

# Hash Tables

---

# Abstract data type

---

## Data Structure

a way to store and organize data to facilitate access and modifications.

*Ex. Linked list, array, heap, ...*

## Abstract Data Type (ADT)

a set of data values and associated operations that are precisely specified independent of any particular implementation.

*Ex. Dictionary, stack, queue, priority queue, set, bag ...*

- ❑ ADT describe the functionality of data structures
- ❑ Data structures **implement** ADT
  - how is the data stored?
  - which algorithms implement the operations?

# Priority queue

---

## Max-priority queue

Stores a set  $S$  of **elements**, each with an associated **key** (integer value).

## Operations

Insert( $S, x$ ): inserts element  $x$  into  $S$ , that is,  $S \leftarrow S \cup \{x\}$

Maximum( $S$ ): returns the element of  $S$  with the largest **key**

Extract-Max( $S$ ): removes and returns the element of  $S$  with the largest **key**

Increase-Key( $S, x, k$ ): give **key**[ $x$ ] the value  $k$

condition:  $k$  is larger than the current value of **key**[ $x$ ]

# Implementing a priority queue

---

	Insert	Maximum	Extract-Max	Increase-Key
sorted list	$\Theta(n)$	$\Theta(1)$	$\Theta(1)$	$\Theta(n)$
sorted array	$\Theta(n)$	$\Theta(1)$	$\Theta(n)$	$\Theta(n)$
heap	$\Theta(\log n)$	$\Theta(1)$	$\Theta(\log n)$	$\Theta(\log n)$

# Dictionary

---

## Dictionary

Stores a set  $S$  of **elements**, each with an associated **key** (integer value).

## Operations

Search( $S, k$ ): return a pointer to an element  $x$  in  $S$  with  $\text{key}[x] = k$ , or **NIL** if such an element does not exist.

Insert( $S, x$ ): inserts element  $x$  into  $S$ , that is,  $S \leftarrow S \cup \{x\}$

Delete( $S, x$ ): remove element  $x$  from  $S$

---

$S$ : personal data

- **key**: burger service number
- name, date of birth, address, ... (**satellite data**)

# Implementing a dictionary

---

	Search	Insert	Delete
linked list	$\Theta(n)$	$\Theta(1)$	$\Theta(1)$
sorted array	$\Theta(\log n)$	$\Theta(n)$	$\Theta(n)$

---

# Hash Tables

---

# Hash tables

---

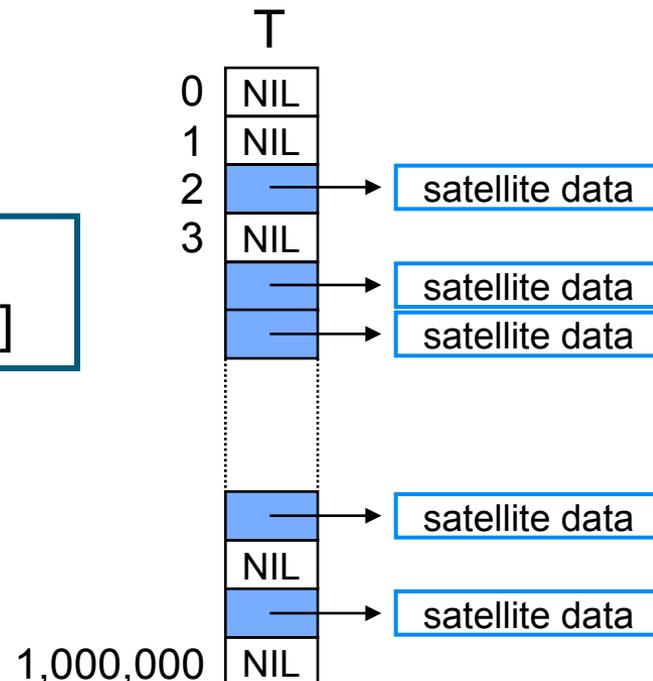
- Hash tables generalize ordinary arrays

# Hash tables

- **S**: personal data in population register
  - **key**: bsn (burgerservicenummer)
  - name, date of birth, address, ... (**satellite data**)

**Assume**: bsn-numbers are integers in the range  $[0 \dots 1,000,000]$

Direct addressing  
use table  $T[0 \dots 1,000,000]$



# Direct-address tables

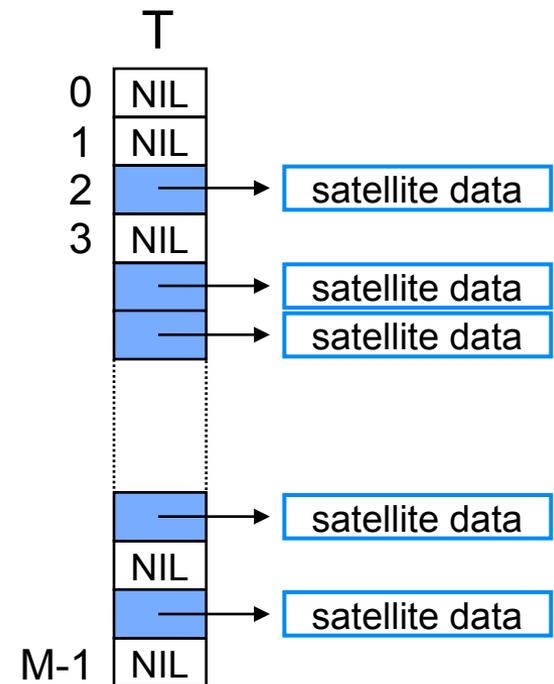
- $S$ : set of elements
  - **key**: unique integer from the universe  $U = \{0, \dots, M-1\}$
  - **satellite data**

- use table (array)  $T[0..M-1]$

$$T[i] = \begin{cases} \text{NIL if there is no element with} \\ \text{key } i \text{ in } S \\ \text{pointer to the satellite data if there} \\ \text{is an element with key } i \text{ in } S \end{cases}$$

## Analysis:

- Search, Insert, Delete:  $O(1)$
- Space requirements:  $O(M)$



# Direct-address tables

---

- S: personal data
    - key: bsn
    - name, date of birth, address, ... (satellite data)
- 

Assume: bsn are integers with 9 digits

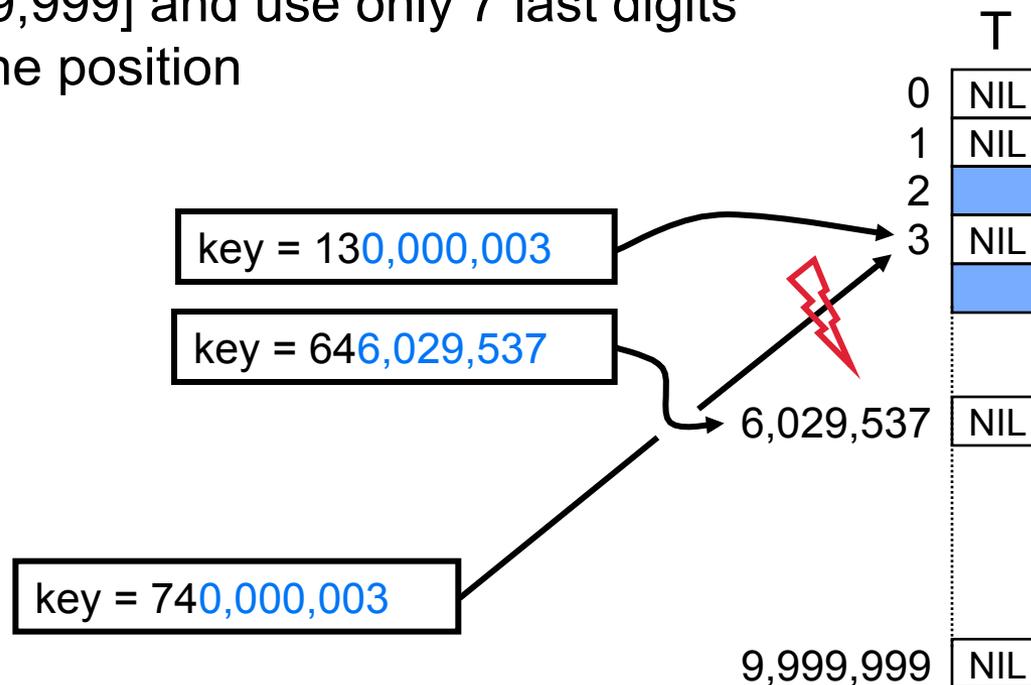
➔ use table  $T[0 \dots 999,999,999]$  **?!?**

- uses too much memory, most entries will be NIL ...
- if the universe  $U$  is large, storing a table of size  $|U|$  may be impractical or impossible
- often the set  $K$  of keys actually stored is small, compared to  $U$ 
  - ➔ most of the space allocated for  $T$  is wasted.

# Hash tables

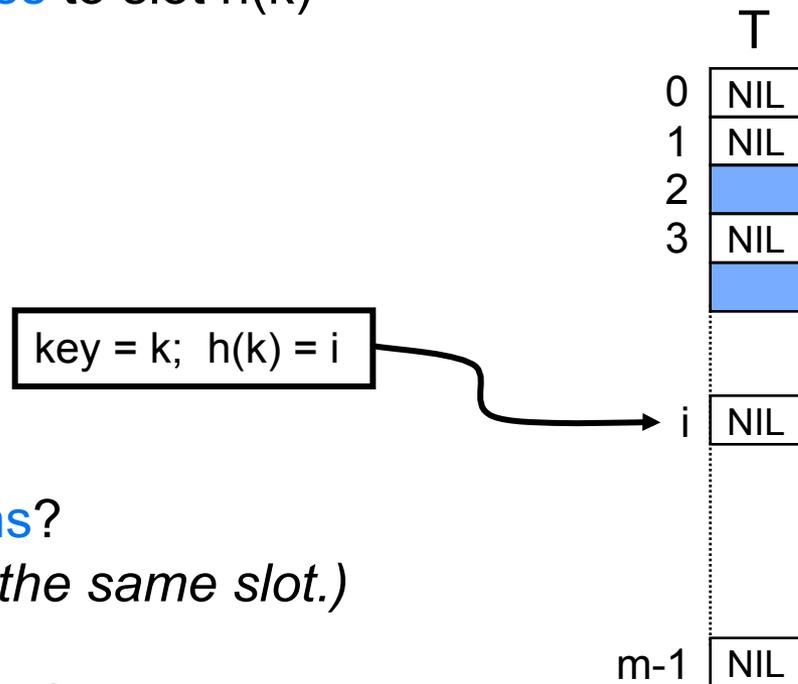
- S: personal data
  - key = bsn = integer from  $U = \{0 \dots 999,999,999\}$

**Idea:** use a smaller table, for example,  
 $T[0 \dots 9,999,999]$  and use only 7 last digits  
to determine position



# Hash tables

- $S$  set of keys from the universe  $U = \{0 \dots M-1\}$ 
  - use a hash table  $T [0..m-1]$  (with  $m \leq M$ )
  - use a hash function  $h : U \rightarrow \{0 \dots m-1\}$  to determine the position of each key: key  $k$  hashes to slot  $h(k)$



- How do we resolve collisions?  
(Two or more keys hash to the same slot.)
- What is a good hash function?

# Resolving collisions: chaining

**Chaining:** put all elements that hash to the same slot into a linked list

Example ( $m=1000$ ):

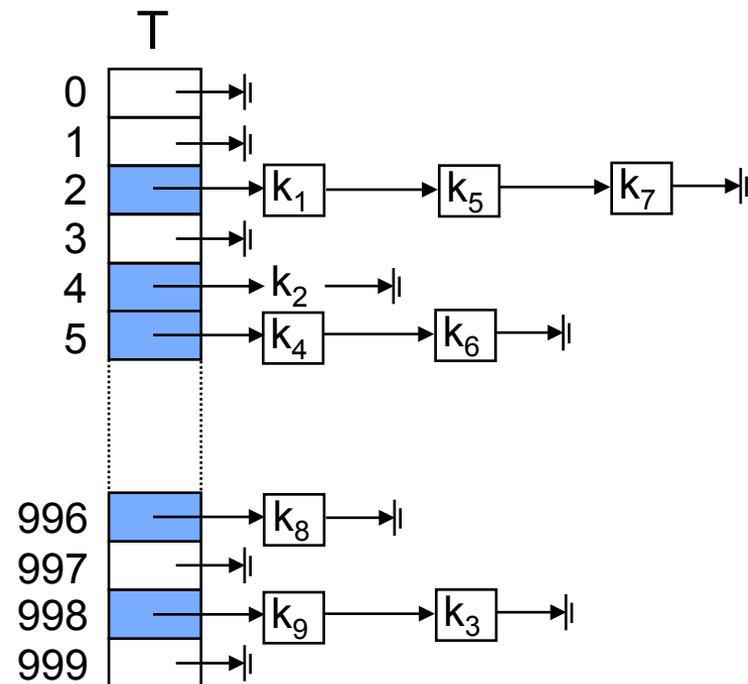
$$h(k_1) = h(k_5) = h(k_7) = 2$$

$$h(k_2) = 4$$

$$h(k_4) = h(k_6) = 5$$

$$h(k_8) = 996$$

$$h(k_9) = h(k_3) = 998$$



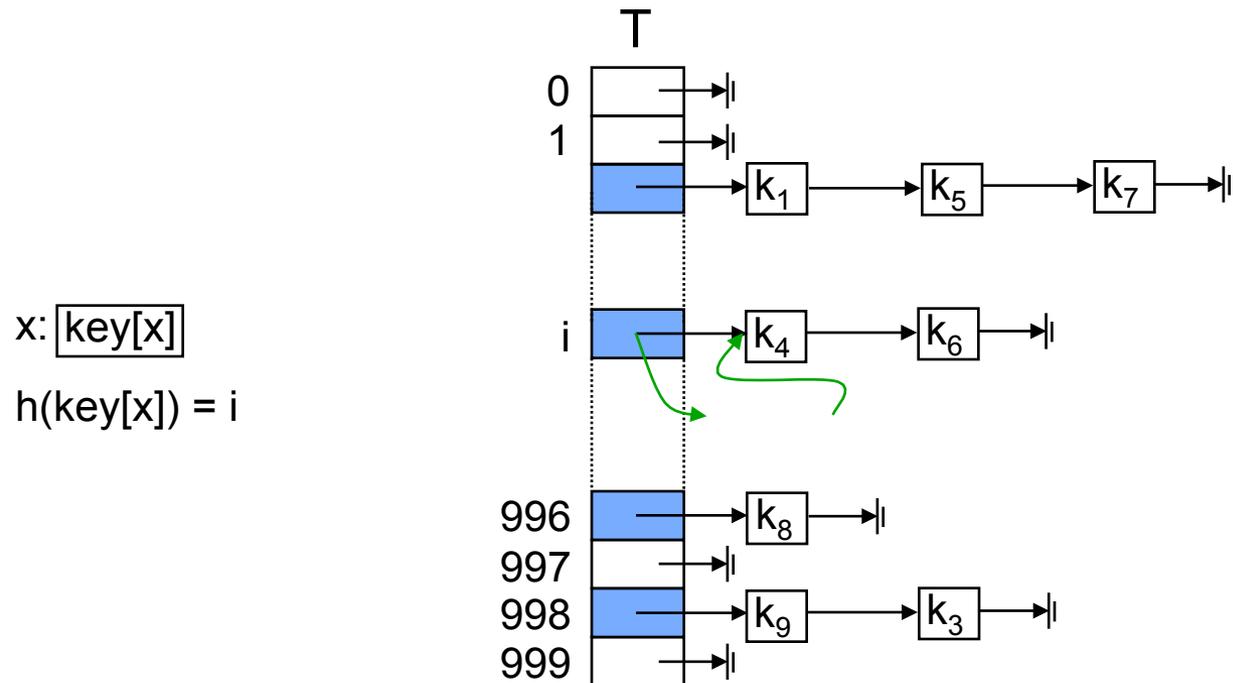
*Pointers to the satellite data also need to be included ...*

# Hashing with chaining: dictionary operations

## Chained-Hash-Insert( $T, x$ )

insert  $x$  at the head of the list  $T[h(\text{key}[x])]$

Time:  $O(1)$



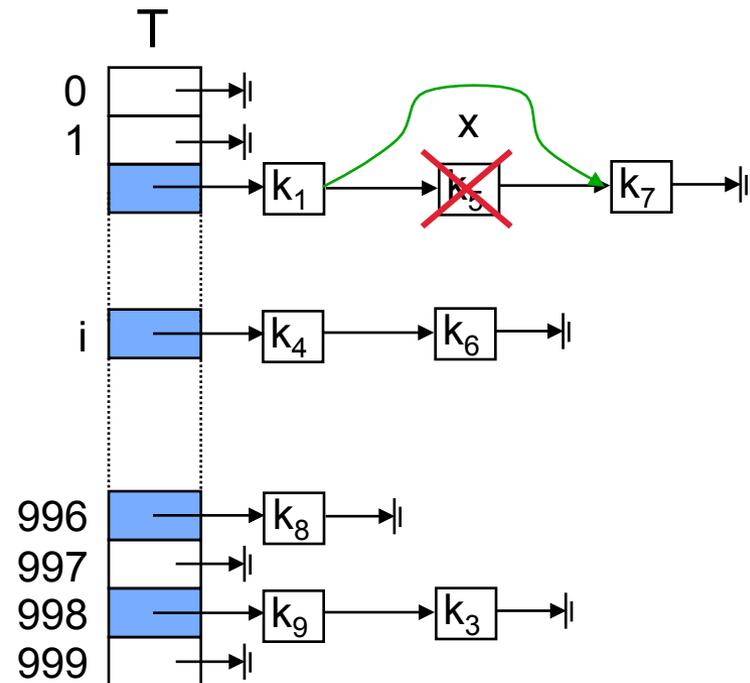
# Hashing with chaining: dictionary operations

**Chained-Hash-Delete**( $T, x$ )  
delete  $x$  from the list  $T[h(\text{key}[x])]$

□  $x$  is a pointer to an element

Time:  $O(1)$

(with doubly-linked lists)



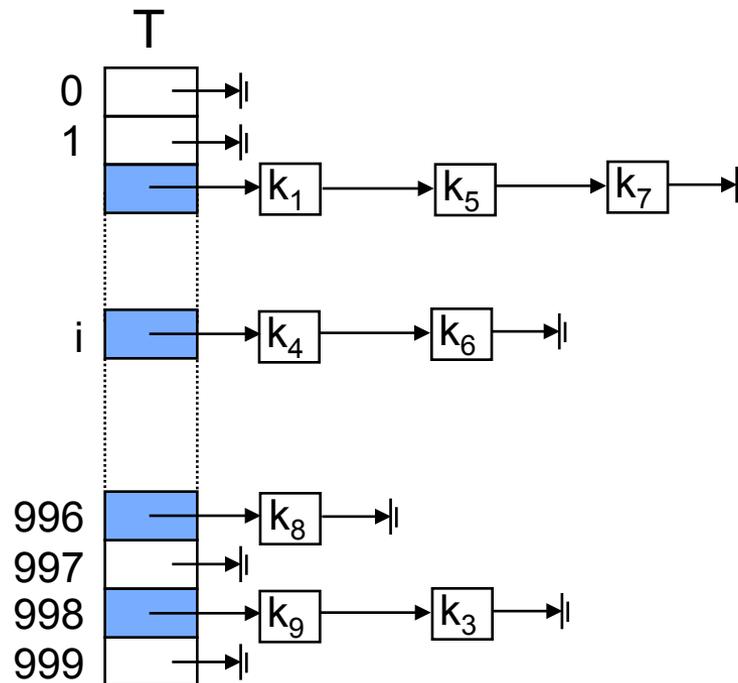
# Hashing with chaining: dictionary operations

## Chained-Hash-Search( $T, k$ )

search for an element with key  $k$  in list  $T[h(k)]$

## Time:

- unsuccessful:  $O(1 + \text{length of } T[h(k)])$
- successful:  $O(1 + \# \text{ elements in } T[h(k)] \text{ ahead of } k)$



# Hashing with chaining: analysis

---

Time:

- unsuccessful:  $O(1 + \text{length of } T[h(k)])$
- successful:  $O(1 + \# \text{ elements in } T[h(k)] \text{ ahead of } k)$

➔ worst case  $O(n)$

Can we say something about the average case?

Simple uniform hashing

any given element is equally likely to hash into any of the  $m$  slots

# Hashing with chaining: analysis

---

## Simple uniform hashing

any given element is equally likely to hash into any of the  $m$  slots

in other words ...

- the hash function distributes the keys from the universe  $U$  uniformly over the  $m$  slots
  - the keys in  $S$ , and the keys with whom we are searching, behave as if they were randomly chosen from  $U$
- ➔ we can analyze the average time it takes to search as a function of the **load factor**  $\alpha = n/m$

*( $m$ : size of table,  $n$ : total number of elements stored)*

# Hashing with chaining: analysis

---

## Theorem

In a hash table in which collisions are resolved by chaining, an unsuccessful search takes time  $\Theta(1+\alpha)$ , on the average, under the assumption of simple uniform hashing.

## Proof (for an arbitrary key)

- the key we are looking for hashes to each of the  $m$  slots with equal probability
  - the average search time corresponds to the average list length
  - average list length = total number of keys / # lists =  $\alpha$
- 
- The  $\Theta(1+\alpha)$  bound also holds for a successful search (although there is a greater chance that the key is part of a long list).
  - If  $m = \Omega(n)$ , then a search takes  $\Theta(1)$  time on average.

---

What is a good hash function?

---

# What is a good hash function?

---

1. as random as possible

*get as close as possible to simple uniform hashing ...*

- the hash function distributes the keys from the universe  $U$  uniformly over the  $m$  slots
- the hash function has to be as independent as possible from patterns that might occur in the input

3. fast to compute

# What is a good hash function?

---

**Example:** hashing performed by a compiler for the symbol table

- keys: variable names which consist of (capital and small) letters and numbers: i, i2, i3, Temp1, Temp2, ...

**Idea:**

- use table of size  $(26+26+10)^2$
- hash variable name according to the first two letters:  
Temp1 → Te

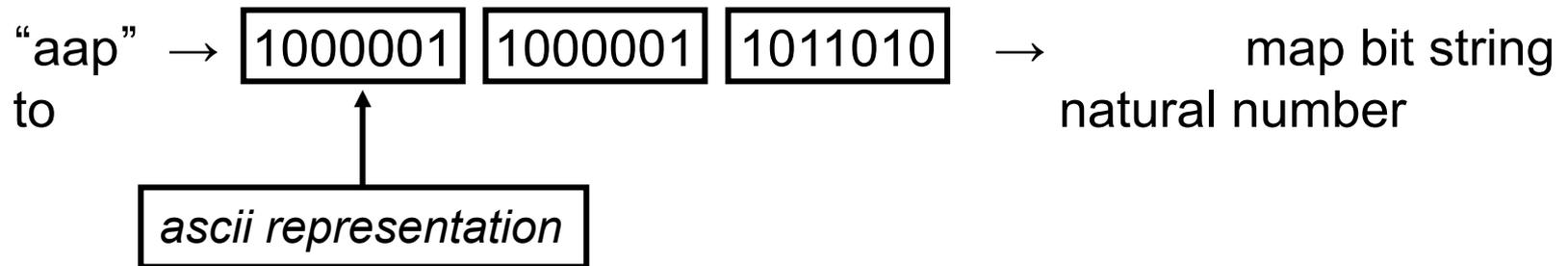


**Bad idea:** too many “clusters”  
*(names that start with the same two letters)*

# What is a good hash function?

**Assume:** keys are natural numbers

*if necessary first map the keys to natural numbers*



➔ the hash function is  $h: \mathbf{N} \rightarrow \{0, \dots, m-1\}$

▣ the hash function always has to depend on all digits of the input

# Common hash functions

**Division method:**  $h(k) = k \bmod m$

Example:  $m=1024$ ,  $k = 2058 \rightarrow h(k) = 10$

- don't use a power of 2  
 $m = 2^p \rightarrow h(k)$  depends only on the  $p$  least significant bits
- use  $m =$  prime number, not near any power of two

**Multiplication method:**  $h(k) = \lfloor m (kA \bmod 1) \rfloor$

1.  $0 < A < 1$  is a constant
  2. compute  $kA$  and extract the **fractional part**
  3. multiply this value with  $m$  and then take the floor of the result
- Advantage: choice of  $m$  is not so important, can choose  $m =$  power of 2

---

# Resolving collisions

---

more options ...

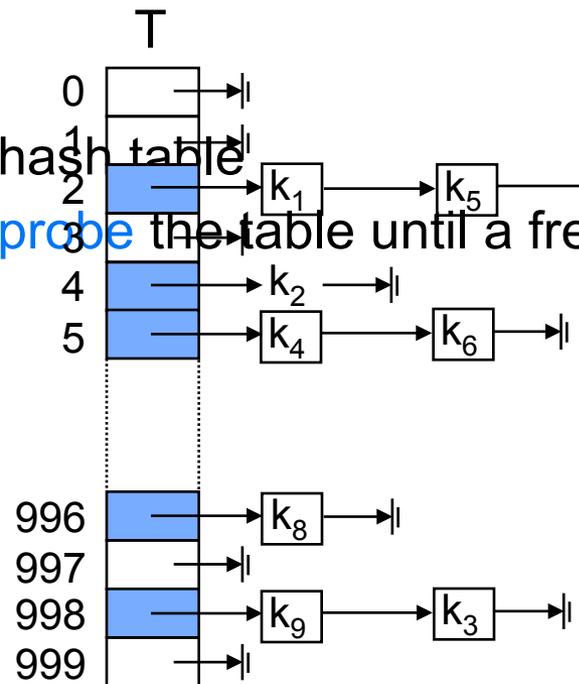
# Resolving collisions

## Resolving collisions

1. **Chaining**: put all elements that hash to the same slot into a linked list

2. **Open addressing**:

- ❑ store all elements in the hash table
- ❑ when a collision occurs, **probe** the table until a free slot is found



# Hashing with open addressing

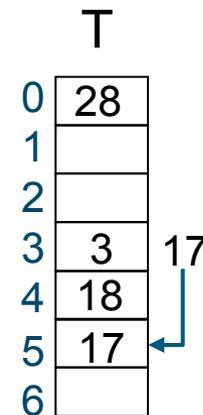
## Open addressing:

- store all elements in the hash table
- when a collision occurs, **probe** the table until a free slot is found

**Example:**  $T[0..6]$  and  $h(k) = k \bmod 7$

1. insert 3
2. insert 18
3. insert 28
4. insert 17

T	
0	28
1	
2	
3	3
4	18
5	17
6	



- no extra storage for pointers necessary
- the hash table can “fill up”
- the load factor  $\alpha$  is always  $\leq 1$

# Hashing with open addressing

---

- there are several variations on open addressing depending on how we search for an open slot
- the hash function has two arguments:  
the key and the number of the current probe
  - ➔ **probe sequence**  $\langle h(k,0), h(k, 1), \dots, h(k, m-1) \rangle$

The probe sequence has to be a permutation of  $\langle 0, 1, \dots, m-1 \rangle$  for every key  $k$ .

# Open addressing: dictionary operations

Hash-Insert( $T, k$ )

we're actually inserting element  $x$  with  $\text{key}[x] = k$

1.  $i = 0$
2. **while** ( $i < m$ ) **and** ( $T[ h(k,i) ] \neq \text{NIL}$ )
3.     **do**  $i = i + 1$
4. **if**  $i < m$
5.     **then**  $T [h(k,i)] = k$
6.     **else** “hash table overflow”

Example: [Linear Probing](#)

- $T[0..m-1]$
- $h'(k)$  ordinary hash function
- $h(k,i) = (h'(k) + i) \bmod m$
- Hash-Insert( $T, 17$ )

T	
0	28
1	
2	
3	3 17
4	18 17
5	17 17
6	

# Open addressing: dictionary operations

## Hash-Search(T,k)

1.  $i = 0$
2. **while**  $(i < m)$  **and**  $(T[h(k,i)] \neq \text{NIL})$
3.     **do if**  $T[h(k,i)] = k$
4.         **then return** “k is stored in slot  $h(k,i)$ ”
5.         **else**  $i = i + 1$
6. **return** “k is not stored in the table”

## Example: Linear Probing

- $h'(k) = k \bmod 7$   
 $h(k,i) = (h'(k) + i) \bmod m$
- Hash-Search(T,17)

T	
0	28
1	
2	
3	3 17
4	18 17
5	17 17
6	

# Open addressing: dictionary operations

## Hash-Search(T,k)

1.  $i = 0$
2. **while**  $(i < m)$  **and**  $(T[h(k,i)] \neq \text{NIL})$
3.     **do if**  $T[h(k,i)] = k$
4.         **then return** “k is stored in slot  $h(k,i)$ ”
5.         **else**  $i = i + 1$
6. **return** “k is not stored in the table”

## Example: Linear Probing

- $h'(k) = k \bmod 7$   
 $h(k,i) = (h'(k) + i) \bmod m$
- Hash-Search(T,17)
- Hash-Search(T,25)

T	
0	28
1	
2	
3	3
4	18
5	17
6	

25

25

25

# Open addressing: dictionary operations

## Hash-Delete( $T, k$ )

1. remove  $k$  from its slot
2. mark the slot with the special value DEL

Example: delete 18

	T
0	28
1	
2	
3	3
4	DEL
5	17
6	

- ❑ Hash-Search passes over DEL values when searching
- ❑ Hash-Insert treats a slot marked DEL as empty
  - ➔ search times no longer depend on load factor
  - ➔ use chaining when keys must be deleted

# Open addressing: probe sequences

---

- $h'(k)$  = ordinary hash function

**Linear probing:**  $h(k,i) = (h'(k) + i) \bmod m$

- $h'(k_1) = h'(k_2) \Rightarrow k_1$  and  $k_2$  have the same probe sequence
  - the initial probe determines the entire sequence
    - ➔ there are only  $m$  distinct probe sequences
  - all keys that test the same slot follow the same sequence afterwards
- 
- Linear probing suffers from **primary clustering**: long runs of occupied slots build up and tend to get longer
    - ➔ the average search time increases

# Open addressing: probe sequences

- $h'(k)$  = ordinary hash function

**Quadratic probing:**  $h(k,i) = (h'(k) + c_1i + c_2i^2) \bmod m$

- $h'(k_1) = h'(k_2) \Rightarrow k_1$  and  $k_2$  have the same probe sequence
  - the initial probe determines the entire sequence
    - ➔ there are only  $m$  distinct probe sequences
  - but keys that test the same slot do not necessarily follow the same sequence afterwards
- 
- quadratic probing suffers from **secondary clustering**: if two distinct keys have the same  $h'$  value, then they have the same probe sequence

*Note:  $c_1$ ,  $c_2$ , and  $m$  have to be chosen carefully, to ensure that the whole table is tested.*

# Open addressing: probe sequences

---

□  $h'(k)$  = ordinary hash function

**Double hashing:**  $h(k,i) = (h'(k) + i h''(k)) \bmod m$ ,  
 $h''(k)$  is a second hash function

- keys that test the same slot do not necessarily follow the same sequence afterwards
- $h''$  must be relatively prime to  $m$  to ensure that the whole table is tested.
- $O(m^2)$  different probe sequences

# Open addressing: analysis

---

## Uniform hashing

each key is equally likely to have any of the  $m!$  permutations of  $\langle 0, 1, \dots, m-1 \rangle$  as its probe sequence

**Assume:** load factor  $\alpha = n/m < 1$ , no deletions

## Theorem

The average number of probes is

- $\Theta(1/(1-\alpha))$  for an unsuccessful search
- $\Theta((1/\alpha) \log (1/(1-\alpha)))$  for a successful search

# Open addressing: analysis

## Theorem

The average number of probes is

- $\Theta(1/(1-\alpha))$  for an unsuccessful search
- $\Theta((1/\alpha) \log(1/(1-\alpha)))$  for a successful search

$$\begin{aligned}\text{Proof: } E[\#\text{probes}] &= \sum_{1 \leq i \leq n} i \cdot \Pr[\#\text{ probes} = i] \\ &= \sum_{1 \leq i \leq n} \Pr[\#\text{ probes} \geq i]\end{aligned}$$

$$\begin{aligned}\Pr[\#\text{probes} \geq i] &= \frac{n}{m} \cdot \frac{n-1}{m-1} \cdot \frac{n-2}{m-2} \cdot \dots \cdot \frac{n-i+2}{m-i+2} \\ &\leq \left(\frac{n}{m}\right)^{i-1} = \alpha^{i-1}\end{aligned}$$

$$E[\#\text{probes}] \leq \sum_{1 \leq i \leq n} \alpha^{i-1} \leq \sum_{0 \leq i \leq \infty} \alpha^i = \frac{1}{1-\alpha} \quad \blacksquare$$

*Check the book for details!*

# Hash tables

---

- Hash tables generalize ordinary arrays
  - map a large universe to a small table
- How do we resolve collisions?
  - Chaining
  - Open addressing: linear and quadratic probing, double hashing
- What is a good hash function?
  - Division method
  - Multiplication method

# Implementing a dictionary

	Search	Insert	Delete
linked list	$\Theta(n)$	$\Theta(1)$	$\Theta(1)$
sorted array	$\Theta(\log n)$	$\Theta(n)$	$\Theta(n)$
hash table	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$

- Running times are average times and assume (simple) uniform hashing and a large enough table (for example, of size  $2n$ ).

**Drawbacks** of hash tables: operations such as finding the min or the successor of an element are inefficient.