

College of Engineering

Dr. Loke Yuan Ren  
Lecturer  
yrloke@ntu.edu.sg

# Overview

- Exhaustive Algorithm: Sequential Search
- Decrease-and-conquer Algorithm: Binary Search
- Data Structures:
  - Hashing
    - Open Hashing
    - Closed Hashing

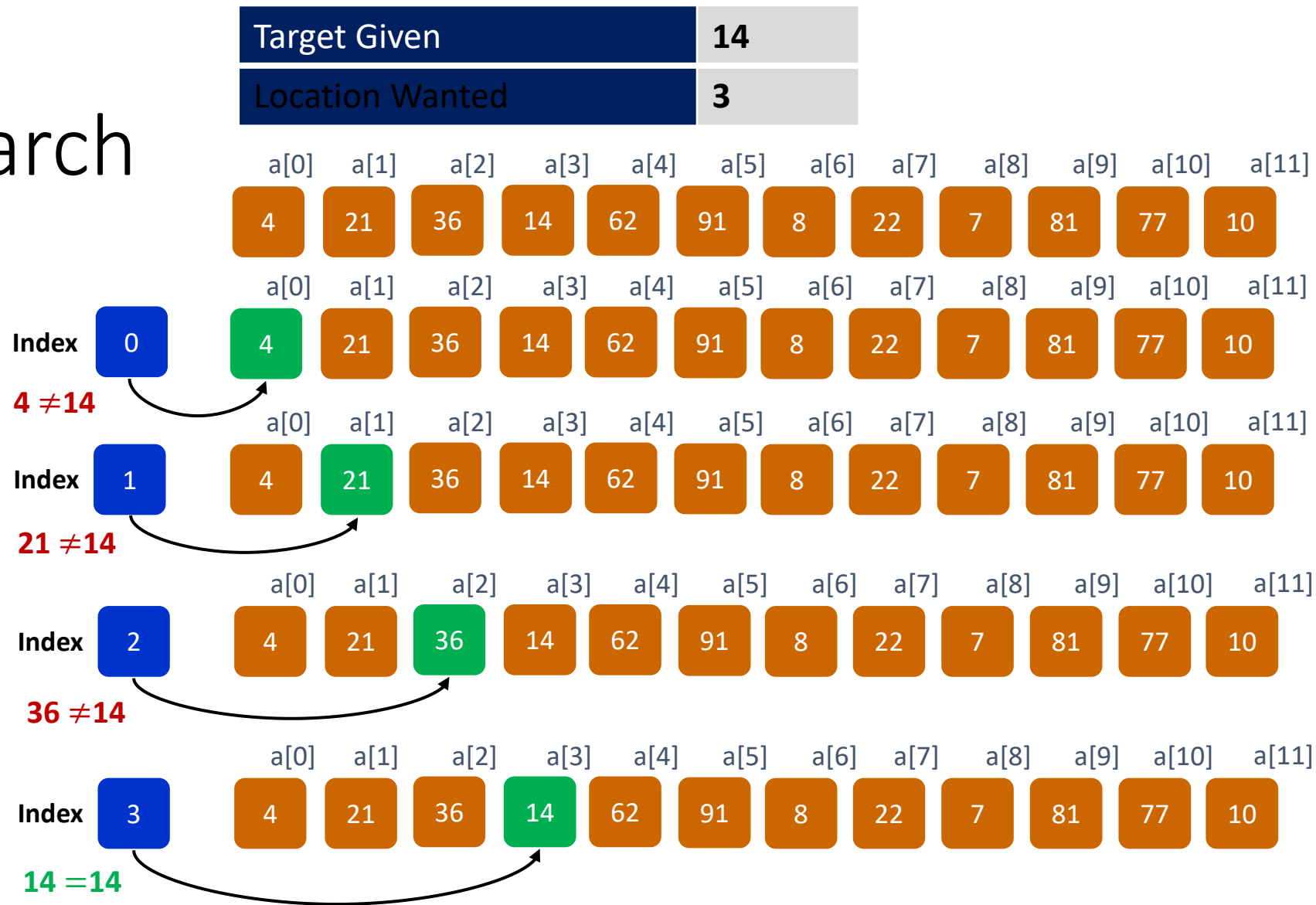
# Sequential Search

- The array is unordered
- The approach is a brute-force approach or a naïve algorithm
- Every element in the array is required to be read and compare

# Sequential Search

## 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
```



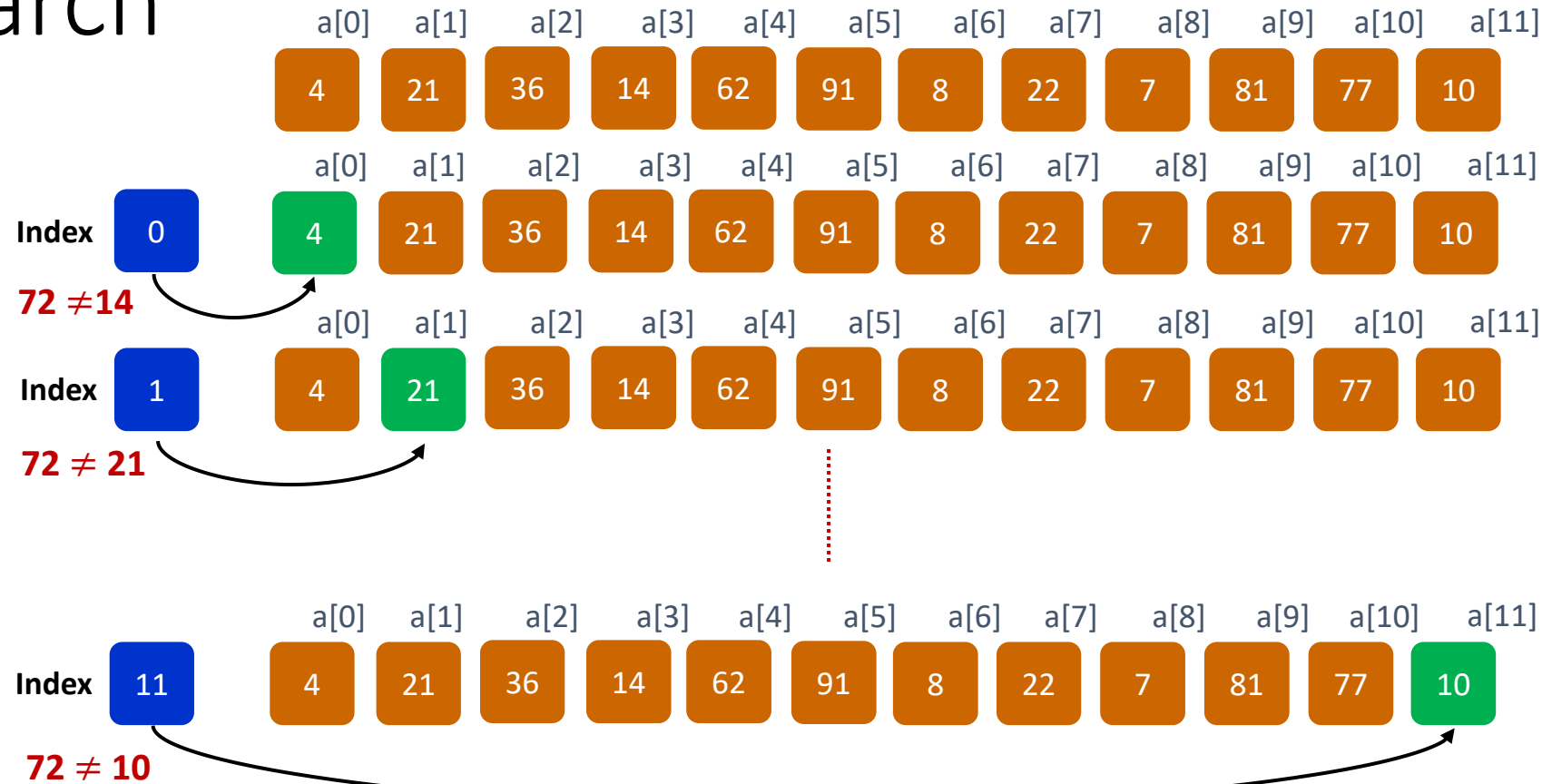
Target Given

72

# Sequential Search

## 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
```



# Time Complexity of Sequential Search

- Best-case complexity:  $\Theta(1)$ , 1 comparison against key (the first item is the search key)
- Worst-case complexity:  $\Theta(n)$ , n comparison against key (Either the last item or no item is the search key)
- Average-case complexity:

## Key is in the search array:

- $e_i$  represents the event that the key appears in  $i^{\text{th}}$  position of array, so its probability  $P(e_i) = 1/n$
- $T(e_i)$  is the no. of comparisons done
- Average complexity:  $A_s(n) = \sum_{i=1}^n \left(\frac{1}{n}\right)(i)$   
$$= \frac{1}{n} \left( \frac{n(n+1)}{2} \right) = \frac{n+1}{2} = \Theta(n)$$

## Key is not in the search array:

- n comparisons (if it is an linked list, you may need to take an extra comparison)
- $A_s(n) = n = \Theta(n)$

$$\Pr(\text{succ}) A_s(n) + \Pr(\text{fail}) A_f(n) = q \frac{n+1}{2} + (1 - q)n$$

Both **worst** and **average complexity** are  $\Theta(n)$

# Binary Search

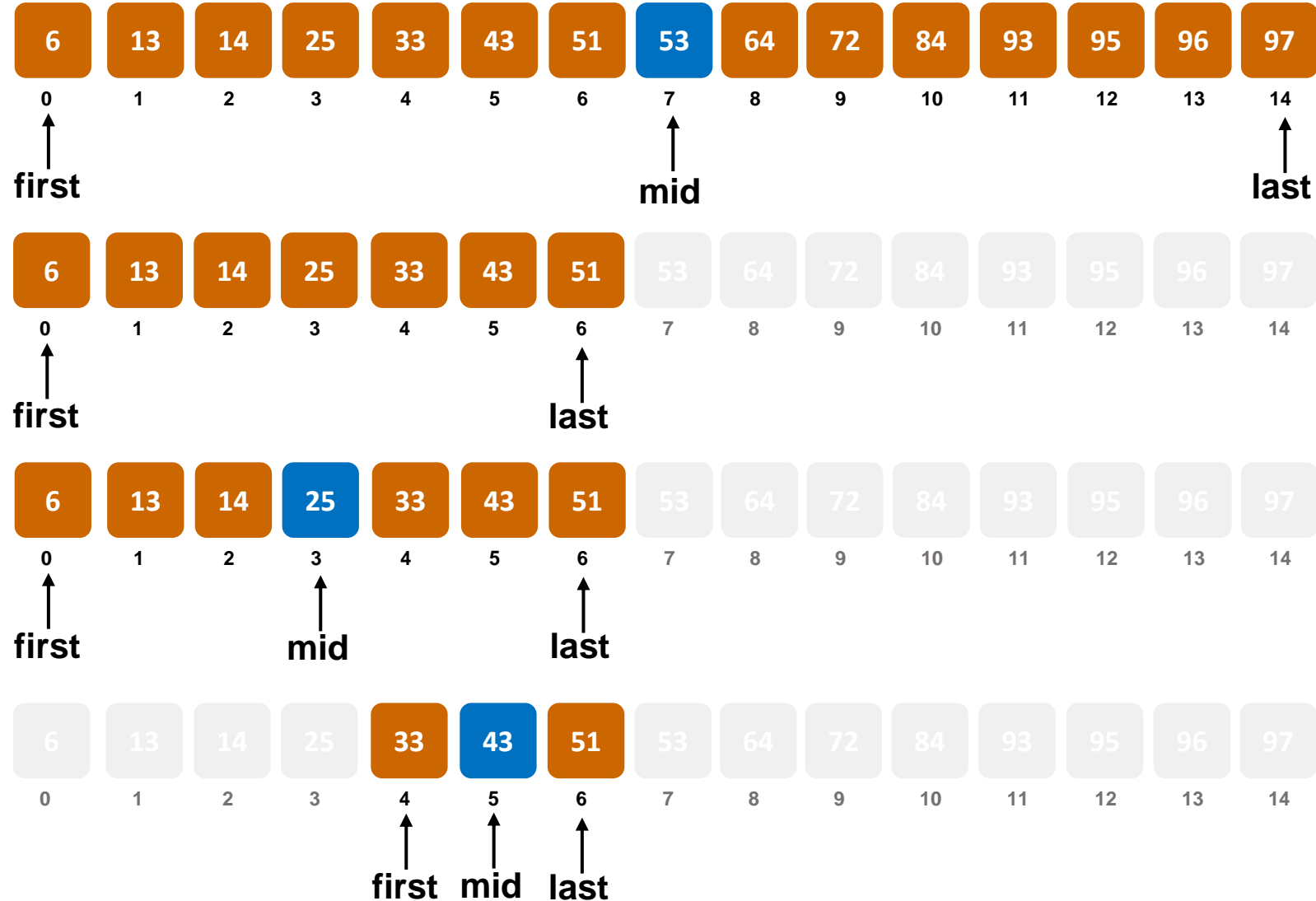
- The array is ordered
- The approach is a **decrease-and-conquer** approach
- A problem is divided into two smaller and similar sub-problem, one of which does not even have to be solved
- The method uses the information of the order to reduce the search space.

# Binary Search

```
1 int binarySearch (int E[], int first, int last, int k)
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);
11        else
12            return binarySearch(E, mid+1, last, k);
13    }
14 }
```

Recursive Version

Binary search for target 33



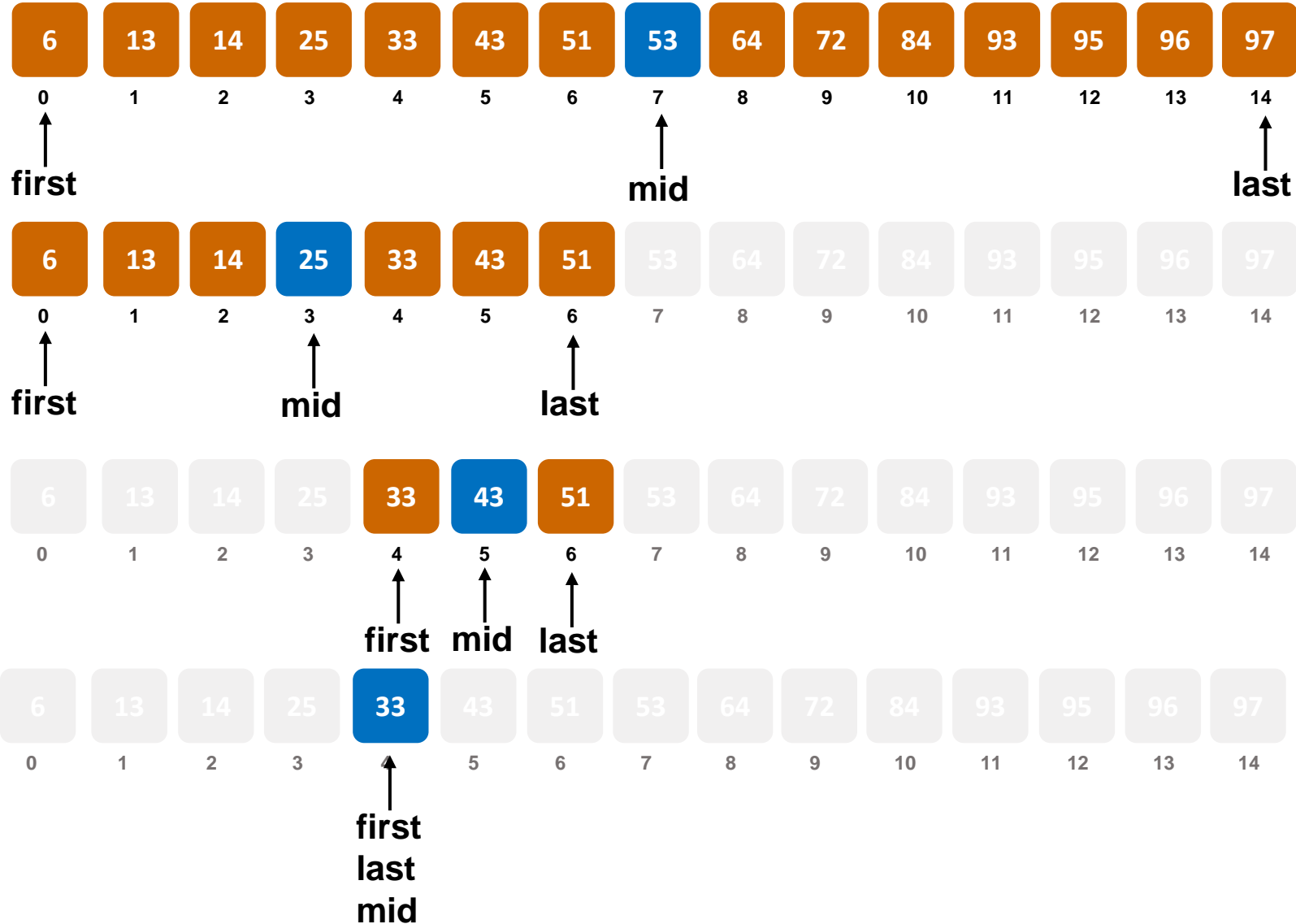


# Binary Search

```
1 int binarySearch (int E[], int first, int last, int k)
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);
11        else
12            return binarySearch(E, mid + 1, last, k);
13    }
14 }
```

Recursive Version

Binary search for target 33



# Binary Search

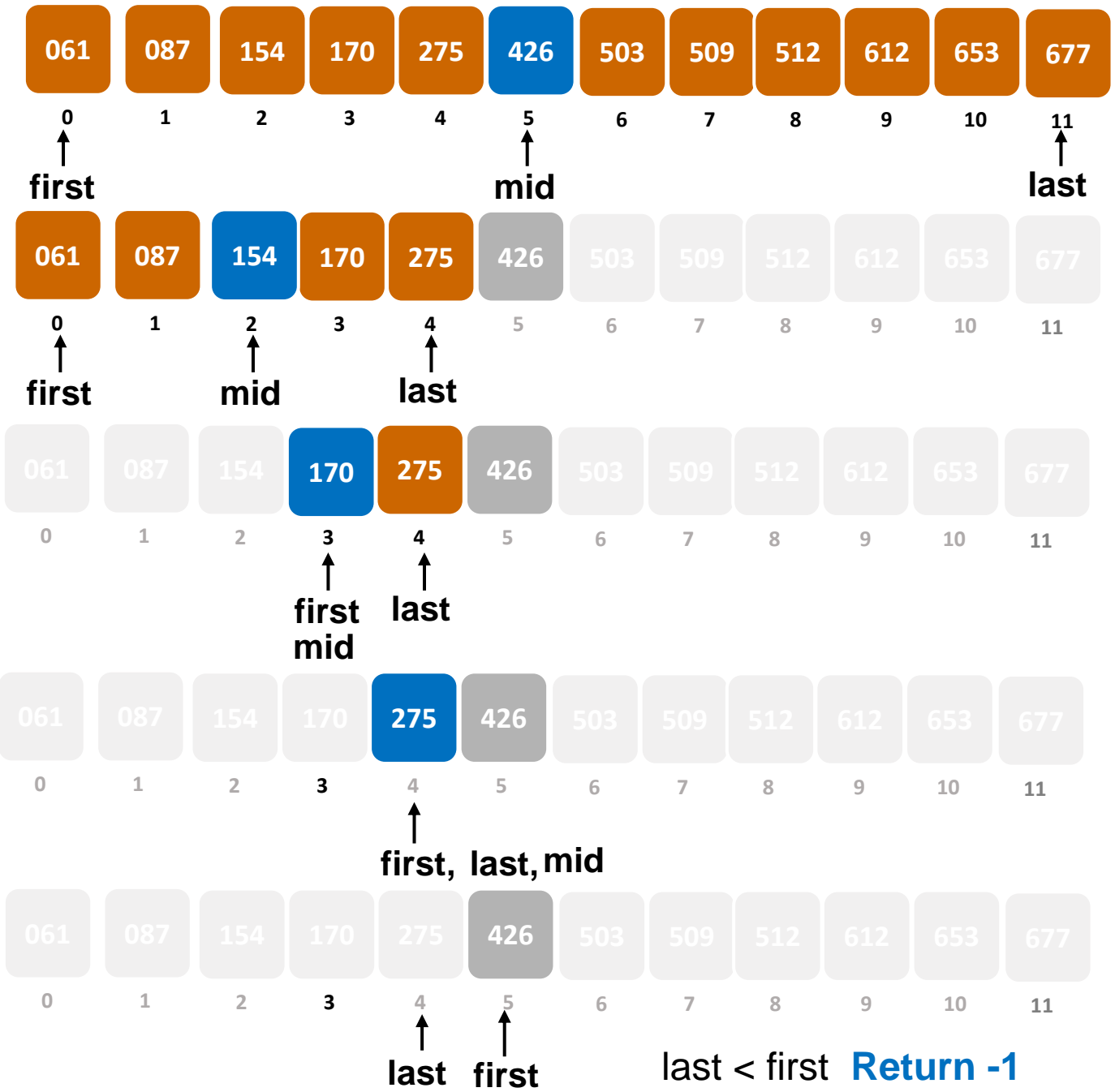
```
1 int binarySearch (int E[], int first, int last, int k)
2 {
3     if (last < first)
4         return -1;
5     else {
6         int mid = (first + last)/2;
7         if (k == E[mid])
8             return mid;
9         else if (k < E[mid])
10            return binarySearch(E, first, mid-1, k);
11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

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 }
```

Iterative Version

Binary search for target **400**



# Time Complexity of Binary Search: Worst Case

```

1 int binarySearch (int E[], int first, int last, int k)
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);
11        else
12            return binarySearch(E, mid+1, last, k);
13    }
14 }

```

Recursive Version

Constant  $c$

$T(n/2)$

$T(n/2)$

$$\begin{aligned}
 T(n) &= T\left(\frac{n}{2}\right) + c \\
 &= T\left(\frac{n}{2^2}\right) + 2c \\
 &= T\left(\frac{n}{2^3}\right) + 3c
 \end{aligned}$$

...

$$\begin{aligned}
 &= T\left(\frac{n}{2^k}\right) + kc \\
 &= T(1) + kc \\
 &= T(0) + (k+1)c \\
 &= (k+1)c \\
 &= (\lfloor \log_2 n \rfloor + 1)c \\
 &= (\lfloor \log_2(n+1) \rfloor)c \\
 &= \Theta(\log_2 n)
 \end{aligned}$$

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

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

Eg.  $k=4$  and  $n$  is integer

$$\begin{aligned}
 2^4 &\leq n < 2^5 \\
 2^4 &< n+1 \leq 2^5
 \end{aligned}$$

$$\begin{aligned}
 \lfloor \log_2(n+1) \rfloor &= k+1 \\
 \lfloor \log_2(n+1) \rfloor &= \lfloor \log_2 n \rfloor + 1
 \end{aligned}$$

# Time Complexity of Binary Search: Average Case

$$A(n) = qA_s(n) + (1 - q)A_f(n)$$

- $A_s(n)$ : # of comparisons for successful search
- $A_f(n)$ : # of comparisons for unsuccessful search (worst case):  $\Theta(\log_2 n)$
- For  $A_s(n)$ , we assume  $n = 2^k - 1$  first

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



1 comparison

# Time Complexity of Binary Search: Average Case

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



↑  
1 comparison



↑                      ↑  
2 comparisons      2 comparisons



↑                      ↑                      ↑                      ↑  
3 comparisons    3 comparisons    3 comparisons    3 comparisons

- We can observe that:
  - 1 position requires 1 comparison
  - 2 positions requires 2 comparisons
  - 4 positions requires 3 comparisons
  - ...
  - $2^{t-1}$  positions requires  $t$  comparisons

•  $n=2^k-1$ , we have

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

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

# Time Complexity of Binary Search: Average Case

- The average complexity is

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

- Binary search does approximately  $\log_2(n + 1)$  comparisons on average for  $n$  entries.
  - $q$  is probability which is always  $\leq 1$
  - $\frac{\log_2(n+1)}{n}$  is very small especially when  $n \gg 1$

# Hashing

- A typical space and time trade-off in algorithm
- To achieve search time in  $O(1)$ , memory usage will be increased
- What is hashing?
  - Hash functions
  - Collision and its resolutions
    - Closed Address Hashing
    - Open Address Hashing
      - Linear Probing
      - Quadratic Probing
      - Double Hashing
  - Delete a key from a hash table
  - Dynamic hash tables

# What is hashing?

## Direct-Address Table

- Assume that the keys of elements  $K$  drawn from the universe of possible keys  $U$
- No two elements have the same key
- Search time is  $O(1)$  but ...
  - The array size is enormous
  - $|U| \gg |K|$

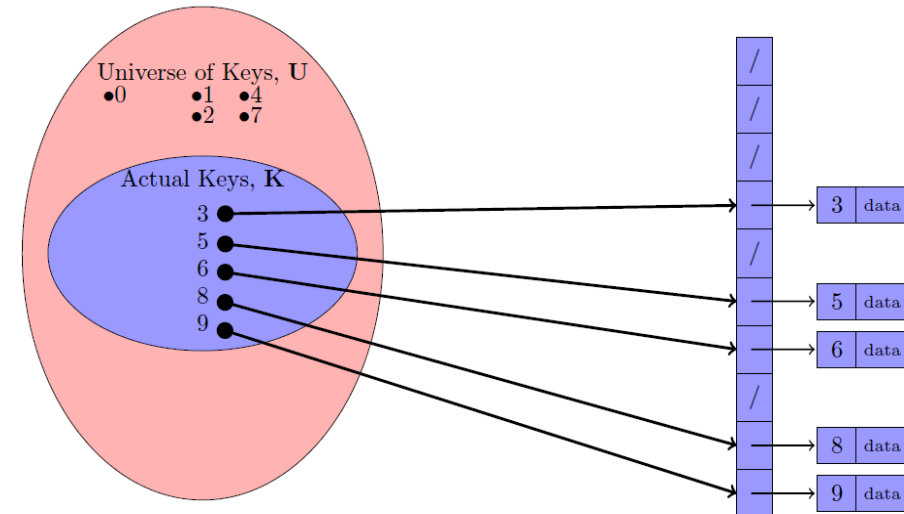


Figure 4.1: Direct Address Table



# What is hashing?

- To reduce the key space to a reasonable size
- Each key is mapped to a unique index (**hash value/code/address**)
- Search time remains  $O(1)$  on the average

**hash function**:  $\{\text{all possible keys}\} \rightarrow \{0, 1, 2, \dots, h-1\}$

- The array is called a **hash table**
- Each entry in the hash table is called a **hash slot**
- When multiple keys are mapped to the same hash value, a **collision** occurs
- If there are  $n$  records stored in a hash table with  $h$  slots, its **load factor** is  $\alpha = \frac{n}{h}$

# Hash Functions

- Must map all possible value within the range of the hash table uniquely
- Mapping should achieve an even distribution of the keys
- Easy and fast to compute
- Minimize collision

1. Modulo Arithmetic
2. Folding
3. Mid-square
4. Multiplicative Congruential Method
5. Etc.

# Hash Functions

## 1. Modulo Arithmetic: $H(k) = k \bmod h$

For keys are chosen from

- decimal number  $\rightarrow h$  avoid to use powers of 10
- unknown lower p-bit patterns  $\rightarrow h$  avoid to use powers of 2
- “real” data  $\rightarrow h$  should be a prime number but not too close to any power of 2

## 2. Folding

- Partition the key into several parts and combine the parts in a convenient way
- Shift folding: Divide the key into a few parts and added up these parts

# Hash Functions

## 3. Mid-square

- The key is squared and the middle part of the result is used as the hash address
  - Eg.  $k=3121$ ,  $k^2 = 3121^2 = 9740641 \rightarrow H(k) = 406$

## 4. Multiplicative Congruential Method

- Pseudo-random number generator
  - $a = 8 \left\lfloor \frac{h}{23} \right\rfloor + 5$
  - $H(k) = (a \times k) \bmod h$

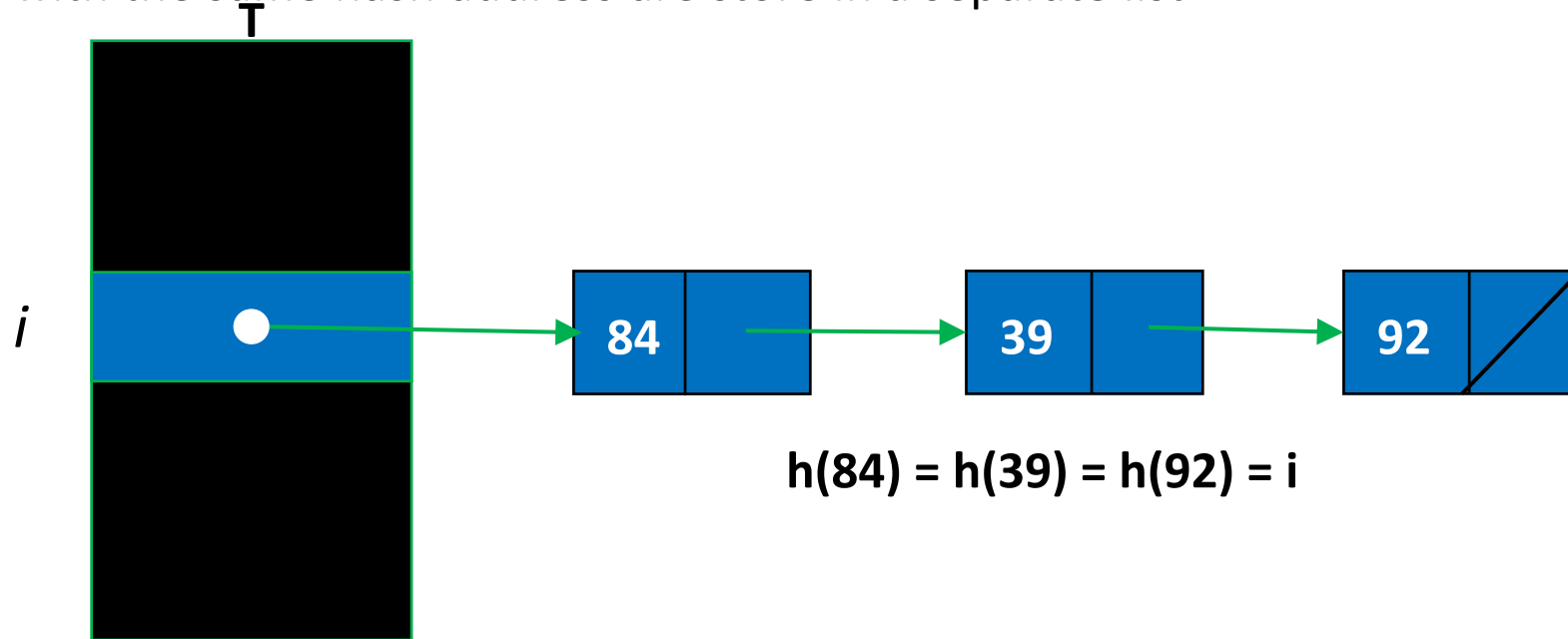
```
1 >> h = 31
2 >> a = 8*floor(h/23) + 5
3
4 a =
5
6     13
7
8 >> k = 1:15
9
10 k =
11
12     1     2     3     4     5     6     7     8     9    10    11    12    13    14    15
13
14 >> fk = mod((a*k),h)
15
16 fk =
17
18     13     26     8     21     3     16     29    11     24     6     19 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

# Collision Resolutions

- Closed Addressing Hashing – a.k.a separate chaining
- Open Addressing Hashing
  - Linear Probing
  - Quadratic Probing
  - Double Probing

# Closed Addressing: Separate Chaining

- Keys are not stored in the table itself
- All the keys with the same hash address are store in a separate list



- During searching, the searched element with hash address  $i$  is compared with elements in linked list  $H[i]$  sequentially
- In closed address hashing, there will be  $\alpha$  number of elements in each linked list on average.

# Closed Addressing: Separate Chaining

Time complexity in the **worst-case analysis**:

- When all elements are hashed to the same slot
- A linked list contains all  $n$  elements
- Its **unsuccessful search** takes  $n$  key comparisons,  $\Theta(n)$
- Its **successful search**, assuming the probability of searching for each item is  $\frac{1}{n}$   
$$\frac{1}{n} \sum_{i=1}^n i = \frac{n+1}{2} = \Theta(n)$$
  - It is just like a sequential search

# Closed Address Hashing: Separate Chaining

Time complexity in the **average-case analysis**:

- All elements are equally likely hashed into  $h$  slots.
- Its **unsuccessful search** takes  $\frac{n}{h}$  key comparisons,  $\Theta(\alpha)$
- Its **successful search** takes 1 more than the number of comparisons done when the sought after item was inserted into the hash table

- Before the  $i^{\text{th}}$  item is inserted into the hash table, the average length of all lists is  $\frac{i-1}{h}$

- When the  $i^{\text{th}}$  item is sought for, the no. of comparisons is  $(1 + \frac{i-1}{h})$

- The average number of key comparisons over  $n$  items

- If  $\alpha$  is constant ( $n$  is proportional to  $h$ ), then  $\Theta(1)$

- Searching takes constant time averagely

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



# Open Addressing

- Keys are stored in the table itself
- $\alpha$  cannot be greater than 1
- When collision occurs, probe is required for the alternate slot

1. Linear Probing: probe the next slot

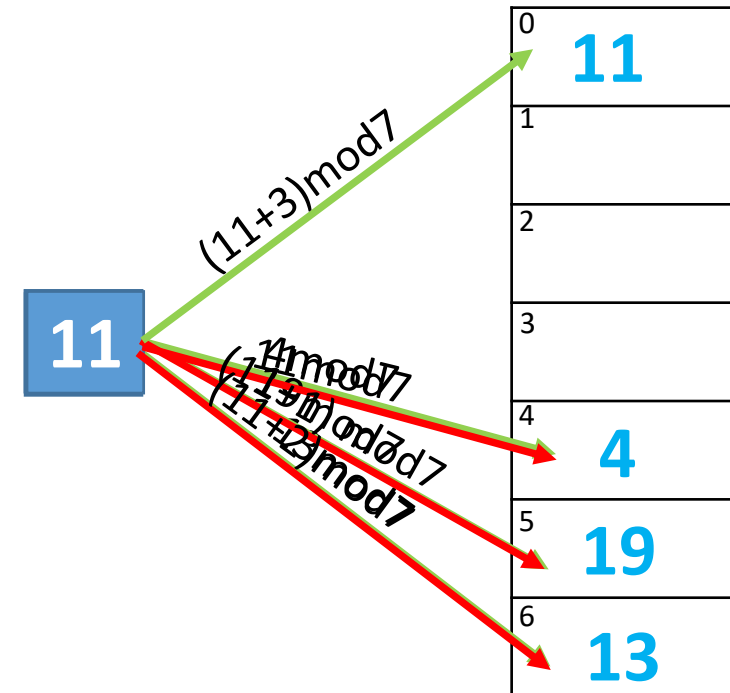
$$H(k, i) = (k + i) \bmod h \text{ where } i \in [0, h - 1]$$

$$\text{eg. } H(k, i) = (k + i) \bmod 7$$

$$k \in \{4, 13, 19, 11\}$$

Primary clustering:

- A long runs of occupied slots
- Average search time is increased



# Open Addressing

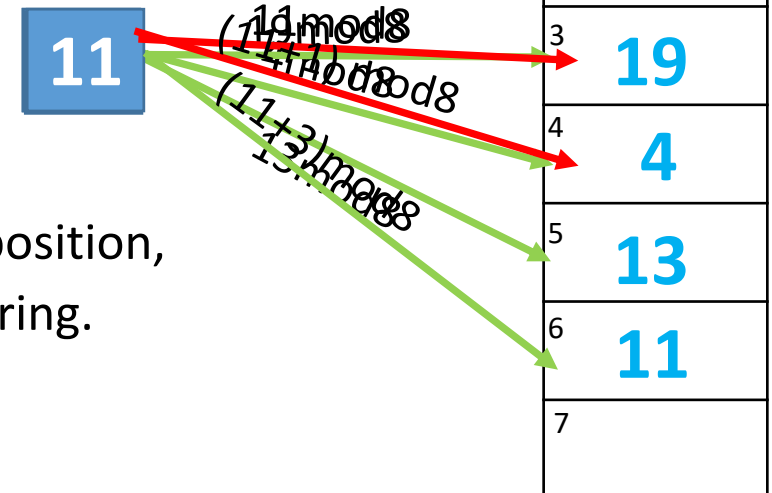
## 2. Quadratic Probing

$$H(k, i) = (k + c_1 i + c_2 i^2) \bmod h \quad \text{where } c_1 \text{ and } c_2 \text{ are constants, } c_2 \neq 0$$

- May not all hash table slots be on the probe sequence (selection of  $c_1, c_2, h$  are important)
- For  $h = 2^n$ , a good choice for the constants are  $c_1 = c_2 = \frac{1}{2}$

eg.  $H(k, i) = \left(k + \frac{1}{2}i + \frac{1}{2}i^2\right) \bmod 8$   
 $k \in \{4, 13, 19, 11\}$

$(\frac{1}{2}i + \frac{1}{2}i^2) \bmod 8$
1
3
6
2
7
5
4



- **Secondary Clustering**: if two keys have the same initial probe position, their probe sequences will be the same. This will form a clustering.
  - Inserting  $k=3$  in the previous example.

# Open Addressing

## 3. Double Hashing: a random probing method

$H(k, i) = (k + iD(k)) \bmod h$  where  $i \in [0, h - 1]$  and  $D(k)$  is another hash function

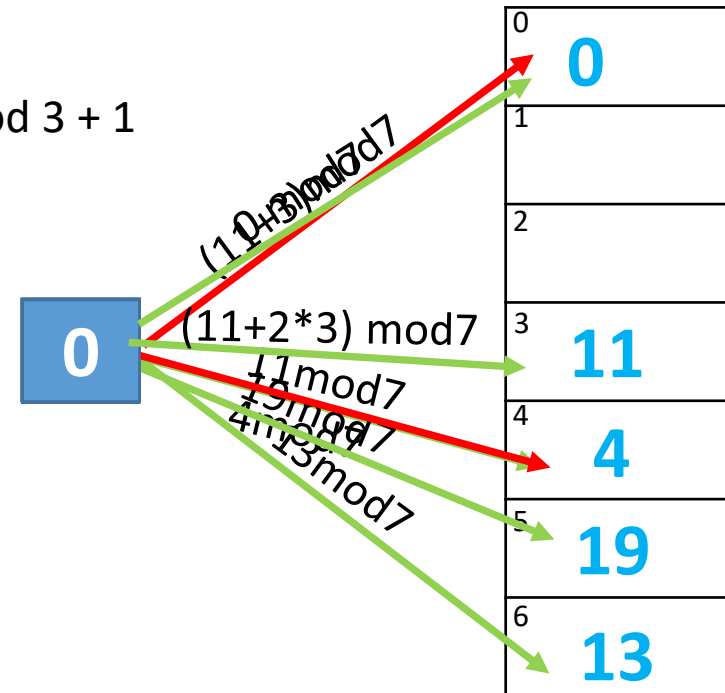
- The hash table size  $h$  should be a prime number

eg.  $H(k, i) = (k + iD(k)) \bmod 7$

$$D(k) = (k) \bmod 3 + 1$$

$$k \in \{0, 4, 13, 19, 11\}$$

$$D(11) = 11 \bmod 3 + 1$$



# Time Complexity

## Linear Probing

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

## Double Hashing

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

# Delete A Key Under Open Addressing

- Leave the deleted key in the table
- Make a marker indicating that it is deleted
- Overwrite it when a new key is inserted to the slot
- May need to do a “garbage collection” when a large number of deletions are done
  - To improve the search time

# Rehashing: Expanding the Hash Table

- As  $\alpha$  increases, the time complexity also increases

Solution:

- Increase the size of hash table (doubled)
- Rehash all keys into new larger hash table

# Summary

- Exhaustive Algorithm: Sequential Search:  $O(n)$
- Decrease-and-conquer Algorithm: Binary Search:  $O(\log_2 n)$
- Data Structures: Hashing
  - Closed Hashing: Separate Chaining:  $O(\alpha)$  on average
  - Open Hashing
    - Linear Probing
    - Quadratic Probing
    - Double Probing
  - Delete keys
  - Rehashing