

#### James 1:19-21 (NASB)

This you know, my beloved brethren: but everyone must be quick to hear, slow to speak and slow to anger;

for the **anger of man**does not achieve the **righteousness of God**.

Therefore, putting aside all filthiness and all that remains of wickedness, in humility receive the word implanted, which is able to save your souls.

- Heap sort (ch5)
  - Intro to trees (more in ch12)
  - Binary heaps and max-heaps
  - Heap sort
  - Max-heaps for priority queue
- Quicksort (ch6)
  - Lomuto partitioning and complexity analysis
  - Randomised Quicksort and analysis
- Monte-Carlo matrix multiply checking

cmpt231.seanho.com/lec3

# Summary of sorting algorithms

- Comparison sorts (ch2, 6, 7):
  - Insertion sort:  $\Theta(n^2)$ , easy to program, slow
  - Merge sort: ⊖(n lg n), out-of-place copy (slow)
  - Heap sort: ⊕(n lg n), in-place, max-heap
  - Quicksort: ⊖(n^2) worst-case, ⊖(n lg n) average
    - and small (fast) constant factors
- Linear-time non-comparison sorts (ch8):
  - Counting sort: k distinct values: ⊖(k)
  - Radix sort: d digits, k values: Θ(d(n+k))
  - Bucket sort: uniform distribution: ⊖(n)
- A sort is stable if preserves order of equal items



### **Binary trees**

- Graph: collection of nodes and edges
  - Edges may be directed or undirected
- Tree: directed acyclic graph (DAG)
  - One node designated root
  - Parent: immediate neighbour toward root
  - Leaf: node with no children
  - Degree: maximum number of children per node
  - Height of node: max num edges to leaf descendant
  - Depth of node: num edges to root
  - Level: all nodes of same depth

- Heap sort (ch5)
  - Intro to trees (more in ch12)
  - Binary heaps and max-heaps
  - Heap sort
  - Max-heaps for priority queue
- Quicksort (ch6)
  - Lomuto partitioning and complexity analysis
  - Randomised Quicksort and analysis
- Monte-Carlo matrix multiply checking

### Binary heaps

- Array storage for certain types of binary trees:
  - Put node i's two children at 2i and 2i+1
  - Fill tree left-to-right, one level at a time
  - e.g.: [ 2, 8, 4, 7, 5, 3, 1, 6 ]
- The max-heap property: every node's value is ≤ its parent
  - (min-heap: ≥)
- max\_heapify()(O(lg n)):
  - move a node i to satisfy max-heap property
- build\_max\_heap()(O(n)):

turn an unordered array into a max-heap

# max\_heapify() on single node

- Input: binary heap A and node index i
  - Precondition: left and right sub-trees of i are each separate max-heaps
- Postcondition: entire subtree at i is a max-heap
- Algorithm:
  - 1. Find largest of: i, left child of i, or right child of i
  - 2. If i is **not** the largest, then:
    - a. Swap i with the largest
    - b. Recurse (or iterate) on that subtree



### max\_heapify()

```
def max_heapify( A, i ):
 max = i
 if 2i \le length(A) and A[2i] > A[max]:
   max = 2i
 else if 2i+1 \le length(A) and A[2i+1] > A[max]: # right
   max = 2i+1
 if max != i:
   swap(A[i],A[max])
   max_heapify( A, max )
```

- Try it on previous heap at i=1
- Running time?

# Building a max-heap

- Input: array A, in any order
  - Postcondition: A is a max-heap
- Algorithm:
  - 1. Last half of array is all leaves
  - 2. Run max\_heapify() on each item in first half
    - Descending order: subtrees are already max-heaps

```
for i = floor( length(A)/2 ) to 1:
  max_heapify( A, i )
```

### Max-heap: complexity

- Group iterations of for loop by height h of node:
  - Each call to max\_heapify(i) takes O(h)
  - Num of nodes with height h is  $\leq \left\lfloor \frac{n}{2^{h+1}} \right\rfloor$ 
    - Reaches that bound when tree is full
- Total running time T(n):  $T(n) = \sum_{h=0}^{\lg n} \left(\frac{n}{2^{h+1}}\right) O(h)$

$$0 \leq n \sum_{h=0}^{\infty} \left( rac{1}{2^{h+1}} 
ight) O(h) \ = n \sum_{h=1}^{\infty} \left( rac{1}{2^h} 
ight) O(h) \ = O(n)$$

- → Can build a max-heap in linear time!
  - But it's not quite a sorting algorithm....

#### Heap sort

- Algorithm:
  - 1. Make array a max-heap
  - 2. Repeat, working backwards from end of array:
    - Swap root with last leaf of heap
    - Shrink heap by 1 and re-apply max\_heapify()
- Loop invariant:

cmpt231.seanho.com/lec3

- First portion of array is still a max-heap
- Last portion of array is sorted (largest items)
- Complexity:  $T(n) = \Theta(n \lg n)$ 
  - $\Theta(n)$  calls to max\_heapify() ( $\Theta(\lg n)$ )

7 /

- Heap sort (ch5)
  - Intro to trees (more in ch12)
  - Binary heaps and max-heaps
  - Heap sort
  - Max-heaps for priority queue
- Quicksort (ch6)
  - Lomuto partitioning and complexity analysis
  - Randomised Quicksort and analysis
- Monte-Carlo matrix multiply checking

## **Priority queue**

- We can use binary heaps to make a priority queue:
  - Set of items with attached priorities
- Interface (set of operations) for priority queue:
  - insert(A, item, pri):add anitem
  - find\_max(A): get item with highest priority
  - pop\_max(A): same but also delete the item
  - set\_pri(A, item, pri):change
    (increase) priority of item
- Initialise queue by building a max-heap:
  - find\_max() is easy: just return A[1]
- pop max ( ) also easy: remove A[1] and

14/35

### Insert into queue

 set\_pri(): start at i and bubble up to proper place:

```
A[ i ] = pri
while i > 1 and A[ i/2 ] < A[ i ]:
```

```
swap(A[i/2], A[i]) i = i/2
```

- Complexity: num iterations =  $\Theta(\lg n)$
- insert(): make new node, then set its priority:

```
A.size++
A[ size ] = item
set_pri( A, A.size, pri )
```

■ Complexity: same as set\_pri(): ⊖(lg n)

cmpt23f.seafhod: speed, often pre-allocate A as fixed-length array.135

# Priority queue operations

- Build queue (with max-heap): ⊖(n)
- Fetch highest priority item: ⊖(1)
- Fetch and remove highest priority item: ⊖(lg n)
- Change priority of an item: ⊖(lg n)
- Insert new item: ⊖(lg n)

- Heap sort (ch5)
  - Intro to trees (more in ch12)
  - Binary heaps and max-heaps
  - Heap sort
  - Max-heaps for priority queue
- Quicksort (ch6)
  - Lomuto partitioning and complexity analysis
  - Randomised Quicksort and analysis
- Monte-Carlo matrix multiply checking

cmpt231.seanho.com/lec3

# Randomised algorithms

- Vegas-style: always correct, fast on average
  - But still slow in worst-case
- Monte Carlo: always fast
  - But not always correct!
  - Approximate: margin of error ε
  - Stochastic: probability P of being correct
  - Estimate improves with more computation time (iterations)

### Quicksort

- **Divide**: partition A[lo..hi] such that:
  - max(A[lo..piv-1]) ≤ A[piv] ≤ min(A[piv+1..
     hi])
  - Not always balanced; this is the "magic sauce"
- Conquer: recurse on each part:
  - quicksort(A, lo, piv-1) and quicksort(A, piv+1, hi)
  - No combine / merge step needed
- In-place sort (only uses swaps) (unlike merge sort)
- Worst case still  $\Theta(n^2)$ , but
  - **Average** case is  $\Theta(n \lg n)$ , with small constants
- One of the best sorts for arbitrary inputs

# Partitioning (Lomuto)

- One option: pick last item as the pivot
- Walk through array from left to right:
  - Throw items smaller than pivot to left part of array
  - Items larger than pivot stay in right part of array
- Lastly, swap pivot in-between two parts

# **Quicksort:** complexity

- Worst case: every partition is maximally uneven:
  - Pivot is either largest or smallest in subarray
  - $T(n) = T(n-1) + T(0) + \Theta(n) = \Theta(n^2)$
  - Example inputs that do this?
- Best case: every partition is exactly half:
  - $T(n) = 2T(n/2) + \Theta(n) = \Theta(n \lg n)$
  - Example inputs that give this?
- What about average case, assuming random input?

### Average case complexity

- Intuition: on average, get splits in-between best and worst
  - If say, average split is 90% vs 10%, then:
  - $T(n) = T(0.90n) + T(*0.10n) + *\Theta(n)$
  - Still results in O(n lg n)
- If assume splits alternate between best and worst:
  - Only adds O(n) work to each of O(lg n) levels
  - Still O(n lg n)! (But maybe larger constants)

### Quicksort with constant splits

- (p.178 #7.2-5): All splits are v vs 1-v, with 0 < v < 1/2
  - ⇒ what is min/max depth of leaf in recursion tree?
- Min depth: follow smaller (v) side of each split
  - How many splits m until reach leaf (1 item)

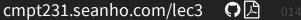
array)? 
$$lpha^m n = 1 \Rightarrow m = -rac{\log(n)}{\log(lpha)}$$

- Max depth: same with 1-v side:  $-\frac{\log(n)}{\log(1-\alpha)}$
- Both are Θ(log n)



# Quicksort with median split

- Best case splits happen when pivot is median:
  - Half of items smaller, half of items larger
  - Not same as average when distribution is skewed
- Median (rank finding) algorithm in O(n): see ch9
  - Partitioning also takes only O(n), so
  - Quicksort T(n) = 2T(n/2) + O(n) = O(n lg n) (always!)
- But, in practise:
  - Extra work, splits are usually already good
  - Benchmarks slower than merge sort



- Heap sort (ch5)
  - Intro to trees (more in ch12)
  - Binary heaps and max-heaps
  - Heap sort
  - Max-heaps for priority queue
- Quicksort (ch6)
  - Lomuto partitioning and complexity analysis
  - Randomised Quicksort and analysis
- Monte-Carlo matrix multiply checking

cmpt231.seanho.com/lec3

# Randomised Quicksort

- With random input, get nice ⊖(n lg n) behaviour
- But presorted input gives worst case behaviour
  - Much real-world data is at least partially presorted
  - Choice of last element (hi) as pivot
- Great candidate for randomised algorithm:
  - Before partitioning, swap hi with a random item
- Still possible to get worst-case behaviour, but unlikely
  - Vegas-style: always correct, and usually fast

### R-Quicksort: complexity

- Assume all items distinct
- Name items according to true order:  $\{z_i\}_{i=1}^n$
- Analyse complexity by counting comparisons performed
  - Worst case: compare all pairs  $(z_i, z_j)$ :  $\Theta(n^2)$
- No comparison can happen multiple times, because
  - Comparisons are only done against pivots, and
  - Each pivot is used only once and not revisited
- So what is the **probability** of a pair  $(z_i, z_j)$  being compared?

# R-Quicksort: pair comparison

- A pair  $(z_i, z_j)$  is compared only if:
  - Either  $z_i$  or  $z_j$  is chosen as a pivot before any other item ordered in-between them:

$$\{z_i, z_{i+1}, ..., z_{j-1}, z_j\}$$

- Otherwise  $z_i$  and  $z_j$  would be on opposite sides of a split, and would never be compared
- Probability of this happening:  $2\left(\frac{1}{j-i+1}\right)$

# R-Quicksort: total comparisons

**Sum** over all pairs  $(z_i, z_j)$ :

$$\sum_{i=1}^{n-1} \sum_{j=i+1}^{n} rac{2}{j-i+1}$$
  $= \sum_{i=1}^{n-1} \sum_{k=1}^{n-i} rac{2}{k+1}$  (let k=j-i)  $< \sum_{i=1}^{n-1} \sum_{k=1}^{n} rac{2}{k}$ 

$$=\sum_{i=1}^{n-1}O(\ln n)$$
 (harmonic series)

$$=O(n {
m lg} n)$$

- Heap sort (ch5)
  - Intro to trees (more in ch12)
  - Binary heaps and max-heaps
  - Heap sort
  - Max-heaps for priority queue
- Quicksort (ch6)
  - Lomuto partitioning and complexity analysis
  - Randomised Quicksort and analysis
- Monte-Carlo matrix multiply checking

cmpt231.seanho.com/lec3

#### Randomised mat-mul check

- Recall matrix multiply: naive  $\Theta(n^3)$ , Strassen  $\overline{\Theta}(n^{2.81})$ 
  - Best-known: Coppersmith-Winograd,  $\Theta(n^{2.376})$
- What if we have 3 n x n matrices A, B, C:
  - Check if A \* B = C, faster than full multiply?
- Frievald's matrix-multiply checker in  $\Theta(n^2)$ :
  - If A \* B = C, always returns True (0% falsenegatives)
  - If A \* B ≠ C, returns False > 50% of the time
- If returns True, run it k times:
  - False-positive rate  $< 2^{-k}$ , in time  $O(kn^2)$

## Frievald's algorithm

- Make a random boolean vector  $\vec{r} = \{r_i\}_1^n$ :
  - $P(r_i = 1) = 0.5$  for all i, independently
  - i.e., flip a fair coin n times
- Return value: check if  $A \cdot (B \cdot \vec{r}) = C \cdot \vec{r}$ 
  - Each multiply is only a (n x n) matrix by a (n x 1) vector
  - lacksquare  $\Rightarrow$  total time still only  $\Theta(n^2)$
- Example of a Monte-Carlo style algorithm
- If A \* B = C, this always returns True
- If A \* B  $\neq$  C, want  $P(A \cdot (B \cdot \vec{r}) \neq C \cdot \vec{r}) > 0.5$



### Frievald: false-positive rate

- Let D = AB C: by assumption,  $D \neq 0$ , so choose  $d_{ij} \neq 0$ 
  - $\Rightarrow$  Want to show  $P(D\vec{r}=0) \leq 0.5$
- $Dec{r}$  is 0 iff all its elts are 0, so  $P(Dec{r}=0) \leq P((Dec{r})_i=0)$
- ullet This is a  $oldsymbol{\mathsf{dot}}$  product:  $\left(Dec{r}
  ight)_i = \sum_{k=1}^n d_{ik} r_k = d_{ij} r_j + y$
- Two possibilities: if y = 0:  $P((D\vec{r})_i = 0) = P(d_{ij}r_j = 0)$   $= P(r_i = 0) = 0.5$
- ullet If y 
  eq 0, then  $P((Dec r)_i = 0) = P(r_j = 1 ext{ and } d_{ij} = -y)$   $\leq P(r_j = 1) = 0.5$

- Heap sort (ch5)
  - Intro to trees (more in ch12)
  - Binary heaps and max-heaps
  - Heap sort
  - Max-heaps for priority queue
- Quicksort (ch6)
  - Lomuto partitioning and complexity analysis
  - Randomised Quicksort and analysis
- Monte-Carlo matrix multiply checking

cmpt231.seanho.com/lec3

