

# Hashing With Chaining

## □ Dictionary

- abstract data type (ADT)
- maintain set of items, each with a key
  - insert item: overwrite any existing key
  - delete item
  - Search key: return item with given key or report doesn't exist.
- $O(lgn)$  time via AVL tree → how to search faster ??

## ○ In python

- $D[key]$  ~ search
- $D[key] = val$ . ~ insert
- $\text{del } D[key]$  ~ delete
- item = (key, val.)

## □ Motivation

- doclist
- database = hashing, search/trees
- compilers & interpreters
- Network: router, server
- substring search
- String commonalities

- file/dir Synchronization
- cryptography

## □ Simple approach

### ◦ Direct-access table

store items in array indexed by key

- Badness:

- ① keys may not be integers
- ② gigantic memory hog

- fix both of above problems

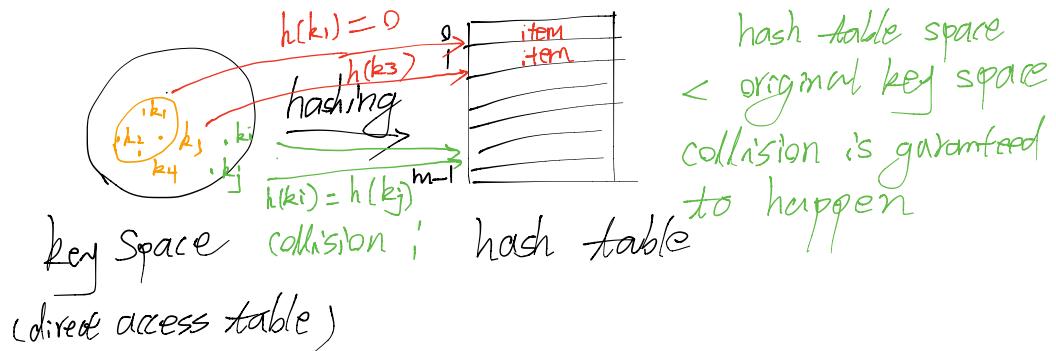
### ◦ Solution to ①: prehash

- maps key → non-negative integers
- In theory, keys are finite & discrete  
(string of bits)
- In python, `hash(x)` is the prehash of  $x$
- $\text{hash('AB')} = \text{hash('DC')} = 64$
- Ideally  $\text{hash}(x) = \text{hash}(y) \Leftrightarrow x = y$
- In python, default "`--hash--`" = `id`

## o Solution to ② : hashing

English meaning of hash = cut into pieces and mix them around.

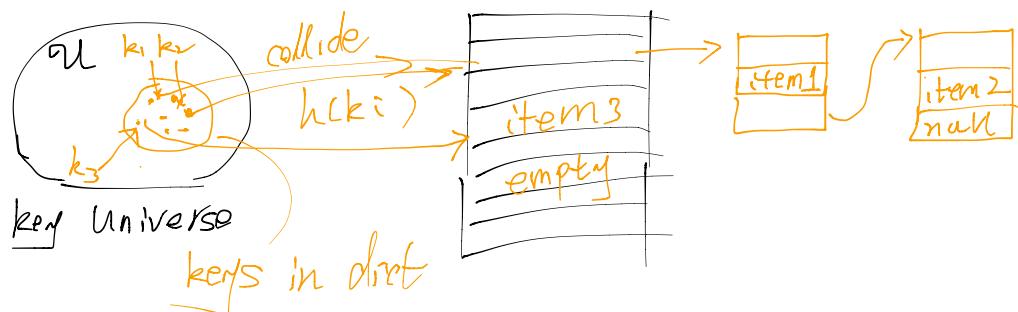
- reduce universe  $U$  of all keys (integers) down to reasonable size  $m$  for table, using hash function  $h()$



- idea:  $m = H(n)$ ,  $n = \#$  keys in the dictionary

## □ chaining — dealing with collisions

chaining = linked list of colliding elements in each slot of hash table



If we use same hashing function, worst case is that  $h(k_i)$  maps all keys in dict to same slot in the hash table and they are thus stored as linked list. In this case search take  $\Theta(n)$ , while we want  $\Theta(1)$

- Simple uniform hashing (assumption, but not true)

Each key is equally likely to be hashed to any slot of the table independent of where other keys hashing.

If above is true, searching takes const. time

- Analysis

- The expected length of chain for  $n$  keys,  $m$  slots in the table  $\rightarrow \frac{n}{m} = \alpha \equiv$  load factor of table as long as  $m = \Theta(n)$ ,  $\alpha = \Theta(1)$

- Running time =  $O(\underbrace{1 + \alpha}_{\text{time computing the hash func.}})$  for searching

## II "Good" hash functions

- Division method

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

- This is a bad method

especially when  $k$  and  $m$  has some common factors

e.g.  $k$  always even and  $m$  is even

↪ only use half of the table

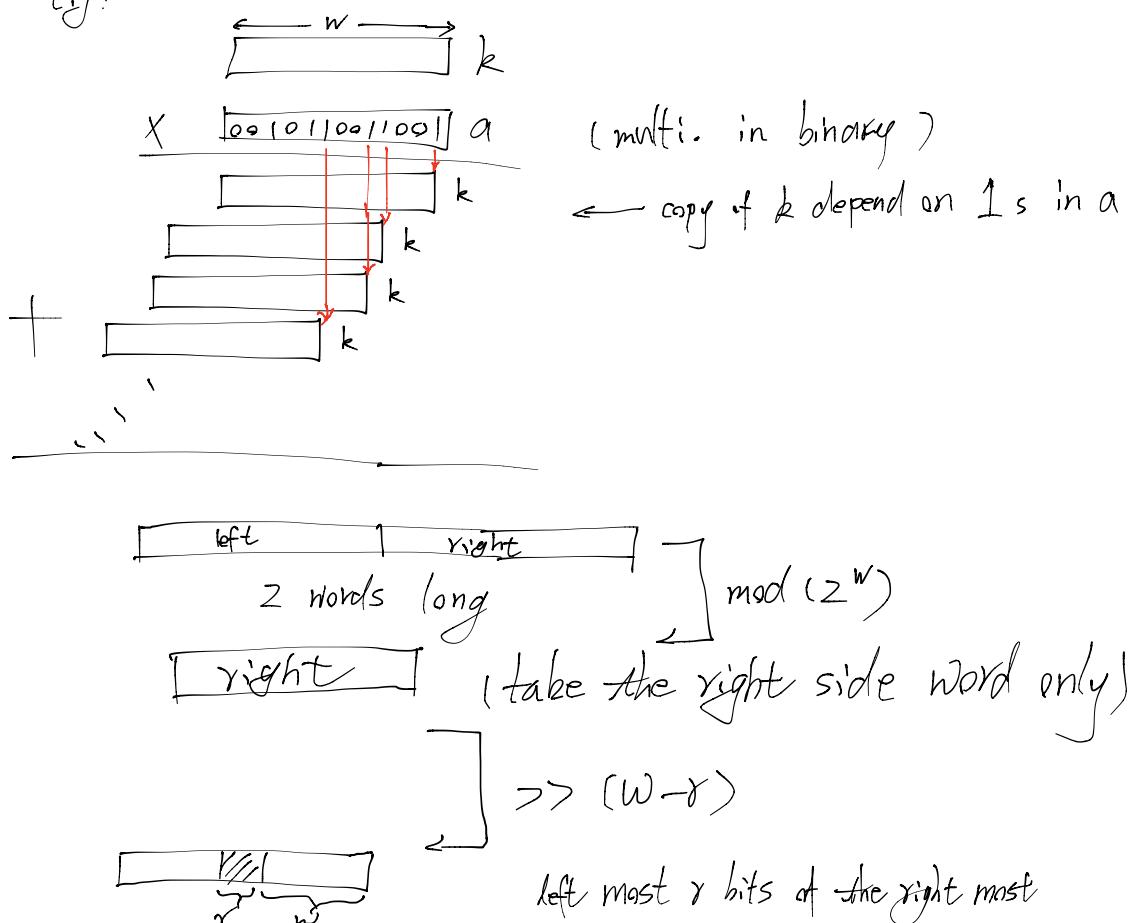
In practice is pretty good when choose  $m$  a prime number

## • Multiplication method

$$h(k) = [(a \cdot k) \bmod 2^w] \gg (w-r)$$

- $w$ :  $w$ -bit machine
- $r$ :  $m = 2^r$
- $a$ : random const. odd

e.g.:





\* not always a good method

### o Universal hashing

$$h(k) = [(ak + b) \bmod p] \bmod m$$

- $a, b$ : random  $\in \{0, 1, \dots, p-1\}$

- $p$ : big prime number  $> |U|$  (bigger than the key universe)

For worst case keys  $k_1 \neq k_2$ , colliding probability

$$\Pr\{h(k_1) = h(k_2) \mid k_1 \neq k_2\} = \frac{1}{m}$$