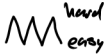# CSE 326 Data Structures

Dave Bacon
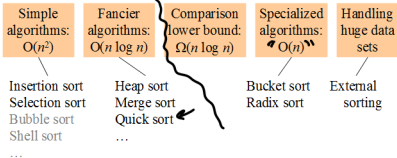
Sorting

# Logistics

- Survey on <u>main web page</u>!

- Homework 6 (due on Friday)

- Project 3, Project 3, Project 3.

- Reading: finish Weiss <u>Chapter 7</u>, start <u>Chapter 9</u>

hard

easy

# Sorting: *The Big Picture*

Given *n* comparable elements in an array,
sort them in an increasing (or decreasing)
order.

| Simple algorithms: $O(n^2)$ | Fancier algorithms: $O(n \log n)$ | Comparison lower bound: $\Omega(n \log n)$ | Specialized algorithms: $O(n)$ | Handling huge data sets |
|---|---|---|---|---|
| Insertion sort Selection sort Bubble sort Shell sort … | Heap sort Merge sort Quick sort … | | Bucket sort Radix sort | External sorting |

# How fast can we sort?

- Heapsort, Mergesort, and Quicksort all run in O(N log N) <u>best</u> case running time
- Can we do any better?
- No, if the basic action is a comparison.

Algorithms

# Sorting Model

- Recall our basic assumption: we can <u>only compare two elements at a time</u>
  - we can only reduce the possible solution space by half each time we make a comparison
- Suppose you are given N elements
  - Assume no duplicates  a, b, c, d, c, ...
- How many possible orderings can you get?
  - Example: a, b, c  (N = 3)

# Permutations

- How many possible orderings can you get?
  - Example: a, b, c  (N = 3)
  - (a b c), (a c b), (b a c), (b c a), (c a b), (c b a)
  - 6 orderings = 3·2·1 = 3!  (ie, "3 factorial")
  - All the possible permutations of a set of 3 elements
- For N elements
  - N choices for the first position, (N-1) choices for the second position, ..., (2) choices, 1 choice
  - N(N-1)(N-2)···(2)(1)= N! possible orderings

$$N(N-1)(N-2)\cdots = N!$$

N!
grows
fast.
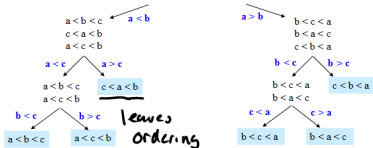
Binary tree    node — set of possible orderings

# Decision Tree

edge = 1-comparison

a < b < c, b < c < a,
c < a < b, a < c < b,
b < a < c, c < b < a

a < b < c
c < a < b          a < b          a > b          b < c < a
a < c < b                                        b < a < c
                                                 c < b < a

a < c          a > c                    b < c          b > c

a < b < c      c < a < b                b < c < a      c < b < a
a < c < b                               b < a < c

b < c          b > c    *leaves*        c < a          c > a
                        *ordering*
a < b < c      a < c < b                b < c < a      b < a < c

The leaves contain all the possible orderings of a, b, c

# Lower bound on Height

- A binary tree of height h has **at most** *how many* leaves?

$$L \leq 2^h$$

$$h = 0 \quad L \leq 1$$

$$L \leq 2^{h+1}$$

$$L \leq 2^{h+1}$$

- A binary tree with L leaves has height **at least**:

$$\log_2 L \leq \log_2 2^h = h \qquad \boxed{h \geq \log_2 L}$$

- The decision tree has how many leaves:

$$\boxed{N!}$$

- So the decision tree has height:

$$h \geq \log \underset{2}{} |N!| \quad \leftarrow \text{worst case}$$

# $\log(N!)$ is $\Omega(N \log N)$

$$\log A \cdot B = \log A + \log B$$

$$\log(N!) = \log(N \cdot (N-1) \cdot (N-2) \cdots (2) \cdot (1))$$

$$= \log N + \log(N-1) + \log(N-2) + \cdots + \underbrace{\log 2 + \log 1}_{\geq 0}$$

select just the first $N/2$ terms

$$\geq \log N + \log(N-1) + \log(N-2) + \cdots + \log \frac{N}{2}$$

each of the selected terms is $\geq \log N/2$

$$\geq \frac{N}{2} \log \frac{N}{2}$$

$$\geq \frac{N}{2}(\log N - \log 2) = \frac{N}{2} \log N - \frac{N}{2}$$

$$= \Omega(N \log N)$$

# $\Omega$(N log N)

- Run time of any comparison-based sorting algorithm is $\Omega$(**N log N**)
- Can we do better if we don't use comparisons?

comparisons.

a < b

# BucketSort (aka BinSort)

If all values to be sorted are *known* to be between **1** and $K$, create an array `count` of size $K$, **increment** counts while traversing the input, and finally output the result.

**Example** $K$=5.  Input = (5,1,3,4,3,2,1,1,5,4,5)

count array

| count | array |
|---|---|
| 1 | 3 ✗ ✗ |
| 2 | 1 |
| 3 | ✗ 2 |
| 4 | ✗ 2 |
| 5 | 3 ✗ ✗ |

1, 1, 1, 2, 3, 3, 4, 4,
5, 5, 5

initial
scan
$N$

read → $K$

**Running time to sort n items?**  no graph

$O(n+K)$

# BucketSort Complexity: O($n+K$)

$K = 1000$   $O(n + 1000)$

- Case 1: $K$ is a constant
  - BinSort is linear time
- Case 2: $K$ is variable ← $K \rightarrow P(n)$
  - Not simply linear time
- Case 3: $K$ is constant but large (e.g. $2^{32}$)
  - ???

$O(n + 2^{32})$

32 bit number

# Fixing impracticality: RadixSort

- Radix = "The base of a number system"
  - We'll use 10 for convenience, but could be anything

- <u>Idea</u>: BucketSort on each **digit**, least significant to most significant (lsd to msd)

# Radix Sort Example (1st pass)



Bucket sort
by 1's digit

Input data

478
537
9
721
3
38
123
67

After 1st pass

721
3
123
537
67
478
38
9

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
|   | 721 |   | 3 123 |   |   |   | 537 67 |   478 38 |   9 |

This example uses B=10 and base 10 digits for simplicity of demonstration. Larger bucket counts should be used in an actual implementation.

# Radix Sort Example (2nd pass)

**After 1st pass**

721
3
123
537
67
478
38
9

Bucket sort
by 10's
digit

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 03 09 | | 721 123 | 537 38 | | | 67 | 478 | | |

**After 2nd pass**

3
9
721
123
537
38
67
478

# Radix Sort Example (3rd pass)

After 2nd pass

| |
|---|
| 3 |
| 9 |
| 721 |
| 123 |
| 537 |
| 38 |
| 67 |
| 478 |

Bucket sort
by 100's
digit

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 003 | 123 | | | 478 | 537 | | 721 | | |
| 009 | | | | | | | | | |
| 038 | | | | | | | | | |
| 067 | | | | | | | | | |

After 3rd pass

| |
|---|
| 3 |
| 9 |
| 38 |
| 67 |
| 123 |
| 478 |
| 537 |
| 721 |

**Invariant**: after k passes the low order k digits are sorted.

# RadixSort

- Input:126, 328, 636, 341, 416, 131, 328

BucketSort on lsd:

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
|   |   |   |   |   |   |   |   |   |   |

BucketSort on next-higher digit:

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
|   |   |   |   |   |   |   |   |   |   |

BucketSort on msd:

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
|   |   |   |   |   |   |   |   |   |   |

# Radixsort: Complexity

- How many passes?

$$p = \log_k(\text{max number})$$

- How much work per pass?

$$O(n + k)$$

- Total time?

$$O(p(n+k))$$

- Conclusion?

$p$ is large

$$\log_2 2^{32} = 32$$

- In practice
  - RadixSort only good for large number of elements with relatively small values
  - Hard on the cache compared to MergeSort/QuickSort

# Internal versus Underline External Sorting

- Need sorting algorithms that minimize disk/tape access time
- **External sorting** – Basic Idea:
  - Load chunk of data into RAM, sort, store this "run" on disk/tape
  - Use the Merge routine from Mergesort to merge runs
  - Repeat until you have only one run (one sorted chunk)
  - Text gives some examples

# Graphs

Chapter 9 in Weiss

# Graph… ADT?

- Not quite an ADT… operations not clear

- A formalism for representing relationships between objects

  Graph $G = (V,E)$

  - Set of *vertices*:
    $V = \{v_1, v_2, \ldots, v_n\}$

  - Set of *edges*:
    $E = \{e_1, e_2, \ldots, e_m\}$
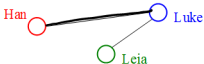    where each $e_i$ connects two
    vertices $(v_{i1}, v_{i2})$



$V = \{Han, Leia, Luke\}$
$E = \{(Luke, Leia),$
$\quad (Han, Leia),$
$\quad (Leia, Han)\}$

# Graph Definitions

In *directed* graphs, edges have a specific direction:



In *undirected* graphs, they don't (edges are two-way):



**v** is *adjacent* to **u** if $\underline{(u, v)} \in \mathbf{E}$

# More Definitions:
# Simple Paths and Cycles

A *simple path* repeats no vertices (except that the first can be the last):

    p = {Seattle, Salt Lake City, San Francisco, Dallas}

    p = {Seattle, Salt Lake City, Dallas, San Francisco, Seattle}

A *cycle* is a path that starts and ends at the same node:

    p = {Seattle, Salt Lake City, Dallas, San Francisco, Seattle}

    p = {Seattle, Salt Lake City, Seattle, San Francisco, Seattle}

A *simple cycle* is a cycle that repeats no vertices except that the first vertex is also the last (in undirected graphs, no edge can be repeated)
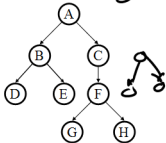
# Trees as Graphs

- Every tree is a graph!
- Not all graphs are trees!
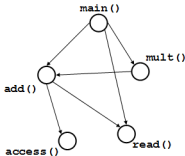
  A graph is a tree if
  - There are *no cycles* (directed or undirected)
  - There is a *path* from the root *to every node*

# Directed Acyclic Graphs (DAGs)

**DAGs** are
directed graphs
with no
(directed)
cycles.

*Aside: if program call-
graph is a DAG, then all
procedure calls can be in-
lined*

# Graph Representations

Han         Luke

Leia

0. List of vertices + list of edges
1. 2-D matrix of vertices (marking edges in the cells)
   "adjacency matrix" ←
2. List of vertices each with a list of adjacent vertices
   "adjacency list" ←

Things we might want to do:
- iterate over vertices
- iterate over edges
- iterate over vertices adj. to a vertex
- check whether an edge exists

Vertices and edges
may be labeled

# Representation 1: Adjacency Matrix

A **|V| x |V|** array in which an element `(u,v)` is true if and only if there is an edge from **u** to **v**



|  | Han | Luke | Leia |
|---|---|---|---|
| Han |  |  |  |
| Luke |  |  |  |
| Leia |  |  |  |

*space requirements:*                    *runtime:*