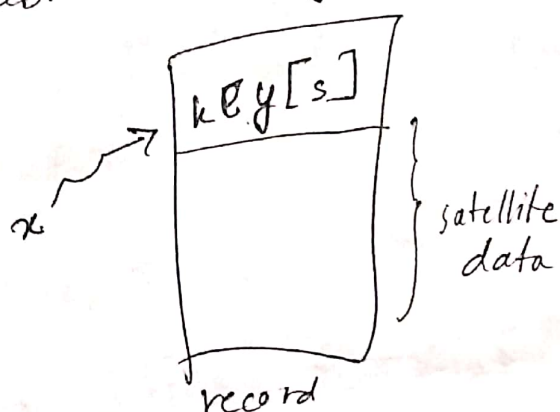


Lecture-7Hashing ISymbol table problem:Table  $S$  holding  $n$  recordsOperations

dynamic

Insert( $S, x$ ) :  $S \leftarrow S \cup \{x\}$ Delete( $S, x$ ) :  $S \leftarrow S - \{x\}$ Search( $S, k$ ) : return  $x$  \$key[ $x$ ] =  $k$  or  
nil, if no  
such  $x$ direct access table

Suppose key are drawn from

$$C = \{0, 1, \dots, m-1\}$$

Assume keys are distinct

Set up array  $T[0, \dots, m-1]$ dynamic set  $S$ 

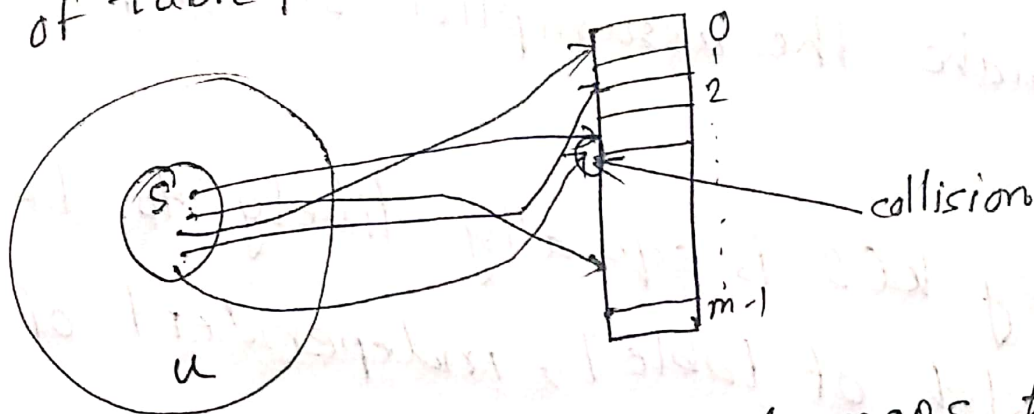
$$T[k] = \begin{cases} x & \text{if } x \in S \text{ and } \text{key}[x] = k \\ \text{nil} & \text{otherwise} \end{cases}$$

The operations take  $\Theta(1)$  timeProblems

The range of keys can be large:

## Hashing:

Hash function  $h$  maps keys "randomly" into slots of table  $T$



When a record to be inserted maps to an already occupied slot  $T$ , a collision occurs.

When a record to be inserted maps to an already occupied slot, a collision occurs.

Resolving collisions by chaining:

Idea: link records in same slot into list

Worst case: Every key hashes to the same slot.

Access time  $O(n)$  if  $|S|=n$



## Average case analysis of chaining

We make the assumption of simple uniform hashing.

Each key  $k \in S$  is equally likely to be hashed at any slot of table  $T$ , independent of where other keys are hashed.

Let  $n$  be the number of key in the table and  $m$  be the number of slots.

Define the load number of  $T$  to be

$$d = n/m = \text{average number of key per slot}$$

The expected time for an unsuccessful search for a record with a given key is

$$O(1 + d)$$

↑  
apply hash function and access the slot

↑  
search the list

Expected search time  $= O(1)$  if  $d = O(1)$  or equivalently  $n = O(m)$

A successful search has same asymptotic bound, but a rigorous argument is a little more complicated.

### Choosing hash function:

- should distributed keys uniformly into slots not
- Regularity in key distribution should affect uniformity

### Division method:

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

Deficiency

Don't pick  $m$  with small divisor  $d$

Ex:  $d=2$  and all keys are even

→ odd slots are never use

Extreme deficiency

Ex:  $m = 2^n$  → hash doesn't depend on all bits of  $k$



Good heuristic:

Pick  $m = \text{prime}$  not too close to  $\text{pow of } 2$  or  $10$

Annoyance:

Sometimes making the table size a prime is convenient.

Multiplication Method:

$m = 2^r$  computer has  $w$  bit words

$$h(k) = (A \cdot k \bmod 2^w) \text{ rsh } (w-r)$$

bit-wise  
right shift

$A$  is an odd integer be in range  $2^{w-1} < A < 2^w$

Don't pick  $A$  too close to  $2^{w-1}$  or  $2^w$

Fast method: module  $2^w$  is fast  
'rsh' operator is fast

$$h(k) = (A \cdot k \bmod 2^w) \text{ rsh } (w-r)$$

suppose that  $m = 8 = 2^3$  and that our comp  
has  $w = 7$ -bit words

$$\begin{array}{r}
 1011001 = A \\
 1101011 = k \\
 \hline
 10010100110011 \\
 \text{ignored} \quad \quad \quad \underline{h(k)}
 \end{array}$$

## Resolving collision by open addressing:

→ No storage ~~for links~~ is used outside of the hash table itself.

- Insertion systematically probes the table until an empty slot is found.
- The hash function depends on both the key and probe number

$$h: U \times \{0, 1, \dots, m-1\} \rightarrow \{0, 1, \dots, m-1\}$$

- The probe sequence  $\langle h(k, 0), h(k, 1), \dots, h(k, m-1) \rangle$  should be a permutation of  $\{0, 1, \dots, m-1\}$
- The table may fill up, and deletion is difficult (but not impossible)

Search: same probe sequence

- successful - find record
- unsuccessful - find nil



## Probing strategies

### Linear probing:

$$h(k, i) = (h(k, 0) + i) \bmod m$$

$$h(k, i) = (h(k) + i) \bmod m$$

This method though simple suffers from primary clustering, where long runs of occupied slots build up, increasing the average search time. Moreover, the long run of occupied slots tend to get longer.

### Double hashing:

Given two ordinary hash function  $h_1(k)$  and  $h_2(k)$  double hashing uses the hash function

$$h(k, i) = (h_1(k) + i \cdot h_2(k)) \bmod m$$

This method generally produces excellent results but  $h_2(k)$  must be relatively prime to  $m$ . One way is to make  $m$  a power of 2 and design  $h_2(k)$  to produce only odd numbers.

## Average case

### Analysis of open addressing:

We make assumptions of uniform hashing

→ Each key is equal to have any one of the  $m!$  permutations as its probe sequence.

Theorem: Given an open addressed hash table with load factor  $\alpha = n/m < 1$ , the expected number of probes in an unsuccessful search is at most  $1/(1-\alpha)$   $E[\# \text{ probes}] = \frac{1}{1-\alpha}$

### Proof of Theorem:

- At least one probe is always necessary
- With probability  $n/m$ , the first probe hits an occupied slot, and a second probe is necessary
- With probability  $(n-1)/(m-1)$  the second probe hits an occupied slot and third probe is necessary

Observe  $\frac{n-i}{m-i} < \frac{n}{m} = \alpha$  for  $1, 2, \dots, n$



Expected number of probes

$$1 + n/m \left( 1 + \frac{n-1}{m-1} \left( 1 + \frac{n-2}{m-2} \left( \dots \left( 1 + \frac{1}{m-n+1} \right) \dots \right) \right) \right)$$

$$\leq 1 + d \left( 1 + d \left( 1 + d \left( \dots (1 + d) \dots \right) \right) \right)$$

$$\leq 1 + d + d^2 + d^3 + \dots$$

$$= \sum_{i=0}^{\infty} d^i$$

$$= \frac{1}{1-d}$$

