



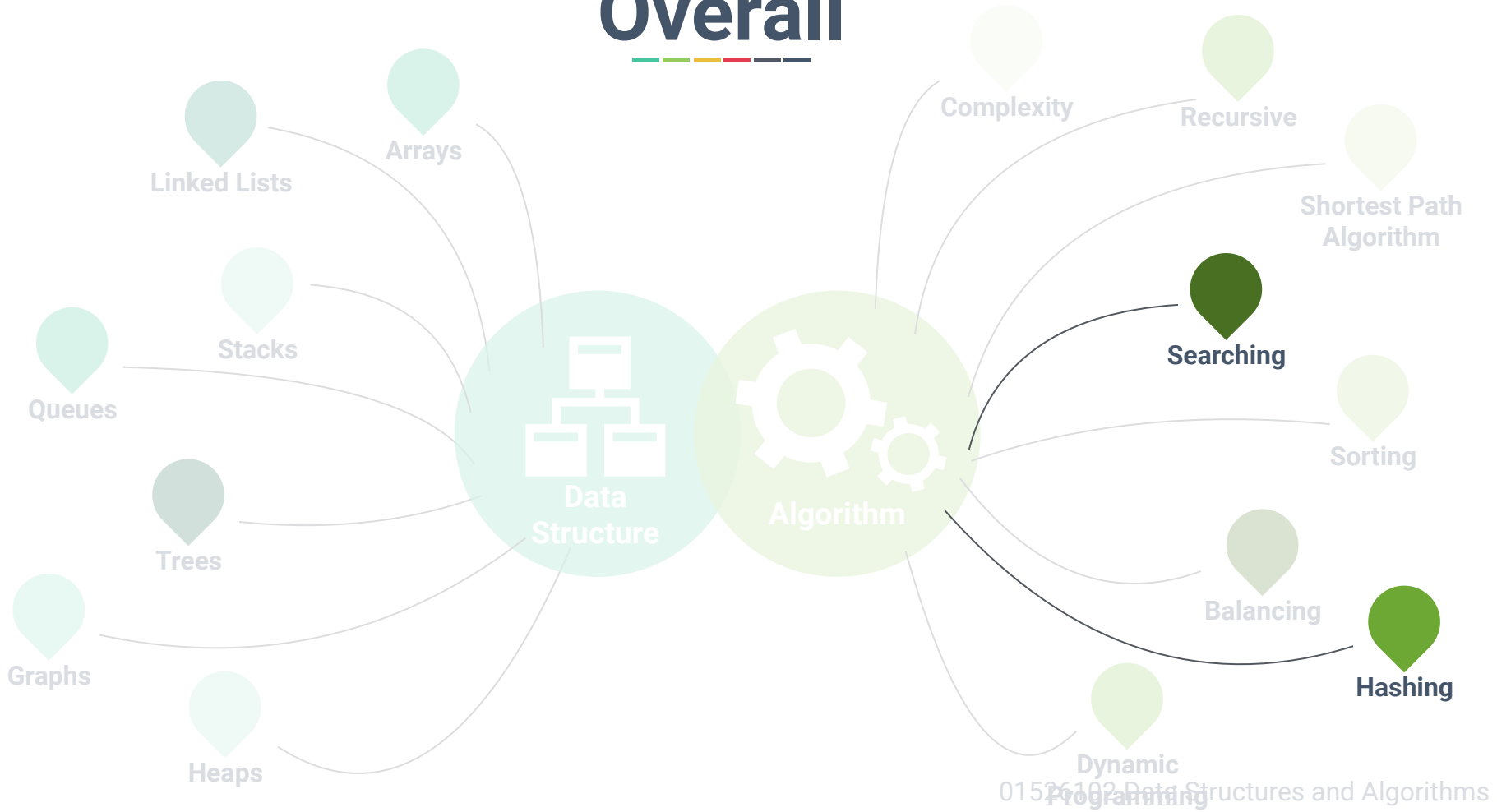
# Chapter 9: Searching



**Dr. Sirasit Lochanachit**



# Overall





# Outline



## Basic Search Algorithms:

- Sequential Search
- Binary Search
- Hashing:
  - Definition and method examples.
  - Collision resolutions



# What is Searching?



- The process of finding a particular item in a collection of items.
- Data management operations, such as insertion, updating, and deletion, usually requires searching algorithm to find an item as fast as possible.



# Searching in an Array

## Sequential



### 1. Search (element by element)





# Searching in an Array

## Sequential

### 1. Search (element by element)



### Pseudocode: Sequential/linear search

```
linear_search (list, target_value)
  for each item in the list
    if item value == target_value
      return the item's location
    end if
  end for
  return 'no match'
END
```



# Sequential Search Algorithm



- Analysis of sequential search algorithm
  - Best Case
  - Worst Case
  - Average Case



# Searching in an Array

## Sequential

### 1. Search (element by element)

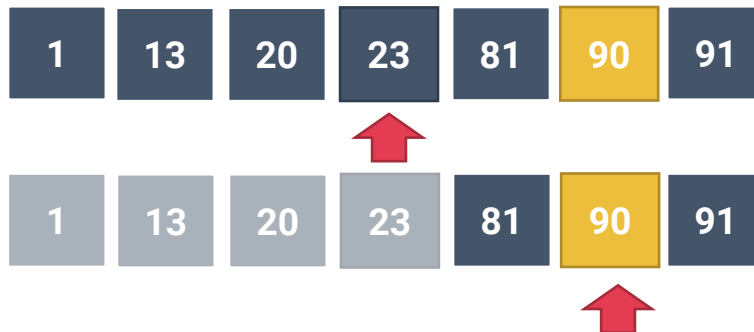


## Binary

### 1. Sort



### 2. Search





# Binary Search



- Implementation
  - Locate a target value in a sequence of  $n$  elements that are sorted.
  - $\text{mid} = (\text{low} + \text{high}) / 2$
  - Initially,  $\text{low} = 0$ ,  $\text{high} = n-1$
- For instance, find number 5.

<b>Data</b>	1	5	7	9	10	11	20
<b>Index</b>	0	1	2	3	4	5	6



# Binary Search



- Implementation
  - If target value < data[mid], next interval is from **low to mid-1**.
  - If target value > data[mid], next interval is from **mid + 1 to high**.

	low			mid			high	
Data	1	5	7	9	10	11	20	mid = (0 + 6) / 2 = 3
Index	0	1	2	3	4	5	6	

	low	mid	high					
Data	1	5	7	9	10	11	20	mid = (0 + 2) / 2 = 1
Index	0	1	2	3	4	5	6	



# Sequential and Binary Search



## Sequential Search

- Easy to implement.
- Suitable for unsorted sequence or small array/list size.
- Search takes linear time -  $O(n)$

## Binary Search

- Suitable for a sorted sequence with large array/list size.
- Search takes logarithmic time -  $O(\log n)$

# What is Hashing?



Index

0	(Data)
1	(Data)
2	(Data)
3	(Data)
4	(Data)
5	(Data)

(1) Array

Bucket

	Hash Table	
0	(Key)	(Data)
1	(Key)	(Data)
2	(Key)	(Data)
3	(Key)	(Data)
4	(Key)	(Data)
5	(Key)	(Data)

(2) Hash table

**Hashing** is a technique that determines the index (bucket) using only a target search key, without the need to search.



# Hashing Components



## 1) Hash Table or Bucket Array

- An array  $A$  of size  $N$ .
- Each cell of  $A$  is thought of as a “Bucket” that is, a collection of key-value pairs.
- Ideally, the keys are well distributed in the range  $[0, N-1]$  by a hash function.
- However, it is possible for two or more distinct keys to get mapped to the same index.

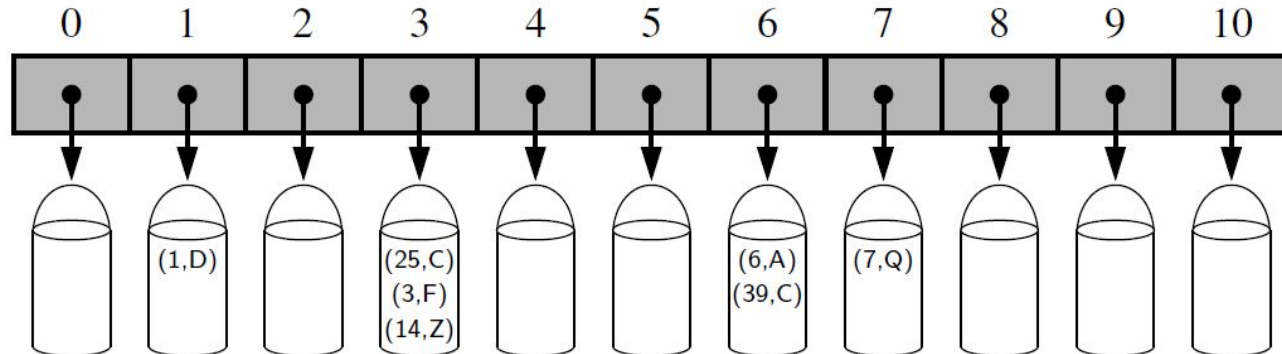
Bucket	Hash Table	
0	(Keys)	(Data)
1	(Keys)	(Data)
2	(Keys)	(Data)
3	(Keys)	(Data)
4	(Keys)	(Data)
5	(Keys)	(Data)

# Hashing Components

## 1) Hash Table or Bucket Array

\* In Python, hashing is implemented as dictionary

Bucket	Hash Table	
0	(Keys)	(Data)
1	(Keys)	(Data)
2	(Keys)	(Data)
3	(Keys)	(Data)
4	(Keys)	(Data)
5	(Keys)	(Data)



# Hashing Components



- a function,  $h$ , that maps each key  $k$  to the corresponding integer index (bucket) in the hash table.

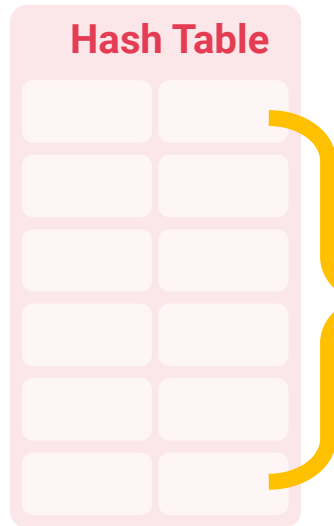
Bucket	Hash Table	
0	(Keys)	(Data)
1	(Keys)	(Data)
2	(Keys)	(Data)
3	(Keys)	(Data)
4	(Keys)	(Data)
5	(Keys)	(Data)



# Hashing Process



**Step 1:** define the capacity of the hash table ( $N$ ).



$N$

**Step 2:** define the hash function.







# Basic Hash Function



**Division method  
(or modulo arithmetic)**

$$\text{Hash value} = k \bmod N$$

**MAD method**

$$\text{Hash value} = [(ak + b) \bmod p] \bmod N$$

**Selection  
or mid-square method**

$$\text{Hash value} = \text{2 digits of } k^2$$

$N$  = size of the bucket array

$k$  = key of the item/element

$p$  = prime number larger than  $N$

$a$  and  $b$  = random integers between  $[0, p-1]$

# Division Method



Bucket	Data
0	
1	
2	
3	
4	
5	
6	
7	
8	
9	
10	

$$\text{Hash value} = k \bmod N$$

Example:

Given a set of integer items {20, 19, 29, 10, 31, 3, 42, 14} and the capacity of the hash table is 11, calculate the hash values for each item.

# Division Method



Bucket	Data
0	
1	
2	
3	
4	
5	
6	
7	
8	
9	
10	

$$\text{Hash value} = k \bmod N$$

To search for an item:

- Simply use the hash function to compute the bucket index for the item.
- Then, check the hash table to see if it exists.

It takes a constant time  $O(1)$  to compute the hash value and retrieve the item.

Highly likely to cause collision of hash values.

# Collision



Bucket	Data
0	
1	
2	
3	
4	
5	
6	
7	
8	
9	
10	

$$\text{Hash value} = k \bmod N$$

If two or more keys produces the same hash value,

- They would need to be in the same bucket.
- This is known as **collision**.



# Division Method Exercise 1



Bucket	Data
0	
1	
2	
3	
4	
5	
6	
7	
8	

$$\text{Hash value} = k \bmod N$$

Exercise 1:

Given a set of integer items {4, 1, 3, 10, 7, 5, 6} and the capacity of the hash table is 9, calculate the hash values for each item.



# Division Method Exercise 2

Bucket	Data
0	
1	
2	
3	
4	
5	
6	
7	
8	

Dec	Hex	Char	Dec	Hex	Char
64	40	@	96	60	`
65	41	A	97	61	a
66	42	B	98	62	b
67	43	C	99	63	c
68	44	D	100	64	d
69	45	E	101	65	e
70	46	F	102	66	f
71	47	G	103	67	g
72	48	H	104	68	h
73	49	I	105	69	i
74	4A	J	106	6A	j
75	4B	K	107	6B	k
76	4C	L	108	6C	l
77	4D	M	109	6D	m
78	4E	N	110	6E	n
79	4F	O	111	6F	o
80	50	P	112	70	p
81	51	Q	113	71	q
82	52	R	114	72	r
83	53	S	115	73	s
84	54	T	116	74	t
85	55	U	117	75	u
86	56	V	118	76	v
87	57	W	119	77	w
88	58	X	120	78	x
89	59	Y	121	79	y
90	5A	Z	122	7A	z
91	5B	[	123	7B	{
92	5C	\	124	7C	
93	5D	]	125	7D	}
94	5E	^	126	7E	~
95	5F	_	127	7F	DEL

$$\text{Hash value} = k \bmod N$$

Exercise 2:

Given a set of character items {A, B, F, G, K, P, R, X, Y, Z} and the capacity of the hash table is set to 9, calculate the hash values for each item.



# Collision Resolution



Two fundamental method for resolving collisions are:

- **First**, use another available spot in the hash table.
- **Second**, change the structure of the hash table so that each array element can store more than one value.



# Collision Resolution



- **First**, use another available spot in the hash table.

Linear Probing

Quadratic Probing

Double Hashing

**Open addressing** is the simple process to find the next empty bucket or address in the hash table.

It starts at the original hash value position and then move to the next hash value until the available bucket is found.



# Open Addressing



Bucket	Data
0	20
1	
2	
3	
4	
5	34
6	19
7	
8	
9	

Hash value = 0

Data = 10



Bucket	Data
0	20
1	
2	
3	
4	
5	34
6	19
7	
8	
9	

10



# Linear Probing



**Linear probing** method resolves a collision during hashing by looking at a sequential location, index by index, in the hash table that is available - starting at the original hash value. This search sequence is **probe sequence**.

Bucket	Data
0	20
1	
2	
3	
4	
5	34
6	19
7	
8	
9	

Hash value = 0

Data = 10

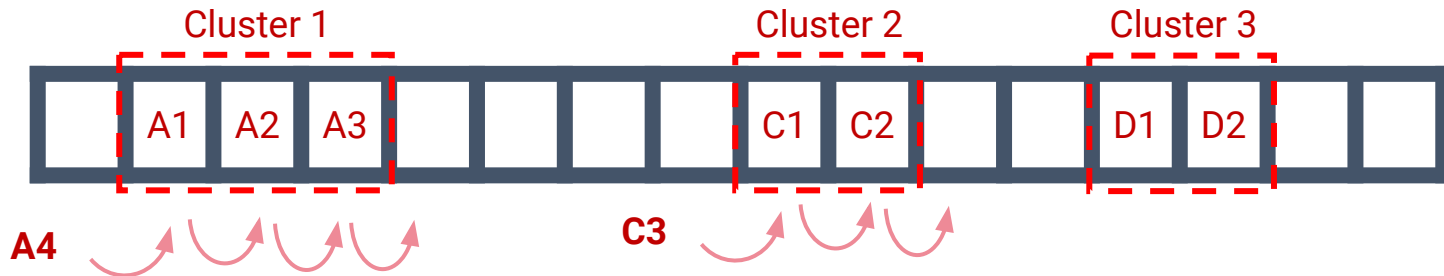


Bucket	Data
0	20
1	10
2	
3	
4	
5	34
6	19
7	
8	
9	





# Linear Probing



Although linear probing is a simple technique, it faces *a primary clustering problem*.

Collisions that are resolved with linear probing cause groups of consecutive locations (as so-called a cluster) in the hash table to be occupied.



# Quadratic Probing



In order to avoid the primary clustering problem, the open addressing of the **quadratic probing** is used. The difference between linear and quadratic probing is that the new location is indicated by *hash+1* for the linear probing whereas it is indicated by *hash+j<sup>2</sup>* for the quadratic probing.

$$Rehash_{new} = (Hash_{old} + j^2) \bmod N$$



The number of the collision

# Quadratic Probing

$$Rehash_{new} = (Hash_{old} + j^2) \bmod N$$

Bucket	Data
0	
1	
2	871
3	
4	
5	
6	
7	
8	415
9	
10	604

Hash value = 8  
Data = 921



Bucket	Data
0	
1	
2	871
3	
4	
5	
6	
7	
8	415
9	921
10	604

Hash value = 9  
Data = 163



Bucket	Data
0	
1	
2	871
3	163
4	
5	
6	
7	
8	415
9	921
10	604

Collision #1  
(8+1<sup>2</sup>) = 9

Collision #2  
(10+2<sup>2</sup>) = 14

Collision #1  
(9+1<sup>2</sup>) = 10



# Quadratic Probing



Although the quadratic probing can avoid the primary clustering problem created by linear probing, it create its own kind of clustering, which is **secondary clustering**.

Moreover, it adds more complexity for the running time. Specifically, it requires more time to compute the indices.



# Collision Resolution



- **First**, use another available spot in the hash table.

Linear Probing

Quadratic Probing

- **Second**, change the structure of the hash table so that each array element can store more than one value.

Chaining

# Chaining



Changing the structure of the hash table so that each array element can represent more than one value by using the **Linked List** structure.

