CONFIDENTIAL

# C Programming Basic – week 14

*Mapping and Hashing*

**Lecturers :**
**Cao Tuan Dung**
**Le Duc Trung**
**Dept of Software Engineering**
**Hanoi University of Technology**
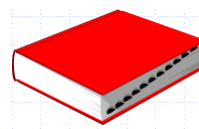
1

# Topics of this week

- Dictionary ADT
- Hash Table
- Hash functions
- Compression maps
- Collision handling
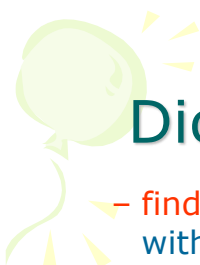- Exercises

2

1

# Dictionary ADT

- The dictionary ADT models a searchable collection of key-element items
- The main operations of a dictionary are searching, inserting, and deleting items
- Multiple items with the same key are allowed
- Applications:
  - address book
  - credit card authorization
  - mapping host names (e.g., csci260.net) to internet addresses (e.g., 128.148.34.101)
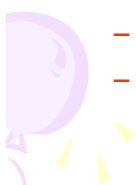
3

# Dictionary ADT methods

- findElement(k): if the dictionary has an item with key k, returns its element, else, returns the special element NO_SUCH_KEY
- insertItem(k, o): inserts item (k, o) into the dictionary
- removeElement(k): if the dictionary has an item with key k, removes it from the dictionary and returns its element, else returns the special element NO_SUCH_KEY
- size(), isEmpty()
- keys(), elements()

4

# Key-Indexed Dictionaries

A[]

| Key | Value |
|-----|-------|
| 1 | Intro to CS 1 |
| 2 | Intro to CS 2 |
| 5 | Theory of Computation |
| 7 | Data Structures |
| 9 | Digital Logic |

| | |
|---|---|
| 0 | |
| 1 | Intro to CS 1 |
| 2 | Intro to CS 2 |
| 3 | |
| 4 | |
| 5 | Theory of Computation |
| 6 | |
| 7 | Data Structures |
| 8 | |
| 9 | Digital Logic |

Space-efficient only if the cardinality of the set is close to N
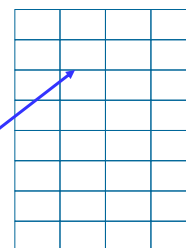
5

5

# Searching without Comparisons

- How could a search algorithm proceed without comparing data elements?
- What if we had some sort of "oracle" that could take the key for a data value and compute, in constant-bounded time, the location at which that key would occur within the data collection?

data key K

O(1) Oracle

$L_i$

location of matching record within the collection

If the container storing the collection supports random access with $\Theta(1)$ cost, as an array does, then we would have a total search cost of $\Theta(1)$.

6

6

# Hash Functions and Hash Tables

- An efficient way of implementing a dictionary is a hash table.
- Use an array (or list) of size N (table)
  - Need to spread keys over range [0,N-1]
  - Collisions occur when elements have same key
- Keys are not always integers, nor in range [0,N-1]
- A hash table for a given key type consists of
  - Hash function h
  - Array (called table) of size N
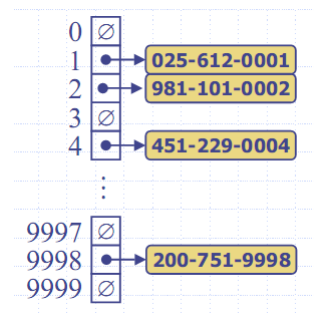- When implementing a dictionary with a hash table, the goal is to store item (k, o) at index $i = h(k)$

# Example

- We design a hash table for a dictionary storing items (SIN, Name), where SIN (social insurance number) is a nine-digit positive integer

- Our hash table uses an array of size N = 10,000 and the hash function

- h(x) = last four digits of x

```
   0 ∅
   1 ●─→ 025-612-0001
   2 ●─→ 981-101-0002
   3 ∅
   4 ●─→ 451-229-0004
     ⋮
9997 ∅
9998 ●─→ 200-751-9998
9999 ∅
```

# Hash functions

- A hash function h maps keys of a given type to integers in a fixed interval [0, N − 1]
- Example:
  h(x) = x mod N is a hash function for integer keys
  The integer h(x) is called the hash value of key x
- A hash function is usually specified as the composition of two functions:
  - Hash code map:
    h1:keys → integers
  - Compression map:
    h2: integers → [0, N − 1]

9

# Hash Code Maps

- Interger cast
  - Bits of the key are interpreted as integer
  - Suitable for keys of length shorter than the number of bits of an integer type
  - Example:
    - 'A' -> 65
    - 'N' ->78

- Component Sum
  - We partition the bits of the key into components of fixed length (e.g., 16 or 32 bits) and we sum the components
  - Suitable for numeric keys of fixed length greater than or equal to the number of bits of the integer type

$$x = \left( x_1, x_2, \ldots, x_{n-1} \right) \Rightarrow h_1(x) = \sum_{i=0}^{n-1} x_i$$

$\underbrace{\phantom{xx}}_{32 \text{ bits}} \underbrace{\phantom{xx}}_{32 \text{ bits}} \underbrace{\phantom{xx}}_{32 \text{ bits}}$

10

# Hash code Maps

- Polynomial accumulation
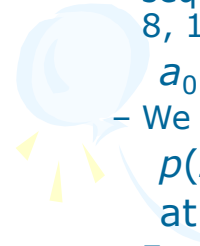  - We partition the bits of the key into a sequence of components of fixed length (e.g., 8, 16 or 32 bits)

    $a_0\ a_1\ \dots\ a_{n-1}$
  - We evaluate the polynomial

    $p(z) = a_0 + a_1\ z + a_2\ z^2 + \dots + a_{n-1}z^{n-1}$

    at a fixed value $z$, ignoring overflows
  - Especially suitable for strings (e.g., the choice $z = 33$ gives at most 6 collisions on a set of 50,000 English words)
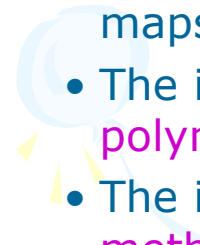
11

11

# Exercise 14.1

- Write three function which implements three type of hash code maps above.
- The input key for integer cast and polynomial is a **string**
- The input key for component sum method is a number of type **long**.

12

12

# Compression Map

- The result of the Hash Code Map needs to be reduced to a value in [0,N-1]
- Division Method:
  - $h_2(y) = |y|$ mod N
  - The size N of the hash table is usually chosen to be a prime

- Multiply, Add and Divide (MAD):
  - h2(y) = |ay + b| mod N
  - a and b are nonnegative integers such that a mod N ≠ 0
  - Otherwise, every integer would map to the same value b

13

13

# Simple implementation of Hash Table

```
#define MAX_CHAR  10
#define TABLE_SIZE  13
typedef struct {
    char key[MAX_CHAR];
    /* other fields */
} element;
element hash_table[TABLE_SIZE];
```

14

14

7

# Hash Algorithm via Division

```
void init_table(element ht[])
{
  int i;
  for (i=0; i<TABLE_SIZE; i++)
    ht[i].key[0]=NULL;
}

int transform(char *key)
{
  int number=0;
  while (*key) number += *key++;
  return number;
}
```

```
int hash(char *key)
{
  return (transform(key)
          % TABLE_SIZE);
}
```

15

# **Conflict Resolution**

- Collisions - occur when $k_1 \neq k_2$ but $h(k_1) = h(k_2)$
- Results in more complex *insertItem*() and *findElement*() operations
- Conflict Resolution Strategies
  - Closed Addressing (Open Hash Table) - i.e. slots other than $h(k)$ are "closed" and can not be used
  - Open Addressing (Closed Hash Table)- look for another open position in the table
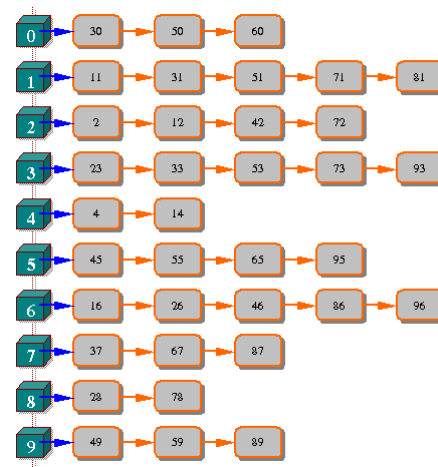


16

# Data structure for Hash Table

- Open Hash Table:
  - Chaining Method

- Closed Hash Table
  - Linear Probing
  - Quadratic Probing
  - Rehashing

17

17

# Data structure for chaining

- Array of pointers
- Each pointer manage a linked list corresponding to a bucket (address).
- This example shows a chaining hash table with hash function
  *N mod 10*



18

18

# Exercise 14.1

- Implement an ADT for chaining hash table providing the following operations:
  - Init
  - Hash function
  - Insert (given key and element)
  - Search, Delete (given key)
  - IsEmpty
  - Clear
  - Traverse

19

19

# Exercise 14-2 Make a hash list

- You assume to make an address book of mobile phone.
- You declare a structure which can hold at least "name," "telephone number," and "e-mail address", and make a program which can manage about 100 these data.
- (1) Read about 10 from an input file, and store them in a hash table which has an "e-mail address" as a key. Then confirm that the hash table is made. In this exercise, the hash function may always return the same value.
- (2)Define the hash function properly, and make the congestion occur as rare as possible

30

30

# Linear Probing
## (linear open addressing)

- Compute f(x) for identifier x
- Examine the buckets
  ht[(f(x)+j)%TABLE_SIZE]
  $0 \leq j \leq$ TABLE_SIZE
  - The bucket contains x.
  - The bucket contains the empty string
  - The bucket contains a nonempty string other than x
  - Return to ht[f(x)]

31

# Linear Probing - example

| | |
|---|---|
| 0 | 49** |
| 1 | 58** |
| 2 | 69** |
| 3 | |
| 4 | |
| 5 | |
| 6 | |
| 7 | |
| 8 | 18 |
| 9 | 89 |

With linear probing *f(i) = i*.

Here is a hash table of size *T = 10*, where the entries 89, 18, 49, 58, and 69 have been inserted. The hash function is *h(key) = key%10*.

Throughout this talk we use a table size T = 10, although in practice it should be prime.

# Exercise 14.3
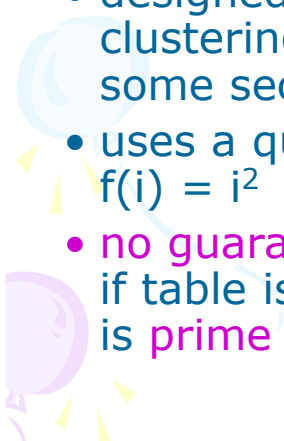
- Implement an ADT Hash Table with linear probing method.

# Quadratic Probing

- Linear probing tends to cluster
  - Slows searches
- designed to eliminate the primary clustering problem of linear (but some secondary clustering)
- uses a quadratic collision function i.e. $f(i) = i^2$
- no guarantee of finding an empty cell if table is > half full unless table size is prime

# Quadratic probing

- Linear probing tends to have clustering problem

➔ Use a quadratic function to calculate the $i^{th}$ probe position for a key $k$:

$$p(k, i) = (h(k) + i^2) \bmod N$$

where

- $N$: array size, better to be prime number
- $h(k)$: hash function for key $k$
- Not guaranteed to succeed when map is half full
- Must use lazy deletion



insert(14)   insert(8)   insert(21)   insert(2)   insert(7)
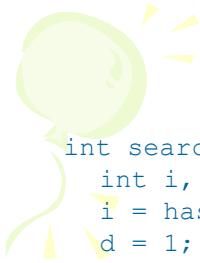14%7 = 0   8%7 = 1   21%7 =0   2%7 = 2   7%7 = 0

41

41

# Exercise 14.4

- Implement an ADT Hash Table with quadratic probing method.
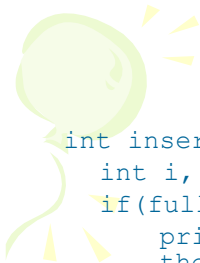
42

42

# Search

```
int search(int k) {
    int i, d;
    i = hashfunc(k);
    d = 1;
    while(hashtable[i].key!=k && hashtable[i].key
    !=NULLKEY){
    //Quadratic probing
    i = (i+d*d) % M;
    d = d+1;
    }
    if(hashtable[i].key==k) // found
        return i;
    else // not found
    return M;
}
```

43

# Insert

```
int insert(int k){
    int i, d;
    if(full()){
        printf("\n Hash table is full. Can not insert
        the   key %d ",k);
        return -1;  // <===
    }
    i=hashfunc(k); d = 1;
    while(hashtable[i].key !=NULLKEY){
        //Quadratic probing
        i = (i+d*d) % M;
        d = d+1;
    }
    hashtable[i].key=k;
    N=N+1;
    return i;
}
```

44

# Double Hashing

- Double hashing uses a secondary hash function $h_2(k)$ and handles collisions by placing an item in the first available cell of the series

  $(i + h_2(k))\ mod\ N$
- The secondary hash function $\boldsymbol{h_2(k)}$ cannot have zero values
- The table size $\boldsymbol{N}$ must be a prime to allow probing of all the cells

- Common choice of compression map for the secondary hash
- function: $h_2(k) = q - k\ mod\ q$
- where
  - $q < N$
  - $q$ is a prime

45

# Double hashing

- Use a second hash function to resolve hash collisions:

  $$p(k, i) = \big(h_1(k) + i \times h_2(k)\big) \bmod N$$

  where

  - $h_2(k)$ should never return 0

Lets say, Hash1 (key) = key % 13

Hash2 (key) = 7 – (key % 7)

Hash1(19) = 19 % 13 = 6

Hash1(27) = 27 % 13 = 1

Hash1(36) = 36 % 13 = 10

Hash1(10) = 10 % 13 = 10

Hash2(10) = 7 – (10%7) = 4

(Hash1(10) + 1*Hash2(10))%13= 1

(Hash1(10) + 2*Hash2(10))%13= 5

Collision

46

# Exercise 14.5

- Implement an ADT Hash Table with rehashing method, using two following hash functions:

- **f1(key)= key % M**
- **f2(key) =(M-2)-key %(M-2)**

# Hash functions

```
int hashfunc(int key)
{
   return(key%M);
}
//Secondary function
int hashfunc2(int key)
{
   return(M-2 – key%(M-2));
}
```