

# Lecture: Analysis of Algorithms (CS583 - 004)

Amarda Shehu

Spring 2019

## 1 Outline of Today's Class

## 2 Order Statistics

- Selection of Order Statistics in Expected Linear Time
  - Randomized Divide and Conquer
- Selection of Order Statistics in Worst-case Linear Time
  - Median of Medians
  - Analysis of Worst-case Running Time
- Order Statistics: Conclusions

# Order Statistics

## Some Order Statistics We Know

Select the  $i^{\text{th}}$  smallest of  $n$  elements (the element with rank  $i$ ):

- $i = 1$ : minimum
- $i = n$ : maximum
- $i = (n + 1)/2$ : median

# Order Statistics

## Some Order Statistics We Know

Select the  $i^{\text{th}}$  smallest of  $n$  elements (the element with rank  $i$ ):

- $i = 1$ : minimum
- $i = n$ : maximum
- $i = (n + 1)/2$ : median

Design a simple algorithm to find the element with rank  $i$

# Order Statistics

## Some Order Statistics We Know

Select the  $i^{\text{th}}$  smallest of  $n$  elements (the element with rank  $i$ ):

- $i = 1$ : minimum
- $i = n$ : maximum
- $i = (n + 1)/2$ : median

Design a simple algorithm to find the element with rank  $i$

# Order Statistics

## Some Order Statistics We Know

Select the  $i^{\text{th}}$  smallest of  $n$  elements (the element with rank  $i$ ):

- $i = 1$ : minimum
- $i = n$ : maximum
- $i = (n + 1)/2$ : median

Design a simple algorithm to find the element with rank  $i$

Naive algorithm: Sort and index  $i^{\text{th}}$  element.

# Order Statistics

## Some Order Statistics We Know

Select the  $i^{\text{th}}$  smallest of  $n$  elements (the element with rank  $i$ ):

- $i = 1$ : minimum
- $i = n$ : maximum
- $i = (n + 1)/2$ : median

Design a simple algorithm to find the element with rank  $i$

**Naive algorithm:** Sort and index  $i^{\text{th}}$  element.

$$\begin{aligned}\text{Worst-case running time} &\in \theta(n \cdot \lg n) + \theta(1) \\ &\in \theta(n \cdot \lg n)\end{aligned}$$

# Order Statistics

## Some Order Statistics We Know

Select the  $i^{\text{th}}$  smallest of  $n$  elements (the element with rank  $i$ ):

- $i = 1$ : minimum
- $i = n$ : maximum
- $i = (n + 1)/2$ : median

Design a simple algorithm to find the element with rank  $i$

**Naive algorithm:** Sort and index  $i^{\text{th}}$  element.

Worst-case running time  $\in \theta(n \cdot \lg n) + \theta(1)$

$\in \theta(n \cdot \lg n)$

...use mergesort or heapsort (not quicksort)



# Order Statistics

## Some Order Statistics We Know

Select the  $i^{\text{th}}$  smallest of  $n$  elements (the element with rank  $i$ ):

- $i = 1$ : minimum
- $i = n$ : maximum
- $i = (n + 1)/2$ : median

Design a simple algorithm to find the element with rank  $i$

**Naive algorithm:** Sort and index  $i^{\text{th}}$  element.

Worst-case running time  $\in \theta(n \cdot \lg n) + \theta(1)$

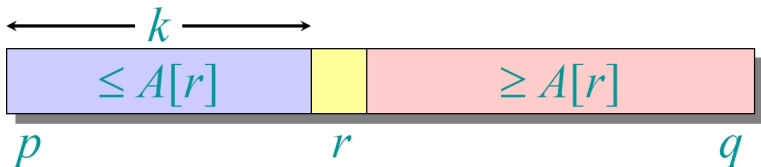
$\in \theta(n \cdot \lg n)$

...use mergesort or heapsort (not quicksort)

## Randomized Divide and Conquer Algorithm

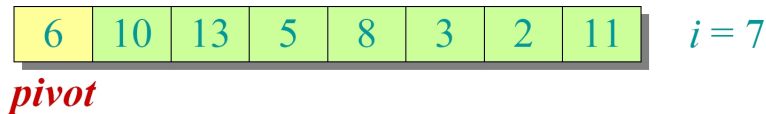
**RAND-SELECT(array A, p, q, i)**     $\triangleright i^{\text{th}}$  smallest of  $A[p \dots q]$

- 1: **if**  $p = q$  **then**
- 2:     **return**  $A[p]$
- 3:  $r \leftarrow$  RAND-PARTITION( $A, p, q$ )
- 4:  $k \leftarrow r - p + 1$      $\triangleright k = \text{rank}(A[r])$
- 5: **if**  $i = k$  **then**
- 6:     **return**  $A[r]$
- 7: **if**  $i < k$  **then**
- 8:     **return** RAND-SELECT( $A, p, r - 1, i$ )
- 9: **else return** RAND-SELECT( $A, r + 1, q, i - k$ )

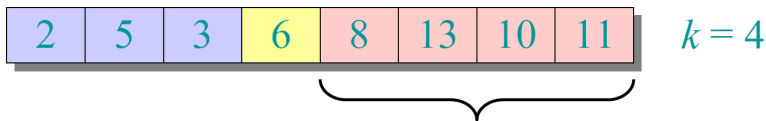


## Randomized Select: Trace

Select the  $i = 7^{\text{th}}$  smallest element from the array below:



Partition:



Now select the  $7 - 4 = 3^{\text{rd}}$  smallest element recursively.

# Randomized Select: Running Time Analysis

- Analysis follows closely that of quicksort
- For simplicity, we will assume that all elements are distinct
- We will first gain intuition through lucky/unlucky scenarios

**Lucky:** [assume a 1 : 9 partition after RAND-PARTITION]

$$\begin{aligned} T(n) &= T(9n/10) + \theta(n) & n^{\log_{10/9}(1)} = n^0 = 1 \\ &= \theta(n) & \text{CASE 3 of master theorem} \end{aligned}$$

**Unlucky:** [assume one side of the partitioned array is empty]

$$\begin{aligned} T(n) &= T(n-1) + \theta(n) & \text{arithmetic series} \\ &= \theta(n^2) & \text{worse than sorting!!!} \end{aligned}$$

# Randomized Select: Analysis of Expected Time

- Analysis similar to randomized quicksort
- Let  $T(n)$  be the random variable for the running time of RAND-SELECT on an input of size  $n$ , assuming random numbers are independent
- To obtain upper bound, assume the  $i^{\text{th}}$  smallest element always falls on the larger side of the partition:

$$T(n) = \begin{cases} T(\max\{0, n-1\}) + \theta(n) & \text{if } 0:n-1 \text{ split} \\ T(\max\{1, n-2\}) + \theta(n) & \text{if } 1:n-2 \text{ split} \\ \dots & \\ T(\max\{n-1, 0\}) + \theta(n) & \text{if } n-1:0 \text{ split} \end{cases}$$

- Summing up we have:

$$E[T(n)] = \frac{1}{n} \sum_{k=0}^{n-1} E[T(\max\{k, n-k-1\}) + \theta(n)]$$

## Randomized Select: Those pesky expectations...

$$\begin{aligned} E[T(n)] &= \frac{1}{n} \sum_{k=0}^{n-1} E[T(\max\{k, n-k-1\})] + \theta(n) \\ &\qquad \qquad \qquad \text{get } \theta(n) \text{ outside} \\ &= \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} E[T(k)] + \theta(n) \\ &\qquad \qquad \qquad \text{upper terms appear twice} \end{aligned}$$

## Randomized Select: Those pesky expectations...

$$\begin{aligned} E[T(n)] &= \frac{1}{n} \sum_{k=0}^{n-1} E[T(\max\{k, n-k-1\})] + \theta(n) \\ &\qquad \qquad \qquad \text{get } \theta(n) \text{ outside} \\ &= \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} E[T(k)] + \theta(n) \\ &\qquad \qquad \qquad \text{upper terms appear twice} \end{aligned}$$

**Prove:**  $E[T(n)] \leq c \cdot n$  for  $c > 0$  ( $c$  large enough for base cases)

## Randomized Select: Those pesky expectations...

$$\begin{aligned} E[T(n)] &= \frac{1}{n} \sum_{k=0}^{n-1} E[T(\max\{k, n-k-1\})] + \theta(n) \\ &\qquad\qquad\qquad \text{get } \theta(n) \text{ outside} \\ &= \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} E[T(k)] + \theta(n) \\ &\qquad\qquad\qquad \text{upper terms appear twice} \end{aligned}$$

**Prove:**  $E[T(n)] \leq c \cdot n$  for  $c > 0$  ( $c$  large enough for base cases)

**Assume:**  $E[T(k)] \leq c \cdot k$ , where  $k < n$



## Randomized Select: Those pesky expectations...

$$\begin{aligned} E[T(n)] &= \frac{1}{n} \sum_{k=0}^{n-1} E[T(\max\{k, n-k-1\})] + \theta(n) \\ &\qquad \qquad \qquad \text{get } \theta(n) \text{ outside} \\ &= \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} E[T(k)] + \theta(n) \\ &\qquad \qquad \qquad \text{upper terms appear twice} \end{aligned}$$

**Prove:**  $E[T(n)] \leq c \cdot n$  for  $c > 0$  ( $c$  large enough for base cases)

**Assume:**  $E[T(k)] \leq c \cdot k$ , where  $k < n$

**Then:**  $E[T(n)] \leq \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} c \cdot k + \theta(n)$

## Randomized Select: Those pesky expectations...

$$\begin{aligned}
 E[T(n)] &= \frac{1}{n} \sum_{k=0}^{n-1} E[T(\max\{k, n-k-1\})] + \theta(n) \\
 &\qquad \qquad \qquad \text{get } \theta(n) \text{ outside} \\
 &= \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} E[T(k)] + \theta(n) \\
 &\qquad \qquad \qquad \text{upper terms appear twice}
 \end{aligned}$$

**Prove:**  $E[T(n)] \leq c \cdot n$  for  $c > 0$  ( $c$  large enough for base cases)

**Assume:**  $E[T(k)] \leq c \cdot k$ , where  $k < n$

**Then:**  $E[T(n)] \leq \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} c \cdot k + \theta(n)$   
 $\sum_{k=\lfloor n/2 \rfloor}^{n-1} k \leq \frac{3}{8} n^2$  (exercise: show it)

## Randomized Select: Those pesky expectations...

$$\begin{aligned}
 E[T(n)] &= \frac{1}{n} \sum_{k=0}^{n-1} E[T(\max\{k, n-k-1\})] + \theta(n) \\
 &\qquad \qquad \qquad \text{get } \theta(n) \text{ outside} \\
 &= \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} E[T(k)] + \theta(n) \\
 &\qquad \qquad \qquad \text{upper terms appear twice}
 \end{aligned}$$

**Prove:**  $E[T(n)] \leq c \cdot n$  for  $c > 0$  ( $c$  large enough for base cases)

**Assume:**  $E[T(k)] \leq c \cdot k$ , where  $k < n$

**Then:**  $E[T(n)] \leq \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} c \cdot k + \theta(n)$   
 $\sum_{k=\lfloor n/2 \rfloor}^{n-1} k \leq \frac{3}{8} n^2$  (exercise: show it)

**So:**  $E[T(n)] \leq cn - (\frac{cn}{4} - \theta(n)) \leq cn$  if  $\frac{cn}{4} - \theta(n) \geq 0$

## Randomized Select: Those pesky expectations...

$$\begin{aligned}
 E[T(n)] &= \frac{1}{n} \sum_{k=0}^{n-1} E[T(\max\{k, n-k-1\})] + \theta(n) \\
 &\qquad \qquad \qquad \text{get } \theta(n) \text{ outside} \\
 &= \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} E[T(k)] + \theta(n) \\
 &\qquad \qquad \qquad \text{upper terms appear twice}
 \end{aligned}$$

**Prove:**  $E[T(n)] \leq c \cdot n$  for  $c > 0$  ( $c$  large enough for base cases)

**Assume:**  $E[T(k)] \leq c \cdot k$ , where  $k < n$

**Then:**  $E[T(n)] \leq \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} c \cdot k + \theta(n)$   
 $\sum_{k=\lfloor n/2 \rfloor}^{n-1} k \leq \frac{3}{8} n^2$  (exercise: show it)

**So:**  $E[T(n)] \leq cn - (\frac{cn}{4} - \theta(n)) \leq cn$  if  $\frac{cn}{4} - \theta(n) \geq 0$

## Randomized Select: Those pesky expectations...

$$\begin{aligned}
 E[T(n)] &= \frac{1}{n} \sum_{k=0}^{n-1} E[T(\max\{k, n-k-1\})] + \theta(n) \\
 &\qquad\qquad\qquad \text{get } \theta(n) \text{ outside} \\
 &= \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} E[T(k)] + \theta(n) \\
 &\qquad\qquad\qquad \text{upper terms appear twice}
 \end{aligned}$$

**Prove:**  $E[T(n)] \leq c \cdot n$  for  $c > 0$  ( $c$  large enough for base cases)

**Assume:**  $E[T(k)] \leq c \cdot k$ , where  $k < n$

**Then:**  $E[T(n)] \leq \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} c \cdot k + \theta(n)$   
 $\sum_{k=\lfloor n/2 \rfloor}^{n-1} k \leq \frac{3}{8} n^2$  (exercise: show it)

**So:**  $E[T(n)] \leq cn - (\frac{cn}{4} - \theta(n)) \leq cn$  if  $\frac{cn}{4} - \theta(n) \geq 0$

Easy to find a large value of  $c$  such that  $\frac{cn}{4}$  dominates  $\theta(n)$

## Randomized Select: Those pesky expectations...

$$\begin{aligned}
 E[T(n)] &= \frac{1}{n} \sum_{k=0}^{n-1} E[T(\max\{k, n-k-1\})] + \theta(n) \\
 &\qquad\qquad\qquad \text{get } \theta(n) \text{ outside} \\
 &= \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} E[T(k)] + \theta(n) \\
 &\qquad\qquad\qquad \text{upper terms appear twice}
 \end{aligned}$$

**Prove:**  $E[T(n)] \leq c \cdot n$  for  $c > 0$  ( $c$  large enough for base cases)

**Assume:**  $E[T(k)] \leq c \cdot k$ , where  $k < n$

**Then:**  $E[T(n)] \leq \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} c \cdot k + \theta(n)$   
 $\sum_{k=\lfloor n/2 \rfloor}^{n-1} k \leq \frac{3}{8} n^2$  (exercise: show it)

**So:**  $E[T(n)] \leq cn - (\frac{cn}{4} - \theta(n)) \leq cn$  if  $\frac{cn}{4} - \theta(n) \geq 0$

Easy to find a large value of  $c$  such that  $\frac{cn}{4}$  dominates  $\theta(n)$

## Randomized Select: Summary

- Works fast in the average case: linear expected time
- Very simple and fast algorithm in practice
- *But*, worst-case behavior is  $\theta(n^2)$
- **Question:** Is there an algorithm that runs in linear time even in the worst case?
- **Answer:** Yes - in 1973, Blum, Floyd, Pratt, and Rivest designed such an algorithm
- **Basic Idea:** Generate good pivots recursively to guarantee a good split

# Worst-case Linear-time Order Statistics

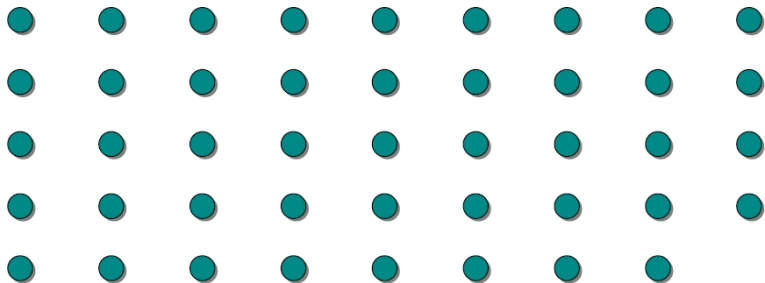
## SELECT(*i*,*n*)

- 1: Divide the  $n$  elements into groups of 5. Find the median of each 5-element group by rote.
- 2: Recursively SELECT the median  $x$  of the  $\lfloor n/5 \rfloor$  group medians to be the pivot
- 3: Partition around the pivot. Let  $k = \text{rank}(x)$
- 4: **if**  $i = k$  **then**
- 5:     **return**  $x$
- 6: **if**  $i < k$  **then**
- 7:     recursively SELECT  $i^{\text{th}}$  smallest element in lower part
- 8: **if**  $i > k$  **then**
- 9:     recursively SELECT  $(i - k)^{\text{th}}$  smallest element in upper part

Note: lines 3.-9. are the same as in RAND-SELECT

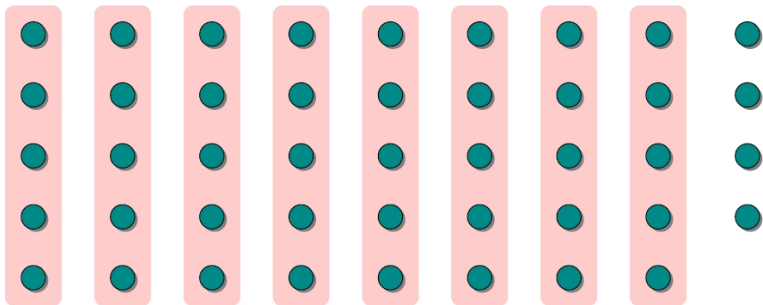


## SELECT: Choosing the Pivot



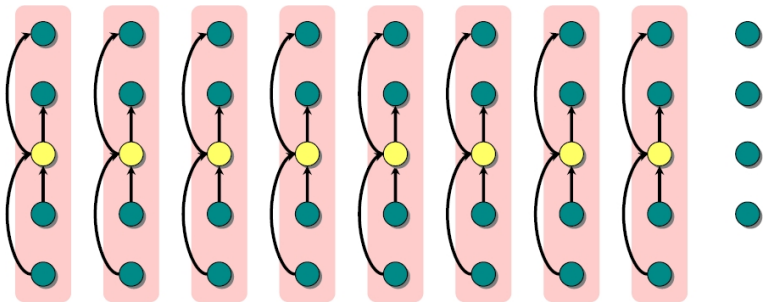
Here is the input:  $n$  elements.

## SELECT: Choosing the Pivot



- 1 Divide the  $n$  elements into groups of 5.

## SELECT: Choosing the Pivot



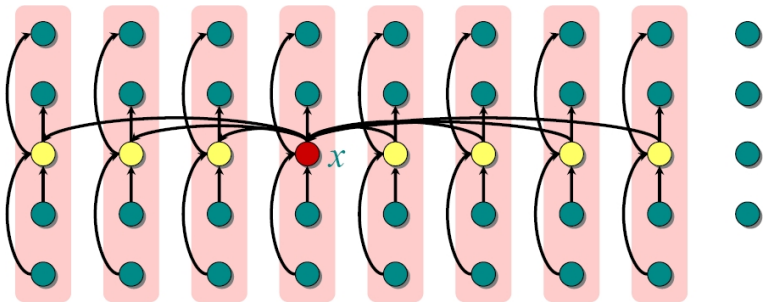
- 1 Divide the  $n$  elements into groups of 5. Find the median of each 5-element group by rote.

*lesser*



*greater*

# SELECT: Choosing the Pivot



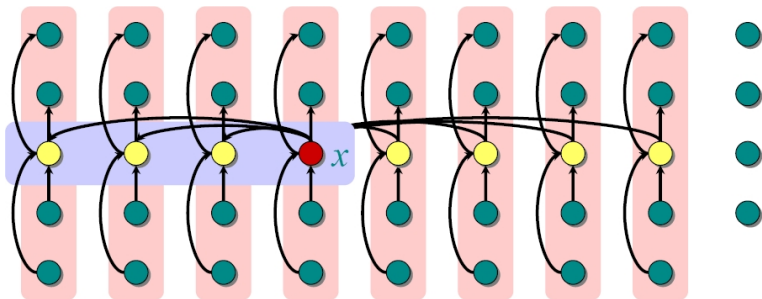
- 1 Divide the  $n$  elements into groups of 5. Find the median of each 5-element group by rote.
- 2 Recursively SELECT the median  $x$  of the  $\lfloor n/5 \rfloor$  group medians to be the pivot.

*lesser*



*greater*

## Select: Running Time Analysis



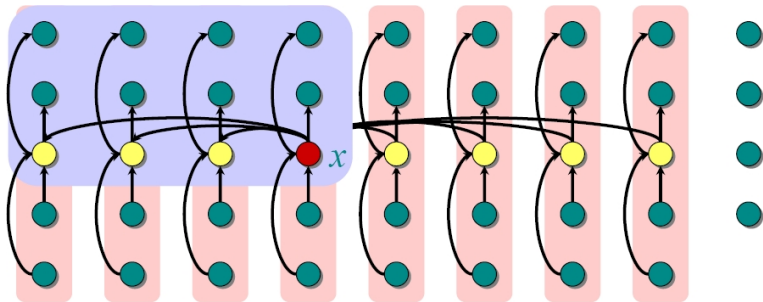
At least half of the group medians are  $\leq x$ ,  
which is at least  $\lfloor \lfloor n/5 \rfloor / 2 \rfloor = \lfloor n/10 \rfloor$  elements.

*lesser*



*greater*

## Select: Running Time Analysis



At least half the group medians are  $\leq x$ , which is at least  $\lfloor \lfloor n/5 \rfloor / 2 \rfloor = \lfloor n/10 \rfloor$  group medians.

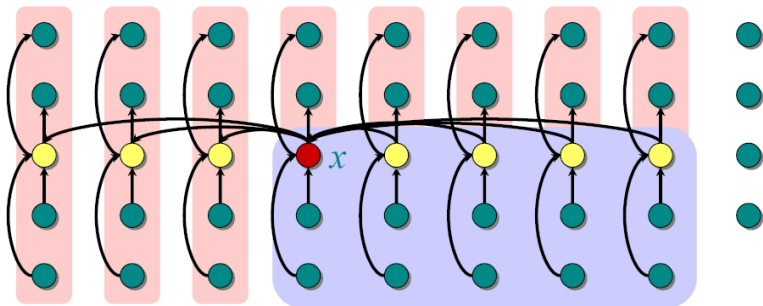
- If we assume that all elements are distinct, then there are  $3\lfloor n/10 \rfloor$  elements  $\leq x$ .

*lesser*



*greater*

## Select: Running Time Analysis



At least half the group medians are  $\leq x$ , which is at least  $\lfloor \lfloor n/5 \rfloor / 2 \rfloor = \lfloor n/10 \rfloor$  group medians.

- If we assume that all elements are distinct, then there are  $3\lfloor n/10 \rfloor$  elements  $\leq x$ .
- Similarly, at least  $3\lfloor n/10 \rfloor$  elements are  $\geq x$

*lesser*



*greater*

# Select: Running Time Analysis

- For  $n \geq 50$ , we have  $3\lfloor n/10 \rfloor \geq n/4$ . So, the call to SELECT in lines 4 and on is executed recursively on at most  $3n/4$  elements
- The recurrence for the running time can assume that lines 4 and on takes  $T(3n/4)$  in the worst case
- For  $n < 50$ , we know that the worst-case time is  $T(n) \in \theta(1)$

The recurrence is:  $T(n) = T(n/5) + \theta(n) + T(3n/4)$

## Breakdown:

- **Line 1:**  $\theta(n)$
- **Line 2:**  $T(n/5)$
- **Line 3:**  $\theta(n)$
- **Lines  $\geq 4$ :**  $T(3n/4)$

## Substitution:

$$\begin{aligned}
 T(n) &\leq \frac{1}{5}c \cdot n + \frac{3}{4}c \cdot n + \theta(n) \\
 &= \frac{19}{20}c \cdot n + \theta(n) \\
 &= c \cdot n - \left(\frac{1}{20}c \cdot n - \theta(n)\right) \\
 &\leq c \cdot n
 \end{aligned}$$

if  $c$  is large enough to dominate  $\theta(n)$



## Order Statistics: Conclusions

- Since the work at each level of recursion is a constant fraction ( $19/20$ ) smaller, the work per level is a geometric series dominated by the linear work at the root
- In practice, this algorithm runs slowly, because the constant in front of  $n$  is large
- The randomized algorithm is far more practical and simpler to implement
- **Exercise:** Why not divide into groups of 3?