# *4* Efficient Sorting

*4 November 2024*

Prof. Dr. Sebastian Wild

# 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

# Why study sorting?

- ▶ fundamental problem of computer science that is still not solved

  <span style="color:darkred">Algorithm with optimal #comparisons in worst case?</span>

- ▶ 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

# Part I

*The Basics*

# Rules of the game

- **Given:**
  - array $A[0..n) = A[0..n-1]$ of $n$ objects
  - a total order relation $\leq$ among $A[0], \ldots, A[n-1]$ $\qquad \leq \quad \text{os} \quad < / = / >$
    (a comparison function)
    *Python:* elements support <= operator (\_\_le\_\_())
    *Java:* Comparable class (x.compareTo(y) <= 0)

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

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

## Clicker Question

running time of fastest solution

**?** What is the <u>complexity of sorting</u>? Type you answer, e. g., as "Theta(sqrt(n))"

(a) algorithm upper bound $O(n \log n)$

(b) lower bound $\Omega(n \log n)$

📱 → *sli.do/cs566*

## 4.1 Mergesort

# Clicker Question

How does mergesort work?

A) Split elements around median, then recurse on small / large elements.

B) Recurse on left / right half, then combine sorted halves.

C) Grow sorted part on left, repeatedly add next element to sorted range.

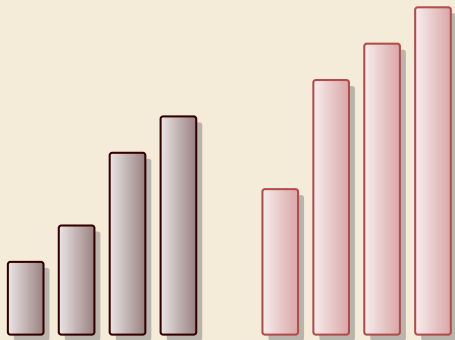D) Repeatedly choose 2 elements and swap them if they are out of order.

E) Don't know.

→ *sli.do/cs566*

# Clicker Question

How does mergesort work?

A. Split elements around median, then recurse on small / large elements.

B. Recurse on left / right half, then combine sorted halves. ✓

C. Grow sorted part on left, repeatedly add next element to sorted range.

D. Repeatedly choose 2 elements and swap them if they are out of order.

E. Don't know.

→ *sli.do/cs566*

# Merging sorted lists

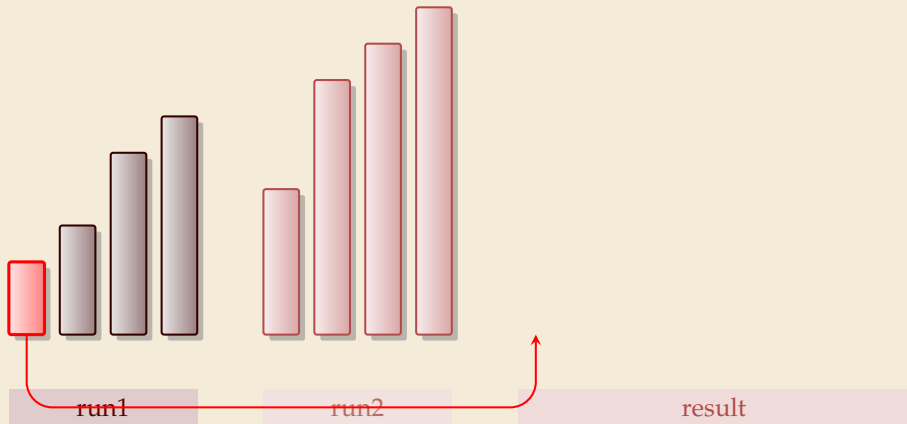# Merging sorted lists



run1          run2                    result

# Merging sorted lists

# Merging sorted lists



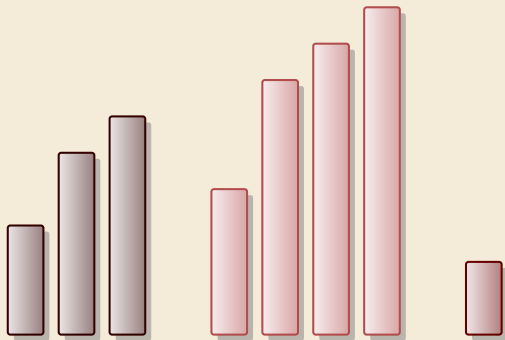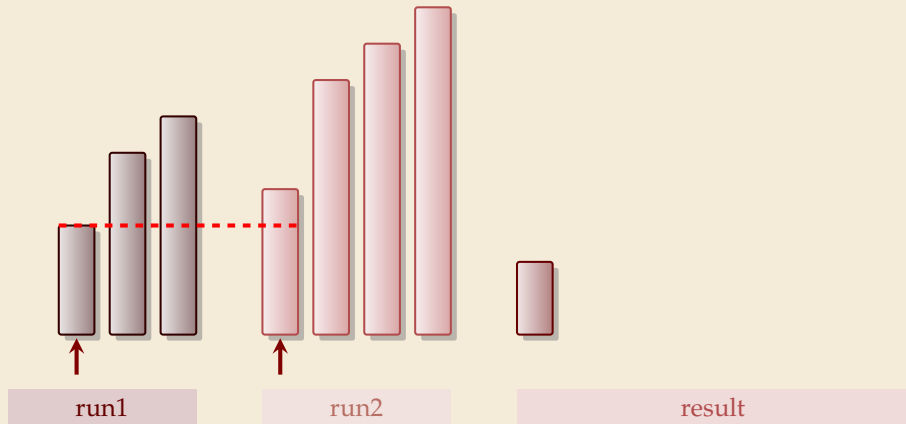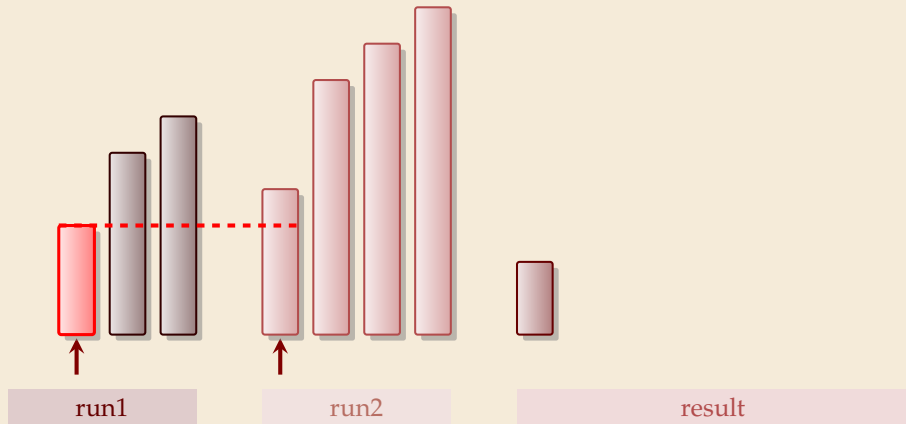run1                    run2                           result

# Merging sorted lists
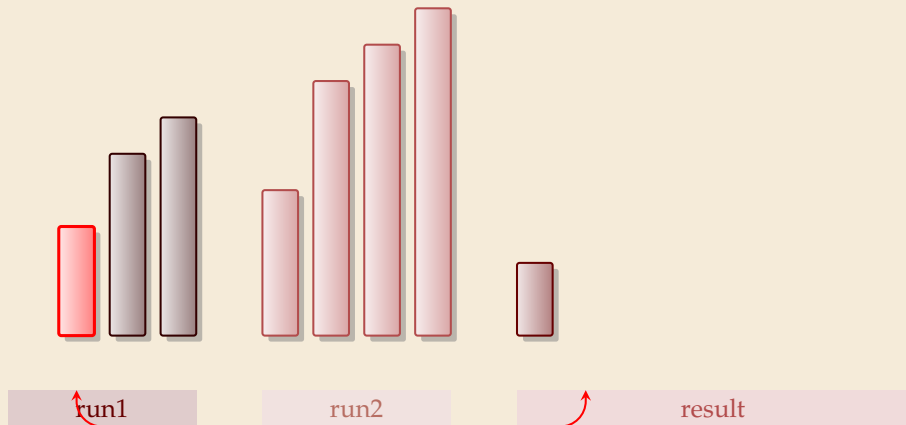
# Merging sorted lists



run1          run2                    result

# Merging sorted lists

# Merging sorted lists



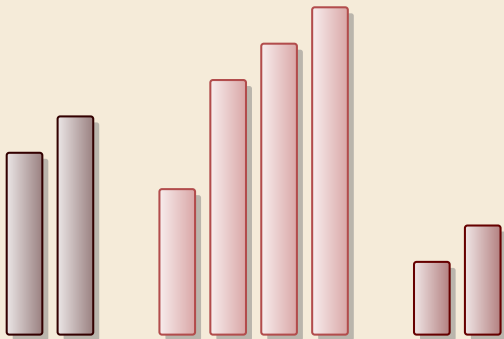run1          run2                    result

# Merging sorted lists

# Merging sorted lists



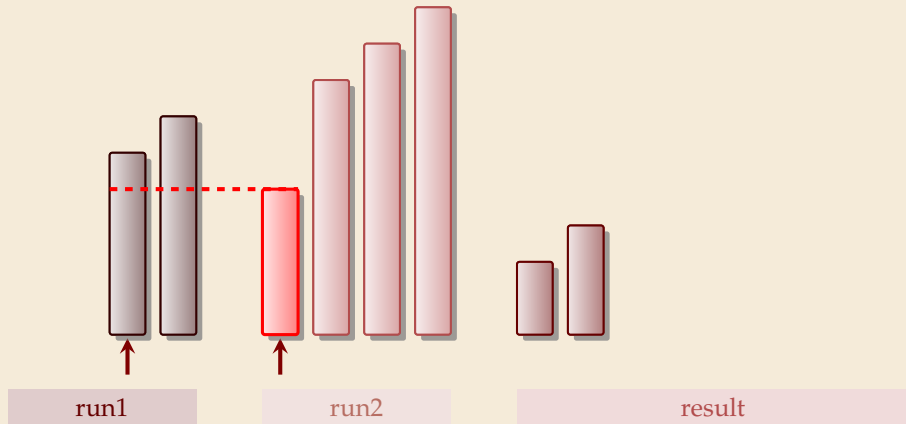run1          run2                    result

# Merging sorted lists

# Merging sorted lists

# Merging sorted lists



run1          run2                    result
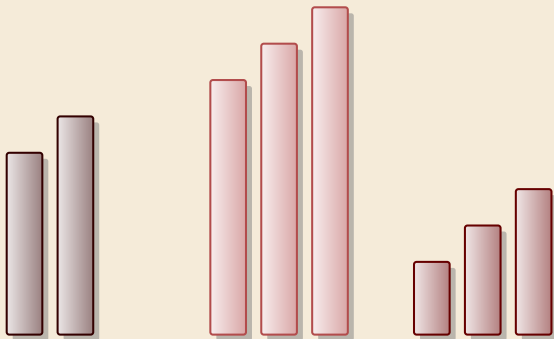
# Merging sorted lists

# Merging sorted lists

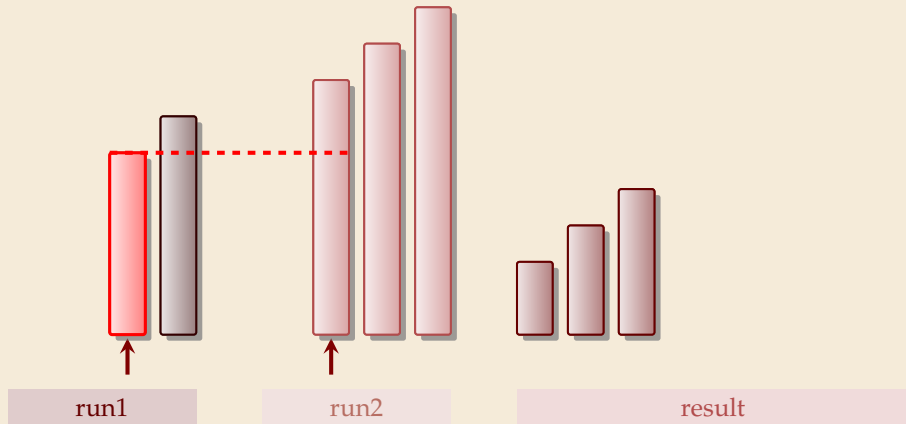# Merging sorted lists



run1          run2                              result

# Merging sorted lists

# Merging sorted lists

# Merging sorted lists



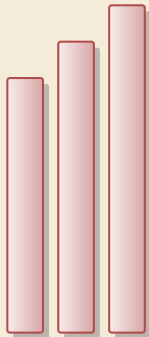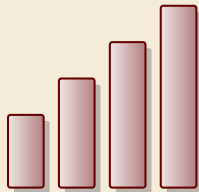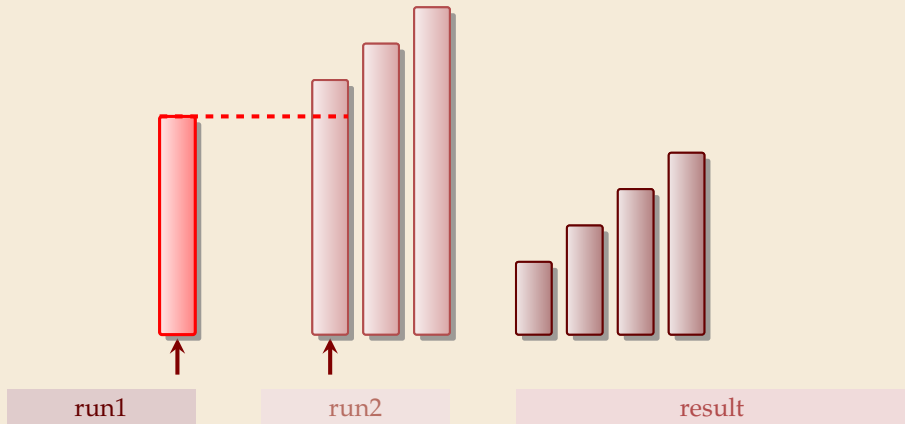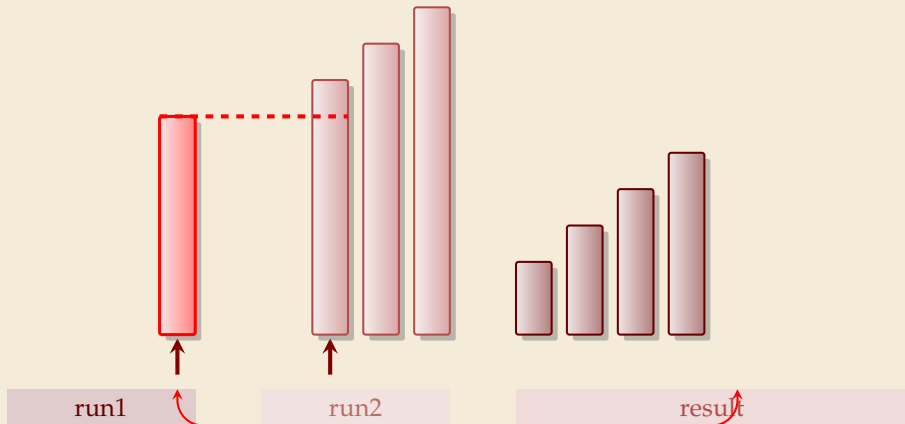run1          run2                    result

# Merging sorted lists



run1     run2     result

# Merging sorted lists



run1    run2    result

# Merging sorted lists



run1        run2        result

# Merging sorted lists



run1   run2   result

# Clicker Question

What is the worst-case running time of mergesort?

A   $\Theta(1)$        G   $\Theta(n \log n)$

B   $\Theta(\log n)$        H   $\Theta(n \log^2 n)$

C   $\Theta(\log \log n)$        I   $\Theta(n^{1+\epsilon})$

D   $\Theta(\sqrt{n})$        J   $\Theta(n^2)$

E   $\Theta(n)$        K   $\Theta(n^3)$

F   $\Theta(n \log \log n)$        L   $\Theta(2^n)$

→ *sli.do/cs566*

# Clicker Question

What is the worst-case running time of mergesort?

- A. ~~$\Theta(1)$~~
- B. ~~$\Theta(\log n)$~~
- C. ~~$\Theta(\log \log n)$~~
- D. ~~$\Theta(\sqrt{n})$~~
- E. ~~$\Theta(n)$~~
- F. ~~$\Theta(n \log \log n)$~~
- G. $\Theta(n \log n)$ ✓
- H. ~~$\Theta(n \log^2 n)$~~
- I. ~~$\Theta(n^{1+\epsilon})$~~
- J. ~~$\Theta(n^2)$~~
- K. ~~$\Theta(n^3)$~~
- L. ~~$\Theta(2^n)$~~

→ `sli.do/cs566`

## Mergesort

```
1  procedure mergesort(A[l..r))
2      n := r − l
3      if n ≤ 1 return
4      m := l + ⌊n/2⌋
5      mergesort(A[l..m))
6      mergesort(A[m..r))
7      merge(A[l..m), A[m..r), buf)
8      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)

## Mergesort

```
1  procedure mergesort(A[l..r])
2      n := r − l
3      if n ≤ 1 return
4      m := l + ⌊n/2⌋
5      mergesort(A[l..m])
6      mergesort(A[m..r])
7      merge(A[l..m], A[m..r], buf)
8      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)

**Analysis:** count "*element visits*" (read and/or write)

$$C(n) = \begin{cases} 0 & n \leq 1 \\ C(\lfloor n/2 \rfloor) + C(\lceil n/2 \rceil) + 2n & n \geq 2 \end{cases}$$

Simplification $\boxed{n = 2^k}$ $= k = \ell_{gn}$

same for best and worst case!

$$C(2^k) = \begin{cases} 0 & k \leq 0 \\ 2 \cdot C(2^{k-1}) + 2 \cdot 2^k & k \geq 1 \end{cases} = 2 \cdot 2^k + 2^2 \cdot 2^{k-1} + 2^3 \cdot 2^{k-2} + \cdots + 2^k \cdot 2^1 = 2k \cdot 2^k$$

$$C(n) = 2n \lg(n) = \Theta(n \log n) \quad \text{(arbitrary } n: \ C(n) \leq C(\text{next larger power of 2}) \leq 4n \lg(n) + 2n = \Theta(n \log n))$$

5

## Mergesort

```
1  procedure mergesort(A[l..r])
2      n := r − l
3      if n ≤ 1 return
4      m := l + ⌊n/2⌋
5      mergesort(A[l..m])
6      mergesort(A[m..r])
7      merge(A[l..m], A[m..r], buf)
8      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)

**Analysis:** count *"element visits"* (read and/or write)

$$C(n) = \begin{cases} 0 & n \leq 1 \\ C(\lfloor n/2 \rfloor) + C(\lceil n/2 \rceil) + 2n & n \geq 2 \end{cases}$$

same for best and worst case!

$$\left( \begin{array}{l} \text{precisely(!) solvable } without \text{ assumption } n = 2^k: \\ C(n) = 2n \lg(n) + (2 - \{\lg(n)\} - 2^{1 - \{\lg(n)\}})2n \\ \text{with } \{x\} := x - \lfloor x \rfloor \end{array} \right)$$

Simplification $\boxed{n = 2^k}$

$$C(2^k) = \begin{cases} 0 & k \leq 0 \\ 2 \cdot C(2^{k-1}) + 2 \cdot 2^k & k \geq 1 \end{cases} = 2 \cdot 2^k + 2^2 \cdot 2^{k-1} + 2^3 \cdot 2^{k-2} + \cdots + 2^k \cdot 2^1 = 2k \cdot 2^k$$

$$C(n) = 2n \lg(n) = \Theta(n \log n) \quad \text{(arbitrary } n: \ C(n) \leq C(\text{next larger power of 2}) \leq 4n \lg(n) + 2n = \Theta(n \log n))$$
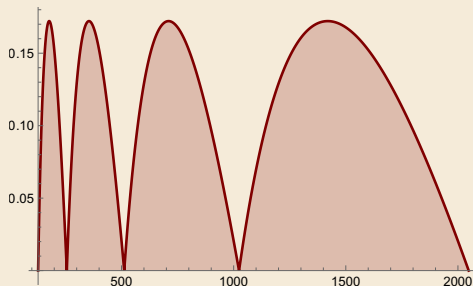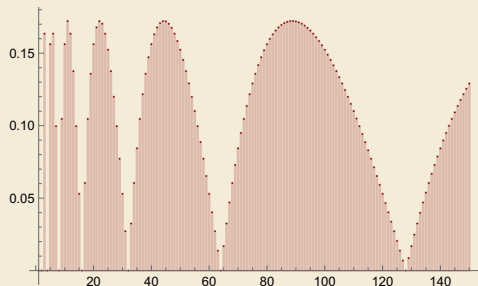
5

# Linear Term of $C(n)$

$n \geq 2$

Recall:
$$C(n) = 2n \lg(n) + \underbrace{(2 - \{\lg(n)\} - 2^{1-\{\lg(n)\}})2n}_{} \qquad \text{with } \{x\} := x - \lfloor x \rfloor$$

Plot of $2(2 - \{\lg(n)\} - 2^{1-\{\lg(n)\}})$



$\rightsquigarrow$ Can prove: $C(n) \leq \underset{\rightleftharpoons}{2n} \lg n + \underline{0.172n}$

6

# Mergesort – Discussion

👍 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

# Clicker Question

How does quicksort work?

A  split elements around median, then recurse on small / large elements.

B  recurse on left / right half, then combine sorted halves.

C  grow sorted part on left, repeatedly add next element to sorted range.

D  repeatedly choose 2 elements and swap them if they are out of order.

E  Don't know.

→ *sli.do/cs566*

# Clicker Question

How does quicksort work?

- A  split elements around median, then recurse on small / large elements. ✓
- B  recurse on left / right half, then combine sorted halves.
- C  grow sorted part on left, repeatedly add next element to sorted range.
- D  repeatedly choose 2 elements and swap them if they are out of order.
- E  Don't know.

→ *sli.do/cs566*

# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot
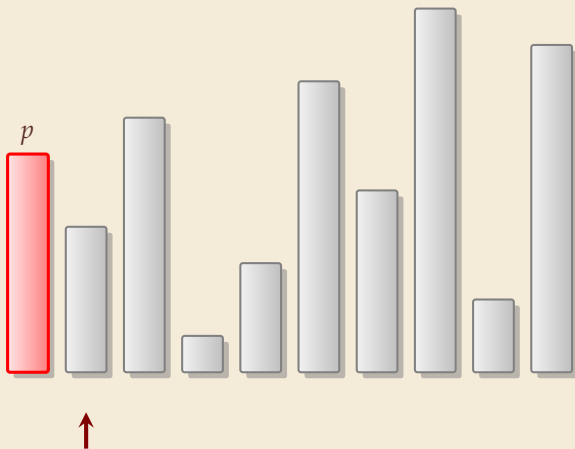
# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot
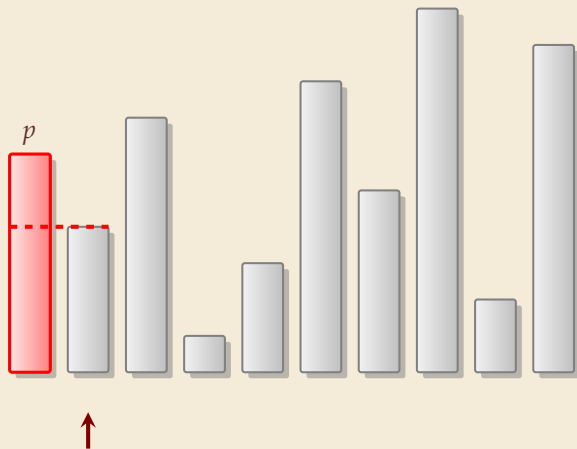
# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot
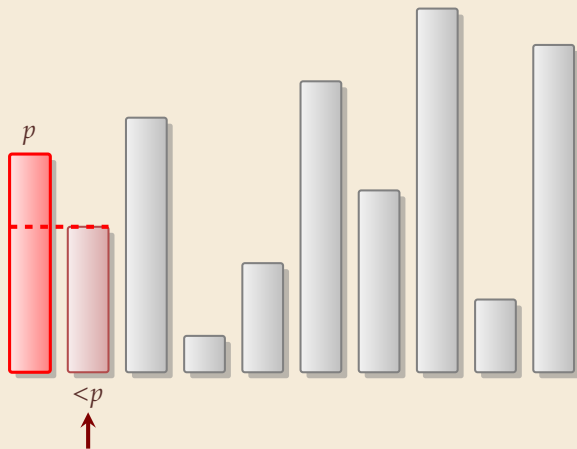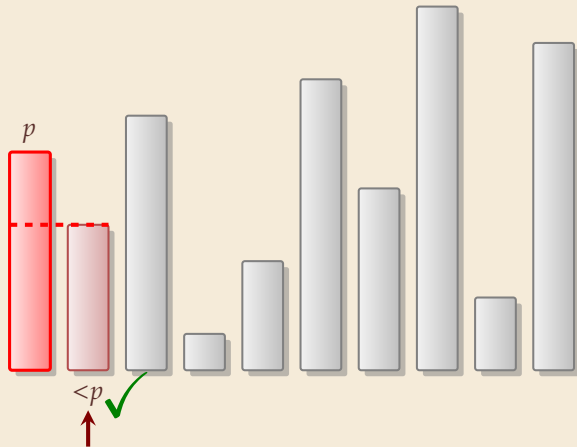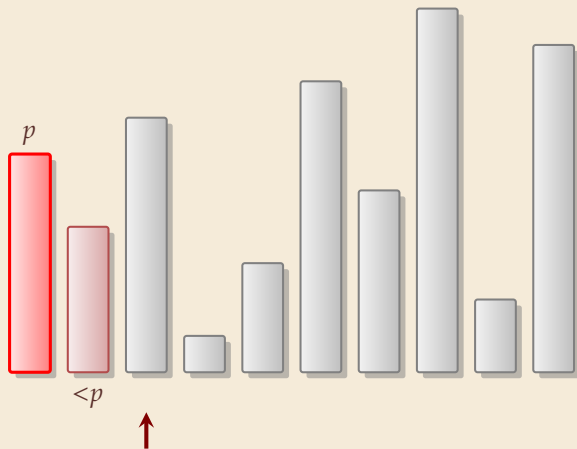
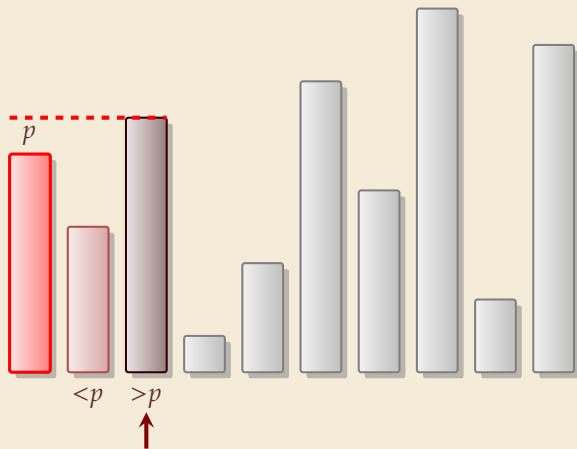# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot

# Partitioning around a pivot

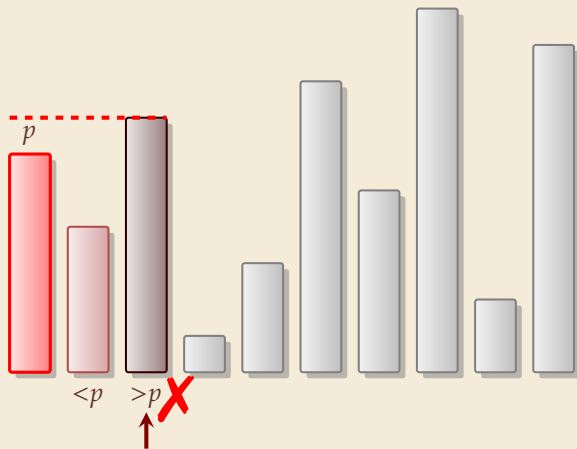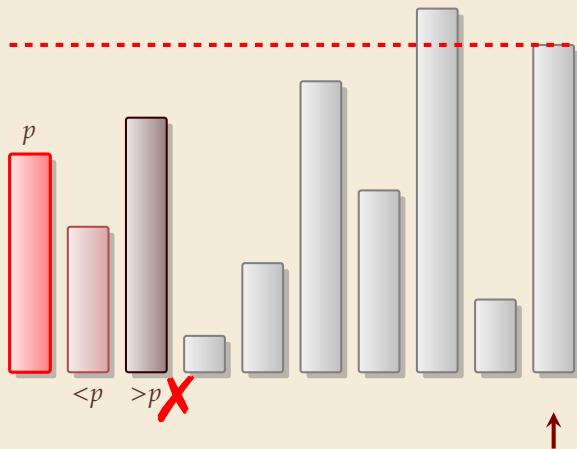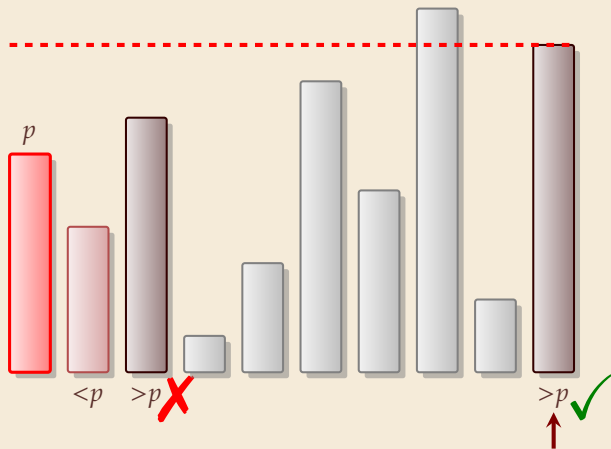# Partitioning around a pivot



- ▶ no extra space needed
- ▶ visits each element once
- ▶ returns rank/position of pivot

## Partitioning – Detailed code

Beware: details easy to get wrong; use this code!   (if you ever have to)

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

**Loop invariant (5–10):**   $A$

| $p$ | $\leq p$ | ? | $\geq p$ |
|-----|----------|---|----------|

## Quicksort

```
1  procedure quicksort(A[l..r))
2      if r − ℓ ≤ 1 then return
3      b := choosePivot(A[l..r))
4      j := partition(A[l..r), b)
5      quicksort(A[l..j))
6      quicksort(A[j + 1..r))
```

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

# Clicker Question

What is the <u>worst-case</u> running time of quicksort?

| A | $\Theta(1)$ | G | $\Theta(n \log n)$ |
|---|---|---|---|
| B | $\Theta(\log n)$ | H | $\Theta(n \log^2 n)$ |
| C | $\Theta(\log \log n)$ | I | $\Theta(n^{1+\epsilon})$ |
| D | $\Theta(\sqrt{n})$ | J | $\Theta(n^2)$ |
| E | $\Theta(n)$ | K | $\Theta(n^3)$ |
| F | $\Theta(n \log \log n)$ | L | $\Theta(2^n)$ |

→ *sli.do/cs566*

# Clicker Question

What is the worst-case running time of quicksort?

A. ~~$\Theta(1)$~~

B. ~~$\Theta(\log n)$~~

C. ~~$\Theta(\log \log n)$~~

D. ~~$\Theta(\sqrt{n})$~~

E. ~~$\Theta(n)$~~

F. ~~$\Theta(n \log \log n)$~~

G. ~~$\Theta(n \log n)$~~

H. ~~$\Theta(n \log^2 n)$~~

I. ~~$\Theta(n^{1+\epsilon})$~~

J. $\Theta(n^2)$ ✓

K. ~~$\Theta(n^3)$~~

L. ~~$\Theta(2^n)$~~

→ *sli.do/cs566*

# Quicksort & Binary Search Trees

**Quicksort**

| 7 | 4 | 2 | 9 | 1 | 3 | 8 | 5 | 6 |
|---|---|---|---|---|---|---|---|---|

# Quicksort & Binary Search Trees

**Quicksort**

| 7 | 4 | 2 | 9 | 1 | 3 | 8 | 5 | 6 |

# Quicksort & Binary Search Trees

**Quicksort**

| 7 | 4 | 2 | 9 | 1 | 3 | 8 | 5 | 6 |
|---|---|---|---|---|---|---|---|---|

| 4 | 2 | 1 | 3 | 5 | 6 |

| 9 | 8 |

# Quicksort & Binary Search Trees

**Quicksort**

| 7 | 4 | 2 | 9 | 1 | 3 | 8 | 5 | 6 |

| 4 | 2 | 1 | 3 | 5 | 6 | 7 | 9 | 8 |

# Quicksort & Binary Search Trees

**Quicksort**

| 7 | 4 | 2 | 9 | 1 | 3 | 8 | 5 | 6 |
|---|---|---|---|---|---|---|---|---|

| 4 | 2 | 1 | 3 | 5 | 6 | 7 | 9 | 8 |
|---|---|---|---|---|---|---|---|---|

# Quicksort & Binary Search Trees

**Quicksort**

| 7 | 4 | 2 | 9 | 1 | 3 | 8 | 5 | 6 |

| 4 | 2 | 1 | 3 | 5 | 6 | 7 | 9 | 8 |

| 2 | 1 | 3 | 4 | 5 | 6 |

# Quicksort & Binary Search Trees

**Quicksort**

| 7 | 4 | 2 | 9 | 1 | 3 | 8 | 5 | 6 |
|---|---|---|---|---|---|---|---|---|

| 4 | 2 | 1 | 3 | 5 | 6 | 7 | 9 | 8 |
|---|---|---|---|---|---|---|---|---|

| 2 | 1 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|

# Quicksort & Binary Search Trees

**Quicksort**

| 7 | 4 | 2 | 9 | 1 | 3 | 8 | 5 | 6 |

| 4 | 2 | 1 | 3 | 5 | 6 | 7 | 9 | 8 |

| 2 | 1 | 3 | 4 | 5 | 6 | 8 | 9 |

# Quicksort & Binary Search Trees

**Quicksort**

| 7 | 4 | 2 | 9 | 1 | 3 | 8 | 5 | 6 |

| 4 | 2 | 1 | 3 | 5 | 6 | 7 | 9 | 8 |

| 2 | 1 | 3 | 4 | 5 | 6 | 8 | 9 |

| 1 | 2 | 3 | 5 | 6 | 8 |

# Quicksort & Binary Search Trees

**Quicksort**

| 7 | 4 | 2 | 9 | 1 | 3 | 8 | 5 | 6 |

| 4 | 2 | 1 | 3 | 5 | 6 | 7 | 9 | 8 |

| 2 | 1 | 3 | 4 | 5 | 6 | 8 | 9 |

| 1 | 2 | 3 | 5 | 6 | 8 |

1    3    6

# Quicksort & Binary Search Trees

**Quicksort**

| 7 | 4 | 2 | 9 | 1 | 3 | 8 | 5 | 6 |

| 4 | 2 | 1 | 3 | 5 | 6 | 7 | 9 | 8 |

| 2 | 1 | 3 | 4 | 5 | 6 | 8 | 9 |

| 1 | 2 | 3 | 5 | 6 | 8 |

1    3    6

# Quicksort & Binary Search Trees

**Quicksort**

| 7 | 4 | 2 | 9 | 1 | 3 | 8 | 5 | 6 |

| 4 | 2 | 1 | 3 | 5 | 6 | 7 | 9 | 8 |

| 2 | 1 | 3 | 4 | 5 | 6 | 8 | 9 |

| 1 | 2 | 3 | 5 | 6 | 8 |

1   3   6

**Binary Search Tree (BST)**

7 4 2 9 1 3 8 5 6

# Quicksort & Binary Search Trees

**Quicksort**

| 7 | 4 | 2 | 9 | 1 | 3 | 8 | 5 | 6 |

| 4 | 2 | 1 | 3 | 5 | 6 | 7 | 9 | 8 |

| 2 | 1 | 3 | 4 | 5 | 6 | 8 | 9 |

| 1 | 2 | 3 | 5 | 6 | 8 |

1    3    6

**Binary Search Tree (BST)**

7 4 2 9 1 3 8 5 6

(7)

# Quicksort & Binary Search Trees

**Quicksort**

| 7 | 4 | 2 | 9 | 1 | 3 | 8 | 5 | 6 |

| 4 | 2 | 1 | 3 | 5 | 6 | 7 | 9 | 8 |

| 2 | 1 | 3 | 4 | 5 | 6 | 8 | 9 |

| 1 | 2 | 3 | 5 | 6 | 8 |

1    3    6

**Binary Search Tree (BST)**

4  2  9  1  3  8  5  6

7

4

# Quicksort & Binary Search Trees

**Quicksort**

| 7 | 4 | 2 | 9 | 1 | 3 | 8 | 5 | 6 |

| 4 | 2 | 1 | 3 | 5 | 6 | 7 | 9 | 8 |

| 2 | 1 | 3 | 4 | 5 | 6 | 8 | 9 |

| 1 | 2 | 3 | 5 | 6 | 8 |

1  3  6

**Binary Search Tree (BST)**

2  9  1  3  8  5  6

# Quicksort & Binary Search Trees

**Quicksort**

| 7 | 4 | 2 | 9 | 1 | 3 | 8 | 5 | 6 |

| 4 | 2 | 1 | 3 | 5 | 6 | 7 | 9 | 8 |

| 2 | 1 | 3 | 4 | 5 | 6 | 8 | 9 |

| 1 | 2 | 3 | 5 | 6 | 8 |

1   3   6

**Binary Search Tree (BST)**

9  1  3  8  5  6

# Quicksort & Binary Search Trees

**Quicksort**

| 7 | 4 | 2 | 9 | 1 | 3 | 8 | 5 | 6 |

| 4 | 2 | 1 | 3 | 5 | 6 | 7 | 9 | 8 |

| 2 | 1 | 3 | 4 | 5 | 6 | 8 | 9 |

| 1 | 2 | 3 | 5 | 6 | 8 |

1    3    6

**Binary Search Tree (BST)**

1  3  8  5  6

# Quicksort & Binary Search Trees

**Quicksort**



**Binary Search Tree (BST)**

# Quicksort & Binary Search Trees

**Quicksort**

| 7 | 4 | 2 | 9 | 1 | 3 | 8 | 5 | 6 |

| 4 | 2 | 1 | 3 | 5 | 6 | 7 | 9 | 8 |

| 2 | 1 | 3 | 4 | 5 | 6 | 8 | 9 |

| 1 | 2 | 3 | 5 | 6 | 8 |

1      3      6

**Binary Search Tree (BST)**

8  5  6

# Quicksort & Binary Search Trees

**Quicksort**



**Binary Search Tree (BST)**

# Quicksort & Binary Search Trees

**Quicksort**



**Binary Search Tree (BST)**

# Quicksort & Binary Search Trees

**Quicksort**

**Binary Search Tree (BST)**



- ▶ recursion tree of quicksort = binary search tree from successive insertion

- ▶ comparisons in quicksort = comparisons to built BST

- ▶ comparisons in quicksort ≈ comparisons to search each element in BST

## Quicksort – Worst Case

- Problem: BSTs can degenerate

- Cost to search for $k$ is $k - 1$

$\rightsquigarrow$ Total cost $\displaystyle\sum_{k=1}^{n}(k - 1) \; = \; \frac{n(n-1)}{2} \; \sim \; \frac{1}{2}n^2$

$\rightsquigarrow$ 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

$\rightsquigarrow$ same as randomly *shuffling* input before sorting

# Randomized Quicksort – Analysis

*mergesort* $2n\lg n$

▶ cost measure: <u>element visits</u> (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 $\sim 2\ln(2) \cdot n\lg n \approx \overset{1.386}{1.39}n\lg n$ *in expectation*
   (see analysis of $C_n$ in Unit 3!)

▶ also:  very unlikely to be much worse:
        e. g., one can prove:  $\Pr[\text{cost} > 10n\lg n] = O(n^{-2.5})$
        distribution of costs is "concentrated around mean"

▶ intuition:  have to be *constantly* unlucky with pivot choice

$\notin$ *exam*



13

# 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)

👎 $\Theta(n^2)$ 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
    ≈ *conservation of energy* in physics

  ▶ **(unique?) feature of computer science:**
    for many problems, solutions are known that (asymptotically) **achieve the lower bound**

  ⇝ can speak of "*optimal* algorithms"

## 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**
  - ↝ 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:
    - <u>comparing</u> two elements
    - moving elements around (e. g. copying, swapping)
    - Cost: number of comparisons.

expert comment

gold standard: cell probe 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.

- This makes very few assumptions on the kind of objects we are sorting.

- Mergesort and Quicksort work in the comparison model.

That's good!
Keeps algorithms general!

## 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  ↝  child subtrees can look different

# Comparison Lower Bound

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

# 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
- ⤳ $\ell$ leaves ⤳ height $\geq \lceil \lg(\ell) \rceil$
- ▶ comparison trees for sorting method must have $\geq n!$ leaves
- ⤳ height $\geq \lg(n!) \sim n \lg n$

more precisely: $\lg(n!) = n \lg n - \lg(e)n + O(\log n)$

Stirling

$$\lg(n!) \sim \lg\left(\sqrt{2\pi n}\left(\frac{n}{e}\right)^n\right)\left(1 + \Theta\left(\frac{1}{n}\right)\right)$$

$n \cdot \lg(n) - n \lg(e)$

$$= \lg(\sqrt{2\pi n}) + \lg\left(1 + \Theta\left(\frac{1}{n}\right)\right) + \lg\left(\left(\frac{n}{e}\right)^n\right)$$

## 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
- ⤳ $\ell$ leaves ⤳ height $\geq \lceil \lg(\ell) \rceil$
- ▶ comparison trees for sorting method must have $\geq n!$ leaves
- ⤳ height $\geq \lg(n!) \sim n \lg n$

more precisely: $\lg(n!) = n \lg n - \lg(e)n + O(\log n)$

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

$\approx 1.4427$

17

# Clicker Question

Does the comparison-tree from the previous slide correspond to a worst-case optimal sorting method?

A Yes

B No

→ *sli.do/cs566*

## Clicker Question

Does the comparison-tree from the previous slide correspond to a worst-case optimal sorting method?

A   Yes ✓

B   ~~No~~

→ *sli.do/cs566*

# 4.4  Integer Sorting

# Clicker Question

Select all **correct formulations** of our **lower bound** from §4.3.

A. Any sorting algorithm requires $O(n \log n)$ running time in the worst case.

B. Every comparison-based sorting algorithm requires $\Omega(n \log n)$ running time in worst case for sorting $n$ elements.

C. Every comparison-based sorting algorithm requires $\Omega(n \log n)$ comparisons in worst case for sorting $n$ elements.

D. Every sorting algorithm requires $\Omega(n \log n)$ comparisons in worst case for sorting $n$ elements.

E. The complexity of sorting $n$ elements in the comparison-model is $\Theta(n \log n)$.

F. The complexity of sorting $n$ elements in the comparison-model is $\Omega(n \log n)$.

→ $sli.do/cs566$

# Clicker Question

Select all **correct formulations** of our **lower bound** from §4.3.

A. ~~Any sorting algorithm requires $O(n \log n)$ running time in the worst case.~~

B. Every comparison-based sorting algorithm requires $\Omega(n \log n)$ running time in worst case for sorting $n$ elements. ✓

C. Every comparison-based sorting algorithm requires $\Omega(n \log n)$ comparisons in worst case for sorting $n$ elements. ✓

D. ~~Every sorting algorithm requires $\Omega(n \log n)$ comparisons in worst case for sorting $n$ elements.~~

E. ~~The complexity of sorting $n$ elements in the comparison-model is $\Theta(n \log n)$.~~

F. The complexity of sorting $n$ elements in the comparison-model is $\Omega(n \log n)$. ✓

→ *sli.do/cs566*

## How to beat a lower bound

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

# 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!*
   ↝ Lower bounds show where to *change* the model!

# 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!*
    - ⤳ 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)
    - ⤳ we are **not** working in the comparison model
    - ⤳ *above lower bound does not apply!*

# 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!*
  ↝ 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)
  ↝ we are **not** working in the comparison model
  ↝ *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 $[0..U) = \{0, \dots, U-1\}$      typically $U = 2^b$ $\rightsquigarrow$ $b$-bit binary numbers

# Counting sort

▶ Important parameter: size/range of numbers
  ▶ numbers in range $[0..U) = \{0, \ldots, U - 1\}$      typically $U = 2^b \rightsquigarrow$ $b$-bit binary numbers

▶ We can sort $n$ integers in $\underbrace{\Theta(n + U)}$ time and $\underbrace{\Theta(U)}$ space when $\boxed{b \leq w}$:

word size

**Counting sort**

```
1  procedure countingSort(A[0..n])
2      // A contains integers in range [0..U).
3      C[0..U] := new integer array, initialized to 0
4      // Count occurrences
5      for i := 0, . . . , n − 1
6          C[A[i]] := C[A[i]] + 1
7      i := 0 // Produce sorted list
8      for k := 0, . . . U − 1
9          for j := 1, . . . C[k]
10             A[i] := k; i := i + 1
```

▶ *count* how often each *possible* value occurs

▶ produce sorted result directly from counts

▶ circumvents lower bound by using integers as array index / pointer offset

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

# Larger Universes: Radix Sort

- ▶ *MSD Radix Sort:*
    - ▶ split numbers into base-$R$ "digits"
    - ▶ Use counting sort on <u>m</u>ost <u>s</u>ignificant <u>d</u>igit
      (with variant of counting sort that moves full number)
    - ⤳ integers sorted with respect to first digit
    - ▶ recurse on sublist for each digit value, using next digit for counting sort
- ⤳ After $\lfloor \log_R(U) \rfloor + 1$ levels of counting sort, fully sorted!
    - ▶ For $R \leq 2^w$, all counting sort calls on same level cost total of $O(n)$ time
      (requires care to avoid reinitialization cost of array $C$)
- ⤳ total time $O(n \log_R(U)) = O\left(n \frac{\log(U)}{\log(R)}\right)$

- ⤳ $O(n)$ time sorting possible for numbers in range $\underline{U = O(n^c)}$ for constant $c$.

# Integer Sorting – State of the art    $\notin$ exam

**Algorithm theory**

- integer sorting on the $w$-bit word-RAM
- suppose $U = 2^w$, but $w$ can be an arbitrary function of $n$

    /usually $w = \Theta(\log n)$

- 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!

# 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!

# Clicker Question

Which statements are correct? Select all that apply.

My computer has 64-bit words, so an `int` has 64 bits. Hence I can sort any `int[]` of length $n$ ...

- (A) in constant time.

- (B) in $O(\log n)$ time.

- (C) in $O(n)$ time.

- (D) in $O(n \log n)$ time.

- (E) some time, but not possible to say from given information.

→ *sli.do/cs566*

# Clicker Question

Which statements are correct? Select all that apply.

My computer has 64-bit words, so an `int` has 64 bits. Hence I can sort any `int[]` of length $n$ ...

A. ~~in constant time.~~

B. ~~in $O(\log n)$ time.~~

C. in $O(n)$ time. ✓

D. in $O(n \log n)$ time. ✓

E. some time, but not possible to say from given information. ✓
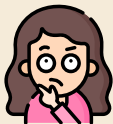
→ *sli.do/cs566*

# Part II
*Exploiting presortedness*

# 4.5 Adaptive Sorting

## Adaptive sorting

▶ Comparison lower bound also holds for the *average case* ⤳ $\lfloor \lg(n!) \rfloor$ cmps necessary

▶ Mergesort and Quicksort from above use $\sim n \lg n$ cmps even in best case

# Adaptive sorting

▶ Comparison lower bound also holds for the *average case* ⇝ $\lfloor \lg(n!) \rfloor$ 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)

⇝ 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 . . .

# 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 $\Omega(n \log n)$ extra comparisons

# 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 $\Omega(n \log n)$ extra comparisons

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

▶ If last element of first run $\leq$ first element of second run, skip merge
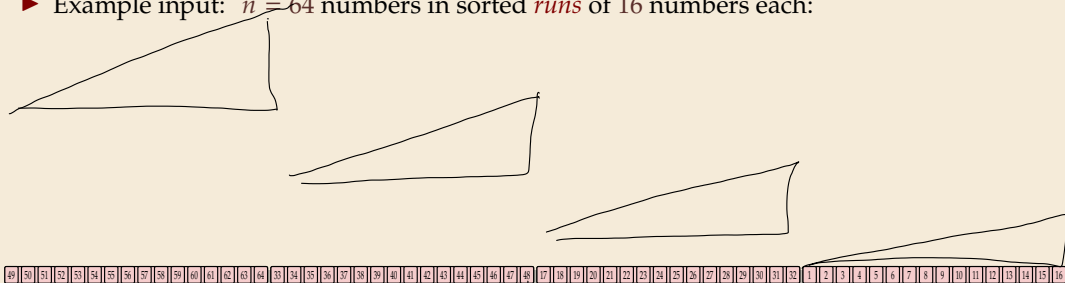
*How effective is this idea?*

```
1 procedure mergesortCheck(A[l..r])
2     n := r - l
3     if n ≤ 1 return
4     m := l + ⌊ n/2 ⌋
5     mergesortCheck(A[l..m])
6     mergesortCheck(A[m..r])
7     if A[m − 1] > A[m]
8         merge(A[l..m], A[m..r], buf)
9         copy buf to A[l..r]
```

- Simplified cost measure: *merge cost* = size of output of merges
  - ≈ number of comparisons
  - ≈ number of memory transfers / cache misses

- Example input: $n = 64$ numbers in sorted *runs* of 16 numbers each:



49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16

## Mergesort with sorted check – Analysis

▶ Simplified cost measure:   *merge cost* = size of output of merges
                                           ≈ number of comparisons
                                           ≈ number of memory transfers / cache misses

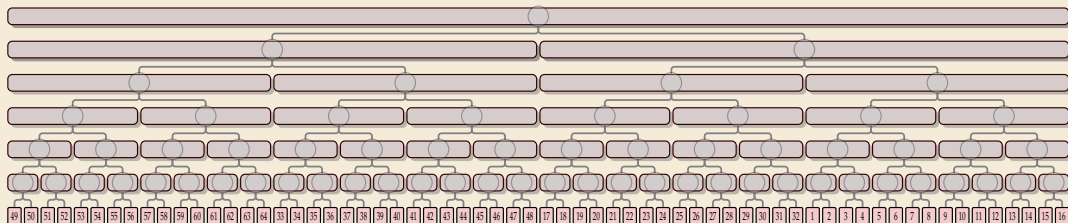▶ Example input:   $n = 64$ numbers in sorted *runs* of $16$ numbers each:

49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16

# Mergesort with sorted check – Analysis

▶ Simplified cost measure:  *merge cost*  = size of output of merges

$$\approx \text{ number of comparisons}$$

$$\approx \text{ number of memory transfers / cache misses}$$

▶ Example input:  $n = 64$ numbers in sorted *runs* of 16 numbers each:
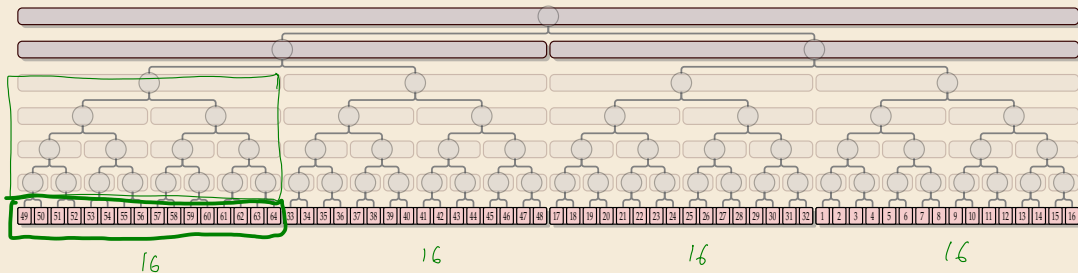


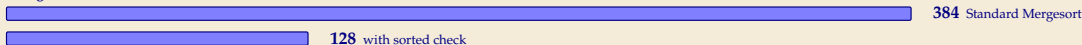**Merge costs:**

**384** Standard Mergesort

# Mergesort with sorted check – Analysis

▶ Simplified cost measure:  *merge cost*  =  size of output of merges
                                          ≈  number of comparisons
                                          ≈  number of memory transfers / cache misses

▶ Example input:  $n = 64$ numbers in sorted *runs* of 16 numbers each:



**Merge costs:**

**384** Standard Mergesort

**128** with sorted check

*Sorted check can help a lot!*

# Alignment issues

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

   = exactly the cost of *creating* this run in mergesort had it not already existed

   ⤳ best savings we can hope for!

⤳ Are overall merge costs $\mathcal{H}(\ell_1, \dots, \ell_r) := \underbrace{n \lg(n)}_{\text{mergesort}} - \underbrace{\sum_{i=1}^{r} \ell_i \lg(\ell_i)}_{\text{savings from runs}}$ ?

      $\ell_i$ = length of $i$th run

run length entropy

alternative intuition about $\mathcal{H}$:

$$= \frac{n!}{\ell_1! \cdots \ell_r!}$$

$$\mathcal{H} = \lg \binom{n}{\ell_1 \ \cdots \ \ell_r}$$

\# bits of information to learn

(previously $\lg(n!)$ )

25

# Alignment issues

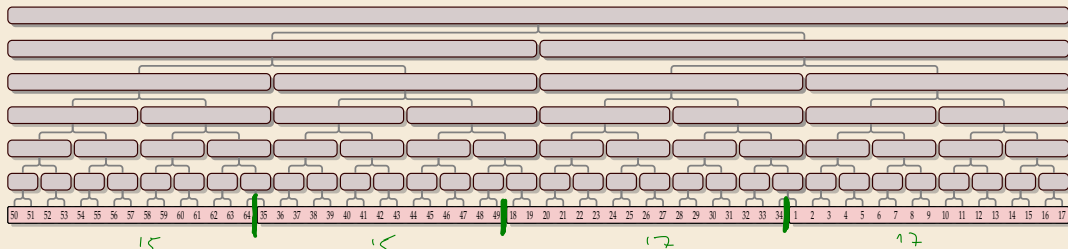- In previous example, each run of length $\ell$ saved us $\ell \lg(\ell)$ in merge cost.

  - = exactly the cost of *creating* this run in mergesort had it not already existed

  - ⤳ best savings we can hope for!

- ⤳ Are overall merge costs $\mathcal{H}(\ell_1, \ldots, \ell_r) := \underbrace{n \lg(n)}_{\text{mergesort}} - \underbrace{\sum_{i=1}^{r} \ell_i \lg(\ell_i)}_{\text{savings from runs}}$ ? $\quad \overset{\ell_i \,=\, \text{length of } i\text{th run}}{}$

**Unfortunately, not quite:**



**Merge costs:**

**384** Standard Mergesort

**127.8** $\mathcal{H}(15, 15, 17, 17)$

# Alignment issues

► In previous example, each run of length $\ell$ saved us $\ell \lg(\ell)$ in merge cost.
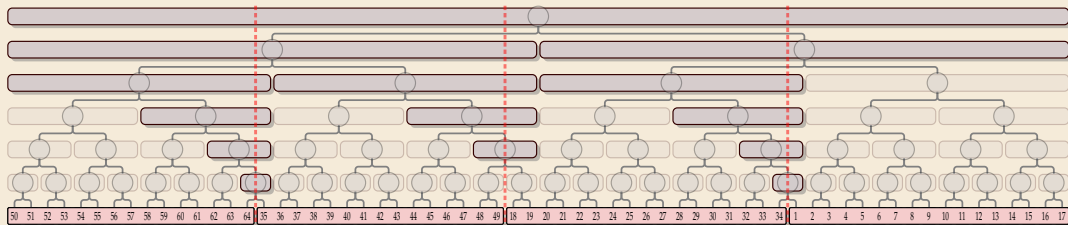
  = exactly the cost of *creating* this run in mergesort had it not already existed
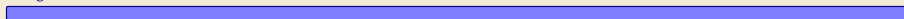
  ⇝ best savings we can hope for!

⇝ Are overall merge costs    $\mathcal{H}(\ell_1, \ldots, \ell_r) := \underbrace{n \lg(n)}_{\text{mergesort}} - \underbrace{\sum_{i=1}^{r} \ell_i \lg(\ell_i)}_{\text{savings from runs}}$ ?   $\ell_i = \text{length of } i\text{th run}$
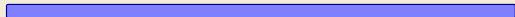
**Unfortunately, not quite:**



Merge costs:

**384** Standard Mergesort

**216** with sorted check

**127.8** $\mathcal{H}(15, 15, 17, 17)$

25

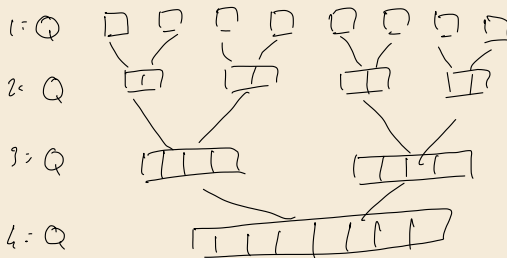## Bottom-Up Mergesort

► Can we do better by explicitly detecting runs?

```
 1  procedure bottomUpMergesort(A[0..n])
 2      Q := new Queue // runs to merge
 3      // Phase 1: Enqueue singleton runs
 4      for i = 0, . . . , n − 1 do
 5          Q.enqueue((i, i + 1))
 6      // Phase 2: Merge runs level −wise
 7      while Q.size() ≥ 2
 8          Q' := new Queue
 9          while Q.size() ≥ 2
10              (i_1, j_1) := Q.dequeue()
11              (i_2, j_2) := Q.dequeue()
12              merge(A[i_1..j_1], A[i_2..j_2], buf)
13              copy buf to A[i_1..j_2]
14              Q'.enqueue((i_1, j_2))
15          if ¬Q.isEmpty() // lonely run
16              Q'.enqueue(Q.dequeue())
17          Q := Q'
```

1: Q

2: Q

3: Q

4: Q

26

## Bottom-Up Mergesort

▶ Can we do better by explicitly detecting runs?

```
1  procedure bottomUpMergesort(A[0..n])
2      Q := new Queue // runs to merge
3      // Phase 1: Enqueue singleton runs
4      for i = 0, . . . , n − 1 do
5          Q.enqueue((i, i + 1))
6      // Phase 2: Merge runs level −wise
7      while Q.size() ≥ 2
8          Q' := new Queue
9          while Q.size() ≥ 2
10             (i₁, j₁) := Q.dequeue()
11             (i₂, j₂) := Q.dequeue()
12             merge(A[i₁..j₁], A[i₂..j₂], buf)
13             copy buf to A[i₁..j₂]
14             Q'.enqueue((i₁, j₂))
15         if ¬Q.isEmpty() // lonely run
16             Q'.enqueue(Q.dequeue())
17         Q := Q'
```

```
1  procedure naturalMergesort(A[0..n])
2      Q := new Queue;  i := 0        find run A[i..j]
3      while i < n                    ↙ starting at i
4          j := i + 1
5          while A[j] ≥ A[j − 1] do j := j + 1
6          Q.enqueue((i, j));  i := j
7      while Q.size() ≥ 2
8          Q' := new Queue
9          while Q.size() ≥ 2
10             (i₁, j₁) := Q.dequeue()
11             (i₂, j₂) := Q.dequeue()
12             merge(A[i₁..j₁], A[i₂..j₂], buf)
13             copy buf to A[i₁..j₂]
14             Q'.enqueue((i₁, j₂))
15         if ¬Q.isEmpty() // lonely run
16             Q'.enqueue(Q.dequeue())
17         Q := Q'
```

# Clicker Question

Suppose we have an input with the 5 elements a, b, c, d, e and we sort them with **bottomUpMergesort**. What sequence of merges are executed?

[A] Policy 1

| a | b | c | d | e |
| a | b | c | d | e |
| a | b | c | d | e |
| a | b | c | d | e |

Policy 1

[B] Policy 2

| a | b | c | d | e |
| a | b | c | d | e |
| a | b | c | d | e |
| a | b | c | d | e |

Policy 2

[C] Policy 3

| a | b | c | d | e |
| a | b | c | d | e |
| a | b | c | d | e |
| a | b | c | d | e |

Policy 3

→ *sli.do/cs566*

# Clicker Question

Suppose we have an input with the 5 elements a, b, c, d, e and we sort them with **bottomUpMergesort**. What sequence of merges are executed?

A  Policy 1 ✓

| a | b | c | d | e |
|---|---|---|---|---|

| a | b | c | d | e |
|---|---|---|---|---|

| a | b | c | d | e |
|---|---|---|---|---|

| a | b | c | d | e |
|---|---|---|---|---|

Policy 1

B  ~~Policy 2~~

| a | b | c | d | e |
|---|---|---|---|---|

| a | b | c | d | e |
|---|---|---|---|---|

| a | b | c | d | e |
|---|---|---|---|---|

| a | b | c | d | e |
|---|---|---|---|---|

Policy 2

C  ~~Policy 3~~

| a | b | c | d | e |
|---|---|---|---|---|

| a | b | c | d | e |
|---|---|---|---|---|

| a | b | c | d | e |
|---|---|---|---|---|

| a | b | c | d | e |
|---|---|---|---|---|

Policy 3

→ sli.do/cs566

# Natural Bottom-Up Mergesort – Analysis

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



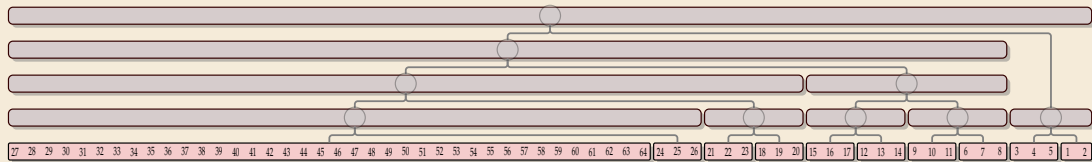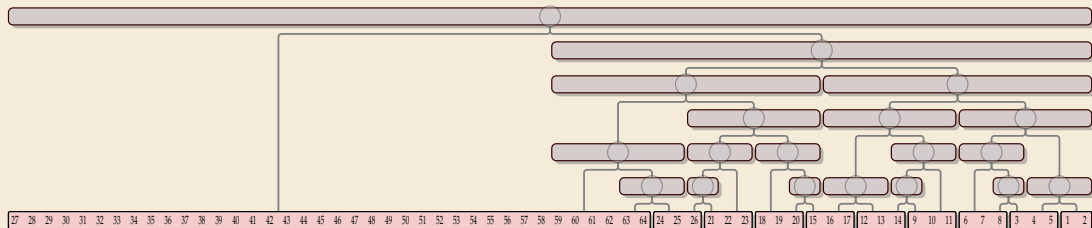Merge costs:

**384** Standard mergesort

**216** Standard mergesort with sorted check

**127.8** $\mathcal{H}(15, 15, 17, 17)$

**128** Natural bottom-up mergesort

# Natural Bottom-Up Mergesort – Analysis [2]

► . . . but less so for widely varying run lengths



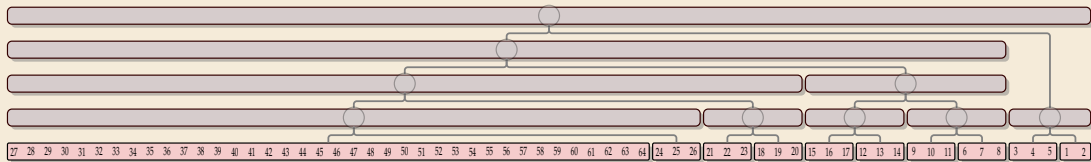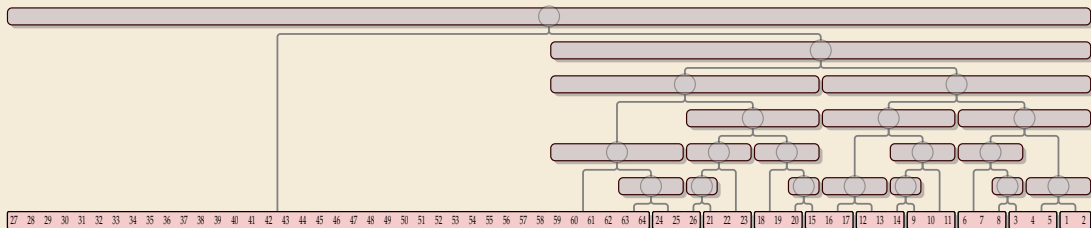**246** Natural bottom-up mergesort



**196** Standard mergesort with sorted check

► ... but less so for widely varying run lengths



**246** Natural bottom-up mergesort



**196** Standard mergesort with sorted check

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

# Good merge orders

*Let's take a step back and breathe.*
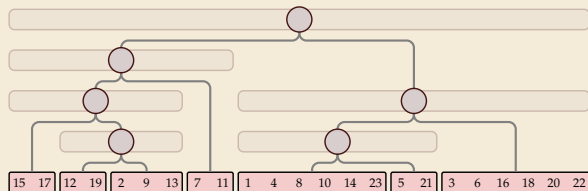
# Good merge orders

◀◀ *Let's take a step back and breathe.*

▶ Conceptually, there are two tasks:

1. Detect and use existing runs in the input $\leadsto \ell_1, \ldots, \ell_r$     (easy)
2. Determine a favorable *order of merges* of runs     ("automatic" in top-down mergesort)

# 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 $\leadsto \ell_1, \ldots, \ell_r$ (easy) ✓
    2. **Determine a favorable *order of merges* of runs** ("automatic" in top-down mergesort)



**Merge cost** = total area of ⬭

# 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 $\rightsquigarrow \ell_1, \ldots, \ell_r$    (easy) ✓
    2. **Determine a favorable *order of merges* of runs**    ("automatic" in top-down mergesort)
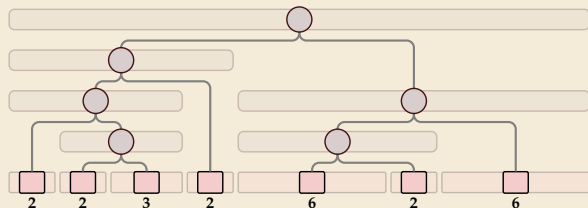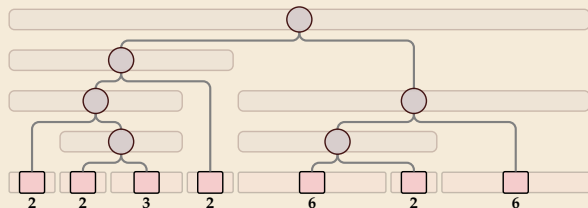


| 15 17 | 12 19 | 2 9 13 | 7 11 | 1 4 8 10 14 23 | 5 21 | 3 6 16 18 20 22 |

**Merge cost** = total area of ▭

             = total length of paths to all array entries

# 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   $\rightsquigarrow$   $\ell_1, \ldots, \ell_r$     (easy) ✓

    *2.* **Determine a favorable *order of merges* of runs**     ("automatic" in top-down mergesort)



**Merge cost** = total area of ⬭

               = total length of paths to all array entries

$$= \sum_{w \text{ leaf}} weight(w) \cdot depth(w)$$

# 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   $\rightsquigarrow$   $\ell_1, \ldots, \ell_r$    (easy) ✓

   *2.* **Determine a favorable *order of merges* of runs**    ("automatic" in top-down mergesort)



**2**    **2**    **3**    **2**    **6**    **2**    **6**

**Merge cost** = total area of ⬭

             = total length of paths to all array entries

$$= \sum_{w \text{ leaf}} weight(w) \cdot depth(w)$$

well-understood problem
with known algorithms

$\rightsquigarrow$ *optimal* merge tree
= optimal *binary search tree*
   for leaf weights $\ell_1, \ldots, \ell_r$
   (optimal expected search cost)

# Nearly-Optimal Mergesort

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

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

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

—— **Abstract** ——
We present two stable mergesort variants, "peeksort" and "powersort", that exploit existing runs and find nearly-optimal merging orders with negligible overhead. Previous methods either require substantial effort for determining the merging order (Takaoka 2009; Barbay & Navarro 2013) or do not have an optimal worst-case guarantee (Peters 2002; Auger, Nicaud & Pivoteau 2015; Buss & Knop 2018). We demonstrate that our methods are competitive in terms of running time with state-of-the-art implementations of stable sorting methods.

## 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 (close-to) optimal comparison-based sorting methods are known, *instance-optimal sorting*, i.e., methods that adapt to the actual input and exploit specific structural properties if present, is still an area of active research. We survey some recent developments 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 Oracle's Java runtime library: adapting to the presence of duplicate keys and using existing sorted segments, called *runs*. The former is achieved by a so-called fat-pivot partitioning variant of quicksort [8], which is also used in the OpenBSD implementation of `qsort` from the C standard library. It is an *unstable* sorting method, though, i.e., the relative order of elements with equal keys might be destroyed in the process. It is hence used in Java solely for primitive-type arrays.

---

▶ 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 $\leq \mathcal{H}(\ell_1, \ldots, \ell_r) + 2n$

  ▶ 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

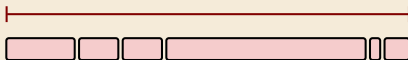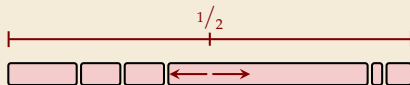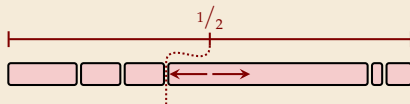- ⇝ split there and recurse
  (instead of at midpoint)

## Peeksort

- based on top-down mergesort

- "peek" at middle of array
  & find closest run boundary

⇝ split there and recurse
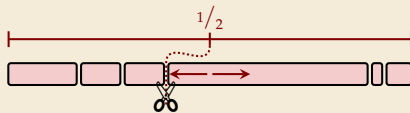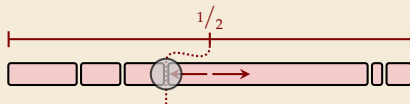  (instead of at midpoint)

## Peeksort



- ▶ based on top-down mergesort

- ▶ "peek" at middle of array
  & find closest run boundary

- ⇝ split there and recurse
  (instead of at midpoint)

## Peeksort



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

# Peeksort

- based on top-down mergesort

- "peek" at middle of array & find closest run boundary

⤳ split there and recurse
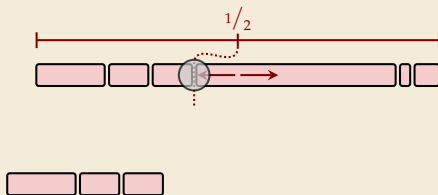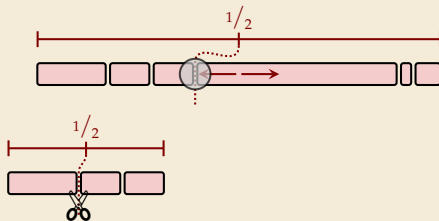   (instead of at midpoint)

# Peeksort



- ▶ based on top-down mergesort

- ▶ "peek" at middle of array
  & find closest run boundary

- ⤳ split there and recurse
  (instead of at midpoint)

# Peeksort

- based on top-down mergesort

- "peek" at middle of array & find closest run boundary
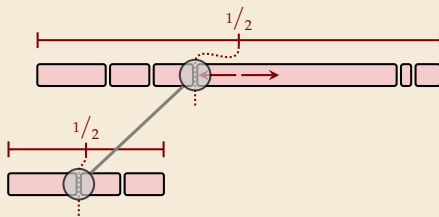
⤳ split there and recurse
  (instead of at midpoint)

# Peeksort

- based on top-down mergesort

- "peek" at middle of array
  & find closest run boundary

⤳ split there and recurse
  (instead of at midpoint)

# Peeksort

- based on top-down mergesort

- "peek" at middle of array & find closest run boundary

⇝ split there and recurse
  (instead of at midpoint)

# Peeksort

- based on top-down mergesort

- "peek" at middle of array
  & find closest run boundary

⤳ split there and recurse
  (instead of at midpoint)

# Peeksort

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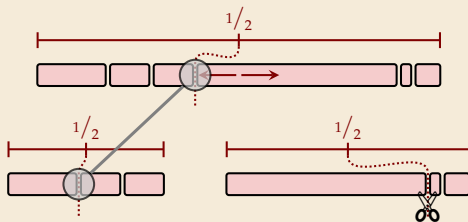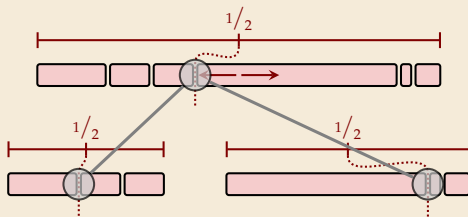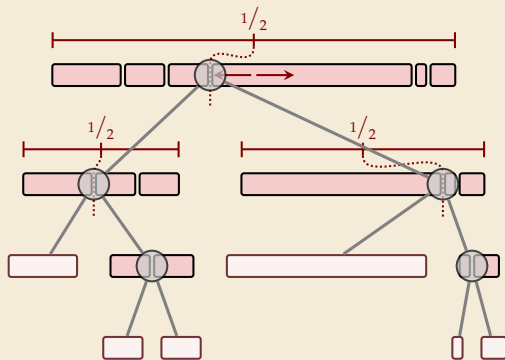- ⇝ split there and recurse
  (instead of at midpoint)

# Peeksort

- based on top-down mergesort

- "peek" at middle of array
  & find closest run boundary

$\rightsquigarrow$ split there and recurse
  (instead of at midpoint)

# Peeksort

- based on top-down mergesort
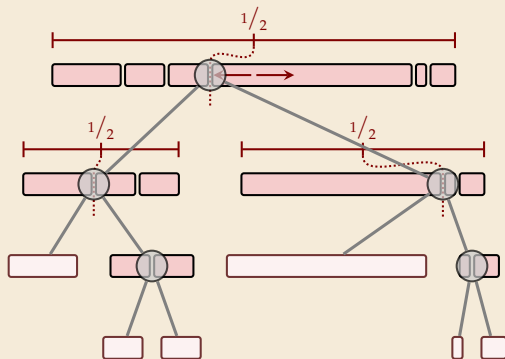
- "peek" at middle of array
  & find closest run boundary

$\rightsquigarrow$ 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



$\ell \qquad \ell+\Delta_\ell \qquad r+\Delta_r \qquad r$

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

## Peeksort – Code

```
1  procedure peeksort(A[ℓ..r], Δ_ℓ, Δ_r)
2      if r − ℓ ≤ 1 then return
3      if ℓ + Δ_ℓ == r ∨ ℓ == r + Δ_r then return
4      m := ℓ + ⌊(r − ℓ)/2⌋
5      i :=  ⎧ ℓ + Δ_ℓ              if ℓ + Δ_ℓ ≥ m
            ⎨
            ⎩ extendRunLeft(A, m)   else
6      j :=  ⎧ r + Δ_r              if r + Δ_r ≤ m ≤ m
            ⎨
            ⎩ extendRunRight(A, m)  else
7      g :=  ⎧ i   if m − i < j − m
            ⎨
            ⎩ j   else
8      Δ_g :=  ⎧ j − i   if m − i < j − m
              ⎨
              ⎩ i − j   else
9      peeksort(A[ℓ..g], Δ_ℓ, Δ_g)
10     peeksort(A[g, r], Δ_g, Δ_r)
11     merge(A[ℓ, g], A[g..r], buf)
12     copy buf to A[ℓ..r]
```

► Parameters:



► initial call:
peeksort($A[0..n], \Delta_0, \Delta_n$) with
$\Delta_0$ = extendRunRight($A$, 0)
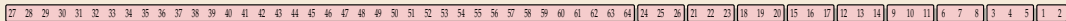$\Delta_n$ = $n$ − extendRunLeft($A$, $n$)

► helper procedure

```
1  procedure extendRunRight(A[0..n], i)
2      j := i + 1
3      while j < n ∧ A[j − 1] ≤ A[j]
4          j := j + 1
5      return j
```

(extendRunLeft similar)

# Peeksort – Analysis

▶ Consider tricky input from before again:

27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 | 24 25 26 | 21 22 23 | 18 19 20 | 15 16 17 | 12 13 14 | 9 10 11 | 6 7 8 | 3 4 5 | 1 2

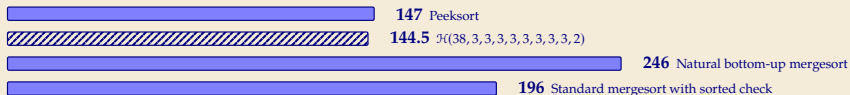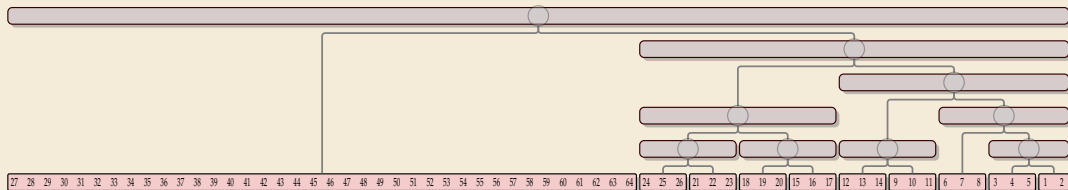**144.5** $\mathcal{H}(38, 3, 3, 3, 3, 3, 3, 3, 3, 2)$

**246** Natural bottom-up mergesort

**196** Standard mergesort with sorted check

# Peeksort – Analysis

▶ Consider tricky input from before again:



▶ One can prove: Mergecost always $\leq \mathcal{H}(\ell_1, \ldots, \ell_r) + 2n$

⤳ We can have the best of both worlds!