \neq j} \Pr[h(k_i) = h(k_j)] + 1$$

$$\leq \sum_{i \neq j} \frac{1}{m} + 1 \quad [\text{From Eqn ①}]$$

$$= \frac{n}{m} + 1$$



Dot product Hash family:

- assume m is prime

- assume $|u| = m^r$ for $\text{int } r^{+ve}$

$$H_a = \{h_a | a \in \{0, \dots, m-1\}^r\}$$

$$h_a(k) = \left(\sum_{i=0}^{r-1} a_i k_i \right) \bmod m = \langle \vec{a} \cdot \vec{k} \rangle \bmod m$$

$$\vec{a} = \langle a_0, a_1, \dots, a_{r-1} \rangle$$

$$\vec{k} = \langle k_0, k_1, \dots, k_{r-1} \rangle$$

base m representation

Space required to store $h_a()$ = space required to store a = $O(\log|u|)$

$m=2$	$K=5$	101
$a_1=4$	$a_2=6$	

$$h_{a_1}(k) = \langle \underline{100}, \underline{101} \rangle \bmod 2$$

$$= 1$$

$$h_{a_2}(k) = \langle \underline{110}, \underline{101} \rangle \bmod 2$$

$$\frac{101}{1+0+0} = 1 \quad \textcircled{1}$$

Theorem Dot product hash family is universal.

Proof: Let $K \neq K' \in \mathcal{U}$

$$|\mathcal{U}| = m^r$$

* $\Pr_a [h_a(K) = h_a(K')] \leq \frac{1}{m}$ (Need to show) !!

$$K = \langle K_0, K_1, \dots, K_d, \dots, K_{r-1} \rangle$$

$$K' = \langle K'_0, K'_1, \dots, K'_d, \dots, K'_{r-1} \rangle$$

$$a = \langle a_0, \dots, \textcircled{a}_d, \dots, a_{r-1} \rangle$$

Since $K \neq K'$, their base m representation differ at say d^{th} digit

assume $K_d \neq K'_d$

$$\text{LHS } \Pr_a [h_a(K) = h_a(K')] = \Pr_a \left[\sum_{i=0}^{r-1} a_i K_i \bmod m = \sum_{i=0}^{r-1} a_i K'_i \bmod m \right] \text{ By def}$$

$$= \Pr_a \left[\underbrace{a_d (K_d - K'_d)} + \sum_{i \neq d} a_i (K_i - K'_i) \equiv 0 \pmod{m} \right]$$

$$= \Pr_a \left[\underbrace{a_d (K_d - K'_d)} = - \sum_{i \neq d} a_i (K_i - K'_i) \pmod{m} \right]$$

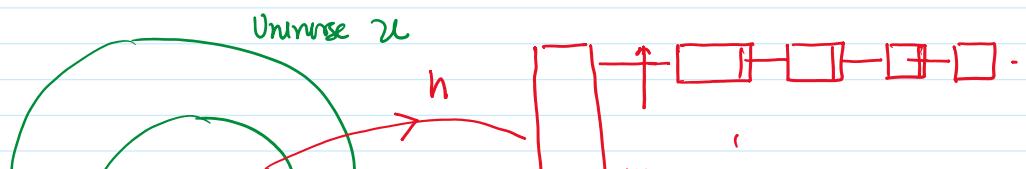
$$= \Pr_a \left[a_d \equiv - (K_d - K'_d)^{-1} \left(\sum_{i \neq d} a_i (K_i - K'_i) \right) \pmod{m} \right]$$

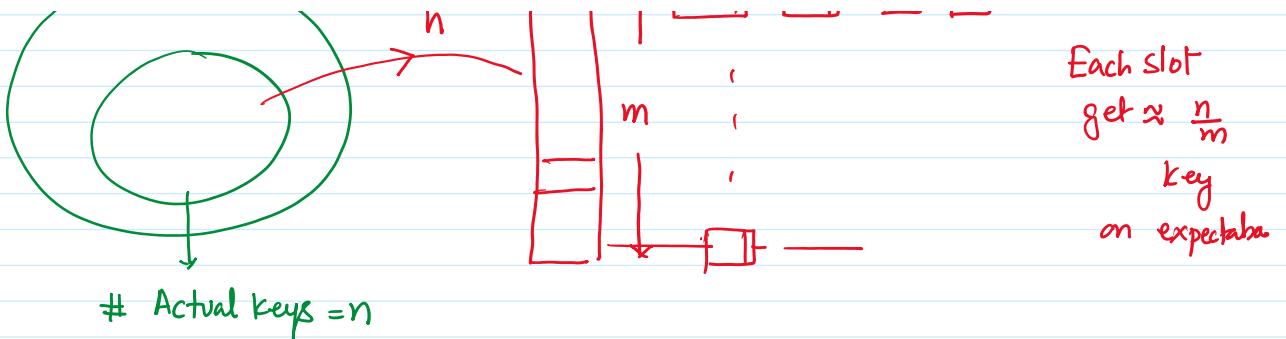
multiplicative inverse of $(K_d - K'_d) \pmod{m}$

$$f(a_0, \dots, a_{d-1}, \textcircled{a}_d, \dots, a_{r-1})$$

$$= \Pr_a [a_d = f(a_0, \dots, a_{d-1}, \textcircled{a}_d, \dots, a_{r-1})]$$

$$= \frac{\# \text{ favourable choices}}{\text{total no of choices}} \leq \frac{m^{r-1}}{m^r} \leq \frac{1}{m}$$





Search (Key) — on expectation n/m

Insert (Key, val) — $O(1)$

delete (Key) — on expectation n/m (due to search)

Static dictionary Problem

Given n keys upfront — prepare a data structure that facilitates fast search operation

(No further insertion/deletion allowed.)

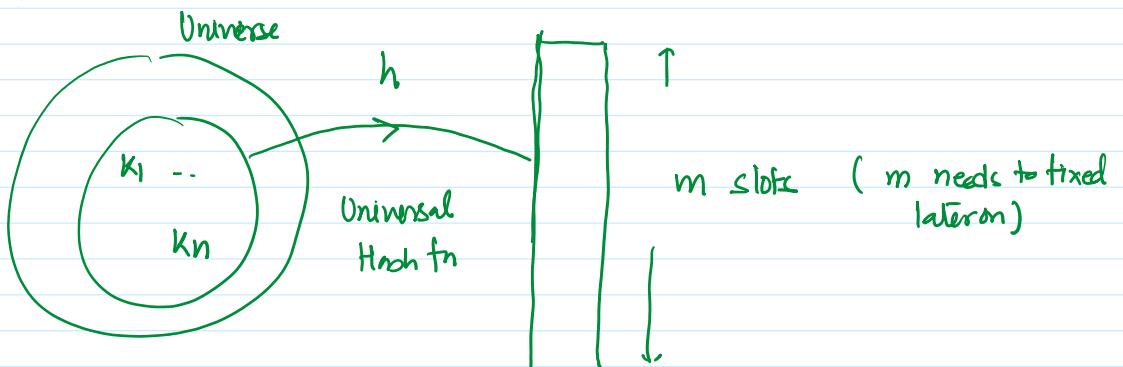
Naive Soln

① Store key in sorted order

② Universal Hashing

Space complexity of preparing data structure	Time complexity (Search)
$O(n)$	$O(\log n)$ — deterministic
$O(n^2)$	$O(1)$ Probabilistic

Simple Soln using Hashing



Lemma 1

$$I_{ij} = \begin{cases} 1 & \text{if } h(k_i) = h(k_j) \\ 0, \quad 0.0 & \text{otherwise} \end{cases}$$

$$\# \text{ of collision} = \frac{1}{2} \sum_{i,j} I_{ij} \quad i, j \in [n]$$

$$E \left[\sum_{i,j} I_{ij} \right] = \sum_{i,j} E[I_{ij}] \leq \sum_{i,j} \frac{1}{m} < \frac{n^2}{m}$$

$$< \frac{1}{10}$$

choose $m = 10n^2$

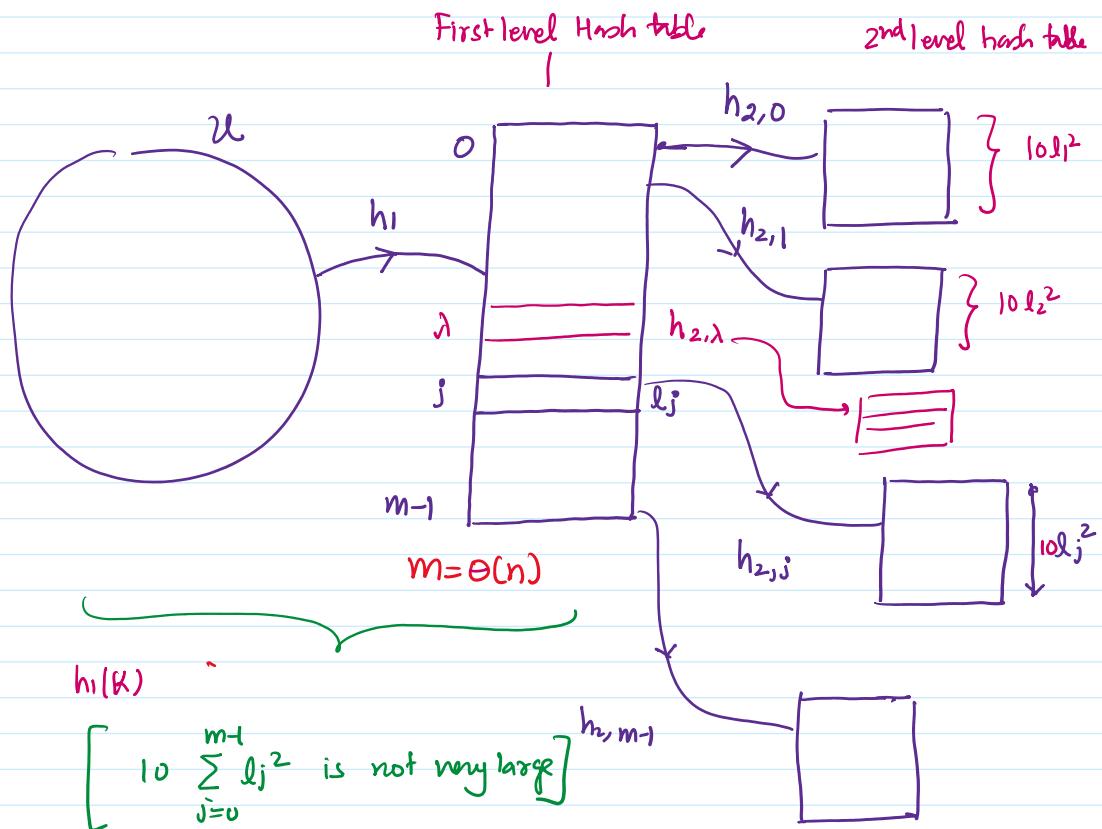
$$\Pr[\# \text{ of collision} > 1] < \frac{E \left[\sum_{i,j} I_{ij} \right]}{1} < \frac{1}{10}$$

\Rightarrow w.p. at least $\frac{9}{10}$ there will not be any collision

Question: coming up with data structure for static dictionary problem

with space complexity $O(n)$

Search $\approx O(1)$.



Step 1: Pick $h_1: \{0 \dots n-1\} \rightarrow \{0 \dots m-1\}$ from universal hash family

Step 1.5 If $\sum_{j=0}^{m-1} l_j^2 > c'n$ then redo step 1

Step 2: For each slot in first hash table $j \in \{0 \dots m-1\}$

$l_j := \# \text{ of key mapped in } j^{\text{th}} \text{ slot}$

Sample $h_{2,j} : \underbrace{\{0 \dots 121-1\}}_{\rightarrow} \rightarrow \{0 \dots (10l_j^2 - 1)\}$

Rehash if there is collision at 2nd level

Step 2.5 Do the following for all $j \in [m]$

while $h_{2,j}(k_i) = h_{2,j}(k_{i'})$ for some j , $k_i \neq k_{i'}$

with $h_1(k_i) = h_1(k_{i'})$

Repick $h_{2,j}$

Similar to proof of lemma 1

$\Rightarrow \Pr[\text{# of collision} > 1] < \frac{1}{10}$

\Rightarrow w.p. $9/10$ there is no collision at 2nd level

Claim Total space complexity is linear in n .

1984
Fredman, Komlos, Szemeredi

Prof:

$$E \left[10 \underbrace{\sum_{j=0}^{m-1} l_j^2}_{\text{Space required for 2nd level of Hash table}} \right] = 10 E \left[\sum_{j=0}^{m-1} l_j^2 \right]$$

FKS Algo

Perfect Hashing

$= 10 E \left[\sum_{i=1}^n \sum_{j=1}^n I_{ij} \right]$

$I_{ij} = \begin{cases} 1 & , h_1(k_i) = h_1(k_j) \\ 0 & . \text{ o.w.} \end{cases}$

$$= 10 \sum_{i=1}^n \sum_{j=1}^n E[I_{ij}]$$

$$\leq \frac{10n^2}{m} \quad \left[\because m = \Theta(n) \right]$$

$$\leq C \cdot n - \textcircled{1}$$

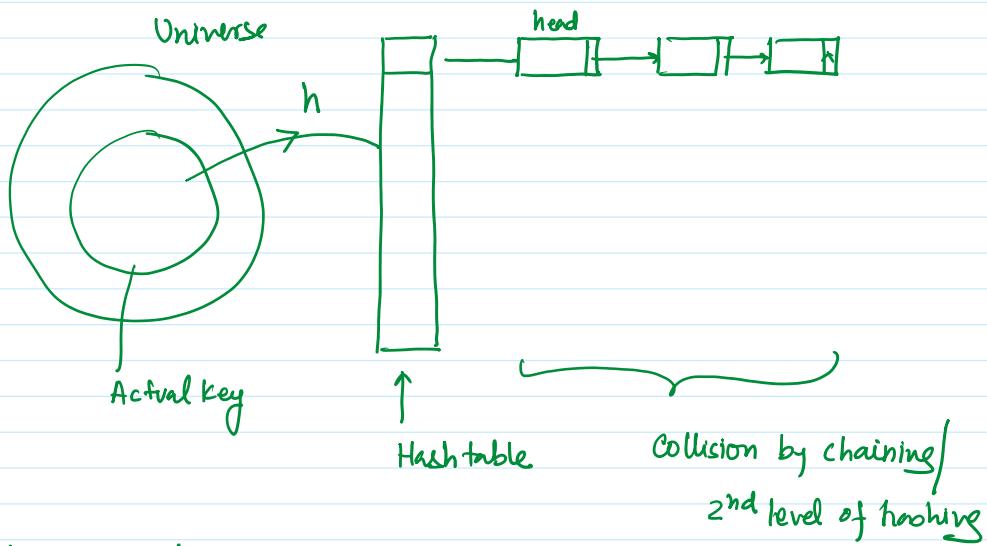
$$\leq c \cdot n - \textcircled{II}$$

Apply Markov on Eqn \textcircled{II}

$$\Pr_{h_1} \left[10 \sum_{j=0}^{m-1} l_j^2 > c' n \right] \leq \frac{E \left[10 \sum_{j=0}^{m-1} l_j^2 \right]}{c' n} \leq \frac{c n}{c' n} < \frac{1}{10} \quad \square$$

choosing c s.t

$$\frac{c}{c'} < \frac{1}{10}$$



Direct level Hashing