

# Introduction and Linear-time Selection

The purpose of this lecture is to give a brief overview of the topic of algorithm analysis and the kind of thinking it involves. As a motivating problem, we consider the famous Quicksort algorithm and the problem of selecting a pivot. We review the complexity of Quicksort under various assumptions and use it to motivate the *selection* problem. This allows us to improve the Quicksort algorithm by deterministically finding the optimal pivot element (the median element) in linear worst-case complexity.

These problems illustrate some of the ideas and tools we will be using (and building upon) in this course. We will practice writing and solving recurrence relations, and bounding the *expected* cost of a randomized algorithm, both of which are key tools for the analysis of algorithms that we will use throughout this class.

## Objectives of this lecture

In this lecture, we cover:

- Formal analysis of algorithms, models of computation,
- The analysis and complexity of the Quicksort algorithm,
- Finding the median: A randomized algorithm in expected linear time,
- Analyzing a randomized recursive algorithm,
- A deterministic linear-time algorithm for medians.

## Recommended study resources

- CLRS, *Introduction to Algorithms*, Chapter 9, Medians and Order Statistics
- DPV, *Algorithms*, Chapter 2.4, Medians

# 1 Goals of the Course

This course is about the design and analysis of algorithms — how to design correct, efficient algorithms, and how to think clearly about analyzing correctness and running time. What is an algorithm? At its most basic, an algorithm is a method for solving a computational problem. A recipe. Along with an algorithm comes a specification that says what the algorithm’s guarantees are. For example, we might be able to say that our algorithm correctly solves the problem in question and runs in time at most  $f(n)$  on any input of size  $n$ . This course is about the whole package: the design of efficient algorithms, *and* proving that they meet desired specifications. For each of these parts, we will examine important techniques that have been developed, and with practice we will build up our ability to think clearly about the key issues that arise.

The main goal of this course is to provide the intellectual tools for designing and analyzing your own algorithms for problems you need to solve in the future. Some tools we will discuss are Dynamic Programming, Divide-and-Conquer, Hashing and other Data Structures, Randomization, Network Flows, and Linear Programming. Some analytical tools we will discuss and use are Recurrences, Probabilistic Analysis, Amortized Analysis, and Potential Functions. We will additionally discuss some approaches for dealing with NP-complete problems, including the notion of approximation algorithms.

Another goal will be to discuss models that go beyond the traditional input-output model. In the traditional model we consider the algorithm to be given the entire input in advance and it just has to perform the computation and give the output. This model is great when it applies, but it is not always the right model. For instance, some problems may be challenging because they require decisions to be made without having full information. Algorithms that solve such problems are called *online* algorithms, which we will also discuss. In other settings, we may have to deal with computing quantities of a “stream” of input data where the space we have is much smaller than the data. In yet other settings, the input is being held by a set of selfish agents who may or may not tell us the correct values.

## 1.1 On guarantees and specifications

One focus of this course is on proving correctness and running-time guarantees for algorithms. Why is having such a guarantee useful? Suppose we are talking about the problem of sorting a list of  $n$  numbers. It is pretty clear why we at least want to know that our algorithm is correct, so we don’t have to worry about whether it has given us the right answer all the time. But, why analyze running time? Why not just code up our algorithm and test it on 100 random inputs and see what happens? Here are a few reasons that motivate our concern with this kind of analysis — you can probably think of more reasons too:

**Composability** A guarantee on running time gives a “clean interface”. It means that we can use the algorithm as a subroutine in another algorithm, without needing to worry whether the inputs on which it is being used now match the kinds of inputs on which it was originally tested.

**Scaling** The types of guarantees we will examine will tell us how the running time scales with the size of the problem instance. This is useful to know for a variety of reasons. For instance,

it tells us roughly how large a problem size we can reasonably expect to handle given some amount of resources.

**Designing better algorithms** Analyzing the asymptotic running time of algorithms is a useful way of thinking about algorithms that often leads to non-obvious improvements.

**Understanding** An analysis can tell us what parts of an algorithm are crucial for what kinds of inputs, and why. If we later get a different but related task, we can often use our analysis to quickly tell us if a small modification to our existing algorithm can be expected to give similar performance to the new problem.

**Complexity-theoretic motivation** In Complexity Theory, we want to know: “how hard is fundamental problem  $X$  really?” For instance, we might know that for a given problem, no algorithm can possibly run in time  $o(n \log n)$  (growing more slowly than  $n \log n$  in the limit) and we have an algorithm that runs in time  $O(n^{3/2})$ . This tells us how well we understand the problem, and also how much room for improvement we have.

It is often helpful when thinking about algorithms to imagine a game where one player is the algorithm designer trying to find an algorithm for the problem, and its opponent, the “adversary”, is trying to find an input that will cause the algorithm to run slowly. An algorithm with good worst-case guarantees is one that performs well no matter what input the adversary chooses.

## 1.2 Formal analysis of algorithms

In this course we are interested in the *formal* analysis of algorithms. In your previous courses on algorithms and programming, you have hopefully studied the notion of the complexity of an algorithm, but you might have done it less formally. Most commonly when first learning about algorithm analysis, you learn to “count the number of operations” done by an algorithm. This is a bit vague and underspecified, but it does the job most of the time.

Our goal is to make this more precise to get a more rigorous view on what constitutes the complexity of an algorithm. To do this, we need the notion of a *model of computation*.

### **Definition: Model of computation**

A model of computation consists of:

1. A set of allowed *operations* that an algorithm may perform, and
2. A cost for each operation, sometimes separately called the *cost model*.

To analyze the complexity of an algorithm, we consider it in a particular model of computation and add up the total costs of all the operations of that algorithm. Sometimes costs will be specified *concretely*, e.g., operation  $X$  costs 1 and operation  $Y$  costs 2. Other times we will only be interested in asymptotic bounds, so we may specify that an operation costs  $O(1)$  time but be uninterested in constant factors in the analysis. We will see many examples in the first lectures.

## 2 Sorting and selection in the comparison model

For the first lecture we will consider algorithm in the *comparison model*. In the comparison model we assume that the input consists of some comparable elements (i.e., we can ask is  $x < y$  for two elements  $x$  and  $y$ ) but we can not assume anything else about their type. For example, we can not assume that they are integers or numbers at all, we can not assume that they are strings, or that we can hash them, etc. Comparing the elements is the *only* information we have about their values. We define the comparison model as follows.

### *Definition: Comparison Model*

In the *comparison model*, we have an input consisting of  $n$  elements in some initial order. An algorithm may compare two elements (asking is  $a_i < a_j$ ?) at a cost of 1. Moving the items, copying them, swapping them, etc., is *free*. No other operations on the items are allowed (using them as indices, adding them, hashing them, etc).

The comparison model is widely used to analyze sorting and selection (e.g., find the max, the second-max, the  $k^{\text{th}}$  largest element, etc.) algorithms. On a practical level, sorting and selection algorithms designed for the comparison model are widely applicable because they make no assumptions about the types of the inputs, so they can be used to sort any types for which the problem of sorting makes sense (the elements must of course be comparable to define what sorting even means!) Conversely however, algorithms designed for the comparison model may be less efficient than algorithms which are specialized for specific types, so we get a tradeoff between generality and performance. We will see an example of this in a couple of lectures.

Since the problems of sorting and selection are fundamentally about ordering, defining the cost of the algorithms in terms of the number of comparisons performed is a natural metric. This number also *usually* (but not always) matches asymptotically the number of operations performed in a less concrete model of computation, so it provides a reasonable prediction of the algorithms' performance. The comparison model is also widely used for studying *lower bounds* on the complexity of sorting and selection problems, which we will study next lecture.

### 2.1 Revisiting a classic: The Quicksort algorithm

Quicksort is one of the most famous and well known algorithms in all of computer science.

#### *Algorithm 1: Quicksort*

- Given array  $A$  of size  $n$ ,
1. Pick an arbitrary pivot element  $p$  from  $A$ .
  2.  $\text{LESS} \leftarrow \{a_i \text{ such that } a_i < p\}$
  3.  $\text{GREATER} \leftarrow \{a_i \text{ such that } a_i > p\}$
  4. **return** Quicksort(LESS) +  $\{p\}$  + Quicksort(GREATER)

The step of splitting  $A$  into LESS and GREATER is called *partitioning*. The pseudocode given above gives the simplest version of the algorithm which copies/moves the elements into a new array rather than partitioning *in place*, which is more complicated, but the fundamental idea of the algorithm and its cost is the same.

We will use Quicksort as a starting point to review what we know about complexity analysis and to ensure that we do so rigorously. Then, we will spend the rest of the lecture figuring out how to improve the algorithm and make its complexity better! First, we need to make sure we remember how to measure complexity. Remember that there isn't a single measure of complexity, there are multiple, so we will review a bunch of them and see how they apply to Quicksort. Note that measures of complexity are orthogonal to the model of computation. We can consider any combination of both to obtain different ways of analyzing the same algorithm!

## 2.2 Worst-case complexity

### *Definition: Worst-case Cost*

The *worst-case* complexity of an algorithm in a given model of computation is the largest cost that the algorithm could incur on any possible input, usually parameterized by the size of the input.

You might remember from your earlier classes that the worst-case cost of Quicksort occurs when the pivot selected happens to always be the smallest or largest element in the array.

### *Theorem: Worst-case cost of Quicksort*

The worst-case cost of Quicksort is  $O(n^2)$  comparisons.

## 2.3 Average-case complexity

The worst-case cost of Quicksort is bad, but it seems to perform well *most* of the time. This is not rigorous, but fortunately there is a formal way to state this using *average-case* complexity.

### *Definition: Average-case Cost*

The *average-case* complexity of an algorithm in a given model of computation is the average of the costs of all possible inputs to the algorithm, usually parameterized by the size of the input.

### *Remark: Random interpretation of average-case cost*

By definition, the average-case cost of an algorithm is equivalent to the *expected value* of its cost over a uniform distribution of possible inputs. You can therefore think of the average cost as the cost of the algorithm on a random input.

Computing average-case complexity tends to be quite a lot more work than computing worst-case complexity since we need a way to enumerate all possible inputs. This can be done for Quicksort using recurrence relations, and you may have seen this in a previous class. We won't repeat the analysis here.

**Theorem: Average-case cost of Quicksort**

The average-case cost of Quicksort is  $O(n \log n)$  comparisons.

## 2.4 Randomized complexity

One interpretation of the average-case cost of Quicksort is that if we run it on a random input then the expected cost is just  $O(n \log n)$ . This is great... if the input to the algorithm is random. However, real life data is almost never random so it is a bad idea to design your algorithms with the assumption that the input is random! Thankfully, there is a simple yet powerful way to overcome this foolish assumption: put the randomness *in the algorithm* instead! This leads to *randomized Quicksort*, a variant of the algorithm that does not perform horribly for any particular input! The difference is just a single word compared to the vanilla variant.

**Algorithm 2: RandomQuicksort**

Given array  $A$  of size  $n$ ,

1. Pick a **random** pivot element  $p$  from  $A$ .
2.  $\text{LESS} \leftarrow \{a_i \text{ such that } a_i < p\}$
3.  $\text{GREATER} \leftarrow \{a_i \text{ such that } a_i > p\}$
4. **return**  $\text{RandomQuicksort}(\text{LESS}) + \{p\} + \text{RandomQuicksort}(\text{GREATER})$

Unlike the previous algorithms, this one now includes randomness, which means that its cost on a particular input is not fixed, it could vary from one run to the next! More formally, the runtime of the algorithm is now a probability distribution rather than a fixed number. When describing the complexity of a randomized algorithm, we usually give some property of the probability distribution, most commonly, we use the *expected value*. By default, unless otherwise specified, we still consider the *worst-case* input, i.e., we want to compute the maximum over all possible inputs, of the expected value of the cost of the algorithm over the random choices made by the algorithm. We refer to this as the *expected complexity* or *expected cost* of the algorithm.

**Definition: Expected Cost**

The *expected* complexity of a randomized algorithm in a given model of computation is the maximum cost over of all possible inputs to the algorithm, of the expected value of the cost of that input, where the expected value is over the distribution of random choices made by the algorithm.

#### *Remark: Misconceptions of expected cost*

The definition of expected complexity is a bit tricky, so it is important to not have misconceptions about it. Here are some important things to keep in mind:

1. Expected complexity by default considers *worst-case inputs*. We are not analyzing the algorithm for a random/average-case input.
2. We are also not considering worst-case random numbers. The input is worst case, and the randomness is well... random. The expected value is computed over the distribution of the random choices of the algorithm.

#### *Theorem: Expected cost of Randomized Quicksort*

The expected cost of Randomized Quicksort is  $O(n \log n)$  comparisons.

The proof is almost identical to that of the average-case cost of (non-randomized) Quicksort since there is essentially a one-to-one mapping between random inputs and randomly selecting pivots. You have probably run across the proof previously so we again won't repeat it here.

## 2.5 Can we do better?

We have watched the analysis of Quicksort go from  $O(n^2)$  worst-case to  $O(n \log n)$  average-case, to expected  $O(n \log n)$  by mixing randomization into the algorithm. Can we (theoretically at least) improve the algorithm even further or is this the end of the journey? Well, we know several other comparison-based sorting algorithms that run in  $O(n \log n)$  comparisons, like Mergesort and Heapsort, and both of them are deterministic! It would be very cool if Quicksort could also be made to take just  $O(n \log n)$  comparisons and still be deterministic too. To achieve this, we would need to find a balanced partition (i.e., pick a good pivot that splits the input into similarly sized halves). That motivates us to consider the **median-finding** problem. If we could find the median of an unsorted array in linear time, we could use that as the pivot and all of our dreams would come true! That will be the remainder of this lecture.

## 3 The Median and Selection Problems

One thing that makes algorithm design “Computer Science” is that solving a problem in the most obvious way from its definitions is often not the best way to get a solution. An example of this is median finding. Recall the concept of the median of a set. For a set of  $n$  elements, this is the “middle” element in the set, i.e., there are exactly  $\lfloor n/2 \rfloor$  elements larger than it. In computer sciencey terms, if the elements are zero-indexed, the median is the  $\lfloor (n-1)/2 \rfloor^{\text{th}}$  element of the set when represented in sorted order.

Given an unsorted array, how quickly can one find the median element? Perhaps the simplest solution, one that can be described in a single sentence and implemented with one or two lines of code in your favorite programming language, is to sort the array and then read off the

element in position  $\lfloor (n-1)/2 \rfloor$ , which takes  $O(n \log n)$  comparisons using your favorite sorting algorithm, such as MergeSort or HeapSort (deterministic) or QuickSort (randomized).<sup>1</sup>

Can one do it more quickly than by sorting? In the remainder of this lecture we describe two linear-time algorithms for this problem: first a simpler randomized algorithm, and then an improvement that makes it deterministic! More generally, we solve the problem of finding the  $k^{\text{th}}$  smallest out of an unsorted array of  $n$  elements.

### 3.1 The problem and a randomized solution

Let's consider a problem that is slightly more general than median-finding:

**Problem: Select- $k$  /  $k^{\text{th}}$  Smallest**

Find the  $k^{\text{th}}$  smallest element in an unsorted array of size  $n$ .

To remove ambiguity, we will assume that our array is zero indexed, so the  $k^{\text{th}}$  smallest element is the element that would be in position  $k$  if the array were sorted. Alternatively, it is the element such that exactly  $k$  other elements are smaller than it. Additionally, let's say all elements are distinct to avoid the question of what we mean by the  $k^{\text{th}}$  smallest when we have duplicates.

**Idea: Eliminate redundancy** The key idea to obtaining a linear time algorithm is to identify and eliminate redundant/wasted work in the “naive” algorithm that just sorts the input and outputs the  $k^{\text{th}}$  element. Note that by sorting the array, we are not just finding the  $k^{\text{th}}$  smallest element for the given value of  $k$ , we are actually finding the answer for *every possible* value of  $k$ . So instead of completely sorting the array, it would suffice to “partially sort” the array such that element  $k$  ends up in the correct position, but the remaining elements might still not be perfectly sorted. This description sounds suspiciously similar to Quicksort, which partitions the array by putting the pivot element into its correctly sorted position in the array (without guaranteeing yet that the rest of the array is sorted). In essence, this is partially sorting the array with respect to the pivot. Recursively sorting both sides then sorts the entire array.

So, what if we were to just run Quicksort, but skip some of the steps that do not matter? Specifically, note that if we run Quicksort and our goal is to output the  $k^{\text{th}}$  element at the end, that after partitioning the array around the pivot, we know which of the two sides must contain the answer. Therefore instead of recursively sorting both sides, we ignore the side that can not contain the answer and recurse just once. That's it!

**Implementation** More specifically, the algorithm chooses a random pivot and then partitions the array into two sets LESS and GREATER consisting of those elements less than and greater than the pivot respectively. After the partitioning step we can tell which of LESS or GREATER has the item we are looking for, just by looking at their sizes. For instance, if we are looking for the 87th-smallest element in our array, and suppose that after choosing the pivot

<sup>1</sup>Of course using a sorting algorithm to find the pivot for Quicksort would be rather useless, but it is a good start to get some intuition on the complexity of the median problem. We know that the problem is solvable in  $O(n \log n)$ , so our goal remains to bring this down to  $O(n)$ .



and partitioning we find that LESS has 200 elements, then we just need to find the 87th smallest element in LESS. On the other hand, if we find LESS has 40 elements, then we just need to find the  $87 - 40 - 1 = 46$ th smallest element in GREATER. (And if LESS has size exactly 86 then we can just return the pivot). One might at first think that allowing the algorithm to only recurse on one subset rather than both would just cut down time by a factor of 2. However, since this is occurring recursively, it compounds the savings and we end up with  $\Theta(n)$  rather than  $\Theta(n \log n)$  time. This algorithm is often called Randomized-Select, or QuickSelect.

### Algorithm 3: QuickSelect

Given array  $A$  of size  $n$  and integer  $0 \leq k \leq n - 1$ ,

1. Pick a pivot element  $p$  at random from  $A$ .
2. Split  $A$  into subarrays LESS and GREATER by comparing each element to  $p$ .
3. (a) If  $|\text{LESS}| > k$ , then return QuickSelect(LESS,  $k$ ).  
 (b) If  $|\text{LESS}| < k$ , then return QuickSelect(GREATER,  $k - |\text{LESS}| - 1$ )  
 (c) If  $|\text{LESS}| = k$ , then return  $p$ . [always happens when  $n = 1$ ]

### Theorem 1

The expected number of comparisons for QuickSelect is at most  $8n$ .

Formally, let  $T(n)$  denote the expected number of comparisons performed by QuickSelect on any (worst-case) input of size  $n$  for any value of  $k$ . What we want is a recurrence relation that looks like

$$T(n) \leq n - 1 + \mathbb{E}_X[T(X)],$$

where  $n - 1$  comparisons come from comparing the pivot to every other element and placing them into LESS and GREATER, and  $X$  is a random variable corresponding to the size of the subproblem that is solved recursively. We can't just go and solve this recurrence because we don't yet know what  $X$  or  $\mathbb{E}[T(X)]$  look like yet.

**A first attempt** Before giving a formal proof, here's some intuition and a slightly incorrect first attempt. First of all, how large is  $X$ , the size of the array given to the recursive call? It depends on two things: the value of  $k$  and the *randomly-chosen pivot*. After partitioning the input into LESS and GREATER, whose size adds up to  $n - 1$ , the algorithm recursively calls QuickSelect on one of them, but which one? Since we are interested in the behavior for a *worst-case input*, we can assume pessimistically that the value of  $k$  will always make us choose the bigger of LESS and GREATER. Therefore the question becomes: if we choose a random pivot and split the input into LESS and GREATER, how large is the larger of the two of them? Well, possible sizes of the splits (ignoring rounding) are

$$(0, n - 1), (1, n - 2), (2, n - 3), \dots, (n/2 - 2, n/2 + 1), (n/2 - 1, n/2),$$

so we can see that the larger of the two is a random number from

$$n-1, n-2, n-3, \dots, n/2+1, n/2.$$

So, the expected size of the larger half is about  $3n/4$ , again, ignoring rounding errors.

Another way to say this is that if we split a candy bar at random into two pieces, the expected size of the larger piece is  $3/4$  of the bar. Using this, we might be tempted to write the following recurrence relation, which is almost correct, but not quite.

$$T(n) \leq n-1 + T(3n/4). \quad \textcolor{red}{(This is wrong!)}$$

If we go through the motions of solving this recurrence, we get  $T(n) = O(n)$ , but unfortunately, this derivation is not quite valid. The reasoning is that  $3n/4$  is only the *expected* size of the larger piece. That is, if  $X$  is the size of the larger piece, we have written a recurrence where the cost of the recursive call is  $T(\mathbb{E}[X])$ , but it was supposed to be  $\mathbb{E}[T(X)]$ , and these are not the same thing! (As an exercise, argue that the two could differ by a lot.)<sup>2</sup> Let's now see this a bit more formally.

**A correct proof** To correct the proof, we need to correctly analyze the *expected value of  $T(X)$* , rather than the value of  $T$  for the expected value of  $X$ . By the definition of expected value, we want to bound the following:

$$\mathbb{E}[T(X)] = \sum_{x=1}^{n-1} \Pr[X = x] \cdot T(x)$$

It is possible to bound the cost by attacking the above expression with induction and a lot of algebra, but it's a bit messy and does not provide very much intuition. Instead, here is a simpler way to go about it. Since we are okay with an upper bound and don't need an exact solution to the recurrence, we can try to upper bound most of the terms in the sum by making them the same. There isn't a single correct way to do this, but the following intuition works. We already know from our slightly incorrect attempt that if the recurrence recurses on  $T(3/4n)$  then we get an  $O(n)$  upper bound, so we can try to make the correct recurrence also depend on  $T(3/4n)$ . To do so, we use the fact that  $T(n)$  is an increasing function.<sup>3</sup>

So, we can upper bound the quantity in question as follows.

$$\mathbb{E}[T(X)] \leq \Pr\left[X \leq \frac{3n}{4}\right] \cdot T\left(\frac{3n}{4}\right) + \Pr\left[X > \frac{3n}{4}\right] \cdot T(n).$$

So, with what probability is  $X$  at most  $3/4$ ? This happens when the smaller of LESS and GREATER is at least one quarter of the elements, i.e., when the pivot is not in the bottom quarter or top quarter. This means the pivot needs to be in the middle half of the data, which happens with probability  $1/2$ . The other half the time, the size of  $X$  will be larger, at most  $n$ . Although this might sound rather loose, this is good enough to write down a good upper bound recurrence!

---

<sup>2</sup>It turns out that these two quantities are equal if  $T$  is linear, which it is in this particular case. However, we can not make this assumption in the proof, because that is we are trying to prove in the first place! Assuming that  $T$  is linear to prove that  $T$  is linear would be circular logic.

<sup>3</sup>If we wanted to be extremely thorough, we could prove this claim separately, but it should be very reasonable to believe that the algorithm does not get faster if we give it a larger input.

*Proof Theorem 1.* We can now use the bound above to show that  $T(n) = O(n)$  as desired. We first write

$$\mathbb{E}[T(X)] \leq \frac{1}{2}T\left(\frac{3n}{4}\right) + \frac{1}{2}T(n).$$

Returning to our original recurrence, we can now correctly assert that

$$T(n) \leq n - 1 + \left(\frac{1}{2}T\left(\frac{3n}{4}\right) + \frac{1}{2}T(n)\right),$$

and multiplying both sides by 2 then subtracting  $T(n)$ , we obtain

$$T(n) \leq 2(n - 1) + T\left(\frac{3n}{4}\right).$$

We can now use induction to prove that this recurrence relation satisfies  $T(n) \leq 8n$ . The base case is simple, when  $n = 1$  there are no comparisons, so  $T(1) = 0$ . Now assume for the sake of induction that  $T(i) \leq 8i$  for all  $i < k$ . We want to show that  $T(k) \leq 8k$ . We have

$$\begin{aligned} T(k) &\leq 2(n - 1) + T\left(\frac{3n}{4}\right), \\ &\leq 2(n - 1) + 8\left(\frac{3n}{4}\right), && \text{Use the inductive hypothesis} \\ &\leq 2n + 6n, \\ &= 8n, \end{aligned}$$

which proves our desired bound. □

## 3.2 A deterministic linear-time selection algorithm

What about a deterministic linear-time algorithm? For a long time it was thought this was impossible, and that there was no method faster than first sorting the array. In the process of trying to prove this formally, it was discovered that this thinking was incorrect, and in 1972 a deterministic linear time algorithm was developed by Manuel Blum, Bob Floyd, Vaughan Pratt, Ron Rivest, and Bob Tarjan.<sup>4</sup>

The idea of the algorithm is that one would like to pick a pivot deterministically in a way that produces a good split. Ideally, we would like the pivot to be the median element so that the two sides are the same size. But, this is the same problem we are trying to solve in the first place! So, instead, we will give ourselves leeway by allowing the pivot to be any element that is “roughly” in the middle (i.e., some kind of “approximate median”). We will use a technique called the *median of medians* which takes the medians of a bunch of small groups of elements, and then finds the median of those medians. It has the wonderful guarantee that the selected element is greater than at least 30% of the elements of the array, and smaller than at least 30% of the elements of the array. This makes it work great as an approximate median and therefore a good pivot choice! The algorithm is as follows:

---

<sup>4</sup>That’s 4 Turing Award winners on that one paper!

#### Algorithm 4: DeterministicSelect

Given array  $A$  of size  $n$  and integer  $k \leq n$ ,

1. Group the array into  $n/5$  groups of size 5 and find the median of each group. (For simplicity, we will ignore integrality issues.)
2. Recursively, find the true median of the medians. Call this  $p$ .
3. Use  $p$  as a pivot to split the array into subarrays LESS and GREATER.
4. Recurse on the appropriate piece in the same way as Quickselect.

#### Theorem 2

DeterministicSelect makes  $O(n)$  comparisons to find the  $k^{\text{th}}$  smallest element in an array of size  $n$ .

*Proof of Theorem 2.* Let  $T(n)$  denote the worst-case number of comparisons performed by the DeterministicSelect algorithm on inputs of size  $n$ .

Step 1 takes time  $O(n)$ , since it takes just constant time to find the median of 5 elements. Step 2 takes time at most  $T(n/5)$ . Step 3 again takes time  $O(n)$ . Now, we claim that at least  $3/10$  of the array is  $\leq p$ , and at least  $3/10$  of the array is  $\geq p$ . Assuming for the moment that this claim is true, Step 4 takes time at most  $T(7n/10)$ , and we have the recurrence:

$$T(n) \leq cn + T\left(\frac{n}{5}\right) + T\left(\frac{7n}{10}\right),$$

for some constant  $c$ . Before solving this recurrence, let's prove the claim we made that the pivot will be roughly near the middle of the array. So, the question is: how bad can the median of medians be? But first, let's do an example. Suppose the array has 15 elements and breaks down into three groups of 5 like this:

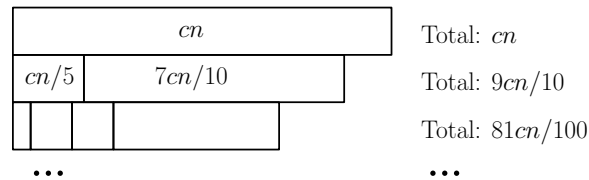
$$\{1, 2, 3, 10, 11\}, \{4, 5, 6, 12, 13\}, \{7, 8, 9, 14, 15\}.$$

In this case, the medians are 3, 6, and 9, and the median of the medians  $p$  is 6. There are five elements less than  $p$  and nine elements greater.

In general, what is the worst case? If there are  $g = n/5$  groups, then we know that in at least  $\lceil g/2 \rceil$  of them (those groups whose median is  $\leq p$ ) at least three of the five elements are  $\leq p$ . Therefore, the total number of elements  $\leq p$  is at least  $3\lceil g/2 \rceil \geq 3n/10$ . Similarly, the total number of elements  $\geq p$  is also at least  $3\lceil g/2 \rceil \geq 3n/10$ .

Now, finally, let's solve the recurrence. We have been solving a lot of recurrences by the “guess and check” method, which works here too, but how could we just stare at this and *know* that the answer is linear in  $n$ ? One way to do that is to consider the “stack of bricks” view of the recursion tree that you might have discussed in your previous classes.

In particular, let's build the recursion tree for the recurrence (1), making each node as wide as the quantity inside it:



Notice that even if this stack-of-bricks continues downward forever, the total sum is at most

$$cn(1 + (9/10) + (9/10)^2 + (9/10)^3 + \dots),$$

which is at most  $10cn$ . This proves the theorem.  $\square$

Notice that in our analysis of the recurrence (1) the key property we used was that  $n/5 + 7n/10 < n$ . More generally, we see here that if we have a problem of size  $n$  that we can solve by performing recursive calls on pieces whose total size is at most  $(1 - \epsilon)n$  for some constant  $\epsilon > 0$  (plus some additional  $O(n)$  work), then the total time spent will be just linear in  $n$ .

#### Lemma 1

For constants  $c$  and  $a_1, \dots, a_k$  such that  $a_1 + \dots + a_k < 1$ , the recurrence

$$T(n) \leq T(a_1 n) + T(a_2 n) + \dots + T(a_k n) + cn$$

solves to  $T(n) = O(n)$ .

### 3.3 Optimal deterministic Quicksort

Armed with a deterministic linear-time median-finding algorithm, we now have the tools to create an  $O(n \log n)$ -cost deterministic Quicksort! Just use the linear-time median algorithm (DeterministicSelect) to select the pivot, then the divide-and-conquer subproblems are at most half the input, which gives us a cost of

$$T(n) = T\left(\frac{n}{2}\right) + T\left(\frac{n}{2}\right) + O(n).$$

Solving this using standard techniques (or remembering the solution from a previous class!) shows us that the cost is  $O(n \log n)$  with no randomness required!

## Exercises: Selection Algorithms & Recurrences

**Problem 1.** Recall the attempted analysis of randomized quickselect where we accidentally assumed that  $\mathbb{E}[T(X)] = T(\mathbb{E}[X])$ . Let  $X$  be a random variable, and find an increasing function  $F : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  such that  $\mathbb{E}[F(i)] \gg F(\mathbb{E}[i])$ .

**Problem 2.** Show that for constants  $c$  and  $a_1, \dots, a_k$  such that  $a_1 + \dots + a_k = 1$  and each  $a_i < 1$ , the recurrence

$$T(n) \leq T(a_1 n) + T(a_2 n) + \dots + T(a_k n) + c n$$

solves to  $T(n) = O(n \log n)$ . Show that this is best possible by observing that  $T(n) = T(n/2) + T(n/2) + n$  solves to  $T(n) = \Theta(n \log n)$ .

**Problem 3.** Consider the median of medians algorithm. What happens if we split the elements into  $n/3$  groups of size 3 instead? Or  $n/k$  groups of size  $k$  for larger odd values of  $k$ ?

**Problem 4.** The rank of an element  $a$  with respect to a list  $A$  of  $n$  distinct elements is  $|\{e \in A \mid e \leq a\}|$  is the number of elements in  $A$  no greater than  $a$ . Hence the rank of the smallest element in  $A$  is 1, and the rank of the median is  $n/2$ . Given an unsorted list  $A$  and indices  $i \leq j$ , give an  $O(n)$  time algorithm to output all elements in  $A$  with ranks lying in the interval  $[i, j]$ .