

CSC 212: Data Structures and Abstractions

Hash Tables

Prof. Marco Alvarez

Department of Computer Science and Statistics
University of Rhode Island

Spring 2025



Data Structure	Worst-case			Average-case			Ordered?
	insert at	delete	search	insert at	delete	search	
sequential (unordered)	$O(n)$	$O(n)$	$O(n)$	$O(n)$	$O(n)$	$O(n)$	No
sequential (ordered) binary search	$O(n)$	$O(n)$	$O(\log n)$	$O(n)$	$O(n)$	$O(\log n)$	Yes
BST	$O(n)$	$O(n)$	$O(n)$	$O(\log n)$	$O(\log n)$	$O(\log n)$	Yes
2-3-4	$O(\log n)$	$O(\log n)$	$O(\log n)$	$O(\log n)$	$O(\log n)$	$O(\log n)$	Yes
Red-Black	$O(\log n)$	$O(\log n)$	$O(\log n)$	$O(\log n)$	$O(\log n)$	$O(\log n)$	Yes

2

Can we do better?

3

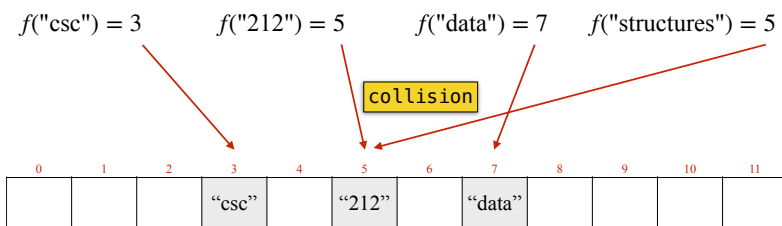
Random access memory

- **Random Access Memory (RAM)** represents a fundamental principle in computer science
 - ✓ it allows the retrieval of any element in constant time $O(1)$, regardless of its position within the memory block
 - ✓ this principle is most commonly observed through arrays
- **Arrays in C++**
 - ✓ contiguous memory allocation
 - ✓ homogeneous elements (same data type)
 - ✓ fixed-length (traditional arrays have predetermined size)
 - ✓ zero-based indexing

4

Hash tables

- A hash table is a data structure that implements an associative array
 - ✓ the array can store keys (set), or key-value pairs (map)
 - ✓ a function is used to compute an index, that can be used to find a desired key in the array
 - ✓ provides an efficient way to implement sets or dictionaries

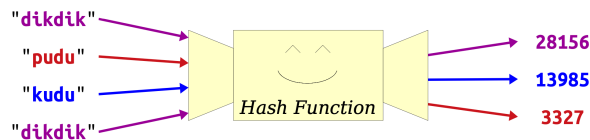


5

Hash function

Hash function

- A hash function is a function that maps an input key to some integer value
 - ✓ must be deterministic (same input produces the same output)
 - ✓ should be well-distributed (the numbers produced are as spread out as possible)



- Hash tables
 - ✓ store keys in an array and use a hash function to map keys to indices
 - ✓ potentials issues:
 - may require a large array
 - array could be very sparse (wasting memory)

Image credit: CS106B @ Stanford

7

Practice

- Which of the following tables is a better choice?
- What is the load factor?

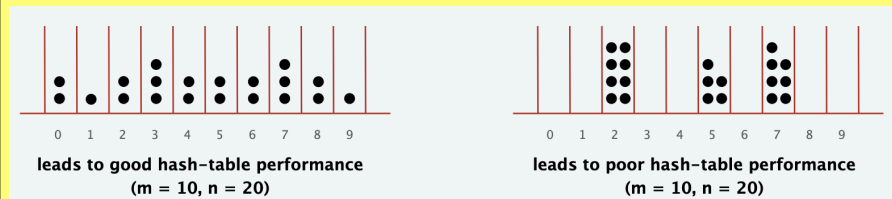


Image credit: COS 226 @ Princeton

8

Hash functions

- A hash function is a function that maps an input key to some integer value
- Properties:
 - ✓ must be **deterministic** (same input produces the same output)
 - ✓ should be **well-distributed** (the numbers produced are as spread out as possible)
 - ✓ should be **fast** to compute
- Hash tables (idea)
 - ✓ store keys in an array and use a hash function to map keys to indices
 - ✓ potentials issues:
 - may require a large array
 - array could be very sparse (wasting memory)

9

Hash functions

- Space efficiency
 - ✓ making all keys equally possible requires a huge array, even if we only have a couple of elements
 - ✓ idea: use a hash function, but modify the result to be within a smaller range (the size of the array)
- ```
// if hash() returns non-negatives
index = hash(key) % capacity

// if hash() returns any integer
index = abs(hash(key)) % capacity
```
- Collision
    - ✓ occurs when two different keys hash to the same index in the hash table
    - ✓ collisions can be resolved using:
      - separate chaining: each slot in the hash table contains a collection of all the keys that hash to that index
      - open addressing: if a collision occurs, the algorithm searches for the next available slot in the hash table
    - ✓ open addressing is more space-efficient than chaining, but it can be slower

10

# Hash functions

- Hash functions can be used on any data type
  - ✓ integers: use the integer value as the hash value
  - ✓ floats: convert to binary and use the integer value as the hash value, or manipulate the bits (e.g. XOR the mantissa and exponent)
  - ✓ strings: use  $31x + y$  rule or apply other variants
  - ✓ compound objects: use  $31x + y$  rule or apply other variants
- Mapping hash values into a smaller “table size”
  - ✓ M is prime: helps distributing keys more uniformly, minimizing collisions
  - ✓ M is a power of two: modulo operation can be replaced with a faster bitwise AND operation

11

# Hash functions (other uses)

- Storing passwords
  - ✓ hash the password and store the hash in the database
- File verification
  - ✓ hash the file (checksum) and compare the hash with the stored hash
  - ✓ e.g. when downloading a file, vendors publish a hash value, client checks whether hash matches, otherwise file is corrupted
- Examples (one-way hashes)
  - ✓ hash functions that are difficult to reverse or to find two keys that map to the same value
  - ✓ MD5, SHA-1, SHA-256, SHA-512, SHA3-512, ...

12

# Separate chaining

## Separate chaining

### • Idea

- ✓ solve collisions by storing a linked list at each index
- ✓ assume duplicated keys are not allowed

### • Operations

- ✓ insert: if a collision occurs, add the new key to the linked list at that index
  - no need to keep the keys on each list in order
- ✓ search: search the linked list at that index
- ✓ delete: search the linked list at that index and removes the key from the list

### • Comments

- ✓ linked list can be replaced by a balanced tree or another data structure to improve access time

14

## Practice

### • Perform the following operations

- ✓ insert(L, 11), delete(D), insert(M, 12), delete(E), search(C)
- assume: insertions occur at front of the lists, hash(L)=3, hash(M)=0

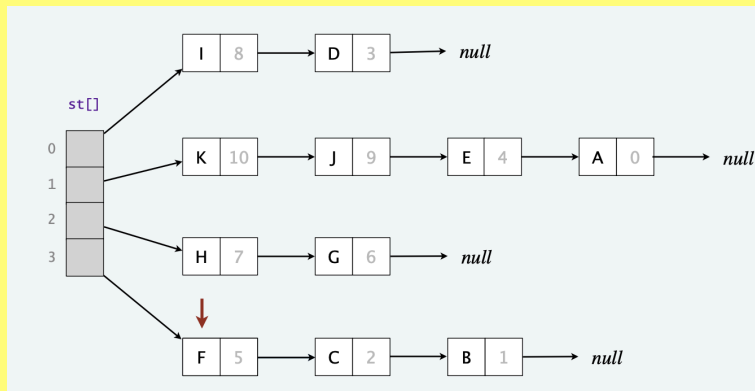


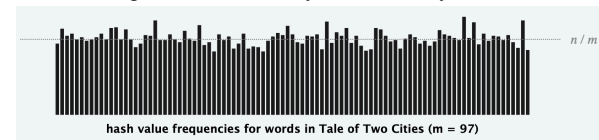
Image credit: COS 226 @ Princeton

15

## Analysis

### • Uniform hashing assumption

- ✓ the hash function is a good one, and all keys are uniformly distributed



### • Load factor ( $\alpha$ )

- ✓ the ratio of the number of keys ( $N$ ) to the number of slots ( $M$ )  $\alpha = \frac{N}{M}$

### • Time complexity

- ✓ average case for search, insert, and delete is  $O(c + \alpha)$ , where  $c$  is the time taken by the hash function
- ✓ worst case for search, insert, and delete is  $O(c + N)$  — all the keys hash to the same index

### • Space complexity

- ✓ the space needed is  $O(N + M)$

16

## Resizing a hash table

- Growing to a larger array when  $\alpha$  exceeds a threshold
  - ✓ create a new table with larger capacity and rehash all the keys

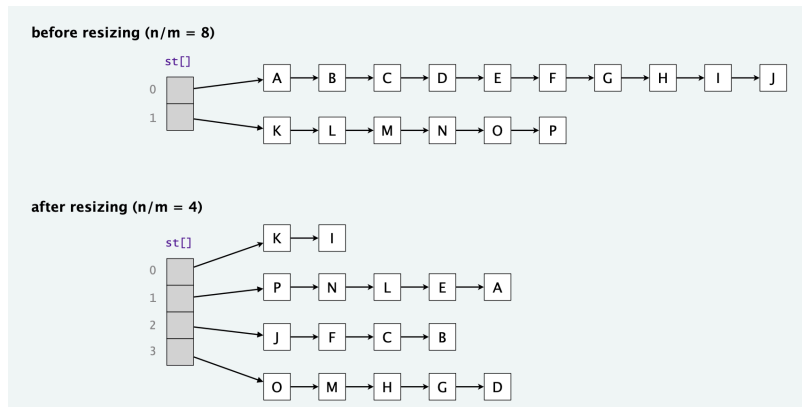


Image credit: COS 226 @ Princeton

17

## Considerations

- Choices for  $\alpha$ 
  - ✓ too small, the hash table will be too large and waste space
  - ✓ too large, the hash table will be too small and cause collisions
- Typical values
  - ✓ between 0.5 and 1.0 often provide a reasonable balance of space efficiency and lookup performance
  - ✓ higher load factors ( $>1.0$ ) remain functional but with degraded performance characteristics
  - ✓ for performance-critical applications, implementers should conduct benchmarks with representative data sets to determine the optimal load factor for their specific use case

18

| Data Structure                        | Worst-case  |             |             | Average-case |             |             | Ordered? |
|---------------------------------------|-------------|-------------|-------------|--------------|-------------|-------------|----------|
|                                       | insert at   | delete      | search      | insert at    | delete      | search      |          |
| sequential (unordered)                | $O(n)$      | $O(n)$      | $O(n)$      | $O(n)$       | $O(n)$      | $O(n)$      | No       |
| sequential (ordered)<br>binary search | $O(n)$      | $O(n)$      | $O(\log n)$ | $O(n)$       | $O(n)$      | $O(\log n)$ | Yes      |
| BST                                   | $O(n)$      | $O(n)$      | $O(n)$      | $O(\log n)$  | $O(\log n)$ | $O(\log n)$ | Yes      |
| 2-3-4                                 | $O(\log n)$ | $O(\log n)$ | $O(\log n)$ | $O(\log n)$  | $O(\log n)$ | $O(\log n)$ | Yes      |
| Red-Black                             | $O(\log n)$ | $O(\log n)$ | $O(\log n)$ | $O(\log n)$  | $O(\log n)$ | $O(\log n)$ | Yes      |
| Hash table (separate chaining)        | $O(n)$      | $O(n)$      | $O(n)$      | $O(1)$       | $O(1)$      | $O(1)$      | No       |

19

Open addressing

# Open addressing

- **Idea**
  - ✓ solve collisions by “*probing*”
    - searching for the next available slot in the hash table
    - each slot holds a single element
    - if using *key-value* pairs, maintain two separate arrays
  - ✓ assume duplicated keys are not allowed and  $M \geq N$
- **Operations**
  - ✓ insert: if a collision occurs, search for the next available slot in the hash table
  - ✓ search: search for the key in the hash table
  - ✓ delete: search for the key in the hash table and mark the slot it as deleted
- **Comments**
  - ✓ open addressing is more space-efficient than chaining, but it can be slower
  - ✓ works better with  $\alpha \approx 0.5$

21

# Probing

- **Linear probing**
  - ✓ moves to the next available index
  - ✓ use index  $i$  if free, otherwise try  $i + 1, i + 2, i + 3, \dots$  (wrap around if necessary)
- **Quadratic probing**
  - ✓ moves to the next available index using a quadratic function
  - ✓ use index  $i$  if free, otherwise try  $i + 1^2, i + 2^2, i + 3^2, \dots$  (wrap around if necessary)
- **Double hashing**
  - ✓ moves to the next available index using a second hash function
  - ✓ use index  $i$  if free, otherwise try  $i + \text{hash2}(\text{key}), i + 2 * \text{hash2}(\text{key}), i + 3 * \text{hash2}(\text{key}), \dots$  (wrap around if necessary)

22

# Practice

- **Perform the following operations (assume linear probing)**
  - ✓ search(w), delete(z), delete(w), search(r), insert(c), insert(d), insert(e)
  - assume:  $\text{hash}(z)=2, \text{hash}(x)=7, \text{hash}(r)=7, \text{hash}(w)=7, \text{hash}(y)=14, \text{hash}(a)=12, \text{hash}(c)=8, \text{hash}(d)=15, \text{hash}(e)=14$

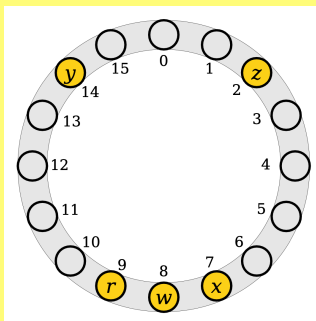


Image credit: CS106B @ Stanford

23

| Data Structure                        | Worst-case  |             |             | Average-case |             |             | Ordered? |
|---------------------------------------|-------------|-------------|-------------|--------------|-------------|-------------|----------|
|                                       | insert at   | delete      | search      | insert at    | delete      | search      |          |
| sequential (unordered)                | $O(n)$      | $O(n)$      | $O(n)$      | $O(n)$       | $O(n)$      | $O(n)$      | No       |
| sequential (ordered)<br>binary search | $O(n)$      | $O(n)$      | $O(\log n)$ | $O(n)$       | $O(n)$      | $O(\log n)$ | Yes      |
| BST                                   | $O(n)$      | $O(n)$      | $O(n)$      | $O(\log n)$  | $O(\log n)$ | $O(\log n)$ | Yes      |
| 2-3-4                                 | $O(\log n)$ | $O(\log n)$ | $O(\log n)$ | $O(\log n)$  | $O(\log n)$ | $O(\log n)$ | Yes      |
| Red-Black                             | $O(\log n)$ | $O(\log n)$ | $O(\log n)$ | $O(\log n)$  | $O(\log n)$ | $O(\log n)$ | Yes      |
| Hash table (separate<br>chaining)     | $O(n)$      | $O(n)$      | $O(n)$      | $O(1)$       | $O(1)$      | $O(1)$      | No       |
| Hash table (open<br>addressing)       | $O(n)$      | $O(n)$      | $O(n)$      | $O(1)$       | $O(1)$      | $O(1)$      | No       |

24