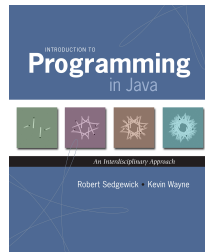


4.2 Sorting and Searching



Sequential Search

Sequential search. Scan through array, looking for key.

- Search hit: return array index.
- Search miss: return -1.

```
public static int search(String key, String[] a) {
    for (int i = 0; i < a.length; i++)
        if (a[i].compareTo(key) == 0)
            return i;
    return -1;
}
```

Search Client: Exception Filter

Exception filter. Read a sorted list of strings from a **whitelist** file, then print out all strings from standard input not in the whitelist.

```
public static void main(String[] args) {
    In in = new In(args[0]);
    String s = in.readLine();
    String[] words = s.split("\\s+");
    while (!StdIn.isEmpty()) {
        String key = StdIn.readString();
        if (search(key, words) == -1)
            StdOut.println(key);
    }
}
```

```
more test.txt           % more whitelist.txt
bob@office             alice@home
carl@beach             bob@office
marv@spam              bob@office
bob@office             bob@office
bob@office             bob@office
m@l@r@y@span          dave@boat
dave@boat             % java BinarySearch whitelist.txt < test.txt
eve@airport           marv@spam
alice@home            m@l@r@y@span
eve@airport           eve@airport
```

Searching Challenge 1

Q. A credit card company needs to whitelist 10 million customer account numbers, processing 1,000 transactions per second.

Using **sequential search**, what kind of computer is needed?

- Toaster
- Cellphone
- Your laptop
- Dual-core server
- Supercomputer
- Google server farm

Binary Search

Twenty Questions

Intuition. Find a hidden integer.

interval	size	Q	A
0 — 128	128	< 64?	false
0 — 64 — 128	64	< 32?	true
0 — 64 — 96	32	< 16?	true
0 — 64 — 80	16	< 8?	false
0 — 72 — 80	8	< 4?	false
0 — 76 — 80	4	< 2?	true
0 — 76 — 78	2	< 1?	false
0 — 77	1	= 77	

Binary Search

Main idea.

- Sort the array (stay tuned).
- Play "20 questions" to determine index with a given key.

Ex. Dictionary, phone book, book index, credit card numbers, ...

Binary search.

- Examine the middle key.
- If it matches, return its index.
- Otherwise, search either the left or right half.

Binary search in a sorted array (one step)

Binary Search: Java Implementation

Invariant. Algorithm maintains $a[lo] \leq key \leq a[hi-1]$.

```

public static int search(String key, String[] a) {
    return search(key, a, 0, a.length);
}

public static int search(String key, String[] a, int lo, int hi) {
    if (hi <= lo) return -1;
    int mid = (hi + lo) / 2;
    int cmp = a[mid].compareTo(key);
    if (cmp > 0) return search(key, a, lo, mid);
    else if (cmp < 0) return search(key, a, mid+1, hi);
    else return mid;
}
    
```

Java library implementation: `Arrays.binarySearch()`

Binary Search: Mathematical Analysis

Analysis. To binary search in an array of size N : do one compare, then binary search in an array of size $N/2$.

$$N \rightarrow N/2 \rightarrow N/4 \rightarrow N/8 \rightarrow \dots \rightarrow 1$$

Q. How many times can you divide a number by 2 until you reach 1?

A.

```

      1
     2 → 1
    4 → 2 → 1
   8 → 4 → 2 → 1
  16 → 8 → 4 → 2 → 1
 32 → 16 → 8 → 4 → 2 → 1
 64 → 32 → 16 → 8 → 4 → 2 → 1
128 → 64 → 32 → 16 → 8 → 4 → 2 → 1
256 → 128 → 64 → 32 → 16 → 8 → 4 → 2 → 1
512 → 256 → 128 → 64 → 32 → 16 → 8 → 4 → 2 → 1
1024 → 512 → 256 → 128 → 64 → 32 → 16 → 8 → 4 → 2 → 1
    
```

Searching Challenge 2

Q. A credit card company needs to whitelist 10 million customer account numbers, processing 1,000 transactions per second.

Using **binary search**, what kind of computer is needed?

- Toaster
- Cell phone
- Your laptop
- Dual-core server
- Supercomputer
- Google server farm

Sorting

Sorting

Sorting problem. Rearrange N items in ascending order.

Applications. Statistics, databases, data compression, bioinformatics, computer graphics, scientific computing, (too numerous to list), ...

Hausser	→	Hanley
Hong		Haskell
Hsu		Hausser
Hayes		Hayes
Haskell		Hong
Hanley		Hornet
Hornet		Hsu

Insertion Sort

13

Insertion Sort

Insertion sort.

- Brute-force sorting solution.
- Move left-to-right through array.
- Exchange next element with larger elements to its left, one-by-one.

		a							
i	j	0	1	2	3	4	5	6	7
6	6	and	had	him	his	was	you	the	but
6	5	and	had	him	his	was	the	you	but
6	4	and	had	him	his	the	was	you	but
		and	had	him	his	the	was	you	but

Inserting a[6] into position by exchanging with larger entries to its left

14

Insertion Sort

Insertion sort.

- Brute-force sorting solution.
- Move left-to-right through array.
- Exchange next element with larger elements to its left, one-by-one.

		a							
i	j	0	1	2	3	4	5	6	7
1	0	was	had	him	and	you	his	the	but
2	1	had	was	him	and	you	his	the	but
3	0	and	had	him	was	you	his	the	but
4	4	and	had	him	was	you	his	the	but
5	3	and	had	him	his	was	you	the	but
6	4	and	had	him	his	the	was	you	but
7	1	and	but	had	him	his	the	was	you
		and	but	had	him	his	the	was	you

Inserting a[1] through a[N-1] into position (insertion sort)

15

Insertion Sort: Java Implementation

```
public class Insertion {
    public static void sort(String[] a) {
        int N = a.length;
        for (int i = 1; i < N; i++)
            for (int j = i; j > 0; j--)
                if (a[j-1].compareTo(a[j]) > 0)
                    exch(a, j-1, j);
                else break;
    }

    private static void exch(String[] a, int i, int j) {
        String swap = a[i];
        a[i] = a[j];
        a[j] = swap;
    }
}
```

16

Insertion Sort: Mathematical Analysis

Worst case. [descending]

- Iteration i requires i comparisons.
- Total = $(0 + 1 + 2 + \dots + N-1) \sim N^2/2$ compares.

E	F	G	H	I	J	D	C	B	A
---	---	---	---	---	---	---	---	---	---

i

Average case. [random]

- Iteration i requires $i/2$ comparisons on average.
- Total = $(0 + 1 + 2 + \dots + N-1)/2 \sim N^2/4$ compares

A	C	D	F	H	J	E	B	I	G
---	---	---	---	---	---	---	---	---	---

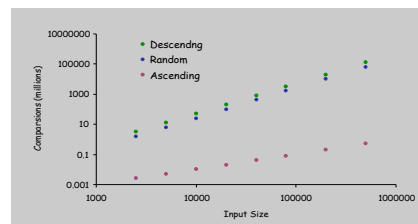
i

17

Insertion Sort: Empirical Analysis

Observation. Number of compares depends on input family.

- Descending: $\sim N^2/2$.
- Random: $\sim N^2/4$.
- Ascending: $\sim N$.



18

Merging

Merging. Combine two pre-sorted lists into a sorted whole.

How to merge efficiently? Use an auxiliary array. ▶

		a									
i	j	k	aux[k]								
			0	1	2	3	4	5	6	7	
			and	had	him	was	but	his	the	you	
0	4	0	and	and	had	him	was	but	his	the	you
1	4	1	but	and	had	him	was	but	his	the	you
1	5	2	had	and	had	him	was	but	his	the	you
2	5	3	him	and	had	him	was	but	his	the	you
3	5	4	his	and	had	him	was	but	his	the	you
3	6	5	the	and	had	him	was	but	his	the	you
3	6	6	was	and	had	him	was	but	his	the	you
4	7	7	you	and	had	him	was	but	his	the	you

Trace of the merge of the sorted left half with the sorted right half

26

Merging

Merging. Combine two pre-sorted lists into a sorted whole.

How to merge efficiently? Use an auxiliary array. ▶

```
String[] aux = new String[N];
// merge into auxiliary array
int i = lo, j = mid;
for (int k = 0; k < N; k++) {
    if (i == mid) aux[k] = a[j++]; //left side of a done already
    else if (j == hi) aux[k] = a[i++]; //right side of a done already
    else if (a[j].compareTo(a[i]) < 0) aux[k] = a[j++];
    else aux[k] = a[i++];
}
// copy back from aux into a
for (int k = 0; k < N; k++) {
    a[lo + k] = aux[k];
}
```

27

Mergesort: Java Implementation

```
public class Merge {
    public static void sort(String[] a) {
        sort(a, 0, a.length);
    }
    // Sort a[lo, hi).
    public static void sort(String[] a, int lo, int hi) {
        int N = hi - lo;
        if (N <= 1) return;
        // recursively sort left and right halves
        int mid = lo + N/2;
        sort(a, lo, mid); //sort a[lo, mid)
        sort(a, mid, hi); //sort a[mid, hi)
        // merge sorted halves (see previous slide)
    }
}
```

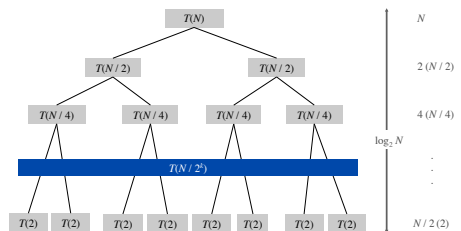


28

Mergesort: Mathematical Analysis

Analysis. To mergesort array of size N , mergesort two subarrays of size $N/2$, and merge them together using $\approx N$ compares.

we assume N is a power of 2 (for simplicity, $N \log_2 N$ bound holds regardless)



$N \log_2 N$

29

Mergesort: Mathematical Analysis

Mathematical analysis.

analysis	comparisons
worst	$N \log_2 N$
average	$N \log_2 N$
best	$1/2 N \log_2 N$

Validation. Theory agrees with observations.

N	actual	predicted
10,000	120 thousand	133 thousand
20 million	460 million	485 million
50 million	1,216 million	1,279 million

30

Sorting Challenge 2

Q. A credit card company sorts 10 million customer account numbers, for use with binary search.

Using mergesort, what kind of computer is needed?

- A. Toaster
- B. Cell phone
- C. Your laptop
- D. Supercomputer
- E. Google server farm

31

Mergesort: Lesson

Lesson. Great algorithms can be more powerful than supercomputers.

Computer	Compares Per Second	Insertion	Mergesort
laptop	10^7	3 centuries	3 hours
super	10^{12}	2 weeks	instant

N = 1 billion

33

Longest Repeated Substring

34

Redundancy Detector

Longest repeated substring. Given a string, find the longest substring that appears at least twice.

a a c a a g t t t a c a a g c

Brute force.

- Try all indices i and j for start of possible match.
- Compute longest common prefix for each pair (quadratic+).

a a c a a g t t t a c a a g c
 i j

Applications. Bioinformatics, data compression, ...

35

Longest Repeated Substring: Brute-Force Solution

Longest repeated substring. Given a string, find the longest substring that appears at least twice.

a a c a a g t t t a c a a g c

Brute force.

- Try all indices i and j for start of possible match.
- Compute longest common prefix (LCP) for each pair.

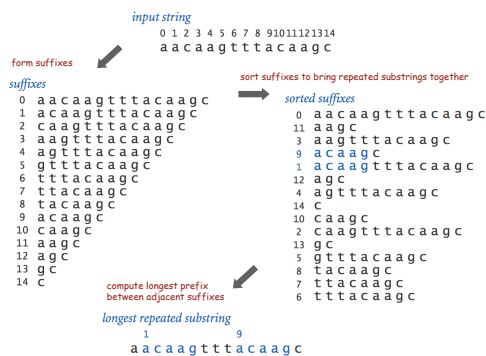
a a c a a g t t t a c a a g c
 i j

Mathematical analysis.

- All pairs: $0 + 1 + 2 + \dots + N-1 \sim N^2/2$ calls on LCP.
- Way too slow for long strings.

37

Longest Repeated Substring: A Sorting Solution



38

Longest Repeated Substring: Java Implementation

Suffix sorting implementation.

```

int N = s.length();
String[] suffixes = new String[N];
for (int i = 0; i < N; i++)
    suffixes[i] = s.substring(i, N);
Arrays.sort(suffixes);
    
```

Longest common prefix. $lcp(s, t)$

- Longest string that is a prefix of both s and t .
 - EX: $lcp("acaagtttac", "acaagc") = "acaag"$.
- Easy to implement (you could write this one).

Longest repeated substring. Search only adjacent suffixes.

```

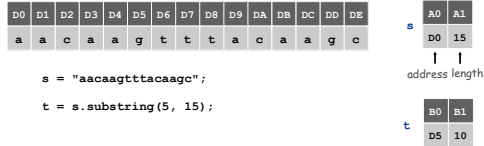
String lrs = "";
for (int i = 0; i < N-1; i++) {
    String x = lcp(suffixes[i], suffixes[i+1]);
    if (x.length() > lrs.length()) lrs = x;
}
    
```

39

OOP Context for Strings

String representation.

- A string is an address and a length.
- Characters can be shared among strings. *does not copy chars*
- `substring()` computes address and length.



Consequences.

- `substring()` is constant-time operation (instead of linear).
- Creating suffixes takes linear space (instead of quadratic).
- Running time of LRS is dominated by the string sort.

40

Sorting Challenge 4

Q. Four researchers A, B, C, and D are looking for long repeated sequences in a genome with over 1 billion characters.

Which one is more likely to find a cure for cancer?

- has a grad student to do it.
- uses brute force (check all pairs) solution.
- uses sorting solution with insertion sort.
- uses sorting solution with mergesort.

41

Longest Repeated Substring: Empirical Analysis

Input File	Characters	Brute	Suffix Sort	Length
LRS.java	2,162	0.6 sec	0.14 sec	73
Amendments	18,369	37 sec	0.25 sec	216
Aesop's Fables	191,945	3958 sec	1.0 sec	58
Moby Dick	1.2 million	43 hours †	7.6 sec	79
Bible	4.0 million	20 days †	34 sec	11
Chromosome 11	7.1 million	2 months †	61 sec	12,567
Pi	10 million	4 months †	84 sec	14

† estimated

Lesson. Sorting to the rescue; enables new research.

42

Summary

Binary search. Efficient algorithm to search a sorted array.

Mergesort. Efficient algorithm to sort an array.

Applications. Many many applications are enabled by fast sorting and searching.

43