

# **Efficient Sorting**

4 November 2024

Prof. Dr. Sebastian Wild

CS566 (Wintersemester 2024/25) Philipps-Universität Marburg version 2024-11-04 18:30 H

#### **Learning Outcomes**

#### Unit 4: Efficient Sorting

- 1. Know principles and implementation of *mergesort* and *quicksort*.
- 2. Know properties and *performance characteristics* of mergesort and quicksort.
- 3. Know the comparison model and understand the corresponding *lower bound*.
- 4. Understand *counting sort* and how it circumvents the comparison lower bound.
- 5. Know ways how to exploit *presorted* inputs.

#### Outline

# **4** Efficient Sorting

- 4.1 Mergesort
- 4.2 Quicksort
- 4.3 Comparison-Based Lower Bound
- 4.4 Integer Sorting
- 4.5 Adaptive Sorting
- 4.6 Python's list sort

## Why study sorting?

- fundamental problem of computer science that is still not solved
- building brick of many more advanced algorithms
  - for preprocessing
  - as subroutine
- playground of manageable complexity to practice algorithmic techniques

#### Here:

- "classic" fast sorting method
- exploit partially sorted inputs
- parallel sorting

- Algorithm with optimal #comparisons in worst case?

# **Part I** *The Basics*

#### Rules of the game

#### Given:

• array A[0..n) = A[0..n-1] of *n* objects

a total order relation ≤ among A[0],..., A[n - 1] (a comparison function) Python: elements support <= operator (\_\_le\_\_()) Java: Comparable class (x.compareTo(y) <= 0)</p>

► **Goal:** rearrange (i. e., permute) elements within A, so that A is *sorted*, i. e.,  $A[0] \le A[1] \le \cdots \le A[n-1]$ 

for now: A stored in main memory (internal sorting) single processor (sequential sorting)

# 4.1 Mergesort

## Merging sorted lists



## Merging sorted lists





run1

run2

result

#### Mergesort

procedure mergesort(A[l..r))

 $_{2}$  n := r - l

 $_3$  if  $n \le 1$  return

$$_{4} \quad m := l + \left| \frac{n}{2} \right|$$

6 mergesort(A[m..r))

<sup>7</sup> merge(
$$A[l..m), A[m..r), buf$$
)

s copy buf to A[l..r)

- recursive procedure
- merging needs
  - temporary storage *buf* for result (of same size as merged runs)
  - to read and write each element twice (once for merging, once for copying back)

#### **Linear Term of** *C*(*n*)

Recall:  $C(n) = 2n \lg(n) + (2 - \{\lg(n)\} - 2^{1 - \{\lg(n)\}}) 2n$ 







 $\rightsquigarrow$  Can prove:  $C(n) \leq 2n \lg n + 0.172n$ 

#### **Mergesort – Discussion**

 $\square$  optimal time complexity of  $\Theta(n \log n)$  in the worst case

*stable* sorting method i.e., retains relative order of equal-key items

memory access is sequential (scans over arrays)

requires  $\Theta(n)$  extra space there are in-place merging methods, but they are substantially more complicated and not (widely) used

# 4.2 Quicksort

#### Partitioning around a pivot



#### Partitioning around a pivot



#### Partitioning – Detailed code

Beware: details easy to get wrong; use this code!

(if you ever have to)

<sup>1</sup> **procedure** partition(*A*, *b*) // input: array A[0..n), position of pivot  $b \in [0..n)$ 2 swap(A[0], A[b])3  $i := 0, \quad j := n$ 4 while true do 5 **do** i := i + 1 while i < n and A[i] < A[0]6 **do** j := j - 1 while  $j \ge 1$  and A[j] > A[0]7 if  $i \ge j$  then break (goto 11) 8 else swap(A[i], A[j])9 end while 10 swap(A[0], A[j])11 return *j* 12



#### Quicksort

procedure quicksort(A[l..r))

- <sup>2</sup> if  $r \ell \leq 1$  then return
- b := choosePivot(A[l..r))
- $_{4}$  j := partition(A[l..r), b)
- 5 quicksort(A[l..j))
- <sup>6</sup> quicksort(A[j+1..r))

- recursive procedure
- choice of pivot can be
  - ► fixed position ~→ dangerous!
  - random
  - more sophisticated, e.g., median of 3

#### **Quicksort & Binary Search Trees**



recursion tree of quicksort = binary search tree from successive insertion

- comparisons in quicksort = comparisons to built BST
- comparisons in quicksort  $\approx$  comparisons to search each element in BST

#### **Quicksort – Worst Case**

- Problem: BSTs can degenerate
- Cost to search for k is k 1

· → Total cost 
$$\sum_{k=1}^{n} (k-1) = \frac{n(n-1)}{2} \sim \frac{1}{2}n^2$$

 $\rightsquigarrow$  quicksort worst-case running time is in  $\Theta(n^2)$ 

terribly slow!

But, we can fix this:

#### **Randomized quicksort:**

- choose a *random pivot* in each step
- → same as randomly *shuffling* input before sorting

#### **Randomized Quicksort – Analysis**

- cost measure: element visits (as for mergesort)
- C(n) = #element visits when sorting n randomly permuted elements
  = cost of searching every element in BST build from input
- → quicksort needs ~  $2 \ln(2) \cdot n \lg n \approx 1.39n \lg n$  in expectation (see analysis of  $C_n$  in Unit 3!)
- also: very unlikely to be much worse: e.g., one can prove: Pr[cost > 10n lg n] = O(n<sup>-2.5</sup>) distribution of costs is "concentrated around mean"
- ▶ intuition: have to be *constantly* unlucky with pivot choice

#### **Quicksort – Discussion**

fastest general-purpose method

- $\Theta(n \log n)$  average case
- works *in-place* (no extra space required)

memory access is sequential (scans over arrays)

- $\bigcirc \Theta(n^2)$  worst case (although extremely unlikely)
- not a *stable* sorting method

Open problem: Simple algorithm that is fast, stable and in-place.

# 4.3 Comparison-Based Lower Bound

#### **Lower Bounds**

- Lower bound: mathematical proof that *no algorithm* can do better.
  - very powerful concept: bulletproof *impossibility* result
    - $\approx$  conservation of energy in physics
  - (unique?) feature of computer science: for many problems, solutions are known that (asymptotically) achieve the lower bound
  - $\rightsquigarrow~$  can speak of "optimal algorithms"

▶ To prove a statement about *all algorithms*, we must precisely define what that is!

- already know one option: the word-RAM model
- Here: use a simpler, more restricted model.

#### The Comparison Model

- ▶ In the *comparison model* data can only be accessed in two ways:
  - comparing two elements
  - moving elements around (e. g. copying, swapping)
  - Cost: number of comparisons.

That's good! /Keeps algorithms general!

- This makes very few assumptions on the kind of objects we are sorting.
- Mergesort and Quicksort work in the comparison model.

→ Every comparison-based sorting algorithm corresponds to a *decision tree*.

- only model comparisons ~>> ignore data movement
- nodes = comparisons the algorithm does
- child links = outcomes of comparison
- leaf = unique initial input permutation compatible with comparison outcomes
- next comparisons can depend on outcomes  $\rightsquigarrow$  child subtrees can look different

#### **Comparison Lower Bound**

**Example:** Comparison tree for a sorting method for *A*[0..2]:



- Execution = follow a path in comparison tree.
- → height of comparison tree = worst-case # comparisons
- comparison trees are *binary* trees

 $\rightsquigarrow \ell \text{ leaves } \rightsquigarrow \text{ height} \geq \lceil \lg(\ell) \rceil$ 

► comparison trees for sorting method must have ≥ n! leaves

- Mergesort achieves ~  $n \lg n$  comparisons  $\rightarrow asymptotically comparison-optimal!$
- ▶ Open (theory) problem: Sorting algorithm with  $n \lg n \lg(e)n + o(n)$  comparisons?

# 4.4 Integer Sorting

#### How to beat a lower bound

• Does the above lower bound mean, sorting always takes time  $\Omega(n \log n)$ ?

▶ Not necessarily; only in the *comparison model*!

 $\rightsquigarrow$  Lower bounds show where to *change* the model!

#### ► Here: sort *n* integers

- ▶ can do *a lot* with integers: add them up, compute averages, ... (full power of word-RAM)
- $\rightsquigarrow~$  we are **not** working in the comparison model
- $\rightsquigarrow$  above lower bound does not apply!
- but: a priori unclear how much arithmetic helps for sorting ...

#### **Counting sort**

- Important parameter: size/range of numbers
  - numbers in range  $[0..U] = \{0, ..., U-1\}$  typically  $U = 2^b \rightsquigarrow b$ -bit binary numbers
- We can sort *n* integers in  $\Theta(n + U)$  time and  $\Theta(U)$  space when  $b \le w$

#### **Counting sort**

- 1 **procedure** countingSort(A[0..n))
- // A contains integers in range [0..U]. 2
- C[0..U] := new integer array, initialized to 0 3
- // Count occurrences 4

5 **for** 
$$i := 0, ..., n - 1$$

- C[A[i]] := C[A[i]] + 16
- i := 0 // Produce sorted list7
- **for** *k* := 8

**for** 
$$j := 1, ..., C[k]$$

9 10

$$0, \dots U - 1 j := 1, \dots C[k] A[i] := k; i := i +$$

count how often each possible value occurs

word size

- produce sorted result directly from counts
- circumvents lower bound by using integers as array index / pointer offset

Can sort *n* integers in range [0..U) with U = O(n) in time and space  $\Theta(n)$ .  $\sim$ 

#### Larger Universes: Radix Sort

#### MSD Radix Sort:

- split numbers into base-R "digits"
- Use counting sort on <u>most significant digit</u> (with variant of counting sort that moves full number)
- $\rightsquigarrow \$  integers sorted with respect to first digit
- recurse on sublist for each digit value, using next digit for counting sort
- $\rightsquigarrow$  After  $\lfloor \log_R(U) \rfloor + 1$  levels of counting sort, fully sorted!
  - For R ≤ 2<sup>w</sup>, all counting sort calls on same level cost total of O(n) time (requires care to avoid reinitialization cost of array C)

 $\rightsquigarrow$  total time  $O(n \log_R(U)) = O\left(n \frac{\log(U)}{\log(R)}\right)$ 

 $\rightarrow O(n)$  time sorting possible for numbers in range  $U = O(n^c)$  for constant *c*.

#### Integer Sorting – State of the art

#### Algorithm theory

- ▶ integer sorting on the *w*-bit word-RAM
- suppose  $U = 2^w$ , but w can be an arbitrary function of n
- ▶ how fast can we sort *n* such *w*-bit integers on a *w*-bit word-RAM?
  - for  $w = O(\log n)$ : linear time (*radix/counting sort*)
  - for  $w = \Omega(\log^{2+\varepsilon} n)$ : linear time (signature sort)
  - ► for *w* in between: can do  $O(n\sqrt{\lg \lg n})$  (very complicated algorithm) don't know if that is best possible!

\* \* \*

... for the rest of this unit: back to the comparisons model!

# **Part II** Exploiting presortedness

# 4.5 Adaptive Sorting

## **Adaptive sorting**

- Comparison lower bound also holds for the *average case*  $\rightarrow \lfloor \lg(n!) \rfloor$  cmps necessary
- Mergesort and Quicksort from above use  $\sim n \lg n$  cmps even in best case



Can we do better if the input is already "almost sorted"?

Scenarios where this may arise naturally:

- Append new data as it arrives, regularly sort entire list (e.g., log files, database tables)
- Compute summary statistics of time series of measurements that change slowly over time (e.g., weather data)
- Merging locally sorted data from different servers (e.g., map-reduce frameworks)
- A Ideally, algorithms should *adapt* to input: *the more sorted the input, the faster the algorithm* ... but how to do that!?

#### Warmup: check for sorted inputs

- Any method could first check if input already completely in order! Best case becomes  $\Theta(n)$  with n - 1 comparisons!
  - Usually n 1 extra comparisons and pass over data "wasted"
  - Only catches a single, extremely special case . . .
- For divide & conquer algorithms, could check in each recursive call!
  Potentially exploits partial sortedness!
  usually adds Ω(n log n) extra comparisons

-<u>\</u>

For Mergesort, can instead check before merge with a **single** comparison

➤ If last element of first run ≤ first element of second run, skip merge

How effective is this idea?

procedure mergesortCheck(A[l..r))

$$n := r - l$$

7

 $if n \le 1$  return

$$m := l + \left\lfloor \frac{n}{2} \right\rfloor$$

$$if A[m-1] > A[m]$$

8 merge(A[l..m), A[m..r), buf)

copy buf to 
$$A[l..r)$$

#### Mergesort with sorted check – Analysis

- ► Simplified cost measure: *merge cost* = size of output of merges
  - $\approx$  number of comparisons
  - $\approx~$  number of memory transfers / cache misses
- Example input: n = 64 numbers in sorted *runs* of 16 numbers each:



Sorted check can help a lot!

#### **Alignment issues**

▶ In previous example, each run of length l saved us  $l \lg(l)$  in merge cost.

- = exactly the cost of *creating* this run in mergesort had it not already existed
- $\rightsquigarrow\,$  best savings we can hope for!
- $\rightsquigarrow$  Are overall merge costs  $\mathcal{H}(\ell_1, \ldots, \ell_r) := n \lg(n)$

$$\underbrace{\binom{n}{\text{sort}}}_{\text{sort}} - \underbrace{\sum_{i=1}^{r} \ell_{i}^{\ell_{i}} \lg(\ell_{i})}_{\text{sort}}?$$

mergesort

savings from runs





#### **Bottom-Up Mergesort**

• Can we do better by explicitly detecting runs?

1	<b>procedure</b> bottomUpMergesort( <i>A</i> [0 <i>n</i> ))
2	<i>Q</i> := new Queue // runs to merge
3	// Phase 1: Enqueue singleton runs
4	<b>for</b> $i = 0,, n - 1$ <b>do</b>
5	Q.enqueue((i, i + 1))
6	// Phase 2: Merge runs level—wise
7	while $Q.size() \ge 2$
8	Q' := new Queue
9	while $Q.size() \ge 2$
10	$(i_1, j_1) := Q.dequeue()$
11	$(i_2, j_2) := Q.dequeue()$
12	$merge(A[i_1j_1), A[i_2j_2), buf)$
13	copy buf to $A[i_1j_2)$
14	$Q'$ .enqueue $((i_1, j_2))$
15	<b>if</b> $\neg Q$ . <i>isEmpty()</i> // <i>lonely run</i>
16	Q'.enqueue( $Q$ .dequeue())
17	Q := Q'

1	<b>procedure</b> naturalMergesort( <i>A</i> [0 <i>n</i> ))
2	Q := new Queue; $i := 0$ find run $A[ij)$
3	while $i < n$
4	j := i + 1
5	<b>while</b> $A[j] \ge A[j-1]$ <b>do</b> $j := j+1$
6	Q.enqueue((i, j)); i := j
7	while $Q.size() \ge 2$
8	Q' := new Queue
9	while $Q.size() \ge 2$
0	$(i_1, j_1) := Q.dequeue()$
1	$(i_2, j_2) := Q.dequeue()$
2	merge( $A[i_1j_1), A[i_2j_2), buf$ )
3	copy buf to $A[i_1j_2)$
4	$Q'$ .enqueue $((i_1, j_2))$
5	<b>if</b> $\neg Q$ .isEmpty() // lonely run
6	Q'.enqueue(Q.dequeue())
7	Q := Q'

#### Natural Bottom-Up Mergesort – Analysis

▶ Works well for runs of roughly equal size, regardless of alignment ...



#### Natural Bottom-Up Mergesort – Analysis [2]

... but less so for widely varying run lengths



... can't we have both at the same time?!

#### Good merge orders

Let's take a step back and breathe.

- Conceptually, there are two tasks:
  - **1.** Detect and use existing runs in the input  $\rightsquigarrow \ell_1, \ldots, \ell_r$
  - 2. Determine a favorable order of merges of runs

("automatic" in top-down mergesort)



#### **Nearly-Optimal Mergesort**

Nearly-Optimal Mergesorts: Fast, Practical Sorting Methods That Optimally Adapt to Existing Runs

J. Ian Munro University of Waterloo, Canada immuro@uwaterloo.ca https://orcid.org/0000-0002-7165-7988

Sebastian Wild University of Waterloo, Canada wild@umsterloo.ca https://orcid.org/0000-0002-6061-9177

#### - Abstract

We present two stable mergener variants, "predoct" and "powersort", that exploit existing runs and find analy-spin-lim merging ender with angighile correctant. Devices methods (wither require substatial effort for determining the merging ender (Takaoka 2009; Barbay & Naurar 2013) or do not have an equiral moved-case guarance (Peters 2002; Agaro, Nixual & F) Processar 2018; Bas & Kaop 2018). We demonstrate that our methods are competitive in terms of running time with state-of-those ring immunities entring methods.

2012 ACM Subject Classification Theory of computation  $\rightarrow$  Sorting and searching

Keywords and phrases adaptive sorting, nearly-optimal binary search trees, Timsort

Digital Object Identifier 10.4230/LIPIcs.ESA.2018.63

Related Version arXiv: 1805.04154 (extended version with appendices)

Supplement Material zenodo: 1241162 (code to reproduce running time study)

Funding This work was supported by the Natural Sciences and Engineering Research Council of Canada and the Canada Research Chairs Programme.

#### 1 Introduction

Sorting is a fundamental building block for numerous tasks and ubiquitous in both the theory and practice of computing. While practical and theoretically (class-to) optimal comparison-based sorting methods are known, instance-optimal sorting, i.e., andthold that adapt to the actual input and exploid specific structural properties if present, is still an areas of active research. We survey some event development in Section 1.1.

Many different structural properties have been investigated in theory. Two of them have also found wide adoption in practice, e.g., in *Grade's*. Jone numling Haryer adapting to the presence of toplicate keys and using existing sorted segments, called reas. The former is disclosed by a so-additional disclosed particular segments of the structure of the method, though, i.e., the relative order of demonstra with equal keys major disclosed in the method, though, i.e., the relative order of demonstra with equal keys major the discrepted in the presents. It is here used in Jam and shift for paintive-type arrays.

 C. J. Lu Manne and Schweiten Weitz Weit Manne Street (1997) (1  In 2018, with Ian Munro, I combined research on nearly-optimal BSTs with mergesort

→ 2 new algorithms: *Peeksort* and *Powersort* 

▶ both adapt provably optimal to existing runs even in worst case: mergecost ≤ H(ℓ<sub>1</sub>,...,ℓ<sub>r</sub>) + 2n

- both are lightweight extensions of existing methods with negligible overhead
- both fast in practice

#### Peeksort

- based on top-down mergesort
- "peek" at middle of array & find closest run boundary
- → split there and recurse (instead of at midpoint)



- can avoid scanning runs repeatedly:
  - find full run straddling midpoint
  - remember length of known runs at boundaries



 $\rightsquigarrow~$  with clever recursion, scan each run only once.

#### **Peeksort – Code**

<sup>1</sup> **procedure** peeksort( $A[\ell..r), \Delta_{\ell}, \Delta_{r}$ ) if  $r - \ell < 1$  then return 2 if  $\ell + \Delta_{\ell} == r \vee \ell == r + \Delta_r$  then return 3  $m := \ell + \lfloor (r - \ell)/2 \rfloor$ 4 5  $i := \begin{cases} \ell + \Delta_{\ell} & \text{if } \ell + \Delta_{\ell} \ge m \\ \text{extendRunLeft}(A, m) & \text{else} \end{cases}$ 6  $j := \begin{cases} r + \Delta_{r} \le m & \text{if } r + \Delta_{r} \le m \le m \\ \text{extendRunRight}(A, m) & \text{else} \end{cases}$ 7  $g := \begin{cases} i & \text{if } m - i < j - m \\ j & \text{else} \end{cases}$ 8  $\Delta_g := \begin{cases} j - i & \text{if } m - i < j - m \\ i - j & \text{else} \end{cases}$ peeksort( $A[\ell ... g), \Delta_{\ell}, \Delta_{\varphi}$ ) 9 peeksort( $A[g, r), \Delta_g, \Delta_r$ ) 10  $merge(A[\ell, g), A[g..r), buf)$ 11 copy buf to  $A[\ell..r)$ 12

Parameters:



• initial call: peeksort( $A[0..n), \Delta_0, \Delta_n$ ) with  $\Delta_0 = \text{extendRunRight}(A, 0)$  $\Delta_n = n - \text{extendRunLeft}(A, n)$ 

#### helper procedure

1 procedure extendRunRight(A[0..n), i) 2 j := i + 13 while  $j < n \land A[j-1] \le A[j]$ 4 j := j + 15 return j

(extendRunLeft similar)

#### **Peeksort – Analysis**

Consider tricky input from before again:



- One can prove: Mergecost always  $\leq \mathcal{H}(\ell_1, \ldots, \ell_r) + 2n$
- $\rightsquigarrow$  We can have the best of both worlds!

# 4.6 Python's list sort

#### Sorting in Python

- CPython
  - Python is only a specification of a programming language
  - The Python Foundation maintains *CPython* as the official reference implementation of the Python programming language
  - If you don't specifically install something else, python will be CPython
- part of Python are list.sort resp. sorted built-in functions
  - implemented in C
  - use *Timsort*, custom Mergesort variant by Tim Peters



Sept 2021: **Python uses** *Powersort*! since CPython 3.11 and PyPy 7.3.6

nsg400864 - Author: Tim Peters (tim.p view)	Date beters) * 🥏 2021-09-01 19:43
I created a PR that implements the powersort m	erge strategy:
https://github.com/python/cpython/pull/28108	
Across all the time this issue report has been to be the top contender. Enough already ;-) It change to make to the code, but that's in rela it's not at all a hard change.	open, that strategy continues 's indeed a more difficult tive terms. In absolute terms,
Laurent, if you find that some variant of Shiv than that, let us know here! I'm a big fan of ' powersort seems to do somewhat better "on aver adaptive ShiversSort (and implementing that to outside of merge_collapse().	ersSort actually runs faster Vincent's innovations too, but age" than even his length- o would require changing code

### **Timsort (original version)**



- above shows the core algorithm; many more algorithm engineering tricks
- Advantages:
  - profits from existing runs
  - locality of reference for merges
- But: not optimally adaptive! (next slide) Reason: Rules A–D (Why exactly these?!)



#### **Timsort bad case**

• On certain inputs, Timsort's merge rules don't work well:



- ▶ As *n* increases, Timsort's cost approach 1.5 · H, i. e., 50% more merge costs than necessary
  - intuitive problem: regularly very unbalanced merges

#### Powersort

--- Timsort's *merge rules* aren't great, but overall algorithm has appeal . . . can we keep that?



25 26 27

19 20

#### **Powersort – Run-Boundary Powers**



- (virtual) perfect balanced binary tree
- midpoint intervals "snap" to closest virtual tree node
  assigns each run boundary a depth = its *power*
- $\rightsquigarrow\,$  merge tree follows virtual tree



#### **Powersort – Run-Boundary Powers are Local**



Computation of powers only depends on two adjacent runs.

#### **Powersort – Computing powers**

- Computing the power of (run boundary between) two runs
  - Image: midpoint interval
  - power = min  $\ell$  s.t. contains  $c \cdot 2^{-\ell}$

<sup>1</sup> **procedure** power( $(i_1, j_1), (i_2, j_2), n$ )  $n_1 := i_1 - i_1$ 2  $n_2 := i_2 - i_2$ 3  $a := \frac{i_1 + \frac{1}{2}n_1 - 1}{2}$ 4  $b := \frac{i_2 + \frac{1}{2}n_2 - 1}{n_1 + \frac{1}{2}n_2 - 1}$  // interval (a, b] 5  $\ell := 0$ while  $|a \cdot 2^{\ell}| = = |b \cdot 2^{\ell}|$ 7  $\ell := \ell + 1$ 8 return ℓ 9

▶ with bitwise trickery *O*(1) time possible





#### **Powersort – Discussion**

Retains all advantages of Timsort

- good locality in memory accesses
- no recursion
- all the tricks in Timsort

optimally adapts to existing runs

minimal overhead for finding merge order