

Design and Analysis  
of Algorithms I

# Linear-Time Selection

---

## Deterministic Selection (Algorithm)

# The Problem

Input: array A with  $n$  **distinct** numbers  
and a number  $i \in \{1, 2, \dots, n\}$ .

for simplicity

Output:  $i^{\text{th}}$  order statistic (i.e.,  $i^{\text{th}}$  smallest element of input array)

Example: median.

( $i = \frac{n+1}{2}$  for  $n$  odd,  
 $i = \frac{n}{2}$  for  $n$  even)

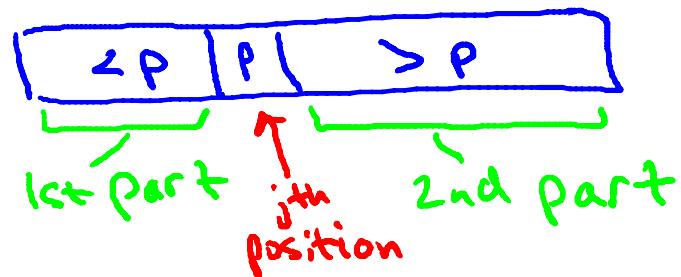
10	8	2	4
----	---	---	---

3<sup>rd</sup> order statistic

# Randomized Selection

RSelect (array A, length n, order statistic i)

- ⑥ if  $n=1$  return  $A[1]$
- ① choose pivot  $p$  from A uniformly at random
- ② partition A around  $p$   
let  $j = \text{new index of } p$
- ③ if  $j=i$  return  $p$
- ④ if  $j > i$  return RSelect (1st part of A,  $j-1$ ,  $i$ )
- ⑤ [if  $j < i$ ] return RSelect (2nd part of A,  $n-j$ ,  $i-j$ )



# Guaranteeing a Good Pivot

Recall: "best" pivot = the median! (seems circular!)

Goal: find pivot guaranteed to be pretty good.

Key idea: use "median of medians"!

# A Deterministic ChoosePivot

ChoosePivot( $A, n$ )

- logically break  $A$  into  $n/5$  groups of size 5 each
- sort each group (e.g., using MergeSort)
- copy  $n/5$  medians (i.e., middle element of each sorted group) into new array  $C$
- recursively compute median of  $C$  (!)
- return this as pivot

# The DSelect Algorithm

DSelect(array A, length n, order statistic i)

1. Break A into groups of 5, sort each group
2. C = the  $n/5$  “middle elements”
3.  $p = \text{DSelect}(C, n/5, n/10)$  [recursively computes median of C]
4. Partition A around p
5. If  $j = i$  return p
6. If  $j < i$  return DSelect( $1^{\text{st}}$  part of A,  $j-1$ ,  $i$ )
7. [else if  $j > i$ ] return DSelect( $2^{\text{nd}}$  part of A,  $n-j$ ,  $i-j$ )

choose pivot

Same  
as  
before

How many recursive calls does DSelect make?

- 0
- 1
- 2
- 3

# Running Time of DSelect

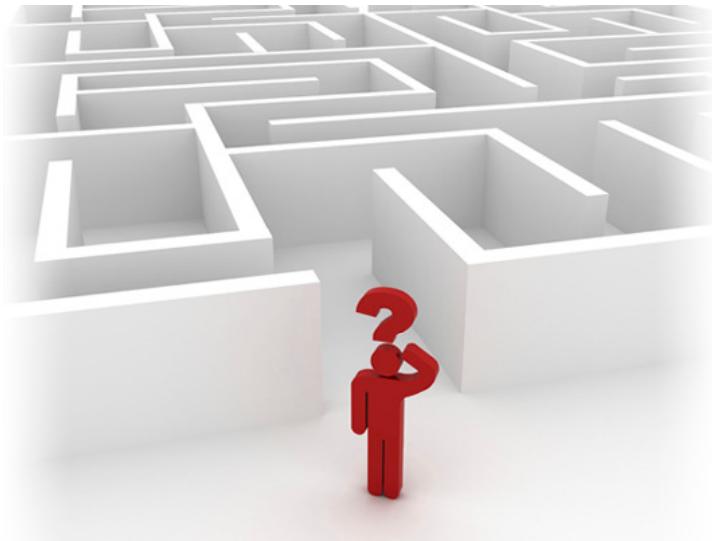
DSelect Theorem: for every input array of length  $n$ , DSelect runs in  $O(n)$  time.

Warning: not as good as QSelect in practice

- ① worse const ants
- ② not-in-place

History: from 1973.

Blum - Floyd - Pratt - Rivest - Tarjan  
1973      1970      1973      1986



Design and Analysis  
of Algorithms I

# Linear-Time Selection

---

## Deterministic Selection (Analysis I)

# The DSelect Algorithm

DSelect(array A, length n, order statistic i)

1. Break A into groups of 5, sort each group

2. C = the  $n/5$  “middle elements”

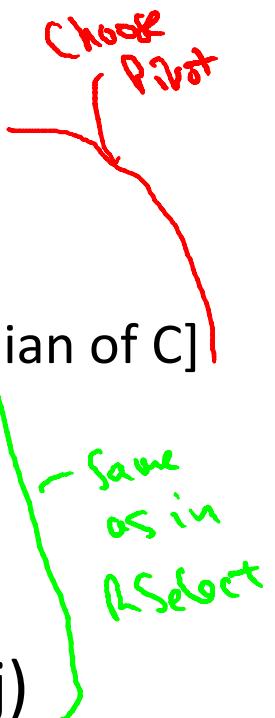
3.  $p = \text{DSelect}(C, n/5, n/10)$  [recursively computes median of C]

4. Partition A around p

5. If  $j = i$  return p

6. If  $j < i$  return DSelect( $1^{\text{st}}$  part of A,  $j-1$ ,  $i$ )

7. [else if  $j > i$ ] return DSelect( $2^{\text{nd}}$  part of A,  $n-j$ ,  $i-j$ )



What is the asymptotic running time of step 1 of the DSelect algorithm?

- $\theta(1)$
- $\theta(\log n)$
- $\theta(n)$
- $\theta(n \log n)$

Note: sorting an array with 5 elements takes  $\leq 120$  operations.

[Why 120? take  $m=5$  in our  $6m(\log_2 m + 1)$  bound for MergeSort]

$$\text{* groups} \quad \xrightarrow{6 \cdot 5 \cdot (\log_2 5 + 1)} \leq 120$$

So:  $\leq \frac{n}{5} \cdot 120 = 24n = O(n)$  for all groups

# The DSelect Algorithm

DSelect(array A, length n, order statistic i)

1. Break A into groups of 5, sort each group  $\rightarrow \Theta(n)$
2. C = the  $n/5$  “middle elements”  $\rightarrow \Theta(n)$
3. p = DSelect(C,  $n/5$ ,  $n/10$ ) [recursively computes median of C]
4. Partition A around p  $\rightarrow \Theta(n) \rightarrow T(\frac{n}{5})$
5. If  $j = i$  return p
6. If  $j < i$  return DSelect(1<sup>st</sup> part of A,  $j-1$ , i)  $\rightarrow T(\cdot)$
7. [else if  $j > i$ ] return DSelect(2nd part of A,  $n-j$ ,  $i-j$ )  $\rightarrow T(\cdot)$

# The Key Lemma

Key lemma: 2nd recursive call (in line 6 or 7) guaranteed to be on an array of size  $\leq \frac{7}{10}n$  (roughly).

Upshot: Can replace "?" by " $\frac{7}{10}n$ ".

Rough Proof: Let  $k = \frac{n}{5} = \# \text{of groups}$ .

Let  $x_i = i^{\text{th}}$  smallest of the  $k$  "middle elements".

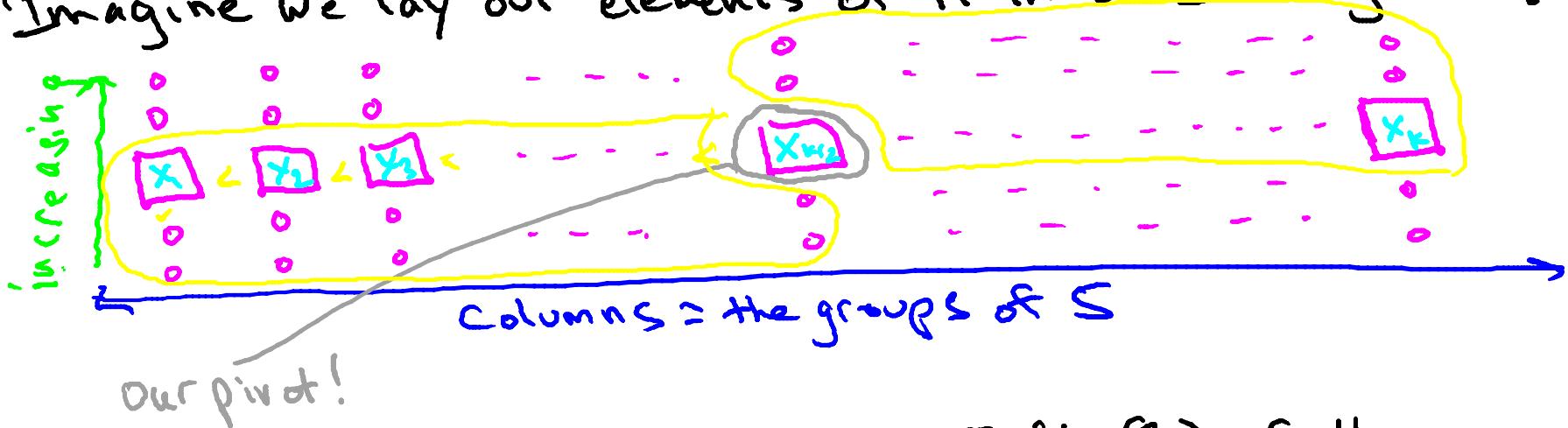
[So  $\text{first} = x_{k+1}$ ]

Goal: 7,30% of input array smaller than  $x_{k+1}$ ,  
7,30% is bigger.

# Rough Proof of Key Lemma

Thought experiment:

Imagine we lay out elements of  $A$  in a 2-D grid:



Key point:  $X_{k,1}$  bigger than 3 out 5 (60%) of the elements in  $\approx 50\%$  of the groups  
 $\Rightarrow$  bigger than 30% of  $A$  (Similarly, smaller than 30% of  $A$ )

Tim Roughgarden

# Example

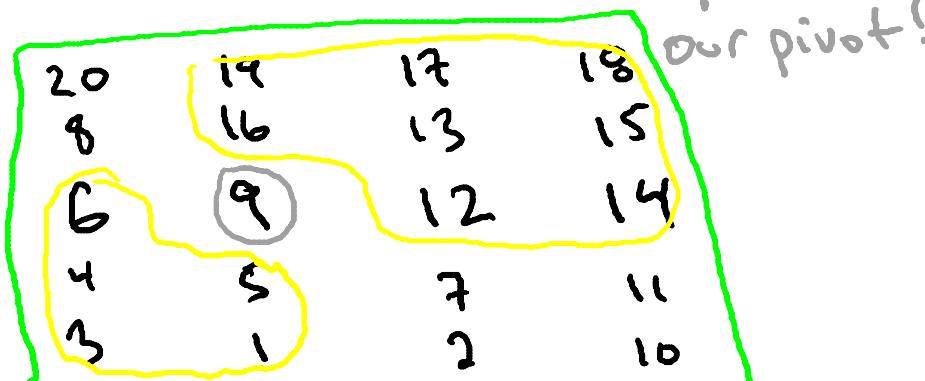
Input

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

After  
sorting  
groups of 5

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

The  
grid:



Tim Roughgarden



Design and Analysis  
of Algorithms I

# Linear-Time Selection

---

## Deterministic Selection (Analysis II)

# Rough Recurrence (Revisited)

Let  $T(n)$  = maximum running time of DSelect on an input array of length  $n$ .

There is a constant  $c \geq 1$  such that :

$$\textcircled{1} \quad T(1) = 1$$

$$\textcircled{2} \quad T(n) \leq \underbrace{cn}_{\text{Sorting the groups, Partition}} + \underbrace{T\left(\frac{n}{5}\right)}_{\text{recursive call in line 3}} + \underbrace{T(?)} \leq \frac{T(n)}{6} \quad \begin{array}{l} \text{by Key Lemma} \\ \text{recursive call in line 6 or 7} \end{array}$$

# Rough Recurrence (Revisited)

$$T(1) = 1, T(n) \leq cn + T(\frac{n}{5}) + T(\frac{7}{10}n)$$

constant  $c > 1$

Note: different-sized subproblems  $\Rightarrow$  can't use Master method!

Strategy: "hope and check".

Hope: there is some constant  $a$  {independent of  $n$ } such that  $T(n) \leq an$   $\forall n \geq 1$ .

[if true, then  $T(n) = O(n)$  and algorithm is linear-time]

# Analysis of Rough Recurrence

Claim: Let  $a = 10c$ .

Then  $T(n) \leq an$  for all  $n \geq 1$ .

Proof: by induction on  $n$ .

Base case:  $T(1) = 1 \leq a \cdot 1$  (since  $a \geq 1$ )

Inductive Step:  $\{n > 1\}$  Inductive Hypothesis:  $T(k) \leq ak$   $\forall k < n$

$$\begin{aligned} \text{We have } T(n) &\stackrel{\text{given}}{\leq} cn + T(\frac{n}{5}) + T(\frac{7}{10}n) \\ &\stackrel{\text{ind hyp}}{\leq} cn + a(\frac{n}{5}) + a(\frac{7}{10}n) \\ &= n(c + \frac{9}{10}a) \stackrel{\text{choose } c}{=} an \end{aligned}$$

QED!

$$\begin{aligned} T(1) &= 1 \quad \text{constant}, 1 \\ T(n) &\leq cn + \\ &T(\frac{n}{5}) + \\ &T(\frac{7}{10}n) \end{aligned}$$