# CS1200: Intro. to Algorithms and their Limitations Lecture 1: Sorting Harvard SEAS - Fall 2025 2025-09-02

#### 1 Announcements

- Watch course overview video (from last year's offering) and read syllabus carefully if you haven't already done so.
- Staff introductions.
- Problem set 0 is available; due next Wednesday. It is required and substantial, so start early!
- Salil's OH this week: TuTh 11:15-12 SEC 3.327; Wed 4-5:30 Zoom (link on Salil's website).
- TF OH and Sections starting; see Ed.
- Ask live questions in class! If you've missed an opportunity to, the TFs will also monitor Ed for questions during class.
- Adam Hesterberg and I have been expanding the course lecture notes into a textbook. The first part of the textbook is posted on Perusall, which is the best place to ask questions about the text or point out typos or confusing text in the book. The lecture notes (like these) are skeleton versions of sections of the textbook, emphasizing the most important points and leaving space for you to take notes. (Note that the numbering in the lecture notes does not match the textbook numbering.) We post later parts of the textbook in Perusall as they are completed. (We may not complete the entire textbook by the end of the semester, in which case we'll revert to making detailed versions of the lecture notes available.)

# 2 Recommended Reading

Throughout the course, the readings other than Hesterberg–Vadhan are optional, meant as supplements if you find them useful for a different presentation of the content or more details/examples/exercises.

- Hesterberg-Vadhan, Chapter 1 (Sorting).
- CS50 Week 3: https://cs50.harvard.edu/x/2022/weeks/3/
- Cormen–Leiserson–Rivest–Stein Chapter 2
- Roughgarden I, Sections 1.4 and 1.5
- We've ordered all of the course texts for purchase at the Coop, for reserve through Harvard libraries (should be available to read as e-books through HOLLIS), and physical copies that can be read in the SEC 2nd floor reading room.

### 3 Motivation

Unit 1: storing and searching data. Vast applicability, and clean setting to develop skills:

- How to mathematically abstract the computational problems that we want to solve with algorithms.
- How to prove that an algorithm correctly solves a given computational problem.
- How to analyze and compare the efficiency of different algorithms.
- How to formalize exactly what algorithms are and what they can and cannot do through precise computational models.

As a concrete motivation for data management algorithms, let us consider the problem of web search, which is one of the most remarkable achievements of algorithms at a massive scale. The World Wide Web consists of billions of web pages and yet search engines are able to answer queries in a fraction of a second. How is this possible?

Let's consider a simplified description of Google's original search algorithm from 1998 (which has evolved substantially since then):

- 1. PAGERANK CALCULATION: For every URL on the entire World Wide Web ( $\forall url \in WWW$ ), calculate its PageRank,  $PR(url) \in [0, 1]$ .
- 2. KEYWORD SEARCH: Given a search keyword w, let  $S_w$  be the set of all webpages containing w. That is, compute  $S_w = \{url \in WWW : w \text{ is contained on the webpage at } url\}$ .
- 3. SORTING: Return the list of URLs in  $S_w$ , sorted in decreasing order of their pagerank.

The definition and calculation of PageRanks (Step 1) was the biggest innovation in Google's search, and is the most computationally intensive of these steps. However, it can be done offline, with periodic updates, rather than needing to be done in real-time response to an individual search query. Pageranks are outside the scope of this textbook, but you can learn about them in more advanced algorithms texts. Keyword Search (Step 2) can be done by creating a *trie* data structure for each webpage, also offline. Tries are covered in introductory programming and data structure texts. Our focus here is Sorting (Step 3), which needs to be extremely fast in response to real-time queries, and operates on a massive scale (e.g. millions of pages).

## 4 The Sorting Problem

Representing data items as key-value pairs (K, V), we can define the sorting problem as follows:

**Input:** An array A of key-value pairs  $((K_0, V_0), \ldots, (K_{n-1}, V_{n-1}))$ , where each key  $K_i \in \mathbb{R}^a$ 

**Output:** An array A' of key-value pairs  $((K'_0, V'_0), \ldots, (K'_{n-1}, V'_{n-1}))$  that is a valid sorting of A. That is, A' should be:

- 1. sorted by key values, i.e.  $K_0' \leq K_1' \leq \cdots \leq K_{n-1}'$ . and
- 2. a reordering of A, i.e.  $\exists$  a permutation  $\pi:[n]\to[n]$  such that  $(K_i',V_i')=(K_{\pi(i)},V_{\pi(i)})$  for  $i=0,\ldots,n-1$ .

#### Computational Problem SORTING

Above and throughout this text, [n] denotes the set of numbers  $\{0, \ldots, n-1\}$ . In pure math (e.g. combinatorics), it is more standard for [n] to be the set  $\{1, \ldots, n\}$ , but being computer scientists, we like to index starting at 0. Similarly, for us, the natural numbers are  $\mathbb{N} = \{0, 1, 2, \ldots\}$ .

To apply the abstract definition of SORTING to the web search problem, we can set:

- $K_i =$
- $\bullet$   $V_i =$

Abstraction  $\rightarrow$  many other applications! Database systems (both Relational and NoSQL), machine learning systems, ranking cat photos by cuteness, ...

Note that some inputs to SORTING have multiple valid outputs. Specifically, if there are repetitions among the input keys (e.g. if the key is a person's age in years in a dataset of 1000 people), then there are multiple reorderings that all yield valid sortings.

In the subsequent =sections, we will see three different sorting algorithms, and compare those algorithms to each other.

## 5 Exhaustive-Search Sort

We begin with a very simple sorting algorithm, which we obtain almost directly from the definition of the Sorting computational problem.

```
ExhaustiveSearchSort(A):

Input
: An array A = ((K_0, V_0), \dots, (K_{n-1}, V_{n-1})), where each K_i \in \mathbb{R}

Output
: A valid sorting of A

o foreach permutation \pi : [n] \to [n] do

1 | if K_{\pi(0)} \le K_{\pi(1)} \le \dots \le K_{\pi(n-1)} then

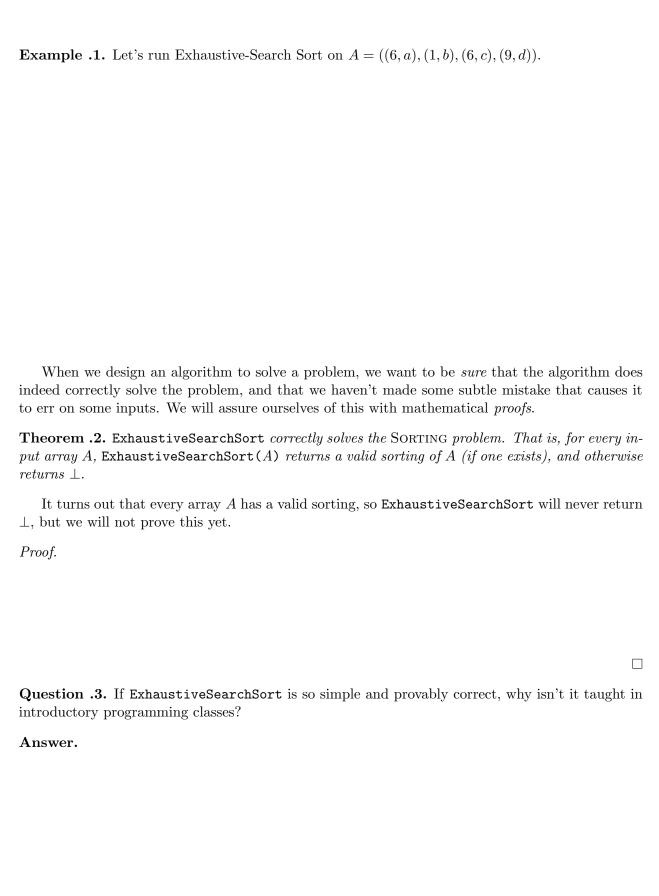
2 | return A' = ((K_{\pi(0)}, V_{\pi(0)}), (K_{\pi(1)}, V_{\pi(1)}), \dots, (K_{\pi(n-1)}, V_{\pi(n-1)}))

3 return \bot
```

Algorithm .1: ExhaustiveSearchSort()

The "bottom" symbol  $\perp$  is one that we will often use to denote failure to find a valid output.

<sup>&</sup>lt;sup>a</sup>When we introduce computational models later in the course, we'll stop working with arbitrary real numbers, and then keys could be restricted to integers or rational numbers.



### 6 Insertion Sort

Now let's see a much more efficient, but still rather simple, sorting algorithm:

```
InsertionSort(A):

Input : An array A = ((K_0, V_0), \dots, (K_{n-1}, V_{n-1})), where each K_i \in \mathbb{R}

Output : A valid sorting of A

0 /* "in-place" sorting algorithm that modifies A until it is sorted */

1 foreach i = 0, \dots, n-1 do

2 | /* Insert A[i] into the correct place among (A[0], \dots, A[i-1]). */

3 | Find the first index j such that A[i][0] \leq A[j][0];

4 | Insert A[i] into position j and shift A[j \dots i-1] to positions j+1, \dots, i

5 return A
```

#### Algorithm .2: InsertionSort()

Above we are using the notation  $A[k \dots \ell]$  as shorthand for the subarray  $(A[k], \dots, A[\ell])$ , for sake of readability.

## **Example .4.** A = ((6, a), (2, b), (1, c), (4, d)).

As our algorithm runs, we produce the following sorted sub-arrays:

Iteration $i$	Array contents after iteration $i$
0	
1	
2	
3	

Above, we only informally described Lines 3 and 4. They are expanded into a more precise implementation in the textbook.

**Theorem .5.** InsertionSort correctly solves the SORTING problem. That is, for every input array A, InsertionSort(A) returns a valid sorting of A.

Note that, as a side-effect, here we will also prove that every array A has a valid sorting. This phenomenon, of proving that a mathematical object exists (e.g. a valid sorting of an arbitrary array A) by exhibiting and analyzing an algorithm to construct it, is a quite common and useful one. See Problem Set 0 for another example.

#### *Proof.* Notation:

<u>Proof strategy:</u> We'll prove by induction on i (from i = 0, ..., n) the following "loop invariant" P(i):

$$P(i) =$$

Notice that the statement P(n) says that  $A^{(n)}$ , which is the output of InsertionSort(A), is a valid sorting of A, as desired.

```
Base Case (P(0)):
Inductive Step (P(i) \Rightarrow P(i+1)):
```

Remarks.

- Not all correctness proofs for algorithms use induction! (cf. Exhaustive-Search Sort.)
- This was not a fully formal proof. It is often necessary to skip steps to make such proofs manageable for humans, but you should be careful when do so. Be sure that (a) you are completely convinced of the correctness of your claims, and (b) you are not omitting the main point or idea of the argument.

# 7 Merge Sort

Finally, you may have already seen the following even more efficient sorting algorithm, MergeSort. The idea is to recursively sort each half of the array and then efficiently "merge" the two sorted halves into a single, sorted array:

We may not cover this in lecture (depending on time), since most of you have already seen it in CS

50. But you should review it from the textbook on your own!

```
MergeSort(A):
Input
: An array A = ((K_0, V_0), \dots, (K_{n-1}, V_{n-1})), where each K_i \in \mathbb{R}
Output
: A valid sorting of A
o if n \leq 1 then return A;
lelse
\begin{vmatrix} i & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & | \\ 0 & | & | & |
```

Algorithm .3: MergeSort()

We omit the implementation of Merge, which you can find in the textbook.

```
Example .6. A = ((7, a), (4.b), (6, c), (9, d), (7, e), (1, f), (2, g), (4, h)).
```

Now, let's confirm that MergeSort works in general:

**Theorem .7.** MergeSort correctly solves the SORTING problem. That is, for every input array A, MergeSort(A) returns a valid sorting of A.

Naturally, the proof of this theorem will rely on the correctness of Merge, whose proof we leave as an exercise:

**Lemma .8.** If B and B' are sorted arrays, then Merge(B, B') is a valid sorting of  $B \circ B'$ .

*Proof Sketch of Theorem .7.* Like with InsertionSort, this is a proof by induction, but we use *strong* induction. (Why?)

Here, the statement we will prove by (strong) induction is simpler than for InsertionSort. It is simply

P(n) = "MergeSort correctly sorts arrays of size n."