

Sorting and Efficiency

Sorting

- Of all the algorithmic problems that computer scientists have studied, the one with the broadest practical impact is certainly the **sorting problem**, which is the problem of arranging the elements of an array or a vector in order.
- The sorting problem comes up, for example, in alphabetizing a telephone directory, arranging library records by catalogue number, and organizing a bulk mailing by ZIP code.
- There are many algorithms that one can use to sort an array. Because these algorithms vary enormously in their efficiency, it is critical to choose a good algorithm, particularly if the application needs to work with large arrays.

The Selection Sort Algorithm

- Of the many sorting algorithms, the easiest one to describe is **selection sort**, which appears in the text like this:

```
void Sort(Comparable* a) {
    int n = a->n();
    for (int k = 0; k < n; k++) {
        int m = k;
        for (int j = k + 1; j < n; j++) {
            if (a[j] < a[m]) m = j;
        }
        int temp = a[k];
        a[k] = a[m];
        a[m] = temp;
    }
}
```

- Coding this algorithm as a single function makes sense for efficiency but complicates the analysis. The next two slides decompose selection sort into a set of functions that make the operation easier to follow.

Decomposition of the Sort Function

```
/*
 * selection: sort
 * -----
 * sorts a Comparable into increasing order. This implementation
 * uses an algorithm called selection sort, which can be described
 * in English as follows: with your left hand (lh), point to each
 * element in the array in turn, starting at index 0. on each
 * step in the cycle:
 *
 * 1. find the smallest element in the range between your left
 *    hand and the end of the array, and point to that element
 *    with your right hand (rh).
 *
 * 2. swap that element into the current position by swapping
 *    the elements indicated by your left and right hands.
 */
void sort(Comparable* a) {
    for (int k = 0; k < a->n(); k++) {
        int m = findSmallest(a, k, a->n() - 1);
        swap(a[k], a[m]);
    }
}
```

Decomposition of the Sort Function

```
/*
 * selection: findSmallest
 * -----
 * returns the index of the smallest value in the range between
 * index positions p1 and p2, inclusive.
 */
int findSmallest(Comparable* a, int p1, int p2) {
    int smallestIndex = p1;
    for (int i = p1 + 1; i <= p2; i++) {
        if (a[i] < a[smallestIndex]) smallestIndex = i;
    }
    return smallestIndex;
}

/*
 * selection: swap
 * -----
 * exchanges the larger values passed by reference.
 */
void swap(int* a, int* b) {
    int temp = *a;
    *a = *b;
    *b = temp;
}
```

Simulating Selection Sort

```
void main() {
    void Sort(Comparable* a) {
        for (int k = 0; k < a->n(); k++) {
            int m = findSmallest(a, k, a->n() - 1);
            swap(a[k], a[m]);
        }
    }
}
```

Efficiency of Selection Sort

- The primary question for today is how one might evaluate the efficiency of an algorithm such as selection sort.
- One strategy is to measure the actual time it takes to run for arrays of different sizes. In C++, you can measure elapsed time by calling the `time` function, which returns the current time in milliseconds. Using this strategy, however, requires some care:
 - The `time` function is often too rough for accurate measurement. It therefore makes sense to measure several runs together and then divide the total time by the number of repetitions.
 - Most algorithms show some variability depending on the data. To avoid distortion, you should run several independent trials with different data and average the results.
 - Some measurements are likely to be wildly off because the computer needs to run some background task. Such data points must be discarded as you work through the analysis.

Measuring Sort Timings

The following table shows the average timing of the selection sort algorithm after removing outlying trials that differ by more than two standard deviations from the mean. The column labeled μ (the Greek letter mu, which is the standard statistical symbol for the mean) is a reasonably good estimate of running time.

	Test 1	Test 2	Test 3	Test 4	Test 5	Test 6	Test 7	Test 8	Test 9	Test 10	μ	σ
N = 10	.0021	.0025	.0022	.0026	.0028	.0030	.0027	.0023	.0025	.0024	.0024	.00029
20	.006	.007	.008	.007	.007		.007	.007	.007	.007	.007	.00036
30	.014	.014	.014	.015	.014	.014	.014	.014	.014	.014	.014	.00013
40	.026	.024	.023	.026	.023	.025	.025	.026	.025	.027	.025	.0014
50	.036	.037	.036	.041	.042	.039		.039	.034	.038	.036	.0025
100	.187	.152	.168	.176	.146	.146	.165	.146	.176	.154	.162	.0151
200	3.94	3.83	4.86	3.76	4.11	3.51	3.48	3.61	3.31	3.45	3.68	0.202
1000	13.40	12.90	13.80	12.80	12.80	14.10	12.70		16.00	15.50	14.32	1.69
10000	332.5	355.9	391.7	321.6	388.3	321.3	321.3	380.7	332.1	321.3	348.4	23.83
100000	1318		1327	1318	1331	1336	1318	1335	1325	1318	1326	7.50

Selection Sort Running Times

- Many algorithms that operate on vectors have running times that are proportional to the size of the array. If you multiply the number of values by ten, you would expect those algorithms to take ten times as long.

N	time
10	.0024
100	0.162
1000	14.32
10000	1322

- As the running times on the preceding slide make clear, the situation for selection sort is very different. The table on the right shows the average running time when selection sort is applied to 10, 100, 1000, and 10000 values.
- As a rough approximation—particularly as you work with larger values of N —it appears that every ten-fold increase in the size of the array means that selection sort takes about 100 times as long.

Counting Operations

- Another way to estimate the running time is to count how many operations are required to sort an array of size N .
- In the selection sort implementation, the section of code that is executed most frequently (and therefore contributes the most to the running time) is the body of the `findSmallest` method. The number of operations involved in each call to `findSmallest` changes as the algorithm proceeds:

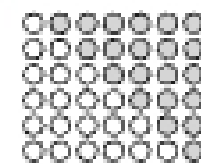
N values are considered on the first call to `findSmallest`.
 $N - 1$ values are considered on the second call.
 $N - 2$ values are considered on the third call, and so on.

- In mathematical notation, the number of values considered in `findSmallest` can be expressed as a summation, which can then be transformed into a simple formula:

$$1 + 2 + 3 + \dots + (N - 1) + N = \sum_{i=1}^N i = \frac{N \cdot (N + 1)}{2}$$

A Geometric Insight

- You can convince yourself that $1 + 2 + 3 + \dots + (N - 2) + (N - 1) + N = \frac{N \cdot (N + 1)}{2}$ by thinking about the problem geometrically.
- The terms on the left side of the formula can be arranged into a triangle, as shown at the bottom of this slide for $N = 6$.
- If you duplicate the triangle and rotate it by 180° , you get a rectangle that in this case contains 6×7 dots, half of which belong to each triangle.



Quadratic Growth

- The reason behind the rapid growth in the running time of selection sort becomes clear if you make a table showing the value of $\frac{N \cdot (N + 1)}{2}$ for various values of N :

N	$\frac{N \cdot (N + 1)}{2}$
10	55
100	5050
1000	500,500
10000	50,005,000

- The growth pattern in the right column is similar to that of the measured running time of the selection sort algorithm. As the value of N increases by a factor of 10, the value of $\frac{N \cdot (N + 1)}{2}$ increases by a factor of around 100, which is 10^2 . Algorithms whose running times increase in proportion to the square of the problem size are said to be **quadratic**.

Big-O Notation

- The most common way to express computational complexity is to use **big-O notation**, which was introduced by the German mathematician Paul Bachmann in 1892.
- Big-O notation consists of the letter O followed by a formula that offers a qualitative assessment of running time as a function of the problem size, traditionally denoted as N . For example, the computational complexity of linear search is

$$O(N)$$

and the computational complexity of radix sort is

$$O(N \log N)$$

- If you read these formulas aloud, you would pronounce them as "big-O of N " and "big-O of $N \log N$ " respectively.

Common Simplifications of Big-O

- Given that big-O notation is designed to provide a qualitative assessment, it is important to make the formula inside the parentheses as simple as possible.
- When you write a big-O expression, you should always make the following simplifications:
 - Eliminate any term whose contribution to the running time ceases to be significant as N becomes large.
 - Eliminate any constant factors.
- The computational complexity of selection sort is therefore

$$O(N^2)$$

and not

$$O\left(\frac{N \cdot (N + 1)}{2}\right)$$


Deducing Complexity from the Code

- In many cases, you can deduce the computational complexity of a program directly from the structure of the code.
- The standard approach to doing this type of analysis begins with looking for any section of code that is executed more often than other parts of the program. As long as the individual operations involved in an algorithm take roughly the same amount of time, the operations that are executed most often will come to dominate the overall running time.
- In the selection sort implementation, for example, the most commonly executed statement is the `if` statement inside the `FindSmallest` method. This statement is part of two `for` loops, one in `FindSmallest` itself and one in `Sort`. The total number of executions is

$$1 + 2 + 3 + \dots + (N - 1) + N$$

which is $O(N^2)$.

Exercise: Computational Complexity

Assuming that none of the steps in the body of the following `for` loops depend on the problem size stored in the variable `n`, what is the computational complexity of each of the following examples:

- a)

```
for (int k = 0; k < n; k++) {
    for (int j = 0; j < n; j++) {
        ...loop body...
    }
}
```
- b)

```
for (int k = 0; k < n; k++) {
    ...loop body...
}
```
- c)

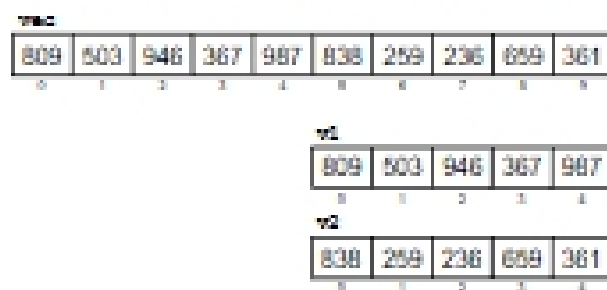
```
for (int k = 0; k < 100; k++) {
    for (int j = 0; j < n; j++) {
        ...loop body...
    }
}
```

Finding a More Efficient Strategy

- As long as arrays are small, selection sort is a perfectly workable strategy. Even for 10,000 elements, the average running time of selection sort is just over a second.
- The quadratic behavior of selection sort, however, makes it less attractive for the very large arrays that one encounters in commercial applications. Assuming that the quadratic growth pattern continues beyond the timings reported in the table, sorting 100,000 values would require two minutes, and sorting 1,000,000 values would take more than three hours.
- The computational complexity of the selection sort algorithm, however, holds out some hope:
 - Sorting twice as many elements takes four times as long.
 - Sorting half as many elements takes only one fourth the time.
 - Is there any way to use sorting half an array as a subtask in a recursive solution to the sorting problem?

The Merge Sort Idea

- Divide the vector into two halves: `v1` and `v2`.
- Sort each of `v1` and `v2` recursively.
- Clear the original vector.
- Merge elements into the original vector by choosing the smallest element from `v1` or `v2` on each cycle.



The Merge Sort Implementation

```

//
// the merge sort algorithm consists of the following steps:
//
// 1. divide the vector into two halves.
// 2. sort each of these smaller vectors recursively.
// 3. merge the two vectors back into the original one.
//
void mergeSort(int a[], int n) {
    int m = n / 2;
    if (n <= 1) return;
    mergeSort(a, m);
    mergeSort(a + m, n - m);
    for (int k = 0; k < n; k++) {
        int i = k / 2;
        int j = (k + 1) / 2;
        int min = i;
        while (i < m && j < n - m) {
            if (a[i] < a[j]) min = i;
            else min = j;
        }
        a[k] = a[min];
    }
}

```