

# **Introduction to data structures**

**William Hendrix**

# Outline

---

- Overview
- Stacks and Queues
- Dictionaries
  - Arrays vs. linked-lists
  - Sorted vs. unsorted
  - Binary search trees
  - Hash tables
  - Bit vectors
- Heaps
- Union-find
- Prefix and suffix trees

# Why are we studying data structures?

- Data structures are a fundamental component of algorithms
- Different data structures support different operations
- Some have multiple implementations
- Strongly affect complexity

- **Example:** SelectionSort

```
Input: data: an array of integers to sort
Input: n: the number of values in data
Output: permutation of data such that
         $data[1] \leq \dots \leq data[n]$ 
1 Algorithm: SelectionSort
2 for  $i = 1$  to  $n$  do
3   |   Let  $m$  be the location of the min value in the
4   |   array  $data[i..n]$ ;
5   |   Swap  $data[i]$  and  $data[m]$ ;
6 end
7 return data;
```

- Heap: data structure that takes  $O(n)$  time to build and  $O(\lg n)$  time to delete the min
- Transforms  $O(n^2)$  SelectionSort into  $O(n \lg n)$  HeapSort

# Data structures in memory

---

- Contiguous
  - Allocate single chunk of memory
  - Retrieve elements by locating data's index in chunk
  - Very fast if elements are accessed consecutively (caching)
  - No memory wasted on storing pointers
  - Following  $k$  links takes  $O(k)$  time vs.  $O(1)$  for pointer arithmetic
- Link-based
  - Data are stored in small "islands" connected via pointers
  - Retrieve elements by starting at the head/tail/root and traversing links
  - Supports data with irregular structure (graphs, trees)
  - Easier for memory manager to allocate
  - Easy to modify structure by changing links (vs. copying data)
- Hybrid

# Stacks and queues

- Abstract data structures
  - Define a set of operations
  - Generally implemented as arrays
- **Stacks**
  - Support *push()* and *pop()* operations
  - Last-In, First-Out (LIFO) order
  - Clears out "new" work quickly
  - DFS generally requires less space than BFS
- **Queues**
  - Support *enqueue()* and *dequeue()* operations
  - First-In, First-Out (FIFO) order
  - More "fair": delay between *enqueue()* and *dequeue()* balanced
- **Dequeues** ("decks")
  - Support all 4 operations



# Operation analysis

- **push(x)**
  - Depends on *capacity*
    - Allocated memory
  - If capacity > size, just add
  - Otherwise, enlarge first
    - Copy all elements over
    - $O(n)$  time

```
1 Algorithm: Push(x)
2 if  $n \geq c$  then
3   | Enlarge();
4 end
5  $n = n + 1$ ;
6  $\text{stack}[n] = x$ ;
```

```
1 Algorithm: Enlarge()
2  $c = ???$ ;
3  $\text{newstack} = \text{allocate}(c)$ ;
4 for  $i = 1$  to  $n$  do
5   |  $\text{newstack}[i] = \text{stack}[i]$ 
6 end
7 free stack;
8  $\text{stack} = \text{newstack}$ ;
```

# Array enlargement policy

- **Bad policy:** increment by 1
    - Cost:  $O(n)$  every time
    - $n$  pushes:  $O(1 + 2 + \dots + n) = O(n^2)$
    - *Average cost per push:*  $O(n)$
  - **Better policy:** increment by  $k$  ( $k=10, 100$ , etc.)
    - Cost:  $O(n)$  every  $k^{\text{th}}$  time
    - $n$  pushes:  $O(1 + (1 + k) + \dots + (1 + k\lfloor \frac{n}{k} \rfloor)) = O\left(\frac{n^2}{k}\right)$
    - *Average cost per push:*  $O\left(\frac{n}{k}\right)$
    - *Caution:* space trade-off
      - Increment by 1B: small stacks will be mostly empty space
  - **Best policy:** double array size
    - Cost:  $O(n)$  after powers of two
    - $n$  pushes:  $O(1 + 2 + 4 + \dots + 2^{\lfloor \lg n \rfloor}) = O(n)$
    - *Average cost per push:*  $O(1)$
    - Array will be at least half-full if not deleting
- } Amortized analysis

# Operation analysis

- **push(x)**
  - Depends on *capacity*
    - Allocated memory
  - If capacity > size, just add
  - Otherwise, enlarge first
  - $O(1)$  time, amortized
- **pop()**
  - Remove last element
  - $O(1)$  time
- **enqueue(x)**
  - Similar to push()
  - $O(1)$  time, amortized
- **dequeue(x)**
  - Advance "head" pointer
  - $O(1)$  time

```
1 Algorithm: Push(x)
2 if  $n \geq c$  then
3   | Enlarge();
4 end
5  $n = n + 1$ ;
6  $\text{stack}[n] = x$ ;
```

```
1 Algorithm: Enlarge()
2  $c = 2c$ ;
3  $\text{newstack} = \text{allocate}(c)$ ;
4 for  $i = 1$  to  $n$  do
5   |  $\text{newstack}[i] = \text{stack}[i]$ 
6 end
7 free stack;
8  $\text{stack} = \text{newstack}$ ;
```

```
1 Algorithm: Pop()
2 if  $n = 1$  then
3   | error "Stack is empty";
4 end
5  $n = n - 1$ ;
6 return  $\text{data}[n]$ ;
```



# Dictionary

---

- Abstract data structure for storing and retrieving values
- **Primary operations**
  - *Search( $x$ )*: returns the location of  $x$  in the dictionary, or NIL if not contained
  - *Insert( $x$ )*: adds  $x$  to the dictionary
  - *Delete( $x$ )*: removes  $x$  from the dictionary
- **Additional operations**
  - *Max()*, *Min()*: return the location of the largest/smallest element
  - *Successor( $x$ )*, *Predecessor( $x$ )*: return the next largest/smallest element than  $x$

# Array-based dictionary operations

- *Unsorted* array

- **Search(x)**
  - Linear scan
  - $O(n)$  time

```
1 Algorithm: Search(x)
2 for i = 1 to n do
3   |   if dict[i] = x then
4   |   |   return &dict[i];
5 end
6 return NIL;
```

- **Delete(x)**
  - Swap "victim" with last node
  - $O(1)$  time

```
1 Algorithm: Delete(x)
2 *x = dict[n];
3 n = n - 1;
```

# Array-based insertion

- *Unsorted* array
- **Insert(x)**
  - Depends on capacity ( $c$ )
  - If capacity is available:
    - Add element to next position
    - $O(1)$  time
  - If space is not available:
    - Allocate larger array
    - Copy elements from old array
    - Deallocate old array
    - Insert as normal
    - $O(1)$  time, amortized

```
1 Algorithm: Insert(x)
2 if  $n \geq c$  then
3   | Enlarge();
4 end
5  $n = n + 1$ ;
6  $\text{dict}[n] = x$ ;
```

```
1 Algorithm: Enlarge()
2  $c = 2c$ ;
3  $\text{newdict} = \text{allocate}(c)$ ;
4 for  $i = 1$  to  $n$  do
5   |  $\text{newdict}[i] = \text{dict}[i]$ 
6 end
7 free dict;
8  $\text{dict} = \text{newdict}$ ;
```

# Sorted array operations

- **Search(x)**
  - Binary search
  - $O(\lg n)$  time

```
1  Algorithm: Search(x)
2  lo = 1;
3  hi = n;
4  while lo < hi do
5      | mid = floor((lo + hi)/2);
6      | if data[mid] = x then
7          |   return &data[mid];
8      | else if data[mid] > x then
9          |   lo = mid+1;
10     | else
11     |   hi = mid-1;
12     | end
13 end
14 if lo ≤ n and data[lo] = x then
15     | return &data[lo];
16 else
17     | return NIL;
18 end
```

# Sorted array operations

- **Delete(x)**

- Shift elements left
- $O(n)$  time

```
1 Algorithm: Delete(x)
2 for i = x-data to n-1 do
3   | data[i] = data[i+1];
4 end
5 n = n-1;
```

- **Insert(x)**

- Enlarge (if needed), then shift right
- $O(n)$  time

```
1 Algorithm: Insert(x)
2 if n ≥ c then
3   | Enlarge();
4 end
5 Use binary search to find i such
  that data[i] ≤ x ≤ data[i+1];
6 for j = n to i+1 step -1 do
7   | data[j+1] = data[j];
8 end
9 data[i+1] = x;
10 n = n+1;
```

- ***Conversion from unsorted array***

- Sort
- $O(n \lg n)$  time

# Summary: array-based dictionaries

| Operation      | Unsorted array            | Sorted array |
|----------------|---------------------------|--------------|
| Search(x)      | $O(n)$                    | $O(\lg n)$   |
| Delete(x)      | $O(1)$                    | $O(n)$       |
| Insert(x)      | $O(1)$ , <i>amortized</i> | $O(n)$       |
| Build          | n/a                       | $O(n \lg n)$ |
| Min()          | $O(n)$                    | $O(1)$       |
| Max()          | $O(n)$                    | $O(1)$       |
| Predecessor(x) | $O(n)$                    | $O(1)$       |
| Successor(x)   | $O(n)$                    | $O(1)$       |

# List-based dictionaries

- Unsorted linked list
  - Largely similar to unsorted array
  - **Search(x)**
    - Linear scan
    - $O(n)$  time
  - **Insert(x)**
    - Append element
    - $O(1)$  time, with tail pointer
  - **Delete(x)**
    - Depends!
    - Singly-linked list: linear scan
      - $O(n)$  time

```
1 Algorithm: Delete-SLL(x)
2 if x = head then
3   | head = x.next;
4   | if x = tail then
5   |   | tail = NIL;
6   | free x;
7 else
8   | curr = head;
9   | while curr ≠ NIL and curr.next
10  |   | ≠ x do
11  |   |   | curr = curr.next;
12  |   | end
13  |   | if curr ≠ NIL then
14  |   |   | curr.next = x.next;
15  |   |   | if x = tail then
16  |   |   |   | tail = curr;
17  |   |   | free x;
18  |   | end
19 end
```

# List-based dictionaries

- Unsorted linked list
  - Largely similar to unsorted array
  - **Search(x)**
    - Linear scan
    - $O(n)$  time
  - **Insert(x)**
    - Append element
    - $O(1)$  time, with tail pointer
  - **Delete(x)**
    - Depends!
    - Singly-linked list: linear scan
      - $O(n)$  time
    - Doubly-linked list: pointer ops
      - $O(1)$  time

```
1  Algorithm: Delete-DLL(x)
2  if x = head then
3      head = x.next;
4      if head  $\neq$  tail then
5          | x.next.prev = NIL;
6          free x;
7  else if x = tail then
8      tail = x.prev;
9      x.prev.next = NIL;
10     free x;
11 else
12     x.prev.next = x.next;
13     x.next.prev = x.prev;
14     free x;
15 end
```



# Sorted linked lists

---

- **Search(x)**
  - No link to midpoint, so no binary search!
  - Linear scan
  - $O(n)$  time
- **Insert(x)**
  - Linear scan
  - $O(n)$  time
- **Delete(x)**
  - Same as unsorted
  - Singly-linked list:  $O(n)$  time
  - Doubly-linked list:  $O(1)$  time
- **Building a sorted list**
  - Possible in  $O(n \lg n)$  time

# Summary: link-based dictionaries

| Operation      | Unsorted SLL | Unsorted DLL | Sorted SLL   | Sorted DLL   |
|----------------|--------------|--------------|--------------|--------------|
| Search(x)      | $O(n)$       | $O(n)$       | $O(n)$       | $O(n)$       |
| Delete(x)      | $O(n)$       | $O(1)$       | $O(n)$       | $O(1)$       |
| Insert(x)      | $O(1)$       | $O(1)$       | $O(n)$       | $O(n)$       |
| Build          | n/a          | n/a          | $O(n \lg n)$ | $O(n \lg n)$ |
| Min()          | $O(n)$       | $O(n)$       | $O(1)$       | $O(1)$       |
| Max()          | $O(n)$       | $O(n)$       | $O(1)$       | $O(1)$       |
| Predecessor(x) | $O(n)$       | $O(n)$       | $O(n)$       | $O(1)$       |
| Successor(x)   | $O(n)$       | $O(n)$       | $O(1)$       | $O(1)$       |

- **Note:** DLL time is strictly better, asymptotically
  - Trade-off: more space, more pointer manipulation

# Coming up

---

- More data structures!
- **Homework 6** (posted tonight) will be due Thursday
- **Project 1** will be posted this weekend
  - Search algorithms, Big-Oh analysis
- **Exam 1** will be returned next week
- **Recommended readings (today):** Sections 3.1-3.4
- **Recommended readings (Tuesday):** Sections 3.5-3.9
- **Practice problems:** 1-2 problems from "Stacks, Queues, and Lists"