

CMP302: Algorithms Design and Analysis



Lecture 02: Sorting & Order Statistics

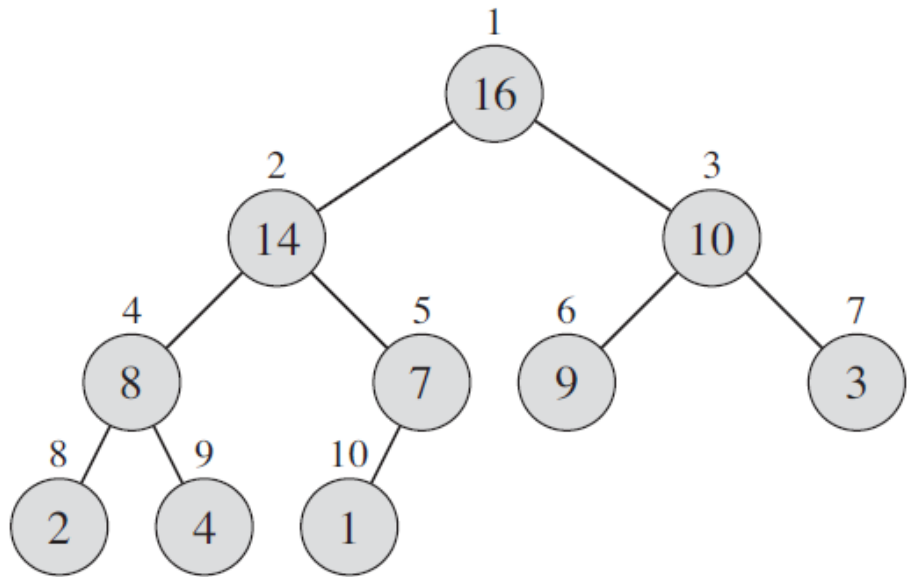
Ahmed Hamdy

Computer Engineering Department

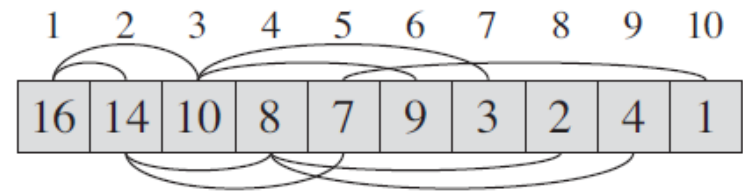
Cairo University

Fall 2017

Heaps



(a)



(b)

PARENT(i)

1 **return** $\lfloor i/2 \rfloor$

LEFT(i)

1 **return** $2i$

RIGHT(i)

1 **return** $2i + 1$

Heaps

Optimizations of heap operations

- $2i$ computed as shift left
- $2i + 1$ computed as shift left and adding 1 / ORing 1
- $\lfloor i/2 \rfloor$ computed as shift right
- Implement heap operations (parent, left, right) as macros or inline functions

Heaps

- Max-heap:

$$A[\text{PARENT}(i)] \geq A[i]$$

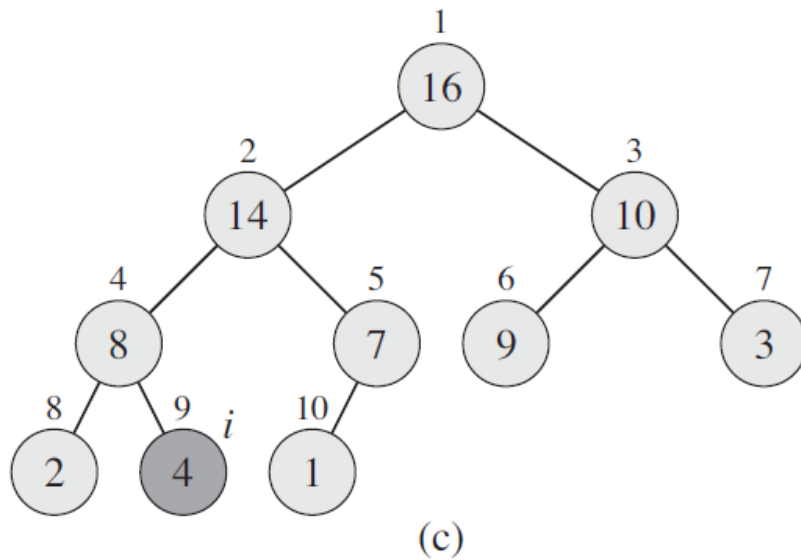
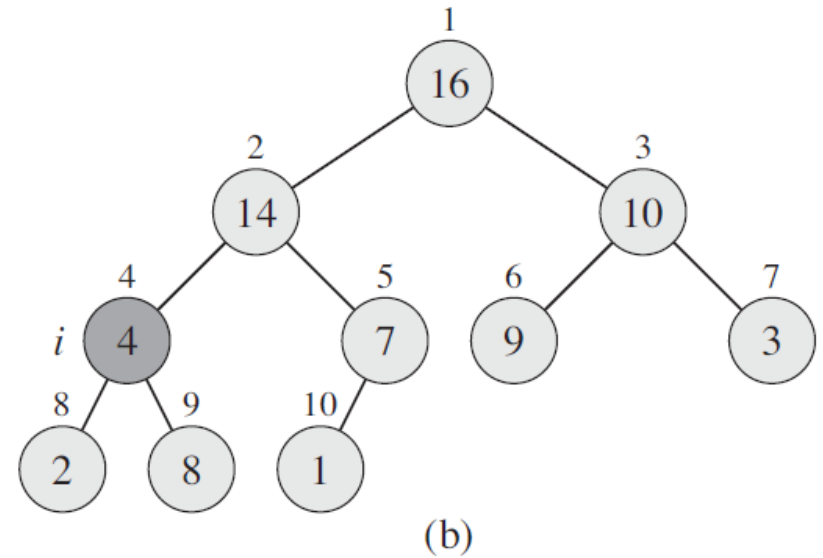
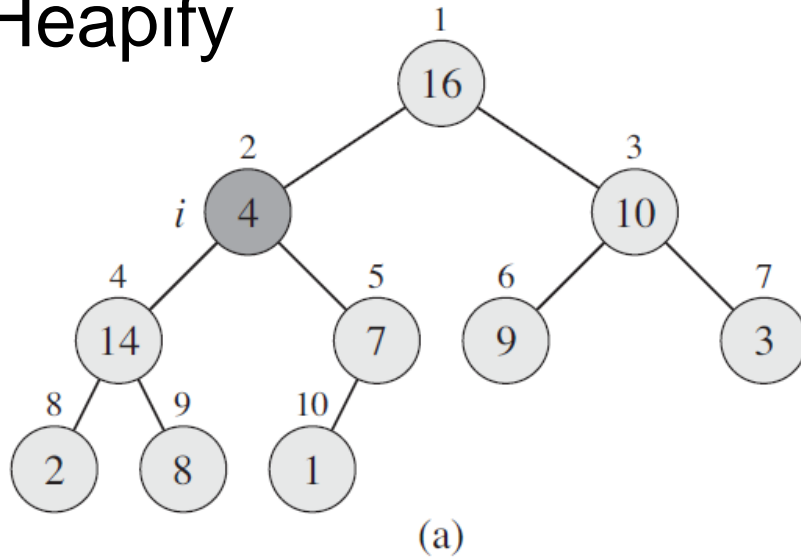
- Min-heap:

$$A[\text{PARENT}(i)] \leq A[i]$$

- Height: $\Theta(\log n)$
- Operations: $O(\log n)$

Heaps

Max-Heapify



Heaps

MAX-HEAPIFY(A, i)

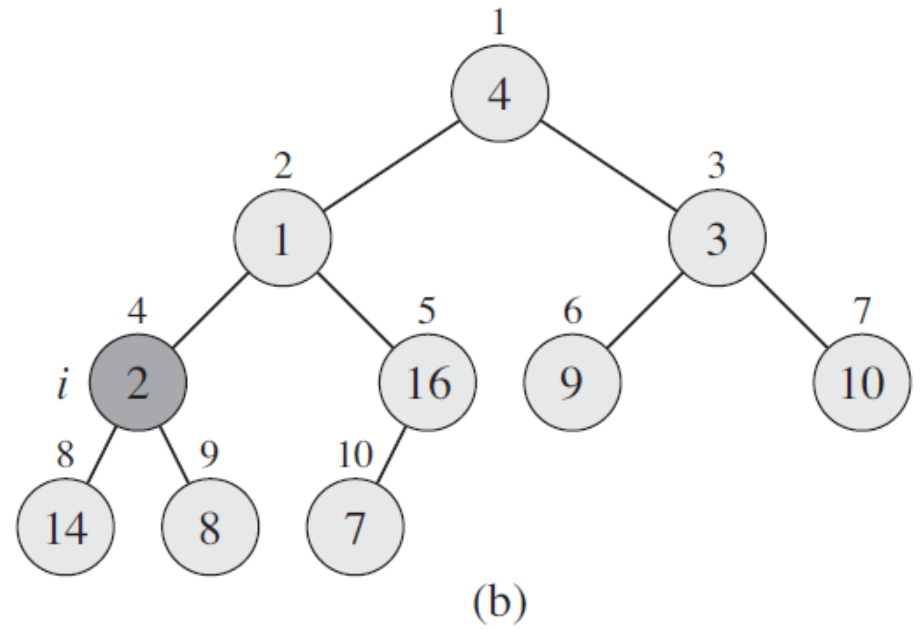
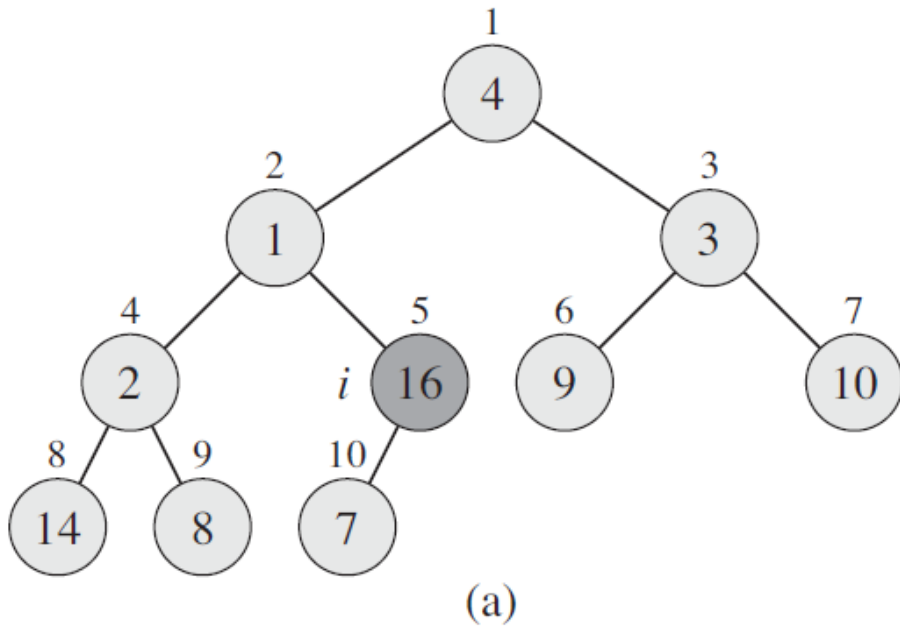
```
1   $l = \text{LEFT}(i)$ 
2   $r = \text{RIGHT}(i)$ 
3  if  $l \leq A.\text{heap-size}$  and  $A[l] > A[i]$ 
4       $\text{largest} = l$ 
5  else  $\text{largest} = i$ 
6  if  $r \leq A.\text{heap-size}$  and  $A[r] > A[\text{largest}]$ 
7       $\text{largest} = r$ 
8  if  $\text{largest} \neq i$ 
9      exchange  $A[i]$  with  $A[\text{largest}]$ 
10     MAX-HEAPIFY( $A, \text{largest}$ )
```

- Max-Heapify complexity: $O(\log n)$

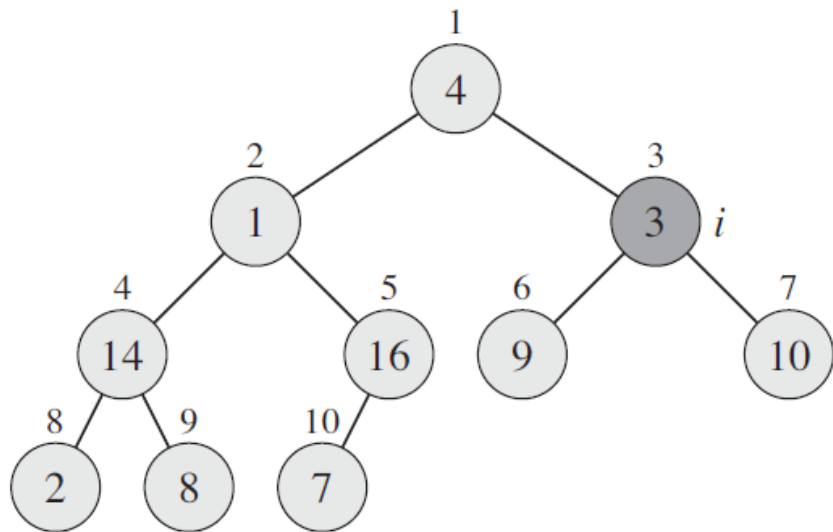
Heaps

A

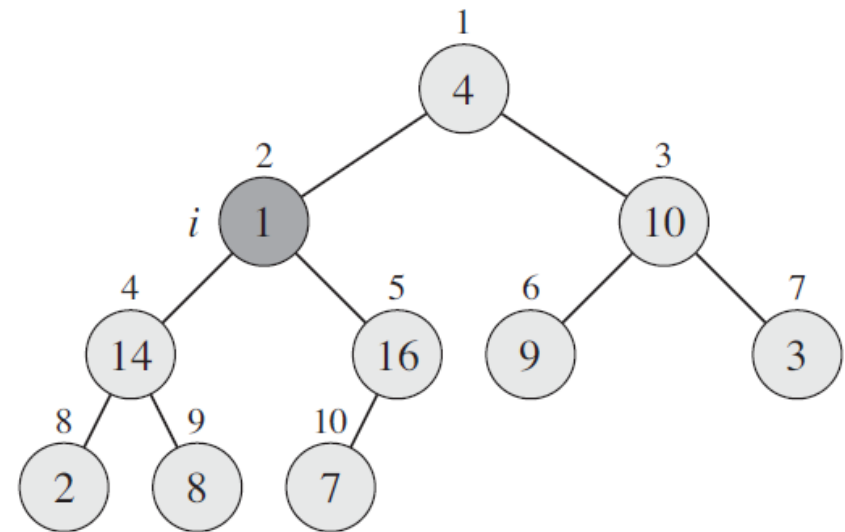
4	1	3	2	16	9	10	14	8	7
---	---	---	---	----	---	----	----	---	---



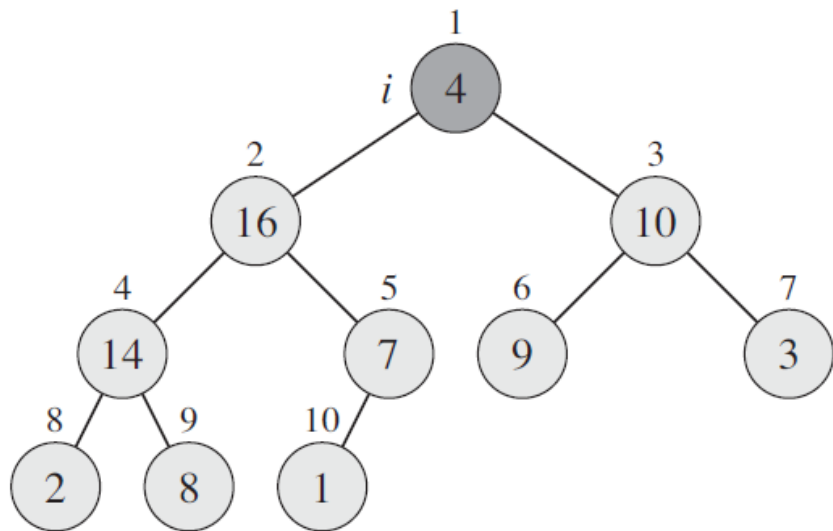
Heaps



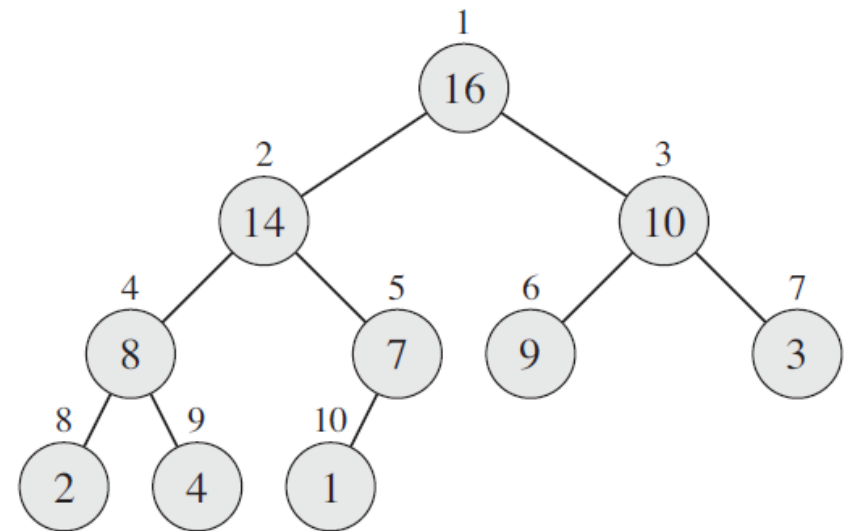
(c)



(d)



(e)



(f)

Heaps

BUILD-MAX-HEAP(A)

```
1   $A.heap-size = A.length$   
2  for  $i = \lfloor A.length/2 \rfloor$  downto 1  
3      MAX-HEAPIFY( $A, i$ )
```

- Build-Max-Heap Complexity:
 - $O(n)$ calls to Max-Heapify
 - Thus overall complexity of Build-Max-Heap: $O(n \log n)$
 - Is this tight bound??

Heaps

- Height of n -element heap: $\lfloor \log n \rfloor$
- Number of nodes at height h : $\left\lceil \frac{n}{2^{h+1}} \right\rceil$
- Build-Max-Heap is bounded by

$$\sum_{h=0}^{\lfloor \lg n \rfloor} \left\lceil \frac{n}{2^{h+1}} \right\rceil O(h) = O \left(n \sum_{h=0}^{\lfloor \lg n \rfloor} \frac{h}{2^h} \right)$$

- Using

$$\sum_{k=0}^{\infty} k x^k = \frac{x}{(1-x)^2} \quad \Rightarrow \quad \sum_{h=0}^{\infty} \frac{h}{2^h} = \frac{1/2}{(1-1/2)^2} = 2.$$

for $|x| < 1$.

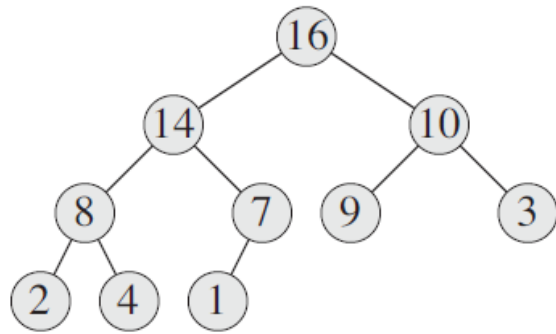
Heaps

$$\sum_{h=0}^{\lfloor \lg n \rfloor} \left\lceil \frac{n}{2^{h+1}} \right\rceil O(h) = O \left(n \sum_{h=0}^{\lfloor \lg n \rfloor} \frac{h}{2^h} \right)$$

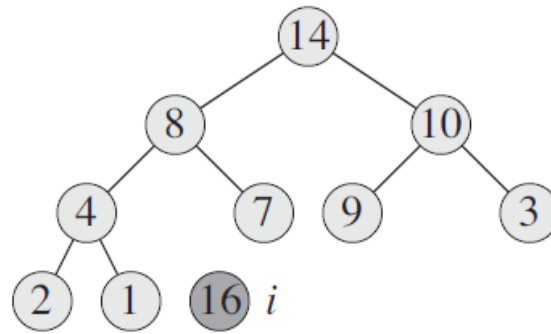
Becomes

$$\begin{aligned} O \left(n \sum_{h=0}^{\lfloor \lg n \rfloor} \frac{h}{2^h} \right) &= O \left(n \sum_{h=0}^{\infty} \frac{h}{2^h} \right) \\ &= O(n) . \end{aligned}$$

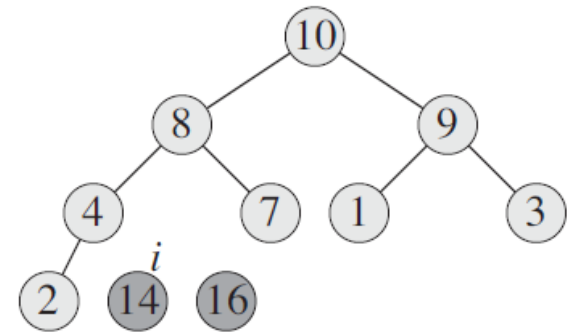
Heapsort



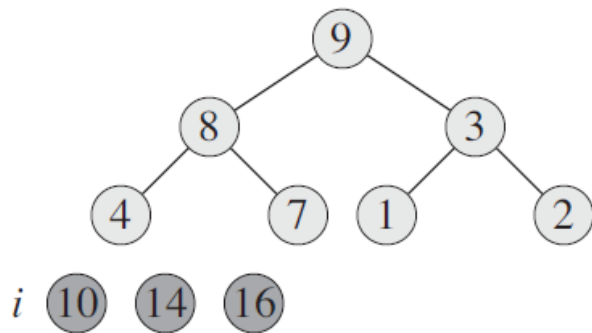
(a)



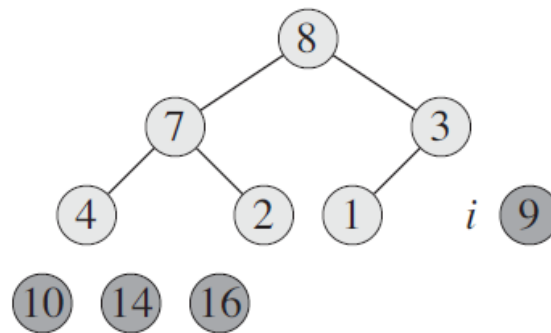
(b)



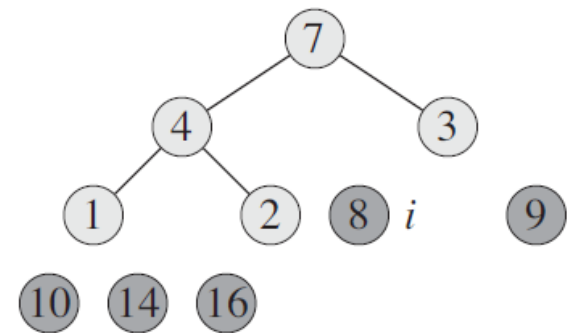
(c)



(d)

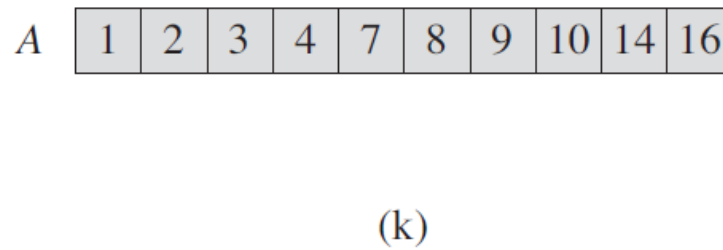
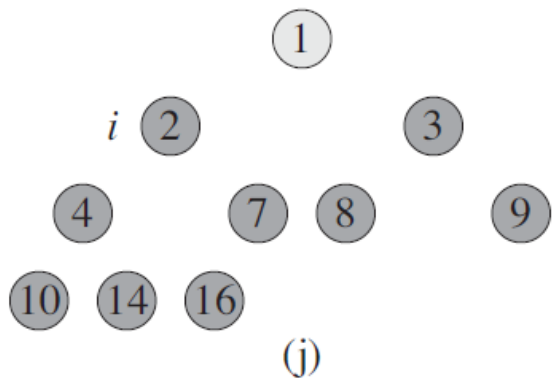
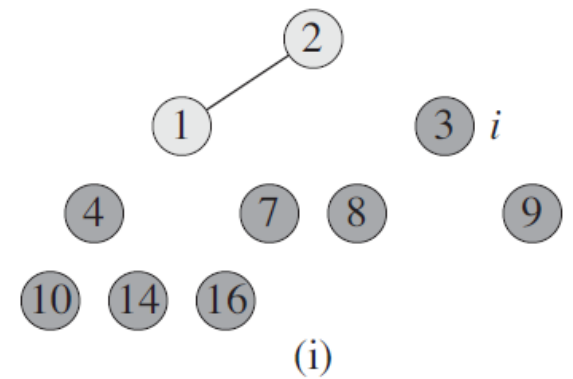
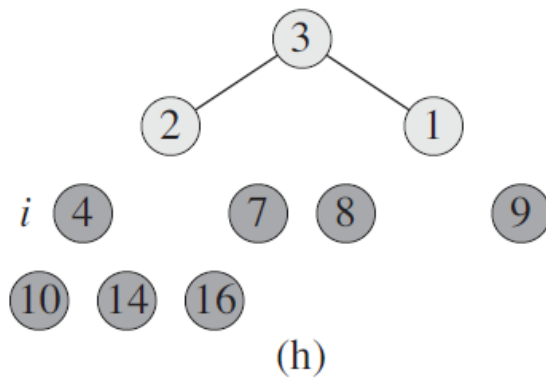
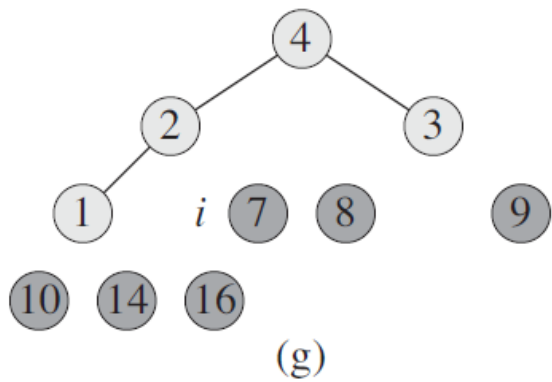


(e)



(f)

Heapsort



Heapsort

HEAPSORT(A)

```
1  BUILD-MAX-HEAP( $A$ )
2  for  $i = A.length$  downto 2
3      exchange  $A[1]$  with  $A[i]$ 
4       $A.heap-size = A.heap-size - 1$ 
5      MAX-HEAPIFY( $A, 1$ )
```

Complexity

- $O(n)$ calls to Max-Heapify
- Thus overall complexity of Heapsort: $O(n \log n)$

Priority queues

- Useful in applications where priority is the criteria for selection
- Example, scheduling jobs
- Max-priority or min-priority queues

Priority queues

HEAP-MAXIMUM(A)

```
1  return  $A[1]$ 
```

HEAP-EXTRACT-MAX(A)

```
1  if  $A.heap-size < 1$   
2      error “heap underflow”  
3   $max = A[1]$   
4   $A[1] = A[A.heap-size]$   
5   $A.heap-size = A.heap-size - 1$   
6  MAX-HEAPIFY( $A, 1$ )  
7  return  $max$ 
```


Priority queues

HEAP-INCREASE-KEY(A, i, key)

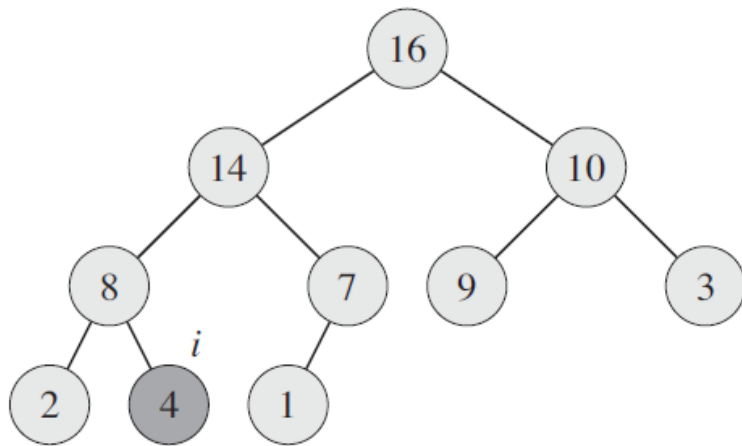
```
1  if  $key < A[i]$ 
2      error “new key is smaller than current key”
3   $A[i] = key$ 
4  while  $i > 1$  and  $A[\text{PARENT}(i)] < A[i]$ 
5      exchange  $A[i]$  with  $A[\text{PARENT}(i)]$ 
6       $i = \text{PARENT}(i)$ 
```

MAX-HEAP-INSERT(A, key)

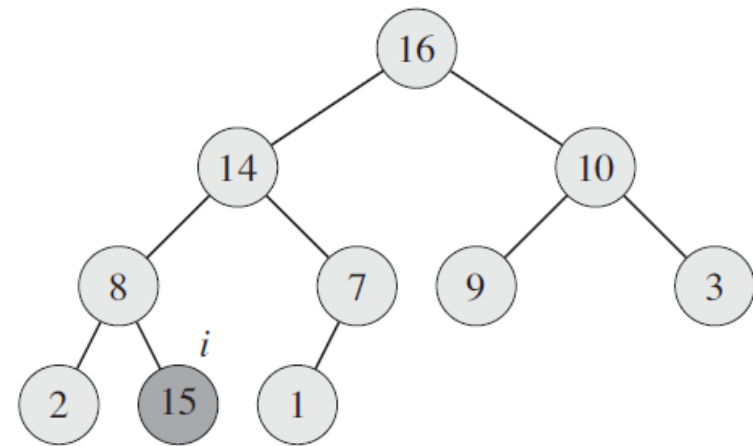
```
1   $A.\text{heap-size} = A.\text{heap-size} + 1$ 
2   $A[A.\text{heap-size}] = -\infty$ 
3  HEAP-INCREASE-KEY( $A, A.\text{heap-size}, key$ )
```

Priority queue

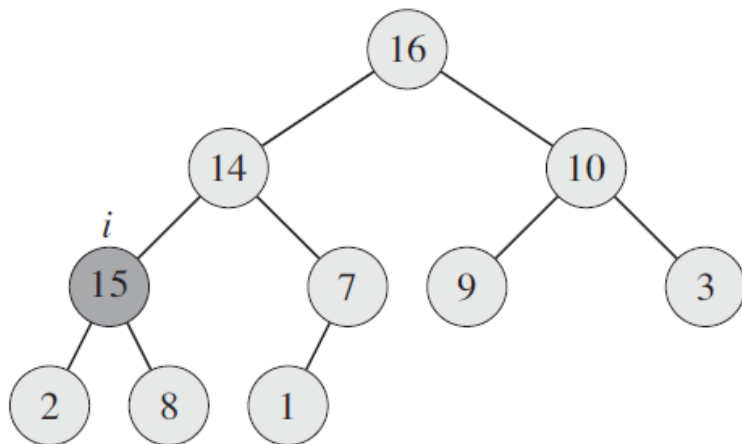
- Heap-Increase-Key



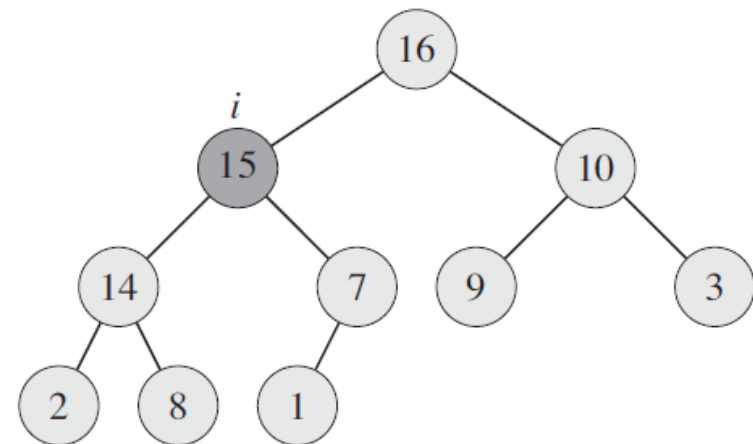
(a)



(b)

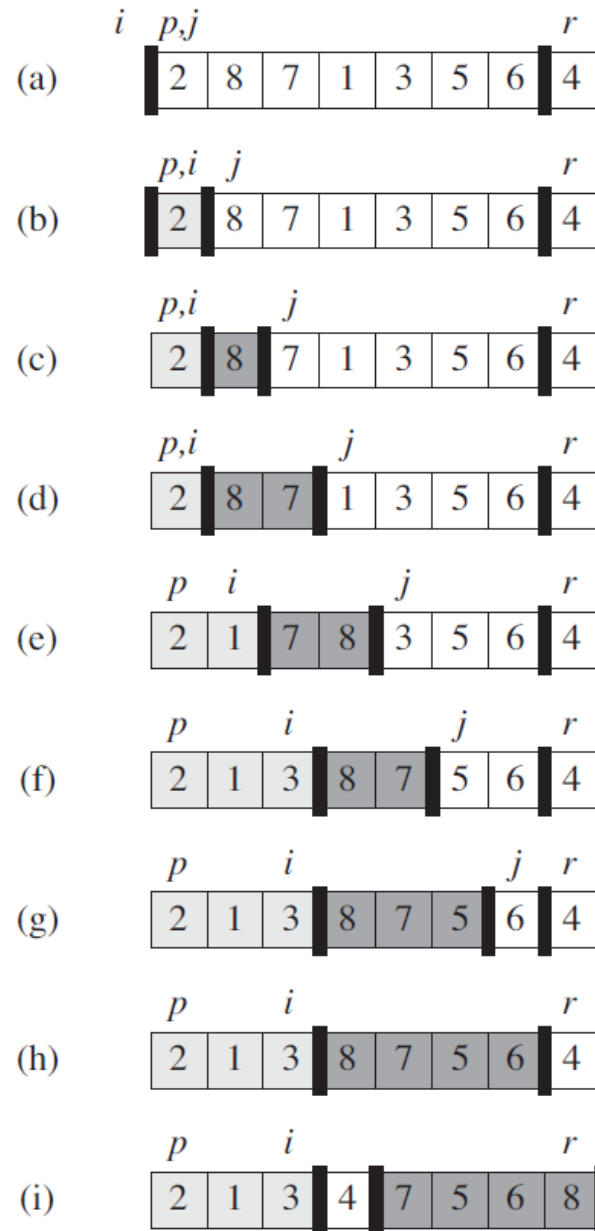


(c)



(d)

Quicksort



Quicksort

QUICKSORT(A, p, r)

```
1  if  $p < r$ 
2       $q = \text{PARTITION}(A, p, r)$ 
3      QUICKSORT( $A, p, q - 1$ )
4      QUICKSORT( $A, q + 1, r$ )
```

PARTITION(A, p, r)

```
1   $x = A[r]$ 
2   $i = p - 1$ 
3  for  $j = p$  to  $r - 1$ 
4      if  $A[j] \leq x$ 
5           $i = i + 1$ 
6          exchange  $A[i]$  with  $A[j]$ 
7  exchange  $A[i + 1]$  with  $A[r]$ 
8  return  $i + 1$ 
```

Quicksort

Performance

- Worst case partitioning:

$$\begin{aligned}T(n) &= T(n-1) + T(0) + \Theta(n) \\ &= T(n-1) + \Theta(n) .\end{aligned}$$

$$T(n) = \Theta(n^2)$$

- Best case partitioning:

$$T(n) = 2T(n/2) + \Theta(n)$$

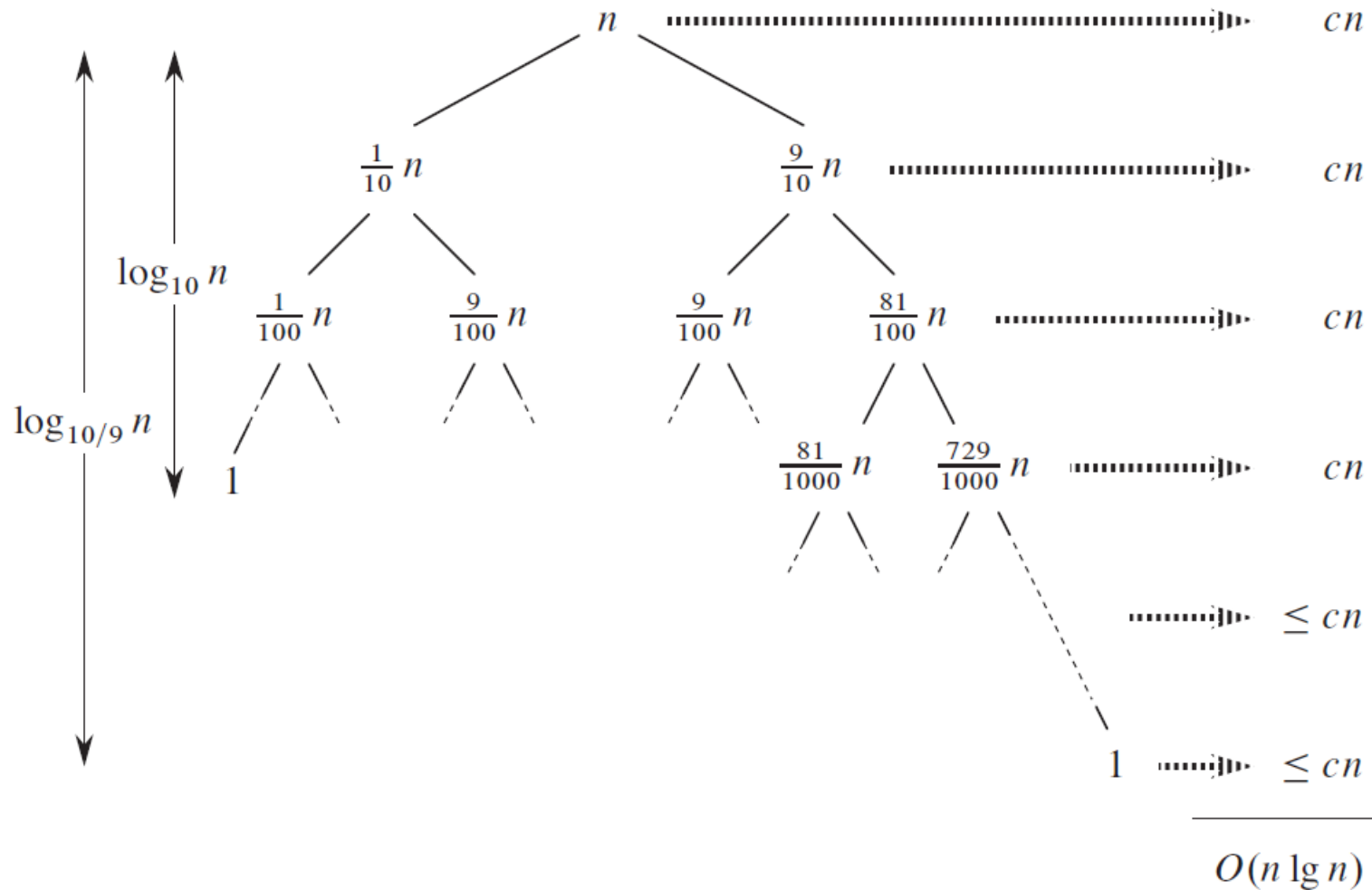
$$T(n) = \Theta(n \lg n)$$

Quicksort

Performance

- Balanced partitioning:

$$T(n) = T(9n/10) + T(n/10) + cn$$



Randomized quicksort

RANDOMIZED-PARTITION(A, p, r)

```
1   $i = \text{RANDOM}(p, r)$   
2  exchange  $A[r]$  with  $A[i]$   
3  return PARTITION( $A, p, r$ )
```

RANDOMIZED-QUICKSORT(A, p, r)

```
1  if  $p < r$   
2       $q = \text{RANDOMIZED-PARTITION}(A, p, r)$   
3      RANDOMIZED-QUICKSORT( $A, p, q - 1$ )  
4      RANDOMIZED-QUICKSORT( $A, q + 1, r$ )
```

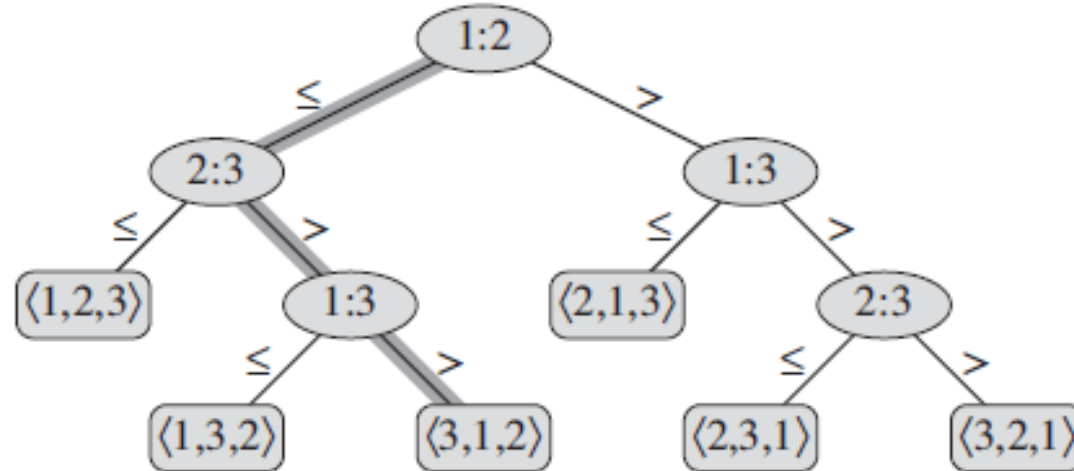
- Expected running time: $O(n \log n)$ when element values are distinct

Comparison sort

- Insertion-sort, mergesort, heapsort and quicksort all use comparisons to gain order
- Can we obtain a lower bound for any comparison sort?

Comparison sort

- Decision tree



- Number of leaves in decision tree: $\geq n!$
- Number of leaves in binary tree: 2^h

$$n! \leq l \leq 2^h$$

$$\begin{aligned} h &\geq \lg(n!) \\ &= \Omega(n \lg n) \end{aligned}$$

Comparison sort

- Mergesort and heapsort are asymptotically optimal comparison sorts
- Both have upper bound of $O(n \log n)$
- Quicksort isn't, why?

Counting sort

	1	2	3	4	5	6	7	8
<i>A</i>	2	5	3	0	2	3	0	3
	0	1	2	3	4	5		
<i>C</i>	2	0	2	3	0	1		

(a)

	0	1	2	3	4	5
<i>C</i>	2	2	4	7	7	8

(b)

	1	2	3	4	5	6	7	8
<i>B</i>							3	
	0	1	2	3	4	5		
<i>C</i>	2	2	4	6	7	8		

(c)

	1	2	3	4	5	6	7	8
<i>B</i>		0					3	
	0	1	2	3	4	5		
<i>C</i>	1	2	4	6	7	8		

(d)

	1	2	3	4	5	6	7	8
<i>B</i>		0				3	3	
	0	1	2	3	4	5		
<i>C</i>	1	2	4	5	7	8		

(e)

	1	2	3	4	5	6	7	8
<i>B</i>	0	0	2	2	3	3	3	5

(f)

Counting sort

COUNTING-SORT(A, B, k)

```
1  let  $C[0..k]$  be a new array
2  for  $i = 0$  to  $k$ 
3       $C[i] = 0$ 
4  for  $j = 1$  to  $A.length$ 
5       $C[A[j]] = C[A[j]] + 1$ 
6  //  $C[i]$  now contains the number of elements equal to  $i$ .
7  for  $i = 1$  to  $k$ 
8       $C[i] = C[i] + C[i - 1]$ 
9  //  $C[i]$  now contains the number of elements less than or equal to  $i$ .
10 for  $j = A.length$  downto 1
11      $B[C[A[j]]] = A[j]$ 
12      $C[A[j]] = C[A[j]] - 1$ 
```

When $k = O(n)$, the sort runs in $\Theta(n)$ time.

Counting sort

COUNTING-SORT(A, B, k)

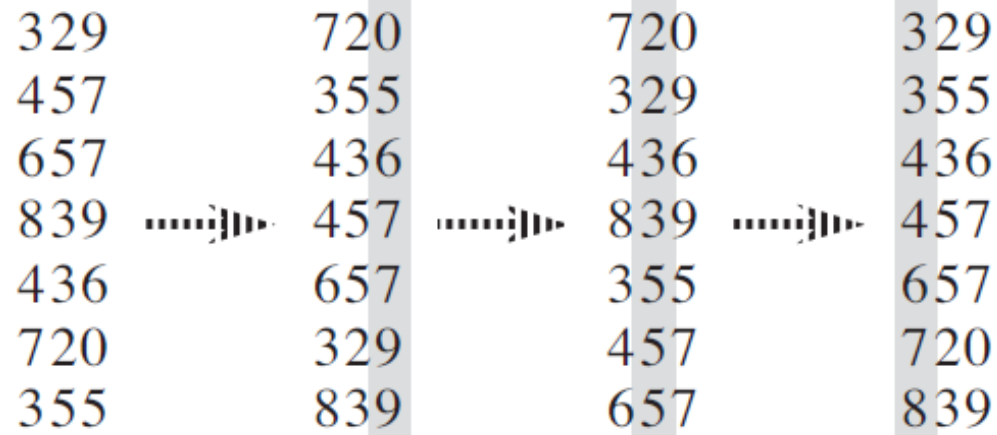
```
1  let  $C[0..k]$  be a new array
2  for  $i = 0$  to  $k$ 
3       $C[i] = 0$ 
4  for  $j = 1$  to  $A.length$ 
5       $C[A[j]] = C[A[j]] + 1$ 
6  //  $C[i]$  now contains the number of elements equal to  $i$ .
7  for  $i = 1$  to  $k$ 
8       $C[i] = C[i] + C[i - 1]$ 
9  //  $C[i]$  now contains the number of elements less than or equal to  $i$ .
10 for  $j = A.length$  downto 1
11      $B[C[A[j]]] = A[j]$ 
12      $C[A[j]] = C[A[j]] - 1$ 
```

- Complexity: $\Theta(n + k)$
- Typically used when $k = O(n)$, complexity $\Theta(n)$
- Stable

Sorting problem

- Sort a set of files according to their timestamps in ascending order??
- Assume just year, month, day.

Radix sort



RADIX-SORT(A, d)

- 1 **for** $i = 1$ **to** d
- 2 use a stable sort to sort array A on digit i

- Use counting sort to sort each digit
- Complexity: $\Theta(d(n + k))$

Radix sort

- What is the optimal d to use given n b -bit numbers?
- Ask it in a different way, what is the optimal r -bits out of the b -bits to use for every digit ($d = \lceil b/r \rceil$)?

$$T(n, b) = \Theta(d(n + k)) = \Theta\left(\frac{b}{r}(n + 2^r)\right)$$

Radix sort

$$T(n, b) = \Theta(d(n + k)) = \Theta\left(\frac{b}{r}(n + 2^r)\right)$$

- As r decreases ($2^r \ll n$), $\frac{b}{r}$ increases and $(n + 2^r)$ stays the same as $\Theta(n)$
- As r increases ($2^r \gg n$), $\frac{b}{r}$ decreases but $(n + 2^r)$ increases exponentially.
- Logically, $(n + 2^r) = \Theta(n)$, then $r = \lfloor \log n \rfloor$

$$T(n, b) = \Theta\left(\frac{b}{r}(n + 2^r)\right) = \Theta\left(\frac{bn}{\log n}\right)$$

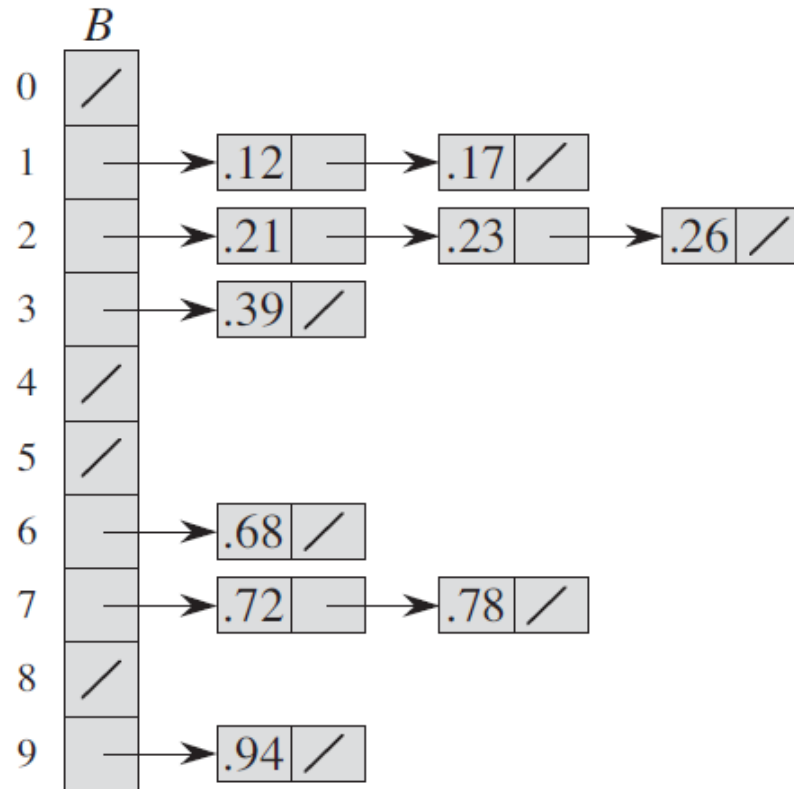
Radix sort

- If $b = O(\log n)$, radix sort becomes $\Theta(n)$
- Counting sort is NOT in-place
- Quicksort is better
 - Better cache utilization
 - In-place

Bucket sort

	<i>A</i>
1	.78
2	.17
3	.39
4	.26
5	.72
6	.94
7	.21
8	.12
9	.23
10	.68

(a)



(b)

Bucket sort

BUCKET-SORT(A)

```
1  let  $B[0 \dots n - 1]$  be a new array
2   $n = A.length$ 
3  for  $i = 0$  to  $n - 1$ 
4      make  $B[i]$  an empty list
5  for  $i = 1$  to  $n$ 
6      insert  $A[i]$  into list  $B[\lfloor nA[i] \rfloor]$ 
7  for  $i = 0$  to  $n - 1$ 
8      sort list  $B[i]$  with insertion sort
9  concatenate the lists  $B[0], B[1], \dots, B[n - 1]$  together in order
```

Bucket sort

- When input is uniformly distributed, the average case is $O(n)$
- It is fast, as counting-sort because it makes an assumption about the input
- **RULE OF THUMB:**
 - Extra assumption/extra information → room for optimization/customization

Order statistics

- The i^{th} *order statistic* of a set of n elements is the i^{th} *smallest element*

Order statistics

- Minimum:

MINIMUM(*A*)

```
1  min = A[1]
2  for i = 2 to A.length
3      if min > A[i]
4          min = A[i]
5  return min
```

- Requires $O(n)$, optimal as $n - 1$ comparisons are needed

Selection algorithm

- Find i^{th} *smallest element*:

RANDOMIZED-SELECT(A, p, r, i)

```
1  if  $p == r$ 
2      return  $A[p]$ 
3   $q = \text{RANDOMIZED-PARTITION}(A, p, r)$ 
4   $k = q - p + 1$ 
5  if  $i == k$            // the pivot value is the answer
6      return  $A[q]$ 
7  elseif  $i < k$ 
8      return RANDOMIZED-SELECT( $A, p, q - 1, i$ )
9  else return RANDOMIZED-SELECT( $A, q + 1, r, i - k$ )
```

- $O(n)$, prove it?

Selection algorithm

- On average, the left or right subarray that will be selected for next iteration will be of size $n/2$
- The recurrence becomes

$$T(n) = T\left(\frac{n}{2}\right) + \Theta(n)$$

- Using the master method with $a = 1$, $b = 2$, $f(n) = \Theta(n)$
- Thus, $O(n)$ on average