

ALGORITHMS & ADVANCED DATA STRUCTURES (#8)



LINEAR TIME SORTING – LOWER BOUNDS

**ADAPTED FROM
CS 146 SJSU (KATERINA POTIKA)**

Sorting So Far – 1st Algorithm

2

- Insertion sort:
 - Easy implementation
 - Fast on small inputs (less than ~50 elements)
 - Fast on nearly-sorted inputs
 - $O(n^2)$ worst case
 - $O(n^2)$ average (equally-likely inputs) case
 - $O(n^2)$ reverse-sorted case

Sorting So Far – 2nd Algorithm

3

- Merge sort:
 - Divide-and-conquer:
 - ◆ Split array in half
 - ◆ Recursively sort subarrays
 - ◆ Linear-time merge step
 - $O(n \lg n)$ worst case
 - Doesn't sort in place

Sorting So Far – 3rd Algorithm

4

- Heap sort:
 - Uses the heap data structure
 - ◆ Complete binary tree
 - ◆ Heap property: parent key $>$ children's keys
 - $O(n \lg n)$ worst case
 - Sorts in place
 - Fair amount of shuffling memory around

Sorting So Far – 4th Algorithm

5

- Quick sort:
 - Divide-and-conquer:
 - ◆ Partition array into two subarrays, recursively sort
 - ◆ All elements of first subarray < all elements of second subarray
 - ◆ No merge step needed!
 - $O(n \lg n)$ average case
 - Fast in practice
 - $O(n^2)$ worst case
 - ◆ Naïve implementation: worst case on sorted input
 - ◆ Address this with randomized quicksort

How Fast Can We Sort?

6

- We will provide a lower bound, then beat it
 - by playing a different game
- First, an observation: all of the sorting algorithms so far are *comparison sorts*
 - The only operation used to gain ordering information about a sequence is the **pairwise comparison of two elements**
 - *Theorem*: all comparison sorts are $\Omega(n \log n)$
 - ◆ A comparison sort must do $\Omega(n)$ comparisons (*why?*)
 - ◆ What about the gap between $\Omega(n)$ and $\Omega(n \log n)$

Decision Trees

7

- *Decision trees* provide an abstraction of comparison sorts
 - A decision tree represents the comparisons made by a comparison sort. Everything else is ignored
 - *What do the leaves represent?*
 - leaf is labeled by the permutation of orders that the algorithm determines
 - *How many leaves must there be?*
 - There are $\geq n!$ leaves, because every permutation appears at least once.

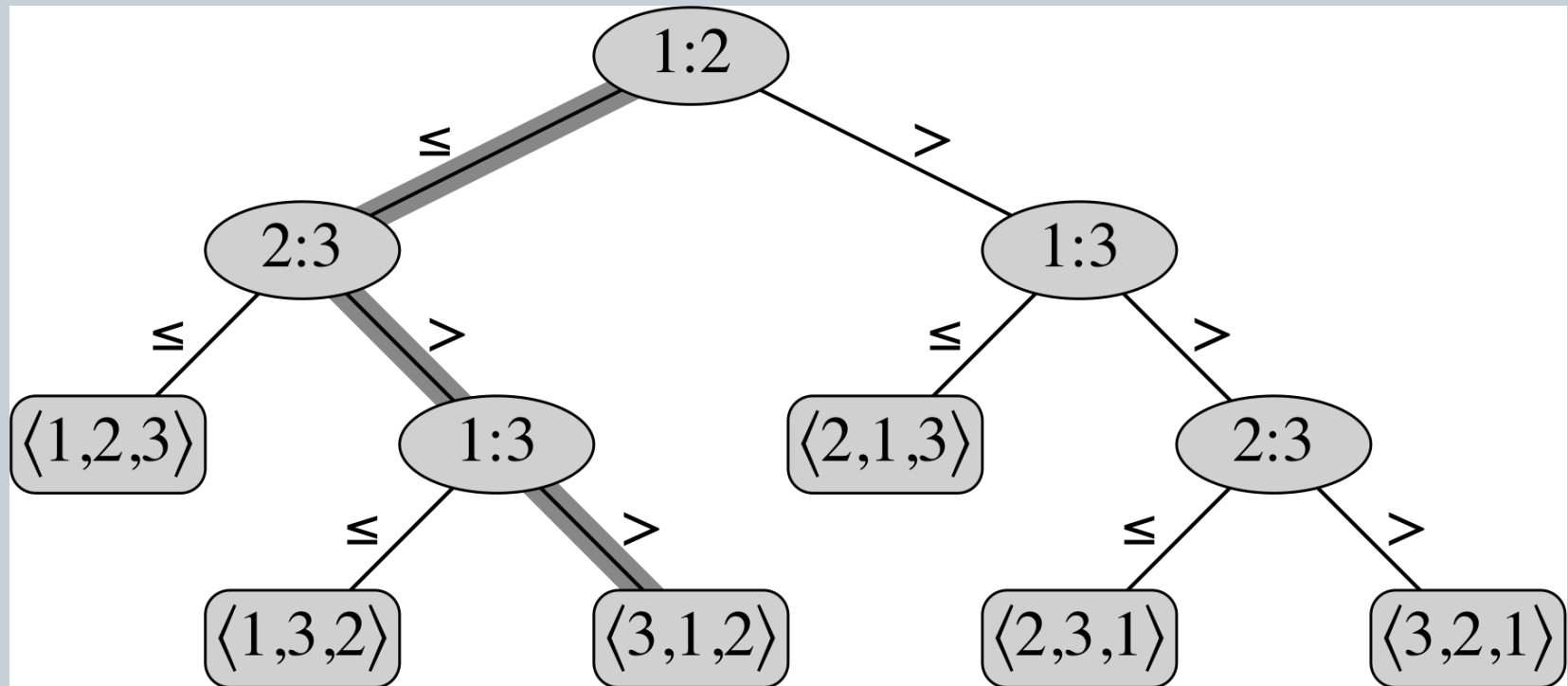
Note: Permutations (Appendix C)

8

- A ***permutation*** of a finite set S is an ordered sequence of all the elements of S , each element appearing exactly once.
-
- For example, if $S = \{a, b, c\}$, then S has 6 permutations: *abc*, *acb*, *bac*, *bca*, *cab*, *cba*
- There are $n!$ permutations of a set of n elements
 - we can choose the first element of the sequence in n ways, the second in $n-1$ ways, the third in $n-2$, etc.

Example: insertion sort for 3 numbers

9



Decision Trees

10

- Decision trees can model comparison sorts. For a given algorithm:
 - One tree for each n
 - Tree paths are all possible execution traces
 - *What's the longest path in a decision tree for insertion sort? For merge sort?*
- *What is the asymptotic height of any decision tree for sorting n elements?*
- Answer: $\Omega(n \log n)$ (proof follows)

Lower Bound For Comparison Sorting

11

- Thm: Any decision tree that sorts n elements has height $\Omega(n \log n)$
- *What's the maximum # of leaves of a binary tree of height h ?*
- Lemma: Any binary tree of height h has $k \leq 2^h$ where k : # of leaves (proof by induction)

Lower Bound For Comparison Sorting

12

- So we have...

$$n! \leq 2^h$$

- Taking logarithms:

$$\lg(n!) \leq h$$

- Stirling's approximation tells us:

$$n! > \left(\frac{n}{e}\right)^n$$

- Thus:

$$h \geq \log\left(\frac{n}{e}\right)^n$$

Lower Bound For Comparison Sorts

13

- So we have

$$h \geq \log \left(\frac{n}{e} \right)^n$$

$$= n \log n - n \log e$$

$$= \Omega(n \log n)$$

- Thus the minimum height of a decision tree is $\Omega(n \log n)$

Lower Bound For Comparison Sorts

14

- Thus the time to comparison sort n elements is $\Omega(n \log n)$
- Corollary: Heapsort and Mergesort are asymptotically optimal comparison sorts
- “Sorting in linear time” (?)
 - *How can we do better than $\Omega(n \log n)$?*

Sorting In Linear Time

15

- Counting sort
 - No comparisons between elements!
 - **But**...depends on assumption about the numbers being sorted
 - ◆ We assume numbers are in the range $1..k$
 - The algorithm:
 - ◆ Input: $A[1..n]$, where $A[j] \in \{1, 2, 3, \dots, k\}$
 - ◆ Output: $B[1..n]$, sorted (notice: not sorting in place)
 - ◆ Also: Array $C[1..k]$ for auxiliary storage

Counting Sort

16

```
1 CountingSort(A, B, k)
2   for i=1 to k
3       C[i]= 0;
4   for j=1 to n
5       C[A[j]] += 1;
6   for i=2 to k
7       C[i] = C[i] + C[i-1];
8   for j=n downto 1
9       B[C[A[j]]] = A[j];
10      C[A[j]] -= 1;
```


Counting Sort

17

```
1 CountingSort(A, B, k)
2   for i=1 to k
3       C[i]= 0;
4   for j=1 to n
5       C[A[j]] += 1;
6   for i=2 to k
7       C[i] = C[i] + C[i-1];
8   for j=n downto 1
9       B[C[A[j]]] = A[j];
10      C[A[j]] -= 1;
```

Takes time $O(k)$



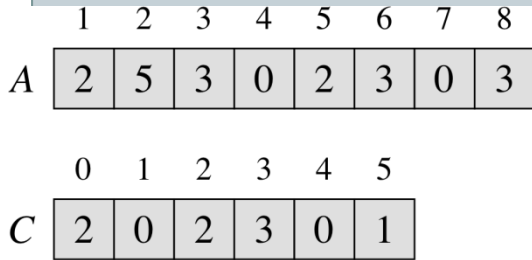
Takes time $O(n)$



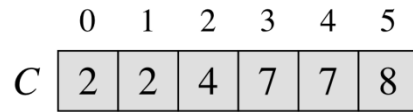
What will be the running time?

Example

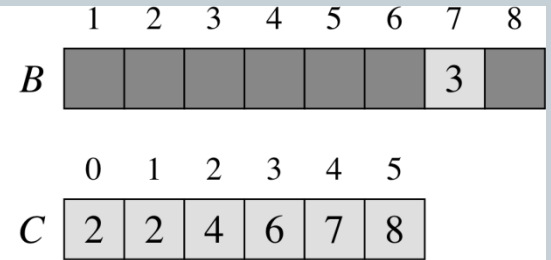
18



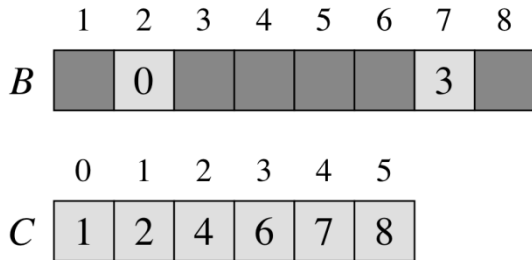
(a)



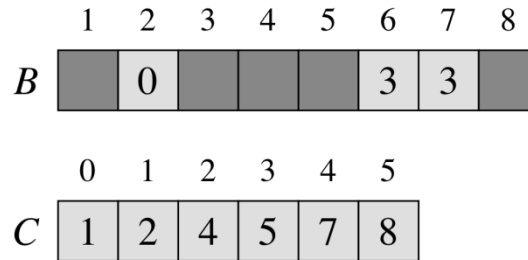
(b)



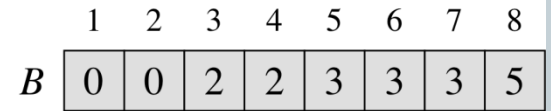
(c)



(d)



(e)



(f)

Counting Sort

19

- Total time: $O(n + k)$
 - Usually, $k = O(n)$
 - Thus counting sort runs in $O(n)$ time
- But sorting is $\Omega(n \log n)$
 - No contradiction--this is **not a comparison** sort (in fact, there are *no* comparisons at all!)
 - Notice that this algorithm is *stable* (what is that?)
- **Stable:** keys with same value appear in same order in output as they did in input

Counting Sort

20

- Cool! *Why don't we always use counting sort?*
- Because it depends on range k of elements
- *Could we use counting sort to sort 32 bit integers? Why or why not? How many possible (distinct) numbers can we have?*
- Answer: no, k too large ($2^{32} = 4,294,967,296$)

Counting Sort

21

- *How did IBM get rich originally?*
- Answer: punched card readers for census tabulation in early 1900's.
 - In particular, a *card sorter* that could sort cards into different bins
 - ◆ Each column can be punched in 12 places
 - ◆ Decimal digits use 10 places
 - Problem: only one column can be sorted on at a time

Clicker Question 9.1

22

Counting sort performs numbers of comparisons between input elements.

- a) 0
- b) n
- c) $n \log n$
- d) n^2

Radix Sort

23

- Intuitively, you might sort on the most significant digit, then the second msd, etc.
- Problem: lots of intermediate piles of cards (read: scratch arrays) to keep track of
- Key idea: sort the *least* significant digit first

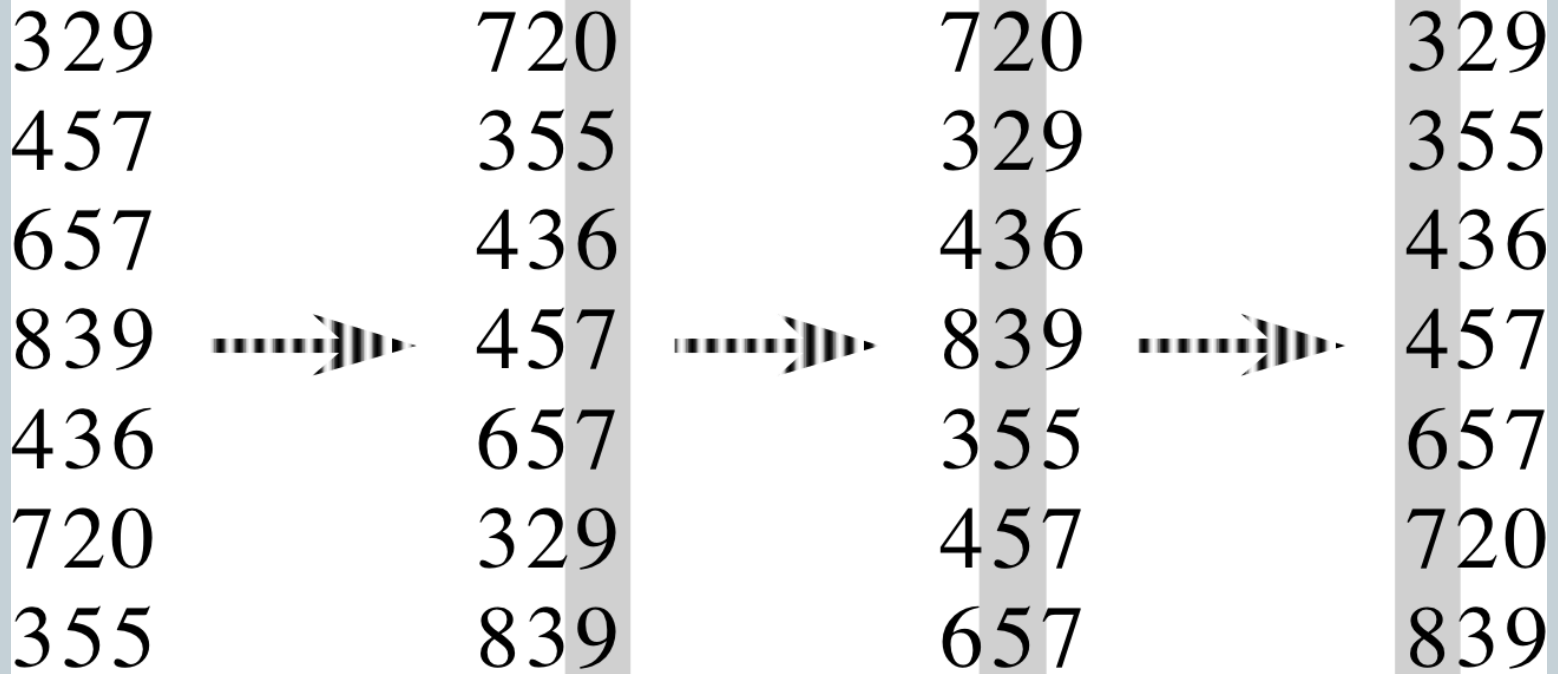
RadixSort(A, d)

for i=1 to d

StableSort(A) on digit i

Example of Radix Sort

24



Radix Sort

25

- *Can we prove it will work?*
- Sketch of an inductive argument (induction on the number of passes):
 - Assume lower-order digits $\{j: j < i\}$ are sorted
 - Show that sorting next digit i leaves array correctly sorted
 - ◆ If two digits at position i are different, ordering numbers by that digit is correct (lower-order digits irrelevant)
 - ◆ If they are the same, numbers are already sorted on the lower-order digits. Since we use a **stable** sort, the numbers stay in the right order

Radix Sort

26

- *What sort will we use to sort on digits?*
- Counting sort is obvious choice:
 - Sort n numbers on digits that range from $1..k$
 - Time: $O(n + k)$
- Each pass over n numbers with d digits takes time $O(n+k)$, so total time $O(dn+dk)$
 - When d is constant and $k=O(n)$, takes $O(n)$ time
- *How many bits in a computer word?*

How to break each key into digits?

27

- n words
- b bits/word
- Break into r -bit digits. Have $d = \lfloor b/r \rfloor$ digits
- Use counting sort with $k = 2^r - 1$
- Example: 32-bit words, 8-bit digits. $b = 32$, $r = 8$, $d = 32/8 = 4$, $k = 2^8 - 1 = 255$.
- • Time = $\Theta(b/r(n + 2^r))$.
- Choose $r \approx \log n$ gives: $\Theta(bn / \log n)$.

Bucket Sort

28

- Assumes the input is generated by a random process that distributes elements uniformly over $[0, 1)$.
- **Idea:**
 - Divide $[0, 1)$ into n equal-sized *buckets*.
 - Distribute the n input values into the buckets.
 - Sort each bucket.
 - Then go through buckets in order, listing elements in each one.
- **Input:** $A[1 \dots n]$, where $0 \leq A[i] < 1$ for all i .
- **Auxiliary array:** $B[0 \dots n - 1]$ of linked lists, each list initially empty.

Bucket sort Code

29

- BUCKET-SORT (A, n)
- **for** $i \leftarrow 1$ **to** n
- **do** insert $A[i]$ into list $B[\lfloor n \cdot A[i] \rfloor]$
- **for** $i \leftarrow 0$ **to** $n - 1$
- **do** sort list $B[i]$ with insertion sort
- concatenate lists $B[0], B[1], \dots, B[n - 1]$ together in order
- **return** the concatenated lists

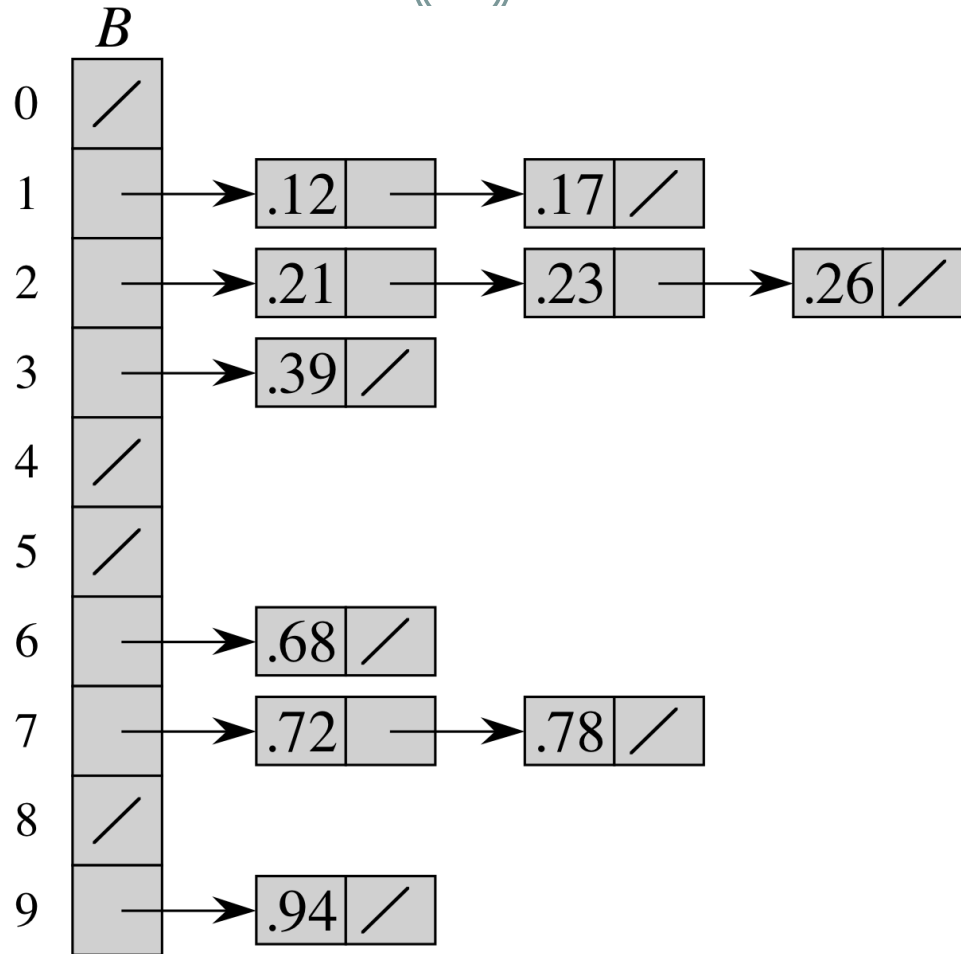
Bucket Sort Example

(30)

A

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

(a)



(b)

Correctness

31

- Consider $A[i]$, $A[j]$. Assume without loss of generality that
- $A[i] \leq A[j]$. Then $\lfloor n \cdot A[i] \rfloor \leq \lfloor n \cdot A[j] \rfloor$. So $A[i]$ is placed into the same bucket as $A[j]$ or into a bucket with a lower index.
- • If same bucket, insertion sort fixes up.
- • If earlier bucket, concatenation of lists fixes up.

Analysis

32

- Relies on *no bucket* getting too many values.
- All lines of algorithm except insertion sorting take $\Theta(n)$ altogether.
- Intuitively, if each bucket gets a constant number of elements, it takes $O(1)$ time to sort each bucket $\Rightarrow O(n)$ sort time for all buckets.
- We “expect” each bucket to have few elements, since the average is 1 element per bucket.

Analysis

33

- Uniform input distribution has $O(1)$ bucket size
- and expected time is $O(n)$

- Later in Hash Tables again the same idea

Clicker Question 9.2

34

Given n numbers of elements in the range $[0 \dots n^3 - 1]$.
which of the following sorting algorithms can sort
them in $O(n)$ time?

- a) Counting sort
- b) Bucket sort
- c) Radix sort
- d) Quick sort

Structures...

35

- Done with sorting and order statistics
- Next part is *data structures Ch 10 (skim)*
- *Ch 11*

Applications

36

- Assume you have an array of objects, instead of integers, and you know that the possible objects are limited (constant). How would you sort?
- Given an array of integers in the range from -5 to 5, write an algorithm that sorts.
- to find the element which appears maximum number of times in the array.
 - Example 4, -1, -5, -2, 1, -5, -2, 2, 0, -5, 3, -2, 4, 1