

ALGORITHMS \$ EFFICIENT
CIENTALGORITHMS \$ EFFI
EFFICIENTALGORITHMS \$
FFICIENTALGORITHMS \$
FICIENTALGORITHMS \$
GORITHMS \$ EFFICIENTAL
HMS \$ EFFICIENTALGORIT

3

Efficient Sorting – The Power of Divide & Conquer

13 October 2023

Sebastian Wild

Learning Outcomes

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.

Unit 3: *Efficient Sorting*



Outline

3 Efficient Sorting

- 3.1 Mergesort
- 3.2 Quicksort
- 3.3 Comparison-Based Lower Bound
- 3.4 Integer Sorting
- 3.5 Adaptive Sorting
- 3.6 Python's list sort
- 3.7 Order Statistics
- 3.8 Further D&C Algorithms

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
- Algorithm with optimal #comparisons in worst case?

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 $A[0], \dots, A[n - 1]$

- (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., $A[0] \leq A[1] \leq \dots \leq A[n - 1]$

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

Clicker Question



What is the complexity of sorting? Type your answer, e. g., as
"Theta(sqrt(n))"

$\Theta(n \log n)$ \leftarrow $O(n \log n)$ & $\Omega(n \log n)$
algorithm lower bound

of that cost impossible to do better



[→ sli.do/comp526](https://sli.do/comp526)

3.1 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/comp526*

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 ~~Crow 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.~~

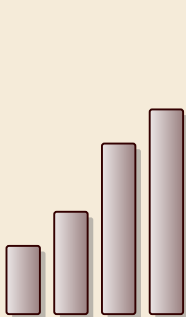


→ sli.do/comp526

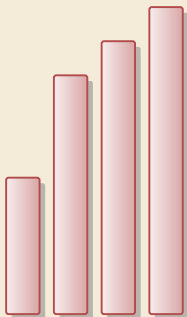
Merging sorted lists



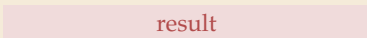
Merging sorted lists



run1

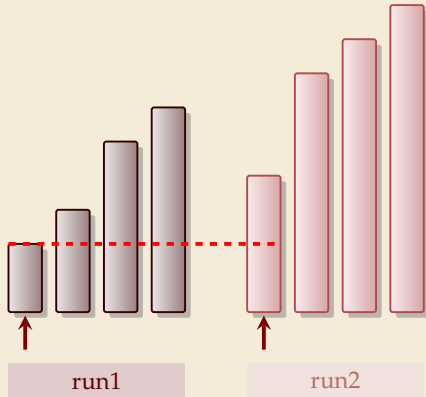


run2

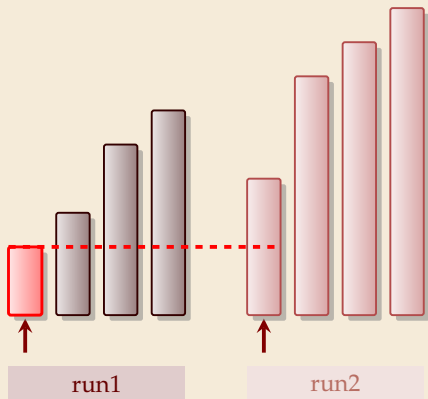


result

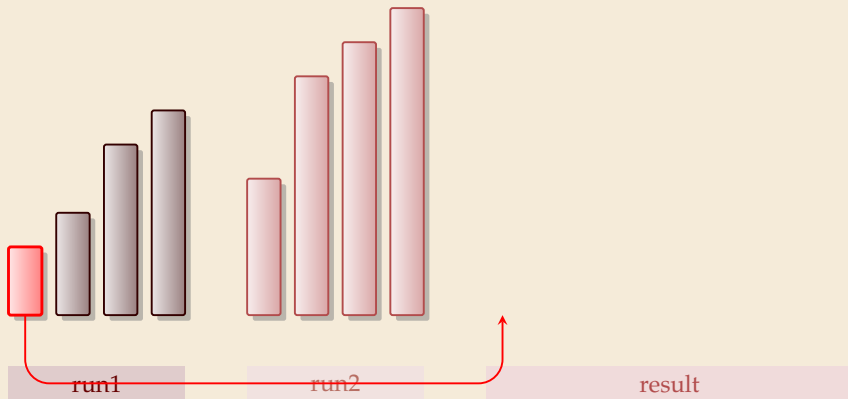
Merging sorted lists



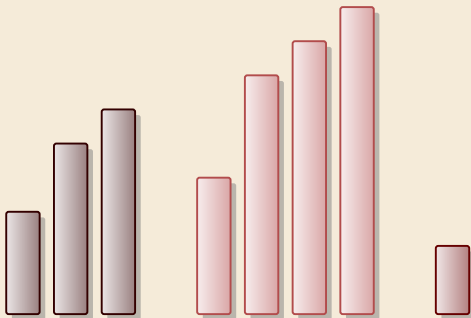
Merging sorted lists



Merging sorted lists



Merging sorted lists

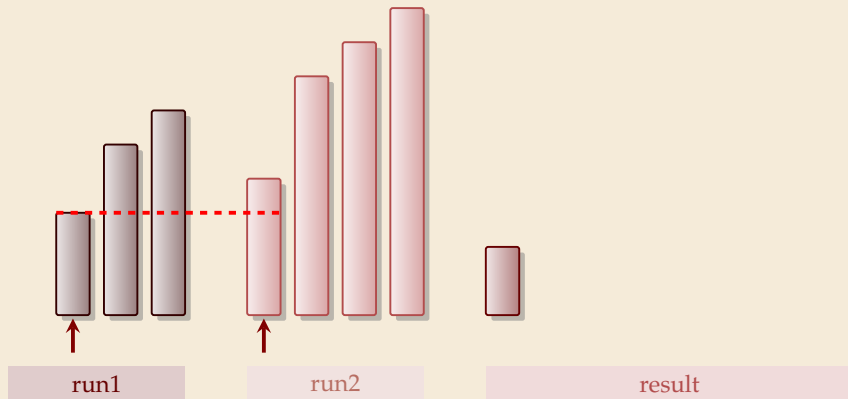


run1

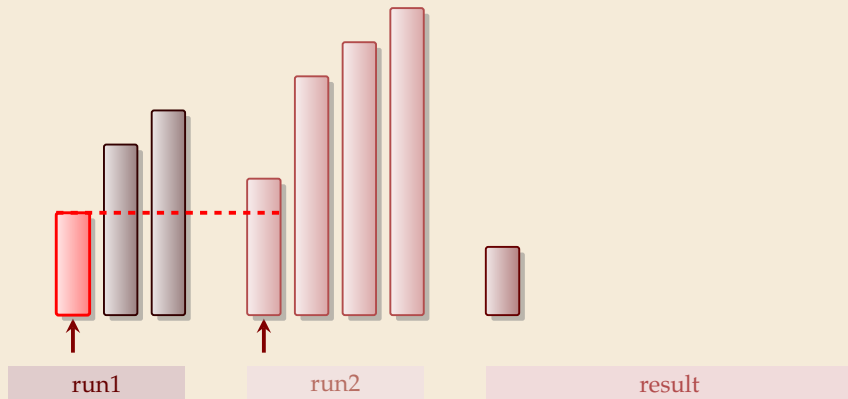
run2

result

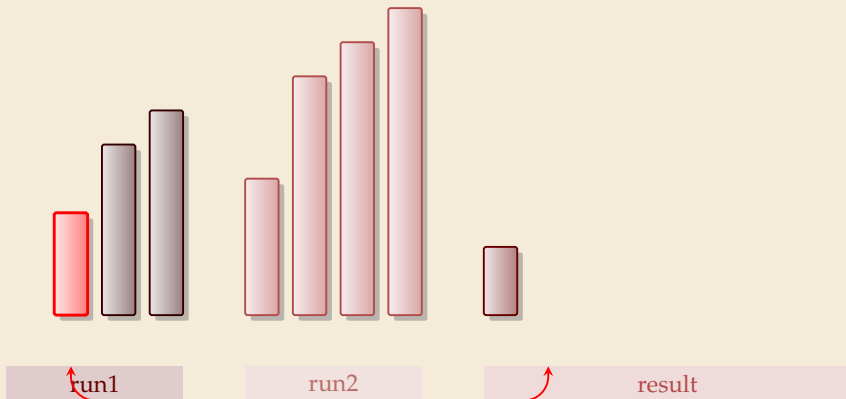
Merging sorted lists



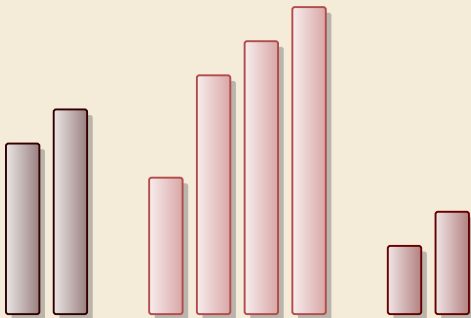
Merging sorted lists



Merging sorted lists



Merging sorted lists

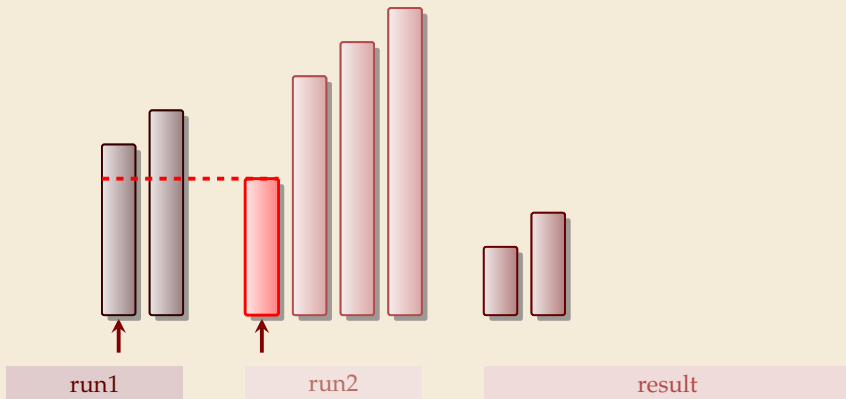


run1

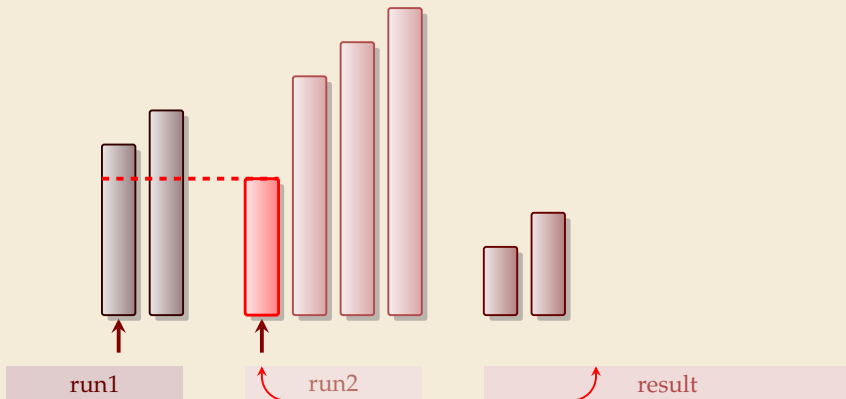
run2

result

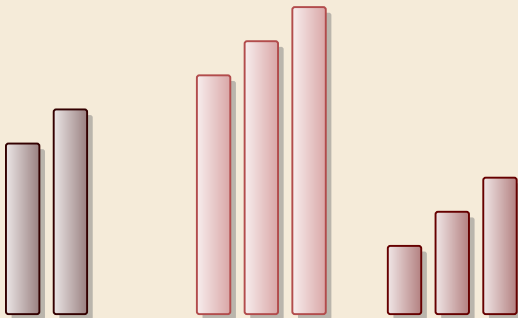
Merging sorted lists



Merging sorted lists



Merging sorted lists

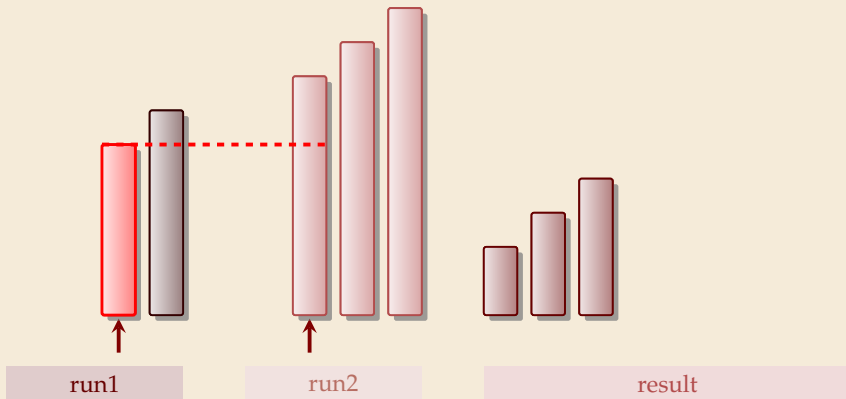


run1

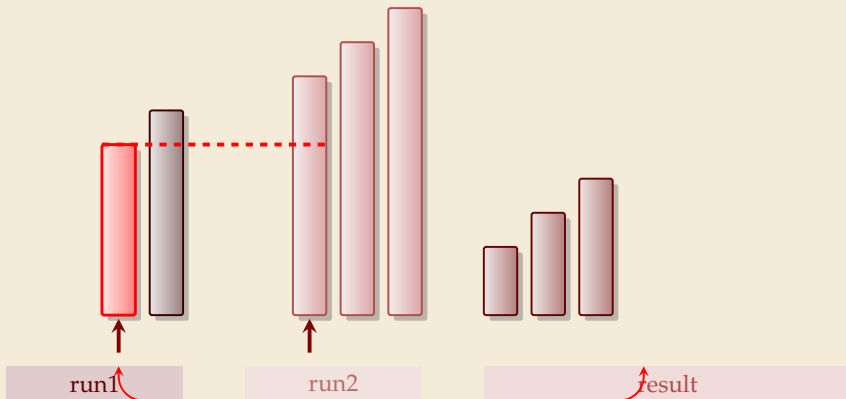
run2

result

Merging sorted lists



Merging sorted lists



Merging sorted lists



run1

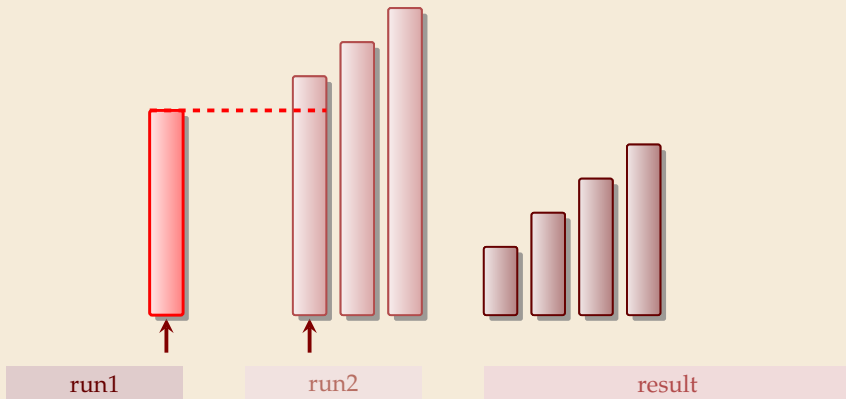


run2

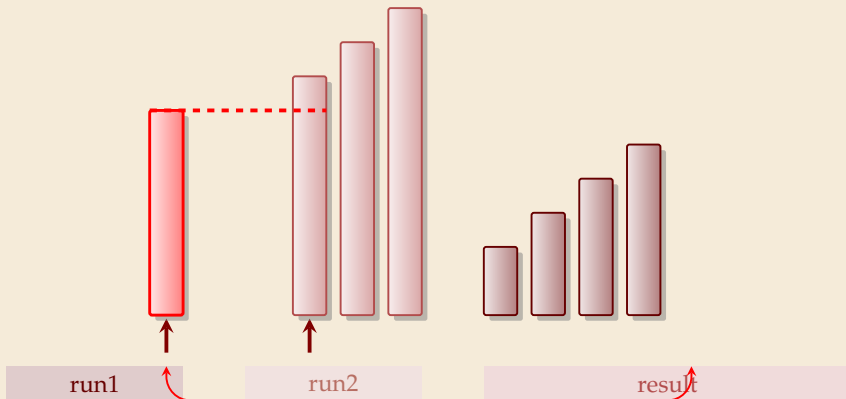


result

Merging sorted lists



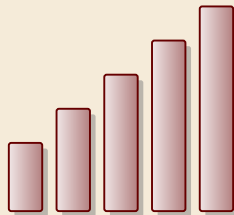
Merging sorted lists



Merging sorted lists



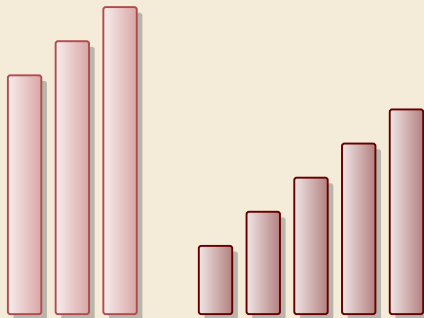
run1



run2

result

Merging sorted lists

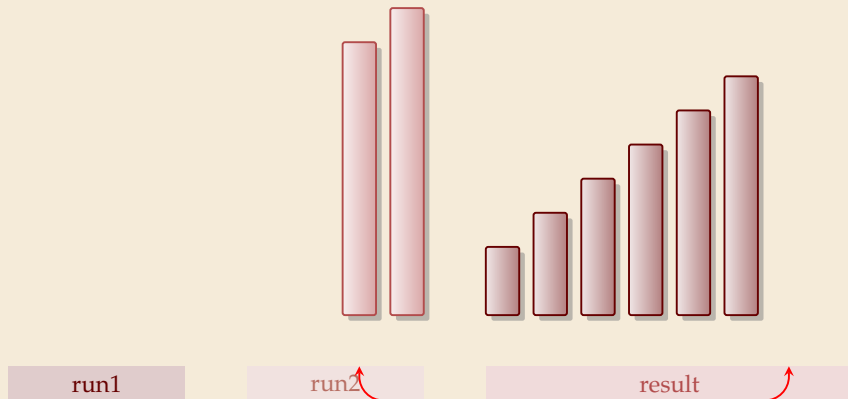


run1

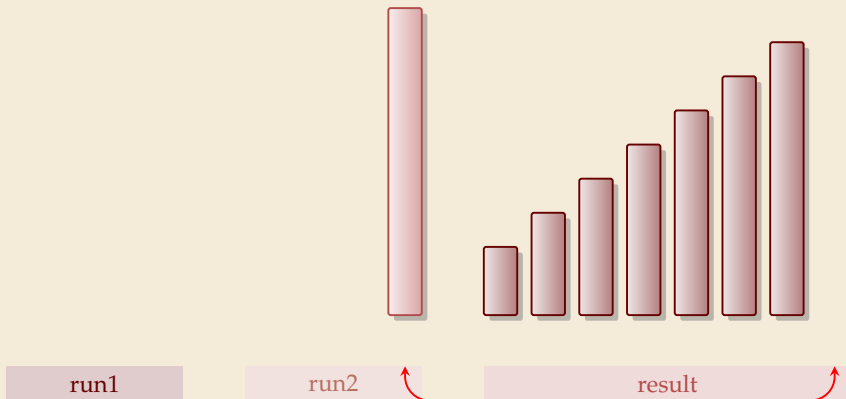
run2

result

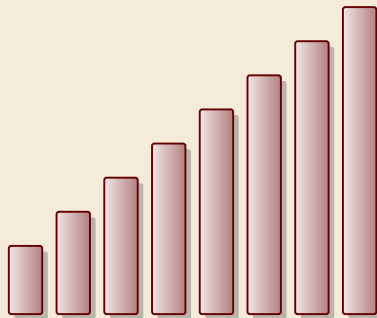
Merging sorted lists



Merging sorted lists



Merging sorted lists



run1

run2

result

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/comp526

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/comp526

Mergesort

```
1 procedure mergesort( $A[l..r]$ )
```

```
2    $n := r - l$ 
```

```
3   if  $n \leq 1$  return
```

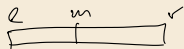
```
4    $m := l + \lfloor \frac{n}{2} \rfloor$ 
```

```
5   mergesort( $A[l..m]$ )
```

```
6   mergesort( $A[m..r]$ )
```

```
7   merge( $A[l..m]$ ,  $A[m..r]$ ,  $buf$ )
```

```
8   copy  $buf$  to  $A[l..r]$ 
```



$mergesort(A[0..n])$

► 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  $n \leq 1$  return  
4    $m := l + \lfloor \frac{n}{2} \rfloor$   
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)

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

comps between $\frac{n}{2}$ and n

Simplification

$$n = 2^k$$

same for best and worst case!

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

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

Mergesort

```
1 procedure mergesort( $A[l..r]$ )  
2    $n := r - l$   
3   if  $n \leq 1$  return  
4    $m := l + \lfloor \frac{n}{2} \rfloor$   
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)

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

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

Simplification $n = 2^k$

same for best and worst case!

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

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

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

3.2 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/comp526*

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.~~

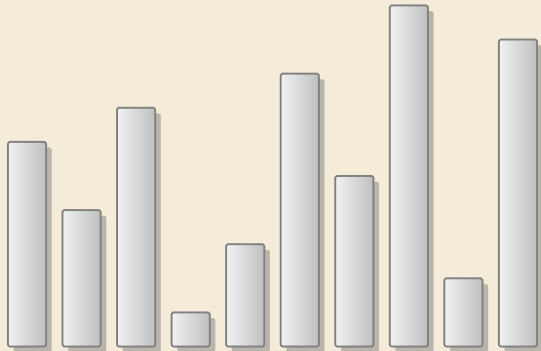


→ sli.do/comp526

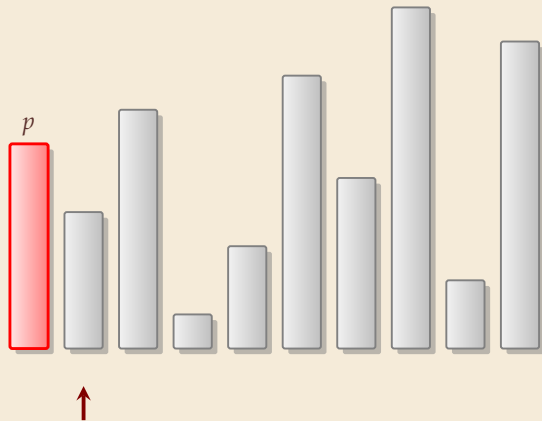
Partitioning around a pivot



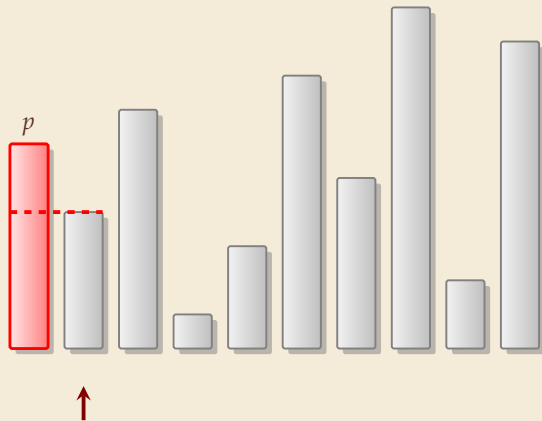
Partitioning around a pivot



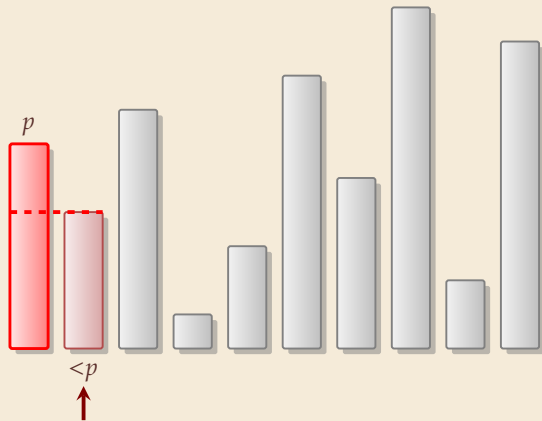
Partitioning around a pivot



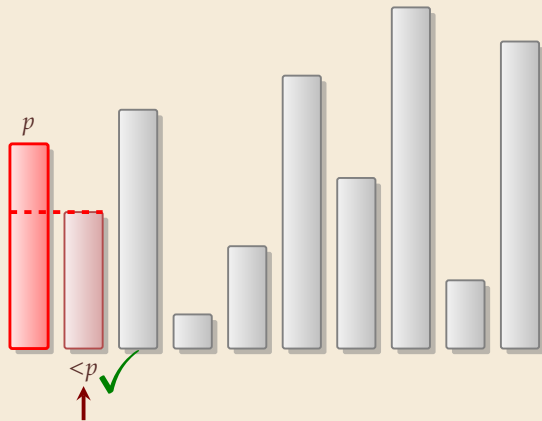
Partitioning around a pivot



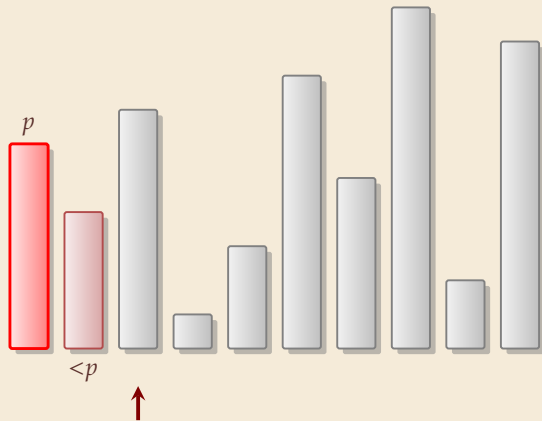
Partitioning around a pivot



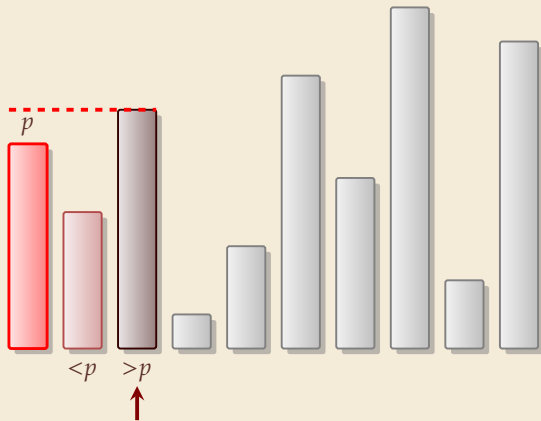
Partitioning around a pivot



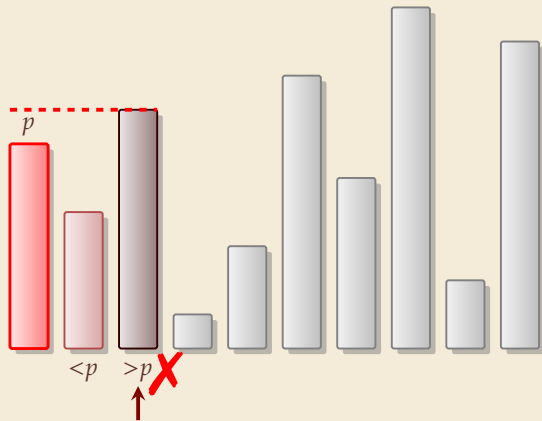
Partitioning around a pivot



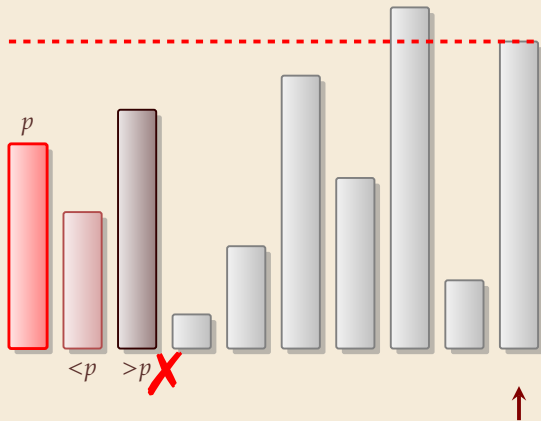
Partitioning around a pivot



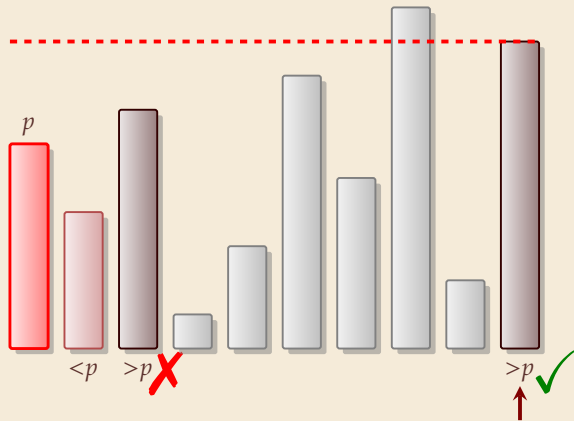
Partitioning around a pivot



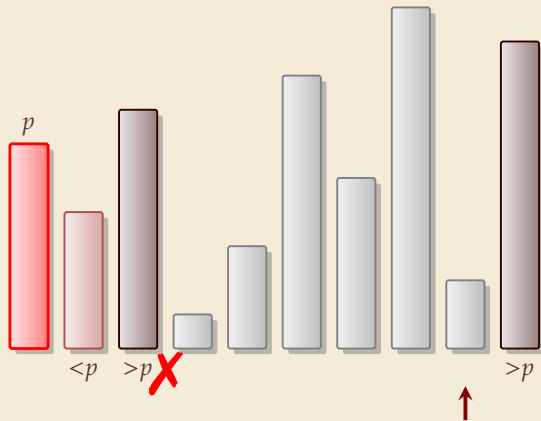
Partitioning around a pivot



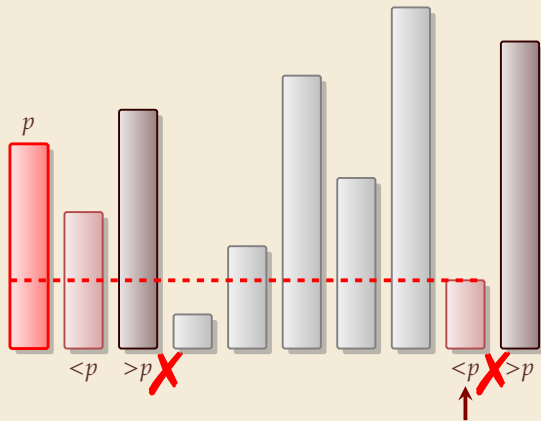
Partitioning around a pivot



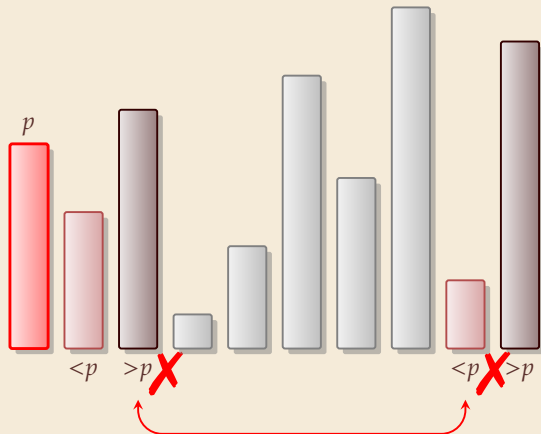
Partitioning around a pivot



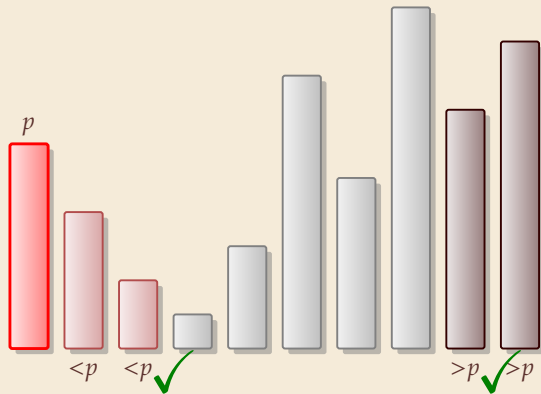
Partitioning around a pivot



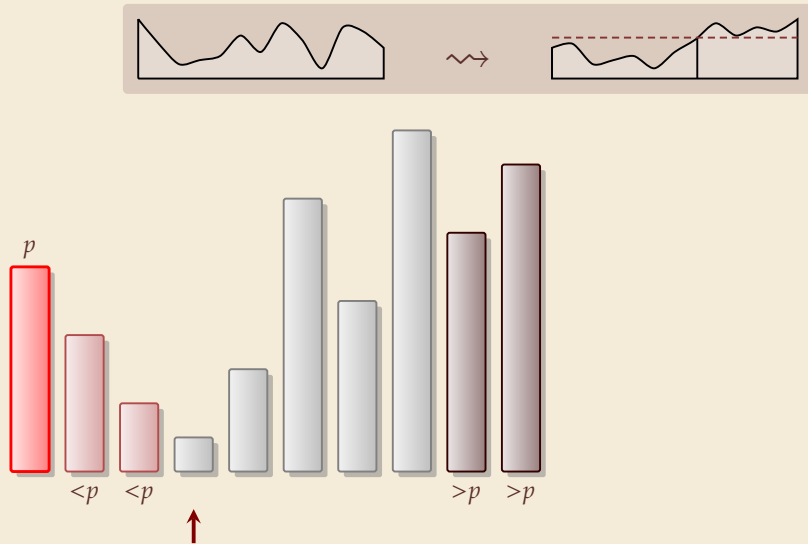
Partitioning around a pivot



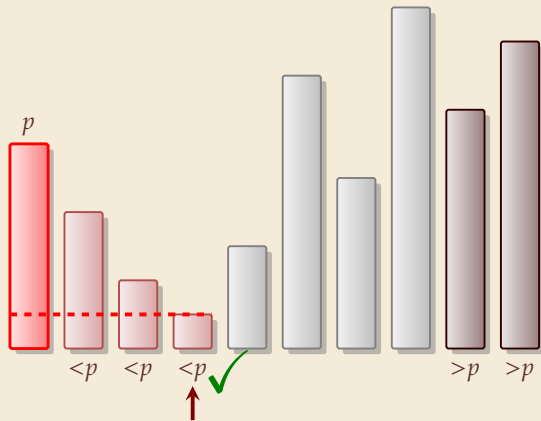
Partitioning around a pivot



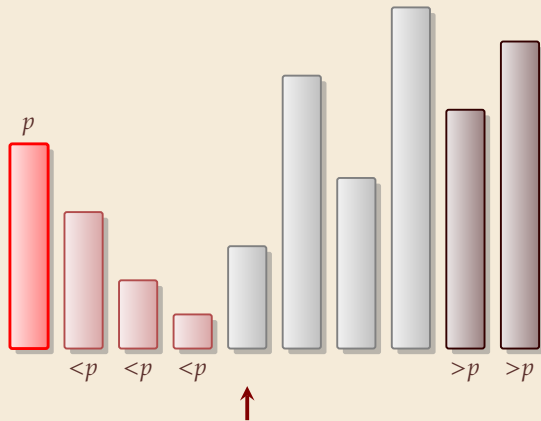
Partitioning around a pivot



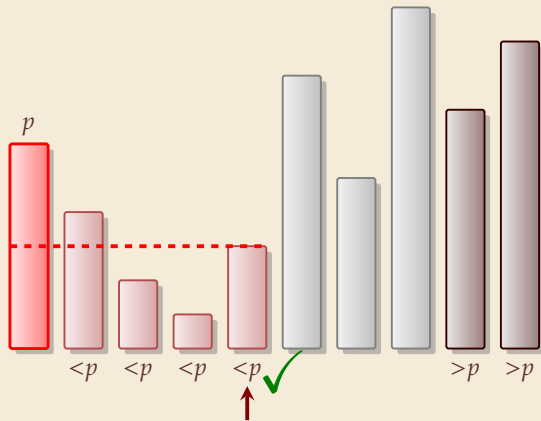
Partitioning around a pivot



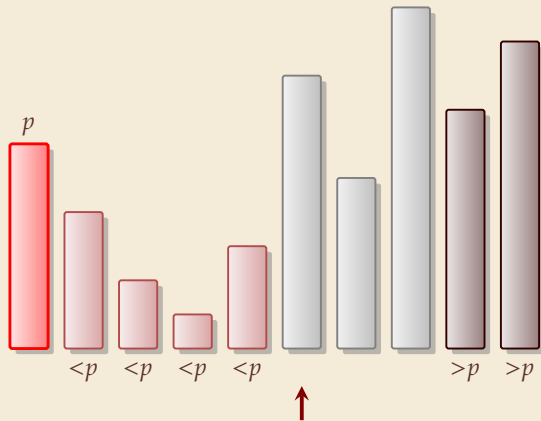
Partitioning around a pivot



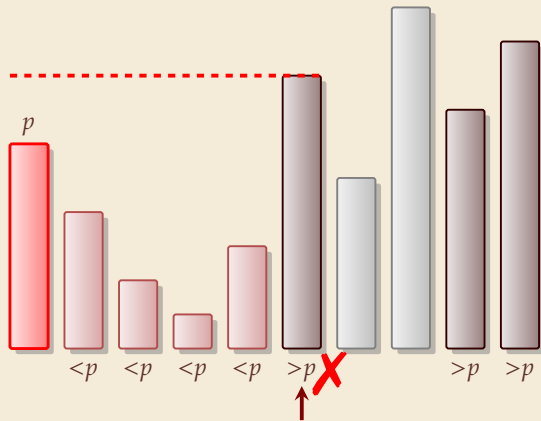
Partitioning around a pivot



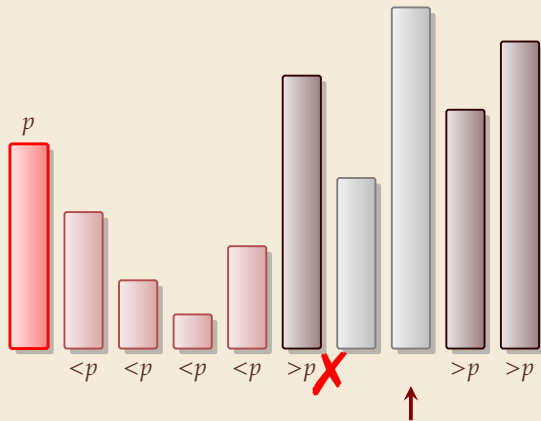
Partitioning around a pivot



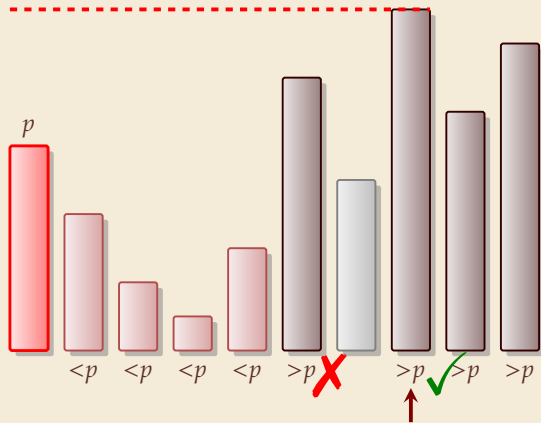
Partitioning around a pivot



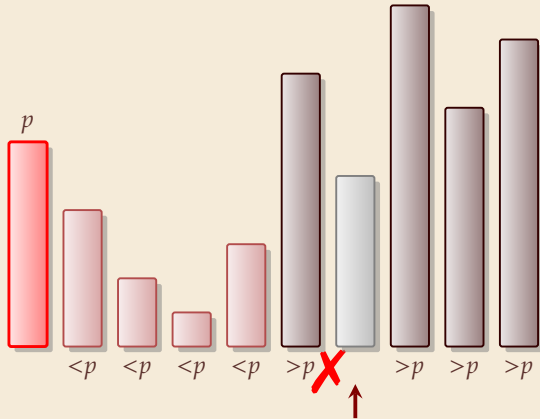
Partitioning around a pivot



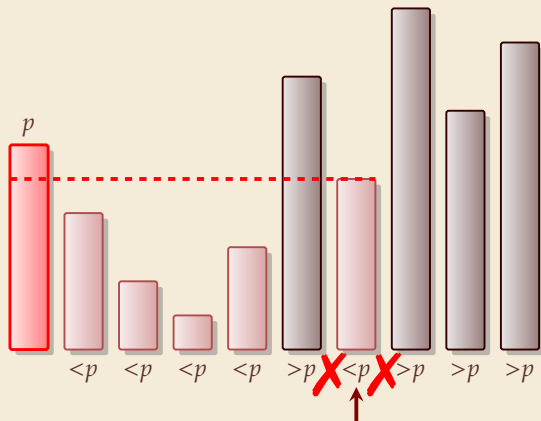
Partitioning around a pivot



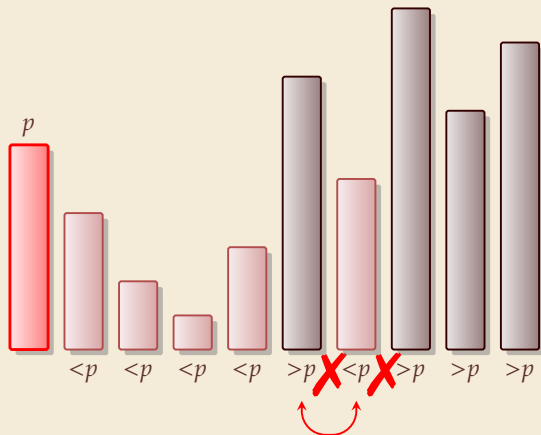
Partitioning around a pivot



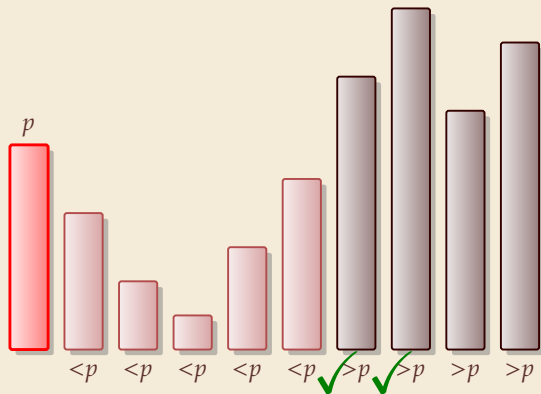
Partitioning around a pivot



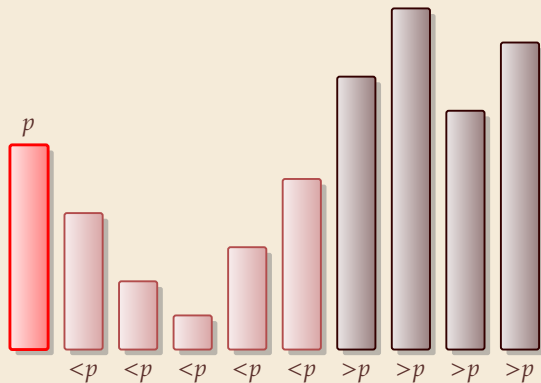
Partitioning around a pivot



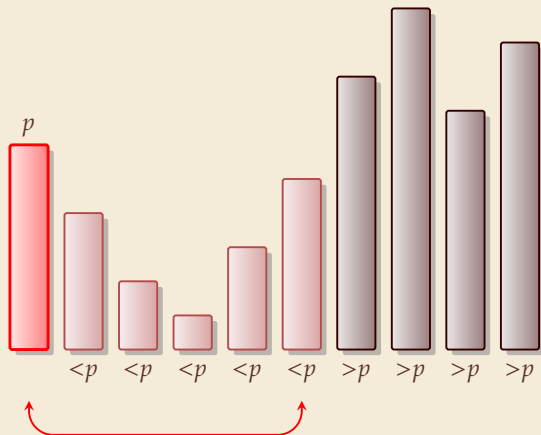
Partitioning around a pivot



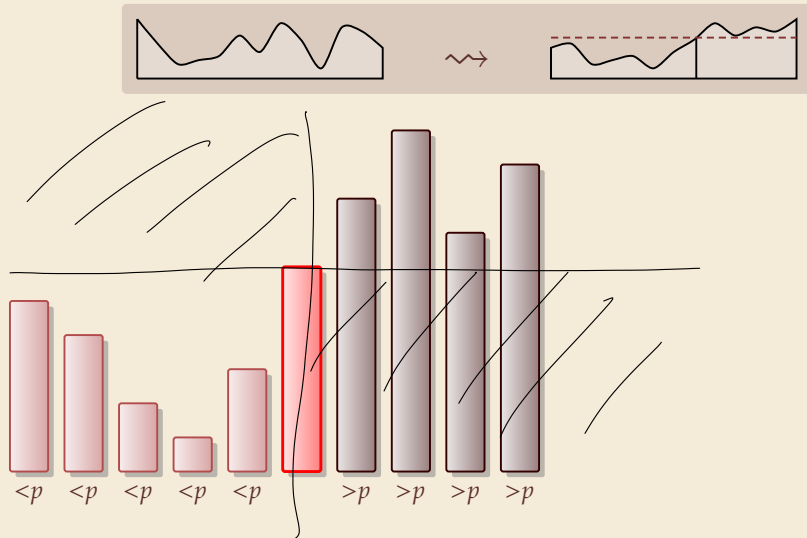
Partitioning around a pivot



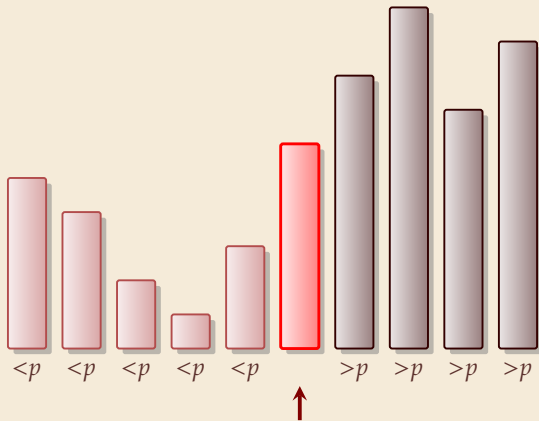
Partitioning around a pivot



Partitioning around a pivot



Partitioning around a pivot



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

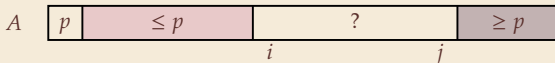
Partitioning – Detailed code

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  $b \in [0..n)$ 
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  $j \geq 1$  and  $A[j] > A[0]$ 
8     if  $i \geq j$  then break (goto 11)
9     else swap( $A[i], A[j]$ )
10  end while
11  swap( $A[0], A[j]$ )
12  return  $j$ 
```

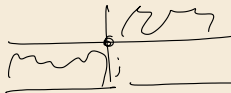
Loop invariant (5–10):



Quicksort

initial call
 $\text{quicksort}(A[0..n])$

```
1 procedure quicksort( $A[l..r]$ )  
2   if  $r - l \leq 1$  then return  
3    $b := \text{choosePivot}(A[l..r])$   
4    $j := \text{partition}(A[l..r], b)$   
5   quicksort( $A[l..j]$ )  
6   quicksort( $A[j + 1..r]$ )
```



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

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/comp526

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/comp526

Quicksort & Binary Search Trees

Quicksort

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

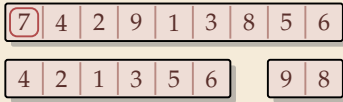
Quicksort & Binary Search Trees

Quicksort

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

Quicksort & Binary Search Trees

Quicksort



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

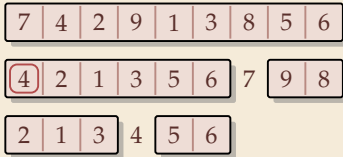
Quicksort

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

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

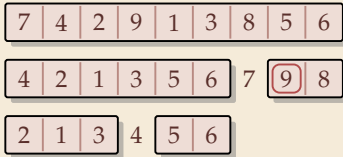
Quicksort & Binary Search Trees

Quicksort



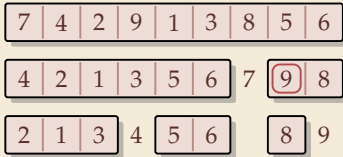
Quicksort & Binary Search Trees

Quicksort



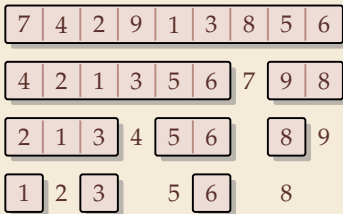
Quicksort & Binary Search Trees

Quicksort



Quicksort & Binary Search Trees

Quicksort



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

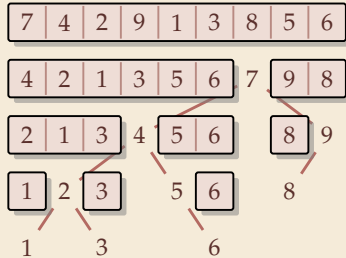
Quicksort & Binary Search Trees

Quicksort



Quicksort & Binary Search Trees

Quicksort

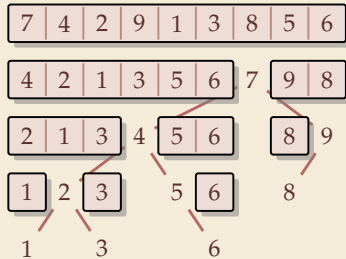


Binary Search Tree (BST)

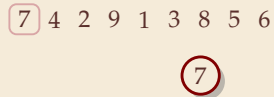
7 4 2 9 1 3 8 5 6

Quicksort & Binary Search Trees

Quicksort



Binary Search Tree (BST)



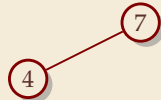
Quicksort & Binary Search Trees

Quicksort



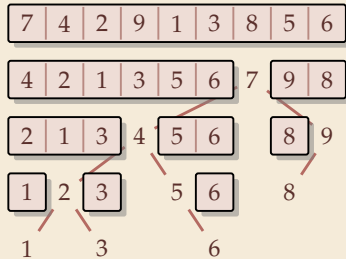
Binary Search Tree (BST)

4 2 9 1 3 8 5 6

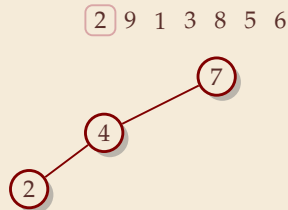


Quicksort & Binary Search Trees

Quicksort



Binary Search Tree (BST)

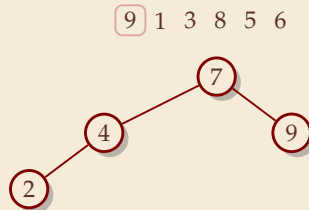


Quicksort & Binary Search Trees

Quicksort

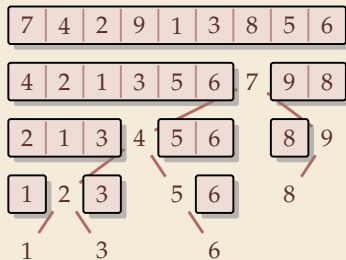


Binary Search Tree (BST)

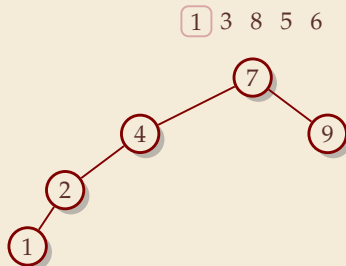


Quicksort & Binary Search Trees

Quicksort

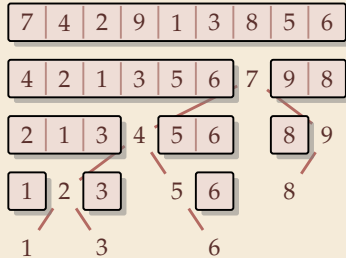


Binary Search Tree (BST)

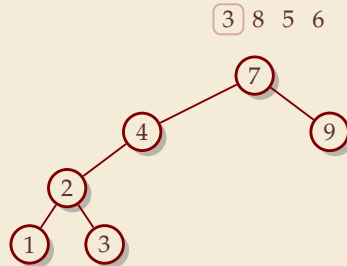


Quicksort & Binary Search Trees

Quicksort



Binary Search Tree (BST)

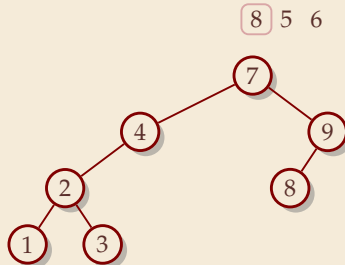


Quicksort & Binary Search Trees

Quicksort

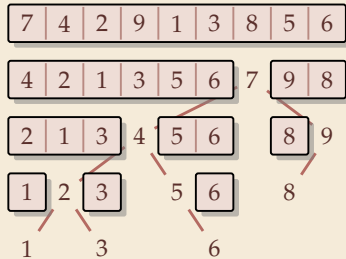


Binary Search Tree (BST)

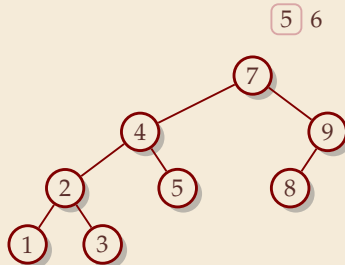


Quicksort & Binary Search Trees

Quicksort

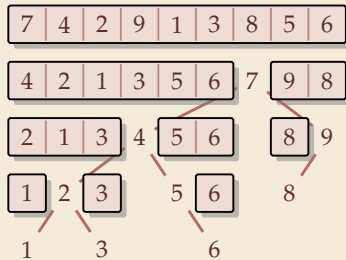


Binary Search Tree (BST)

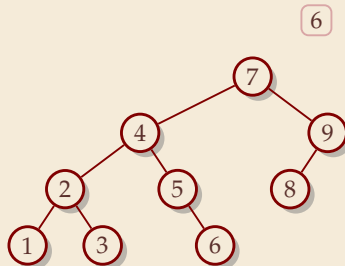


Quicksort & Binary Search Trees

Quicksort

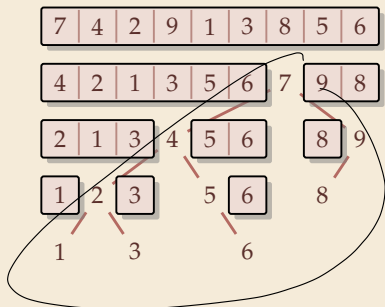


Binary Search Tree (BST)

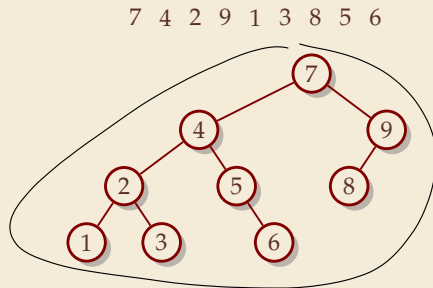


Quicksort & Binary Search Trees

Quicksort



Binary Search Tree (BST)



► recursion tree of quicksort = binary search tree from successive insertion

► comparisons in quicksort = comparisons to build BST

► comparisons in quicksort \approx comparisons to search each element in BST

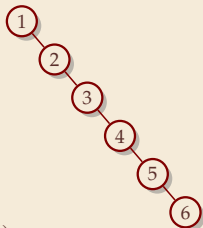
Quicksort – Worst Case

► Problem: BSTs can degenerate

► Cost to search for k is $k - 1$

~> Total cost $\underbrace{\sum_{k=1}^n (k - 1)} = \frac{n(n - 1)}{2} \sim \frac{1}{2}n^2$

~> 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

- ▶ $C(n)$ = element visits (as for mergesort)

↪ quicksort needs $\sim 2 \ln(2) \cdot n \lg n \approx \underline{1.39n \lg n}$ *in expectation*

- ▶ 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



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.

3.3 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
≈ *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:
 - ▶ comparing two elements
 - ▶ moving elements around (e. g. copying, swapping)
 - ▶ Cost: number of these operations.

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 these operations.
 - ▶ 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 $a \leq b$
- ▶ moving elements around (e. g. copying, swapping)
- ▶ Cost: number of these operations.

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.

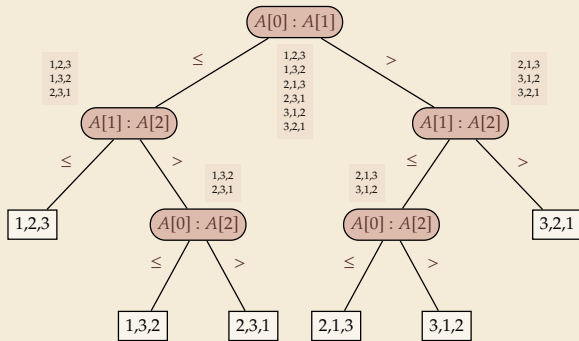


Every comparison-based sorting algorithm corresponds to a *decision tree*.

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

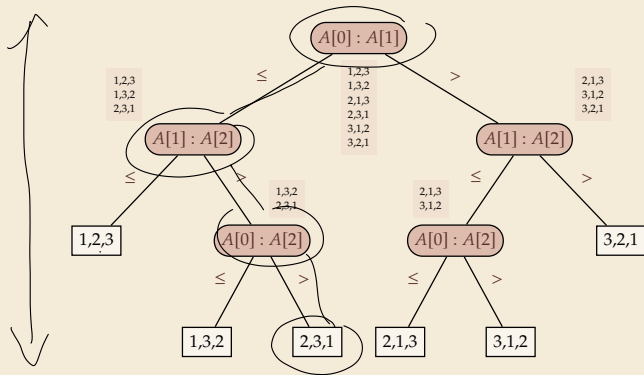
Comparison Lower Bound

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

~ ℓ leaves ~ height $\geq \lceil \lg(\ell) \rceil$

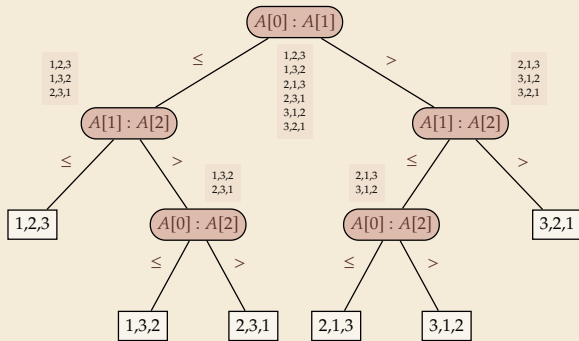
► comparison trees for sorting method must have $\geq n!$ leaves

~ height $\geq \lg(n!) \sim n \lg n$

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

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

↪ ℓ leaves ↪ height $\geq \lceil \lg(\ell) \rceil$

► comparison trees for sorting method must have $\geq n!$ leaves

↪ height $\geq \lg(n!) \sim \underline{n \lg n}$

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?

≈ 1.4427

Clicker Question



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

A Yes

B No



→ *sli.do/comp526*

Clicker Question



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

A Yes ✓

B ~~No~~



→ *sli.do/comp526*