



CX1107 Data Structures and Algorithms 9: Searching

2020/21 Semester 2

School of Computer Science and Engineering

Nanyang Technological University

9.1 Searching

Searching problem is one of the classic computational problem to find a search key in a given data set.

Algorithm 1 Generic Search Algorithm

```

1: begin
2: while there is more possible data to examine do
3:   examine one datum
4:   if the datum = search key then
5:     return succeed
6:   return fail
7: end

```

9.2 Sequential Search

When the given data set is unsorted, we have to check every single element in the data set until either the key is found (successful search) or all elements in the data set have been retrieved once and no match is found (unsuccessful search). The search is called sequential search. It is a brute-force algorithm.

Algorithm 2 Sequential Search

```

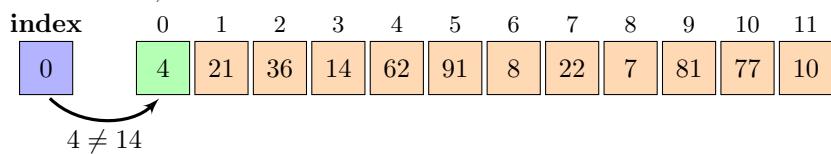
1: function seqSearch(int[] Data, int n, int key)
2: begin
3:   for index = 0 to n - 1 do
4:     begin
5:       if Data[index] == key then
6:         return index;                                {success}
7:       end
8:     return -1;                                  {failure}
9:   end

```

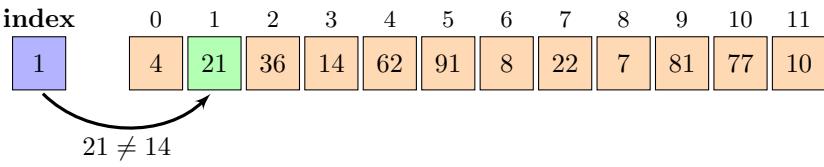
9.2.1 Example: Success Case

Search key: 14

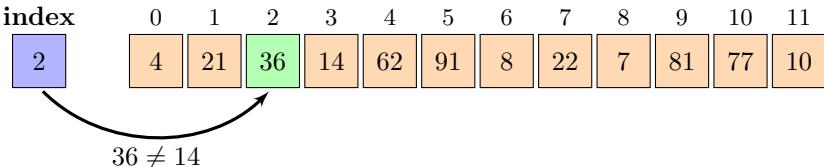
Let index = 0,



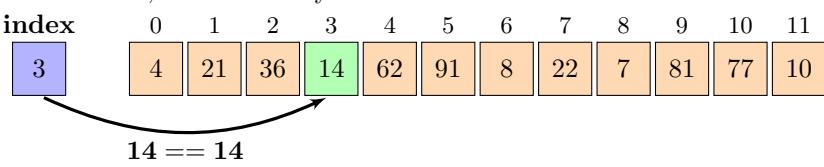
Let index =1,



Let index =2,



Let index =3, the search key is found

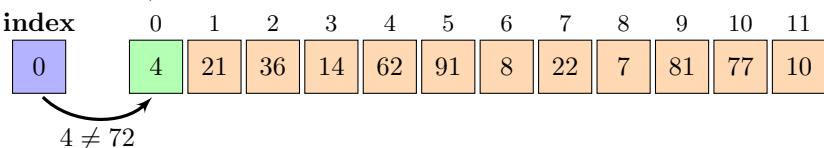


The search key, 14 is found at index =3. **3** return.

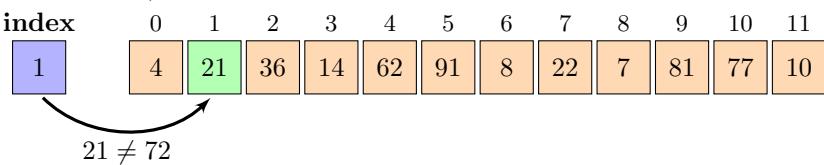
9.2.2 Example: Failure Case

Search key: **72**

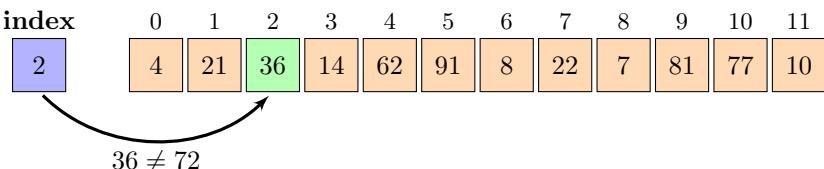
Let index =0,



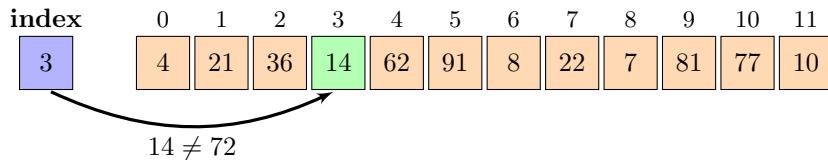
Let index =1,



Let index =2,



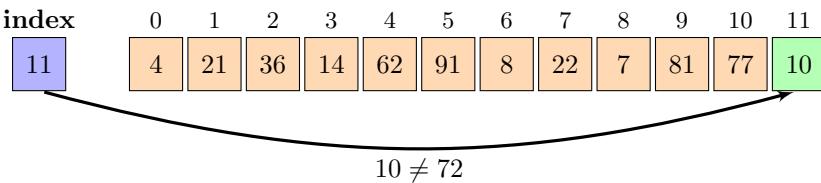
Let index =3,



⋮

⋮

Let index = 11,



After searching throughout the given data set, there is no match with the search key, **72**. It is an unsuccessful search. **-1 return.**

9.2.3 Complexity Analysis

- Best-case analysis: 1 comparison against key (**the first item is the search key**)
- Worst-case analysis: n comparisons against key (Either the last item or no item is the search key)
- Average-case Analysis:
- **Key is always in search array:**
 - Let e_i represents the event that the key appears in i^{th} position of array, its probability $P(e_i) = \frac{1}{n}$
 - $T(e_i)$ is the number of comparisons done
 - $0 \leq i \leq n$ and $T(e_i) = i + 1$
 - Since we assume that key definitely is in the array, the average-case analysis can be done as follow
$$\begin{aligned} A_s(n) &= \sum_{i=0}^{n-1} P(e_i)T(e_i) = \sum_{i=0}^{n-1} \frac{1}{n} \times (i+1) \\ &= \sum_{i=0}^{n-1} \left(\frac{1}{n}\right)(i+1) \\ &= \frac{1}{n} \sum_{i=1}^n i \\ &= \frac{1}{n} \frac{n(n+1)}{2} = \frac{n+1}{2} \end{aligned}$$

- **Key is not always in search array:**

- When the search key is not in the array, the number of comparisons, $A_f(n)$, is n . Refer to 9.2.2.
- Combine success cases and failure cases with their probability, $P(succ) + P(fail) = 1$:

$$P(succ)A_s(n) + P(fail)A_s(n) = q * (n + 1)/2 + (1 - q) * n$$

- If there is a 50-50 chance that key is not in the array, $P(succ) = P(fail) = 0.5$.

- The average no of key comparisons is $\frac{3n}{4} + \frac{1}{4}$ or about $\frac{3}{4}n$ of entries are examined

In conclusion, both worst-case and average-case complexity are $\Theta(n)$.

9.3 Binary Search

Clearly Search an item from an unsorted array will take $\Theta(n)$ in average. To improve the searching performance, the data set need to be sorted out first. It will be discussed in the later lecture. Next, we assume that the array is sorted in order. We can use the information of its order to reduce the search work. Binary search is an example which is using **decrease-and-conquer** approach. In short, this approach divide a problem into two smaller sub-problems, one of which does not even have to be solved.

Binary search first comparing a search key with the sorted array's middle element. If they match, then the algorithm stops; otherwise, the same operation is repeated recursively for either the upper half of the array if the search key is lesser than the middle element or the lower half of the array if the search key is greater.

Binary Search

```

1 int binarySearch (int E[], int first, int last, int k) --> T(n)
2 {
3     if(last < first)
4         return -1;
5     else {
6         int mid = (first + last)/2;
7         if(k == E[mid])
8             return mid;
9         else if(k < E[mid])
10            return binarySearch(E,first, mid-1,k); --> T(n/2)
11        else
12            return binarySearch(E,mid+1, last,k); --> T(n/2)
13    }
14 }
```

Listing 1: Recursive Version

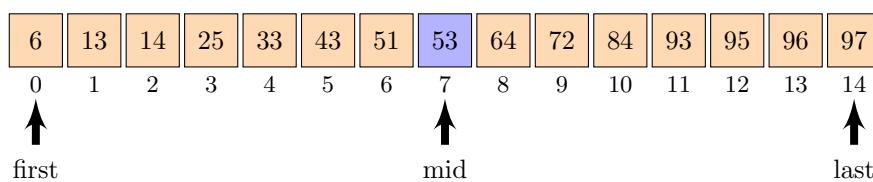
```

1 int binarySearch_iter(int E[], int first, int last, int k)
2 {
3     while (first <= last) {
4         int mid = (first+last) / 2;
5         if (E[mid] == k)
6             return k;
7         else if (k < E[mid] )
8             last = mid - 1;
9         else
10            first = mid + 1;
11    }
12    return -1;
13 }
```

Listing 2: Iterative Version

9.3.1 Example: Success Case

Given a search key, **33**, the **first** index is 0, the **last** index is 14 and the **mid** is $7 \Leftarrow \frac{0+14}{2}$



Since the search key, **33 < 53**, we change **last** index from 14 to 6 (**mid-1**)

6	13	14	25	33	43	51	53	64	72	84	93	95	96	97
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14

A diagram illustrating the relationship between 'first' and 'last'. On the left, an upward-pointing arrow originates from the word 'first'. On the right, another upward-pointing arrow originates from the word 'last'.

The new **mid** is $3 \Leftarrow \frac{0+6}{2}$.

6	13	14	25	33	43	51	53	64	72	84	93	95	96	97
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14

↑
first mid last

Since the search key, **33 > 25**, we change **first** index from 0 to 4 (**mid+1**)

6	13	14	25	33	43	51	53	64	72	84	93	95	96	97
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14

↑ ↑
first last

The new **mid** is 5 $\leftarrow \frac{4+6}{2}$.

6	13	14	25	33	43	51	53	64	72	84	93	95	96	97
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14

first mid last

Since the search key $33 < 43$, we change **last** index from 6 to 4 ($\text{mid}-1$)

6	13	14	25	33	43	51	53	64	72	84	93	95	96	97
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14

```

graph TD
    4[4] --> 2[2]
    4 --> 6[6]
    2 --> 1[1]
    2 --> 3[3]
    6 --> 5[5]
    6 --> 7[7]
    style 1 fill:#fff,stroke:#000
    style 2 fill:#fff,stroke:#000
    style 3 fill:#fff,stroke:#000
    style 4 fill:#fff,stroke:#000
    style 5 fill:#fff,stroke:#000
    style 6 fill:#fff,stroke:#000
    style 7 fill:#fff,stroke:#000

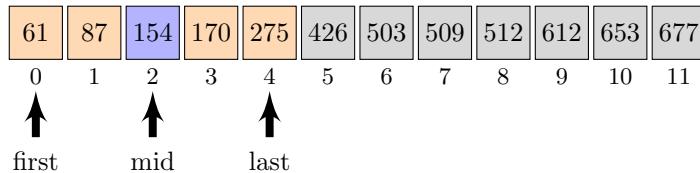
```

first

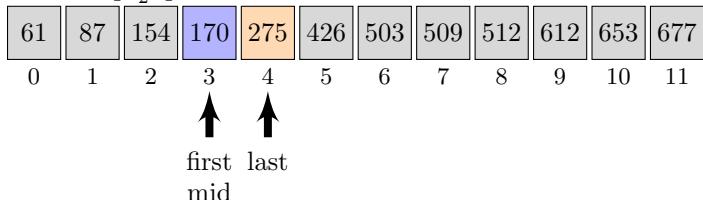
9.3.2 Example: Failure Case

Given a search key, **400**, the **first** index is 0, the **last** index is 11 and the **mid** is $5 \leftarrow \frac{0+11}{2}$

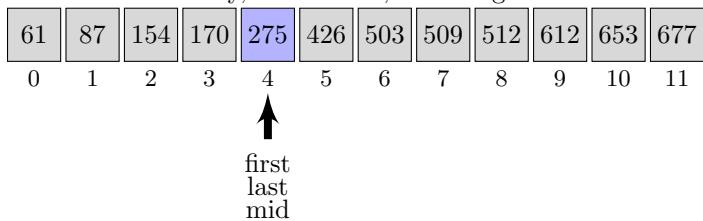
Since the search key, $400 < 426$, we change **last** index from 11 to 4 (**mid-1**)



Since the search key, 400 > 154, we change **first** index from 0 to 3 (**mid+1**). The new **mid** remain unchanged $\lfloor \frac{3+4}{2} \rfloor = 3$



Since the search key, 400 > 170, we change **first** index from 3 to 4 (**mid+1**).



Since the search key, 400 > 275, we change **first** index from 4 to 5 (**mid+1**). Since the **first** index larger than **last**, the search function returns -1 showing that searching fails. 400 does not exist in the given array.

9.3.3 Complexity Analysis

- Best-case analysis: 1 comparison against key (the first item is the search key)
- Worst-case analysis: Refer to Listing 1, each function call will define the new **mid** as $\lfloor \frac{\text{first}+\text{last}}{2} \rfloor$ if $\text{last} \geq \text{first}$. The running time is constant, c . If the **mid** element and search key do not match, there will be a recursive call by passing either right or left section. If the size of array, n is even, there are $\frac{n}{2} - 1$ entries in the left section and $\frac{n}{2}$ in the right section. If n is odd, there are $\frac{n-1}{2}$ entries in both sections. Refer to 9.3.1 and 9.3.2. The worst case is $\frac{n}{2}$. When n is 1, **first** and **last** are the same index. Hence, the new **mid** will be the index. If this only one entry is not the search key, it will make a recursive call. This is the last recursive call which will return -1 (no match). How many recursive call will be made until n is 1? It is noted that when n is 1, we simply let $T(1) = c$ (constant time).

$$\begin{aligned}
 T(n) &= T\left(\frac{n}{2}\right) + c \\
 &= (T\left(\frac{n}{4}\right) + c) + c \\
 &= (T\left(\frac{n}{8}\right) + c) + c + c \\
 &\quad \dots \\
 &= T(1) + d * c \\
 &= (d + 1)c
 \end{aligned}$$

recursive calling

From the equation above, d is the number of recursive call. We divide n by 2 every call. At the d^{th} recursive call, we have

$$\begin{aligned}\frac{n}{2^d} &= 1 \\ n &= 2^d \\ d &= \log_2 n\end{aligned}$$

We can observe that the number of iteration, d is less than $\log_2 n$ due to the number of entries can be even or odd number.

$$\begin{aligned}d &\leq \log_2 n \\ &= \lfloor \log_2 n \rfloor\end{aligned}$$

Prove: $T(n) = \lfloor \log_2 n \rfloor + 1 = \lceil \log_2(n+1) \rceil \quad \forall n \geq 1 \wedge n \in \mathbb{N}$

Let k be some integer number, $k \in \mathbb{N}$ that $k = \lfloor \log_2 n \rfloor$, then

$$\begin{aligned}k &\leq \log_2 n < k + 1 \\ 2^k &\leq n < 2^{k+1} \\ 2^k &< n + 1 \leq 2^{k+1}\end{aligned}$$

For example, if $k = 4$, then we have $16 \leq n < 32$. How about the range of $(n+1)$? it is $16 < n+1 \leq 32$ because n is an integer number.

$$\begin{aligned}2^k &< n + 1 \leq 2^{k+1} \\ k &< \log_2(n+1) \leq k + 1\end{aligned}$$

The last inequality implies that $\lceil \log_2(n+1) \rceil = k + 1$ and we define that $k = \lfloor \log_2 n \rfloor$,

$$\lceil \log_2(n+1) \rceil = \lfloor \log_2 n \rfloor + 1$$

$$\therefore T(n) = c(\lfloor \log_2 n \rfloor + 1) = c\lceil \log_2(n+1) \rceil.$$

The worst-case running time is in $\Theta(\log_2 n)$

- Average-case analysis:

To analysis the average case, we first consider two scenario, successful search and failed search. By Law of Expectations:

$$A_q(n) = q\underline{A_s(n)} + (1-q)\underline{A_f(n)}$$

where q is the probability. Here we consider $n = 2^k - 1$ for simplicity but other values is very close to the following result.

The failed search is the worst case. Thus, its complexity is:

$$A_f(n) = \lceil \log_2(n+1) \rceil = \log_2(n+1) = k$$

The successful search of $n = 2^k - 1$ entries may take from 1 to k comparisons. It depends on the position of the search key. $1, 2, 4, 8, \dots, 2^{k-1}$ positions require to take $1, 2, 3, 4, \dots, k$ comparisons respectively.

复习的时候自己做一个 analysis

For example, $n = 2^3 - 1 = 7$ entries



1 comparison



In this example, we can observe that if the search key is at the 4th position, we just need one comparison. If the search key is at the 2nd or the 6th position, we need two comparisons etc.

We assume that the probability of each position being searched is equal, i.e. $\frac{1}{n}$. The complexity of successful search is:

$$\begin{aligned}
 A_s(n) &= \frac{1}{n} \sum_{t=1}^k t2^{t-1} \quad \text{(?)} \\
 &= \frac{(k-1)2^k + 1}{n} \\
 &= \frac{[\log_2(n+1) - 1](n+1) + 1}{n} \\
 &= \log_2(n+1) - 1 + \frac{\log_2(n+1)}{n}
 \end{aligned}$$

第一次执行
 有这个这样位置.

The derivation can be found in [9.5]

The average-case time complexity can be obtained by substituting $A_s(n)$ and $A_f(n)$ into $A_q(n)$:

$$\begin{aligned}
 A_q(n) &= qA_s(n) + (1-q)A_f(n) \\
 &= q[\log_2(n+1) - 1 + \frac{\log_2(n+1)}{n}] + (1-q)(\log_2(n+1)) \\
 &= \log_2(n+1) - q + q\frac{\log_2(n+1)}{n} \\
 &= \Theta(\log_2(n))
 \end{aligned}$$

q is the probability which is always ≤ 1 and $\frac{\log_2(n+1)}{n}$ is negligible when n is large.

Therefore, binary search does approximately $\log_2(n+1)$ comparisons on average for n entries.

9.3.4 Summary of Binary Search

- Decrease-and-conquer strategy
- Examine data item in the middle and search recursively on single half
- Average and worst time complexities are both $\Theta(\log(n))$

9.4 Hashing

It is a typical space-and-time trade-off strategy in algorithm. To achieve search time in $\mathcal{O}(1)$, memory usage have to be increased. Here we will discuss the following hashing approach:

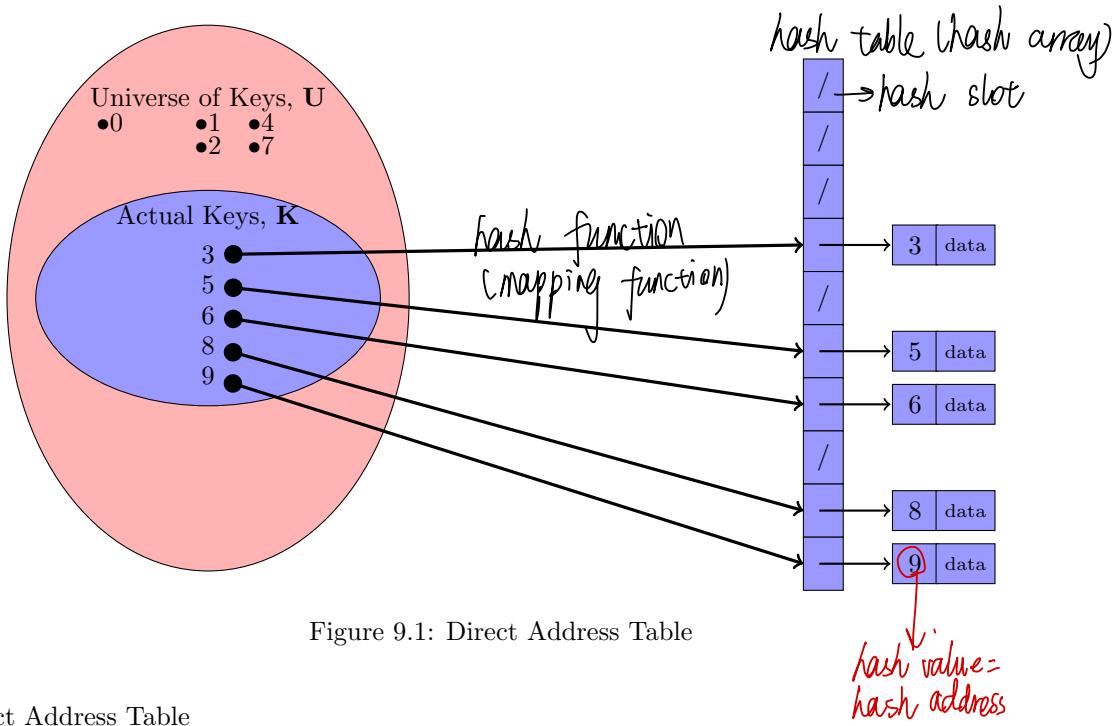


Figure 9.1: Direct Address Table

- Direct Address Table
- Closed Address Hashing
- Open Address Hashing
 - Linear Probing
 - Quadratic Probing
 - Double Hashing

9.4.1 Direct Address Table *(Space - costly)*

Suppose that the set of actual keys is $K \subseteq \{0, 1, 2, \dots, m-1\}$ and keys are distinct. We can define an array, $T[0 \dots m-1]$ or **direct-address table**:

$$T[x.k] = \begin{cases} x, & \text{if } k \in K \wedge x.k = k ; \\ NIL, & \text{otherwise.} \end{cases}$$

Thus, operations take $\mathcal{O}(1)$ time.

The direct address tables are impractical when the range of keys can be very large ($m \gg |K|$, e.g. 64-bit numbers (range of $m \approx 18.45 \times 10^{18}$). To overcome the large range issue, we can use a **hash function**, $h(key)$ to map the universe, U of all keys into hash table slots, $\{0, 1, 2, \dots, m-1\}$. However, **multiple keys may be mapped to the same slots. It is known as collision.** See Figure 9.2. Motivation of Hashing is to be able assign a unique array index to every possible key that could occur in an application.

- Key space may be too large for an array on the computer while only a small fraction of the key values will appear.
- The purpose of hashing is to translate an extremely large key space into a reasonably small set of integers.

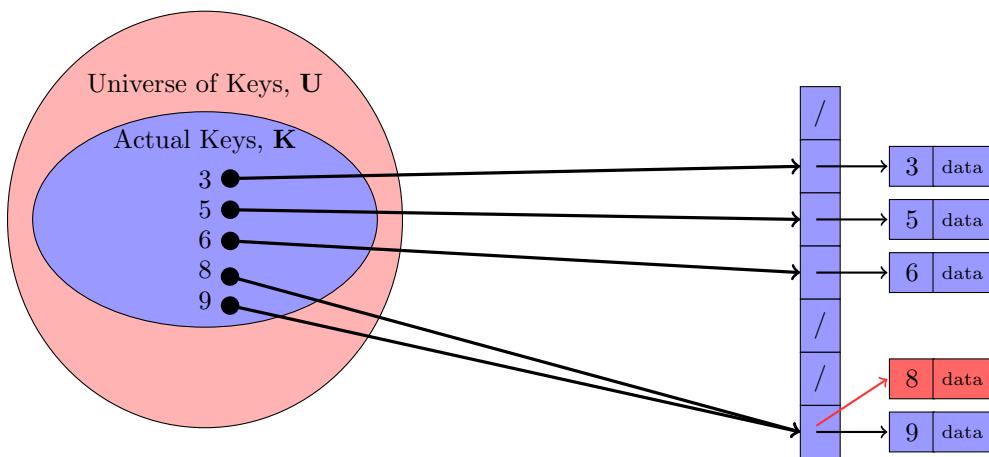


Figure 9.2: Collision Issue in Hash Table

- A hash function f : key space \rightarrow hash codes.

Example: A hash table of 200 entries.

A possible hash function is

$$hash(k) = k \bmod 200$$

When multiple keys are mapped to the same hash code, a collision occurs. e.g. $k = 200$ and $k = 400$ are mapped to $\text{hash}(k) = 0$.

9.4.2 What hash function to use?

- A hash function **MUST** return a value within the hash table range.
 - It should achieve an even distribution of the keys that actually occur across the range of indices.
 - It should be easy and quick to compute.
[Search time should remain $O(1)$]

Here we introduce three types of hash functions:

1. The Division Method: *to even distribute key*

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

- The return value is the last four bits of k , $f(k) = k \bmod 16$
 - Avoid to use power of 10 for decimal numbers as keys
原因：10的倍数进位点
 - The best table size is often a prime number not too close to exact powers of 2 for “real” data.
 - Real data may not always evenly distribute
原因：真实生活中数据不均匀分布。

2. The Folding Method:

- All bits contribute to the result .
- Partition the key into several parts and combine the parts in a convenient way (e.g. addition or multiplication).
- Example 1: Sum the key value, take modulus of the sum. The sum must be large enough compared to the quotient, M .

```

1 int h(char x[10])
2 { int i, sum;
3   for (sum=0, i=0; i < 10; i++)
4     sum += (int) x[i];
5   return(sum % M);
6 }
```

$$\text{Eq. } k = 3121 \quad k^2 = 9740641 \Rightarrow H(k) = 406$$

- Example 2: Mid-square method: Square the key value, take the middle r bits (from the result) for a hash table of 2^r slots

3. Multiplicative Congruential method: (pseudo-number generator)

Step 1: Choose the hash table size, h .

Step 2: Choose the multiplier, a .

$$a = 8 \lfloor \frac{h}{23} \rfloor + 5$$

Step 3: Define the hash function, f

$$f(k) = (a * k) \bmod h$$

```

1 >> h = 31
2 >> a= 8*floor(h/23) +5
3
4 a =
5
6     13
7
8 >> k = 1:15
9
10 k =
11
12     1     2     3     4     5     6     7     8     9     10    11    12    13    14    15
13
14 >> fk = mod((a*k),h)
15
16 fk =
17
18     13     26      8     21      3     16     29      11     24      6     19      1     14     27      9
19
```

9.4.3 How to handle collisions?

To handle collision issue in hash table, we can use

- Closed Address Hashing
- Open Address Hashing

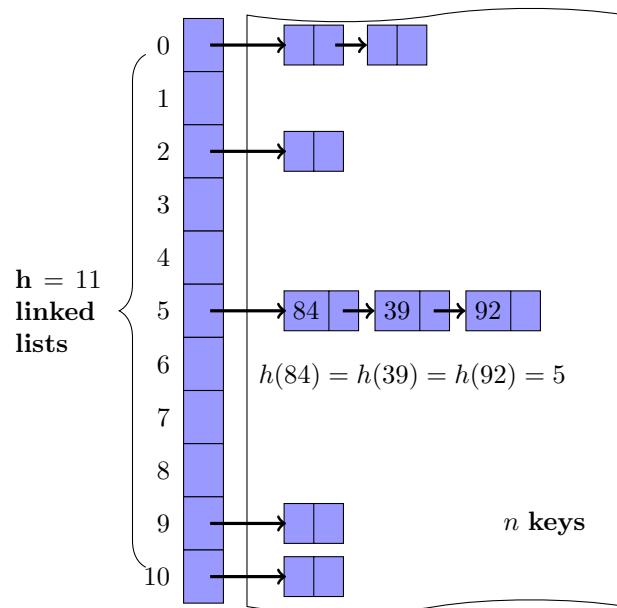


Figure 9.3: Closed Addressing Hash Table

9.4.4 Closed Address Hashing *(位置不变, 存取的东西变成 list)*

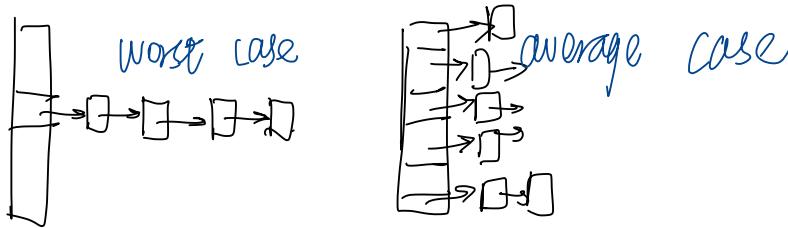
- Maintains the original hashed address
- Records hashed to the same slot are linked into a list
- The address is closed (fixed). Each key has a corresponding fixed address
- Also called **chained hashing**
 - Initially, all entries in the hash table are empty lists.
 - All elements with hash address i will be inserted into the linked list $H[i]$.
 - If there are n records to store in the hash table, then $\alpha = \frac{n}{h}$ is the **load factor** of the hash table. *(on average; how much records per slot; to see how many slots are occupied)*
 - In closed address hashing, there will be α number of elements in each linked list on average.
 - During searching, the searched element with hash address i is compared with elements in linked list $H[i]$ sequentially

9.4.4.1 Analysis of Chained Hashing

The **worst case** behaviour of hashing happens when all elements are hashed to the same slot. In this case

- The linked list contains all n elements
- An unsuccessful search will do n key comparisons

*fixing hash address \Rightarrow
Do linear search*



- A successful search, assuming the probability of searching for each item is $1/n$, will take

$$\frac{1}{n} \sum_{i=1}^n i = \frac{n+1}{2} = O(n)$$

- In this analysis, checking if link list is NULL is not counted as a key comparison

The **average case**: If we assume that any given item is equally likely to hash into any of the h slots, an **unsuccessful search on average does n/h key comparisons.**

$$(\alpha = \frac{n}{h})$$

- An unsuccessful search means searching to the end of the list.
- The number of comparisons is equal to the length of the list.
- The average length of all lists is the load factor, α
- Thus the **expected number of comparisons in an unsuccessful search is α .**

(Search \propto times)

Assume that any given item is equally likely to hash into any of the h slots, a **successful search on average does $\Theta((1 + \alpha))$ comparisons.**

insert \downarrow \downarrow $\alpha \times |$

- We assume that the items are inserted at the end of the list in each slot.
- The expected no. of comparisons in a successful search is 1 more than the no. of comparisons done when the sought after item was inserted into the hash table.
- When the i^{th} item is inserted into the hash table, the average length of all lists is $(i - 1)/h$. So when the i^{th} item is sought for, the no. of comparisons is

当作 $i-1$ 次
Unsuccessful search $\left(1 + \frac{i-1}{h}\right)$
 i item 需要 insert

找第 i 个 item 看作前一个 item 中
的第 i 个. 前面已有 $i-1$ 个.

- So the average number of comparisons over n items is:

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \left(1 + \frac{i-1}{h}\right) &= \frac{1}{n} \sum_{i=1}^n (1) + \frac{1}{nh} \sum_{i=1}^n (i-1) \\ &= 1 + \frac{1}{nh} \sum_{i=0}^{n-1} i \\ &= 1 + \frac{n-1}{2h} \end{aligned}$$

Load 越高，每次搜索时间越长
但每个 item 的搜索时间不变

- Therefore, a successful search on average does $\Theta(1 + \alpha)$ comparisons
- If n is proportional to h , i.e. $n = \mathcal{O}(h)$, then $n/h = \mathcal{O}(h)/h = \mathcal{O}(1)$.
- Thus, each successful search with chained hashing takes **constant time** averagely

9.4.5 Open Address Hashing

- To store all elements in the hash table
- The load factor $\alpha = n/h$ is never greater than 1 ($\frac{1}{h}$ slots $\neq \frac{n}{h} - 1$)
- When collision occurs, probe is required for an alternative slot
- The address is open (not fixed)
 - Linear Probing
 - Quadratic Probing
 - Double Hashing

9.4.5.1 Linear Probing

- Suppose the hash function, $H'(k) = k \bmod h$
 - If $H'(k)$ is an empty slot, key k will be hashed to $H'(k)$.
 - Otherwise, we need to probe the next empty slot by taking $(H'(k) + 1) \bmod h$.
 - Repeating the probing function,
- $$H(k, i) = (k + i) \bmod h$$
- for $i = 1, 2, \dots, h - 1$ until an empty slot is found
- Consider the linear probing policy for storing the following keys:

1055, 1492, 1776, 1812, 1918, 1942

The linear probing hash function, $H(k, i) = (k + i) \bmod 10$

~~Find slot~~

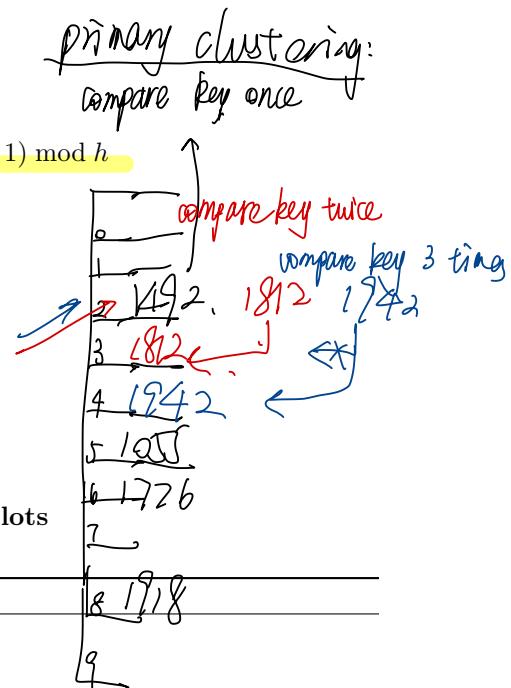
9.4.5.2 Searching A Key In A Hash Table With Probing New Slots

Algorithm 3 Searching a Key

```

1: function Search ( $k$ )
2: begin
3:    $code \leftarrow \text{hash}(k)$ 
4:    $loc \leftarrow code$ 
5:    $ans \leftarrow \emptyset$ 
6:   while  $H[loc] \neq \emptyset$  do
7:     if  $H[loc].key == k$  then
8:       begin
9:          $ans \leftarrow H[loc]$ 
10:        break
11:      end
12:    else
13:      begin
14:         $loc \leftarrow \text{probe}(loc)$ 
15:        if  $loc == code$  then
16:          break
17:        end
18:      end

```



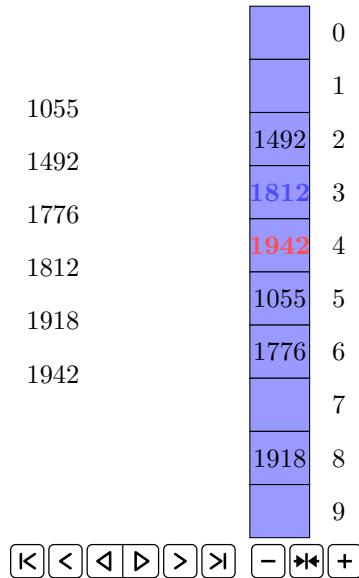


Figure 9.4: Open Addressing via Linear Probe

Three outcomes of searching:

1. key found at $H[loc]$ – Success
2. position empty – Fail
3. all slots are searched

9.4.5.3 Limitations of Linear Probe

Linear probing is a simple probing new slot method. Unfortunately, it has some limitations. As we can imagine, in open address hashing, searching becomes expensive when the load factor α approaches to 1. For linear probing, the problem can happen earlier, if keys are hashed to nearby places in the hash table. We call this phenomenon as **primary clustering**, which is characterized by long runs of occupied slots. The search time and insertion time will be increased.

9.4.5.4 Quadratic Probing

Instead of linearly probing the next empty slot in the hash table, we can probe the slot quadratically by the following probe function:

$$H(k, i) = (k + c_1 i + c_2 i^2) \bmod h$$

However, c_1 , c_2 and h need to be carefully selected to ensure that all hash slots can be probed. For $h = 2^n$, a good choice for the constants are $c_1 = c_2 = \frac{1}{2}$.

Although quadratic probing can avoid primary clustering, it still face another phenomenon known as **secondary clustering**. If two keys have the same initial probe position, their probe sequences will be the same. It implies that the later inserted keys will take longer search time.

(*Modulus same*)

slot empty \Rightarrow no key comparison
 key comparison: compare the key in the slot (not slot).

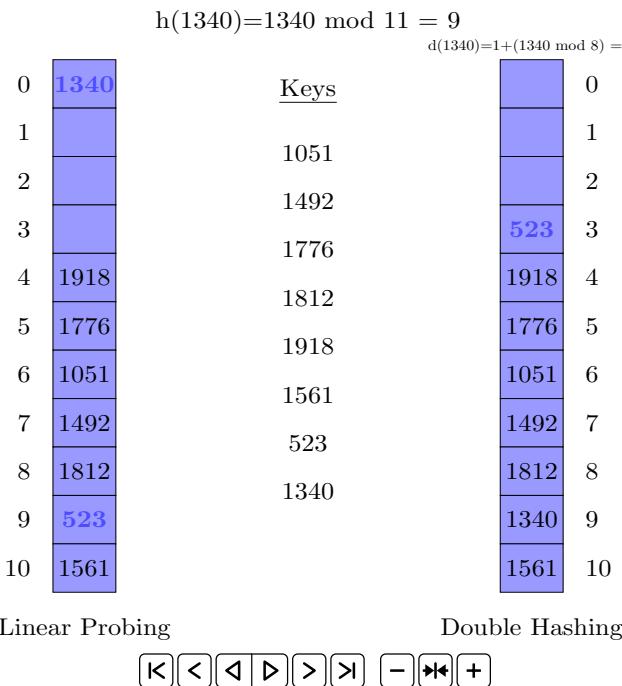


Figure 9.5: Open Addressing via Linear Probe and Double Hashing

9.4.5.5 Double Hashing

$key = search item$

It is another better way to alleviate the collision issue. Its probing to the new slot is a more random method.
 Suppose the hash function, $H(k, i)$

$$H(k, i) = (k + iD(k)) \bmod h$$

where $D(k)$ be another hash function and $i \in [0, h - 1]$.

If $H(k, 0)$ slot is occupied, a probe method will be applied to find the new slot.

The hash table size, h , should be a prime number. Using a prime number makes it impossible for any number to divide it evenly, so the probe sequence will eventually check every slot.

9.4.5.6 Example: Linear Probing VS Double Hashing

Given the following keys:

1051, 1492, 1776, 1812, 1918, 1561, 523, 1340

- Linear Probing: $H(k, i) = (k + i) \bmod 11$ $i = 0, 1, 2, \dots, 10$
- Double Hashing: $H(k, i) = (k + iD(k)) \bmod 11$, where $D(k) = 1 + (k \bmod 8)$ and $i = 0, 1, 2, \dots, 10$

double hashing need to make sure to visit every slot

可以保证在较小的位数上

9.4.5.7 Time Complexity

The time complexity of open addressing methods as follow: Linear Probing:

Function of α

- Successful Search: $\frac{1}{2}(1 + \frac{1}{1-\alpha})$
- Unsuccessful Search: $\frac{1}{2}(1 + (\frac{1}{1-\alpha})^2)$

Quadratic Probing (*approximate any open addressing method which causes secondary clusters)

- Successful Search: $1 - \ln(1 - \alpha) - \frac{\alpha}{2}$
- Unsuccessful Search: $\frac{1}{1-\alpha} - \alpha - \ln(1 - \alpha)$

Double Hashing:

- Successful Search: $\frac{1}{\alpha} \ln \frac{1}{1-\alpha}$
- Unsuccessful Search: $\frac{1}{1-\alpha}$

The proof can be found in [Knuth97]. In this course, you do not need to memorize these time complexities but you need to know that they are functions of $\frac{1}{1-\alpha}$. What does it mean? Think about it.

constant number & bigger load \Rightarrow longer time

9.4.5.8 Rehashing The Hash Table

We can observe that the multiple probing is always required when most of slots are occupied (α is approaching to 1). It is unrelated to the probing methods. Therefore, we need to rehash the hash table when it reaches the load factor threshold.

9.4.5.9 Deletion A Key Under Open Addressing

Removing a key from an open-addressing hash table is a tedious task. Instead of deleting the key, we leave the key in the table but make a marker like ‘obsolete’ or ‘tombstone’ to indicate that it is deleted. It can be overwritten by a new key when it is inserted to the slot.

We may do a “garbage collection” when we have done a large number of deletions. It can improve the insertion and searching time.

9.5 Appendix

$$\begin{aligned}
 \sum_{t=1}^k t2^{t-1} &= 1 \cdot 1 + 2 \cdot 2 + 3 \cdot 4 + 4 \cdot 8 + \dots + k \cdot 2^{k-1} \\
 2 \sum_{t=1}^k t2^{t-1} &= 1 \cdot 2 + 2 \cdot 4 + 3 \cdot 8 + \dots + (k-1) \cdot 2^{k-1} + k \cdot 2^k \\
 (2-1) \sum_{t=1}^k t2^{t-1} &= -1 \cdot 1 - 1 \cdot 2 - 1 \cdot 4 - 1 \cdot 8 - \dots - 1 \cdot 2^{k-1} + k \cdot 2^k \quad \triangleright \text{eq. 2 - eq. 1} \\
 \sum_{t=1}^k t2^{t-1} &= -2^k + 1 + k \cdot 2^k \quad \triangleright \text{geometric series} \\
 &= 2^k(k-1) + 1
 \end{aligned}$$

References

- [Knuth97] KNUTH, DONALD E. "The Art of Computer Programming, Volume 3 (3rd Ed.): Sorting and Searching," *Addison Wesley Longman Publishing Co., Inc.*, 1997.