# Lecture ~~7~~: Order Statistics

# ORDER STATISTICS

Select the $i$th smallest of $n$ elements.

- $i = 1$: *minimum*;
- $i = n$: *maximum*;
- $i = \lfloor (n+1)/2 \rfloor$ or $\lceil (n+1)/2 \rceil$: *median*.

***Naive algorithm***: Sort and index $i$th element.

Worst-case running time $= \Theta(n \lg n) + \Theta(1)$
$$= \Theta(n \lg n),$$

using merge sort or heapsort (*not* quicksort).

# RANK OF AN ELEMENT

Rank of an element is the position of the element in the sorted array

Given Numbers:

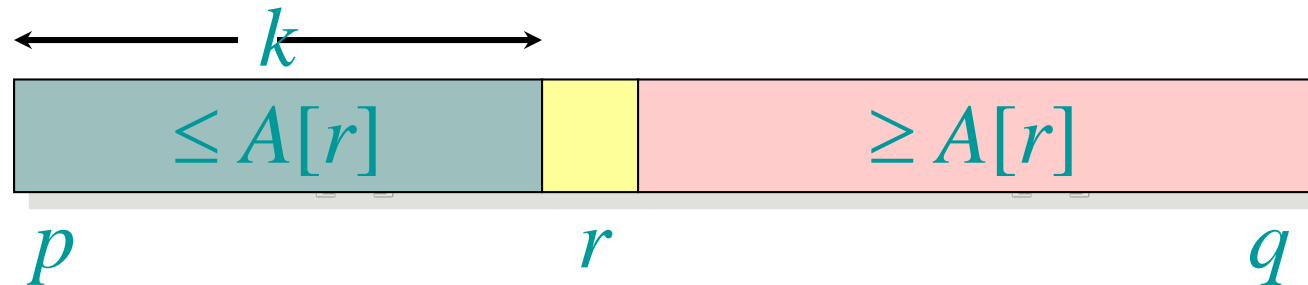| 10 | 13 | 5 | 8 | 3 | 2 | 11 |
|----|----|---|---|---|---|----|

Sorted Array:

| 2 | 3 | 5 | 8 | 10 | 11 | 13 |
|---|---|---|---|----|----|----|

*rank* (8)=4

Finding the $i$th smallest of $n$ elements = Finding an element with *rank i*.

# DIVIDE-AND-CONQUER ALGORITHM

$\text{SELECT}(A, p, q, i)$      ▷ $i$th smallest of $A[p \mathinner{.\,.} q]$

     **if** $p = q$ **then return** $A[p]$

     $r \leftarrow \text{PARTITION}(A, p, q)$

     $k \leftarrow r - p + 1$      \\ k=rank of the pivot

     **if** $i = k$ **then return** $A[r]$

     **if** $i < k$

        **then return** $\text{SELECT}(A, p, r - 1, i)$
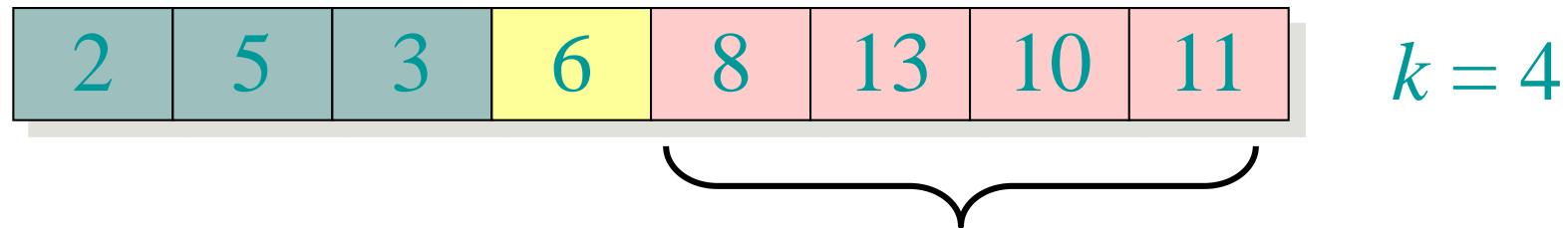
        **else return** $\text{SELECT}(A, r + 1, q, i - k)$

# EXAMPLE

Select the $i = 7$th smallest:

| 6 | 10 | 13 | 5 | 8 | 3 | 2 | 11 |
|---|----|----|---|---|---|---|----|

$i = 7$

*pivot*

Partition:

| 2 | 5 | 3 | 6 | 8 | 13 | 10 | 11 |
|---|---|---|---|---|----|----|----|

$k = 4$

Select the $7 - 4 = 3$rd smallest recursively.

# INTUITION FOR ANALYSIS

(All our analyses today assume that all elements are distinct.)

**Lucky:**

$$T(n) = T(9n/10) + \Theta(n)$$
$$= \Theta(n)$$

$$n^{\log_{10/9} 1} = n^0 = 1$$

CASE 3

**Unlucky:**
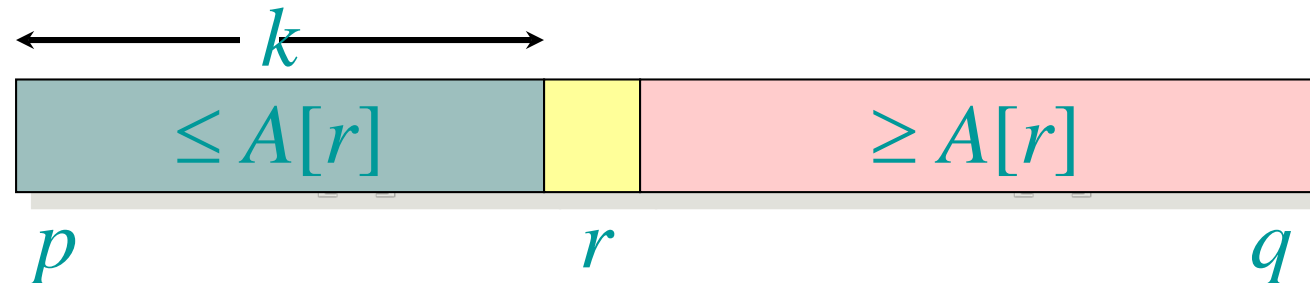
$$T(n) = T(n - 1) + \Theta(n)$$
$$= \Theta(n^2)$$

arithmetic series

*Worse than sorting!*

# RANDOMIZED DIVIDE-AND-CONQUER ALGORITHM

RAND-SELECT($A$, $p$, $q$, $i$)     ▷ $i$th smallest of $A[p \mathinner{.\,.} q]$

   **if** $p = q$ **then return** $A[p]$

   $r \leftarrow$ RAND-PARTITION($A$, $p$, $q$)

   $k \leftarrow r - p + 1$

   **if** $i = k$ **then return** $A[r]$

   **if** $i < k$

      **then return** RAND-SELECT($A$, $p$, $r - 1$, $i$)

      **else return** RAND-SELECT($A$, $r + 1$, $q$, $i - k$)

# ANALYSIS OF EXPECTED TIME

The analysis follows that of randomized quicksort, but it's a little different.

Let $T(n) =$ the random variable for the running time of RAND-SELECT on an input of size $n$, assuming random numbers are independent.

For $k = 0, 1, \ldots, n-1$, define the *indicator random variable*

$$X_k = \begin{cases} 1 & \text{if PARTITION generates a } k : n-k-1 \text{ split,} \\ 0 & \text{otherwise.} \end{cases}$$

# ANALYSIS (CONTINUED)

To obtain an upper bound, assume that the $i$th element always falls in the larger side of the partition:

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

$$= \sum_{k=0}^{n-1} X_k \left( T(\max\{k, n-k-1\}) + \Theta(n) \right).$$

# CALCULATING EXPECTATION

$$E[T(n)] = E\left[\sum_{k=0}^{n-1} X_k \left(T(\max\{k, n-k-1\}) + \Theta(n)\right)\right]$$

Take expectations of both sides.

# CALCULATING EXPECTATION

$$E[T(n)] = E\left[\sum_{k=0}^{n-1} X_k \left(T(\max\{k, n-k-1\}) + \Theta(n)\right)\right]$$

$$= \sum_{k=0}^{n-1} E\left[X_k \left(T(\max\{k, n-k-1\}) + \Theta(n)\right)\right]$$

Linearity of expectation.

# CALCULATING EXPECTATION

$$E[T(n)] = E\left[\sum_{k=0}^{n-1} X_k \left(T(\max\{k, n-k-1\}) + \Theta(n)\right)\right]$$

$$= \sum_{k=0}^{n-1} E\left[X_k \left(T(\max\{k, n-k-1\}) + \Theta(n)\right)\right]$$

$$= \sum_{k=0}^{n-1} E\left[X_k\right] \cdot E\left[T(\max\{k, n-k-1\}) + \Theta(n)\right]$$

Independence of $X_k$ from other random choices.

# CALCULATING EXPECTATION

$$E[T(n)] = E\left[\sum_{k=0}^{n-1} X_k \left(T(\max\{k, n-k-1\}) + \Theta(n)\right)\right]$$

$$= \sum_{k=0}^{n-1} E\left[X_k \left(T(\max\{k, n-k-1\}) + \Theta(n)\right)\right]$$

$$= \sum_{k=0}^{n-1} E\left[X_k\right] \cdot E\left[T(\max\{k, n-k-1\}) + \Theta(n)\right]$$

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

Linearity of expectation; $E[X_k] = 1/n$ .

# CALCULATING EXPECTATION

$$E[T(n)] = E\left[\sum_{k=0}^{n-1} X_k \left(T(\max\{k, n-k-1\}) + \Theta(n)\right)\right]$$

$$= \sum_{k=0}^{n-1} E\left[X_k \left(T(\max\{k, n-k-1\}) + \Theta(n)\right)\right]$$

$$= \sum_{k=0}^{n-1} E[X_k] \cdot E\left[T(\max\{k, n-k-1\}) + \Theta(n)\right]$$

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

$$\leq \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} E[T(k)] + \Theta(n)$$

Upper terms
appear twice.

# HAIRY RECURRENCE

(But not quite as hairy as the quicksort one.)

$$E[T(n)] = \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} E[T(k)] + \Theta(n)$$

**Prove:** $E[T(n)] \leq cn$ for constant $c > 0$.

• The constant $c$ can be chosen large enough so that $E[T(n)] \leq cn$ for the base cases.

**Use fact:** $\sum_{k=\lfloor n/2 \rfloor}^{n-1} k \leq \frac{3}{8}n^2$ (exercise).

# SUBSTITUTION METHOD

$$E[T(n)] \leq \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} ck + \Theta(n)$$

Substitute inductive hypothesis.

# SUBSTITUTION METHOD

$$E[T(n)] \leq \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} ck + \Theta(n)$$

$$\leq \frac{2c}{n} \left( \frac{3}{8} n^2 \right) + \Theta(n)$$

Use fact.

# SUBSTITUTION METHOD

$$E[T(n)] \leq \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} ck + \Theta(n)$$

$$\leq \frac{2c}{n} \left( \frac{3}{8} n^2 \right) + \Theta(n)$$

$$= cn - \left( \frac{cn}{4} - \Theta(n) \right)$$

Express as *desired* – *residual*.

# SUBSTITUTION METHOD

$$E[T(n)] \leq \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} ck + \Theta(n)$$

$$\leq \frac{2c}{n} \left( \frac{3}{8} n^2 \right) + \Theta(n)$$

$$= cn - \left( \frac{cn}{4} - \Theta(n) \right)$$

$$\leq cn,$$

if $c$ is chosen large enough so that $cn/4$ dominates the $\Theta(n)$.

# SUMMARY OF RANDOMIZED ORDER-STATISTIC SELECTION

- Works fast: linear expected time.
- Excellent algorithm in practice.
- But, the worst case is *very* bad: $\Theta(n^2)$.

*Q.* Is there an algorithm that runs in linear time in the worst case?

*A.* Yes, due to Blum, Floyd, Pratt, Rivest, and Tarjan [1973].

**IDEA:** Generate a good pivot recursively.

# 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 $x$. Let $k = \text{rank}(x)$.

4. **if** $i = k$ **then return** $x$
   **elseif** $i < k$
      **then** recursively SELECT the $i$th
         smallest element in the lower part
      **else** recursively SELECT the $(i–k)$th
         smallest element in the upper part

Same as
RAND-
SELECT

# CHOOSING THE PIVOT
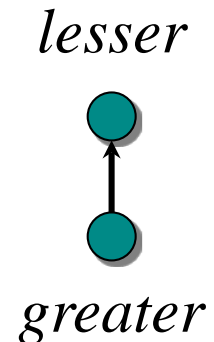
# CHOOSING THE PIVOT



1. Divide the $n$ elements into groups of 5.
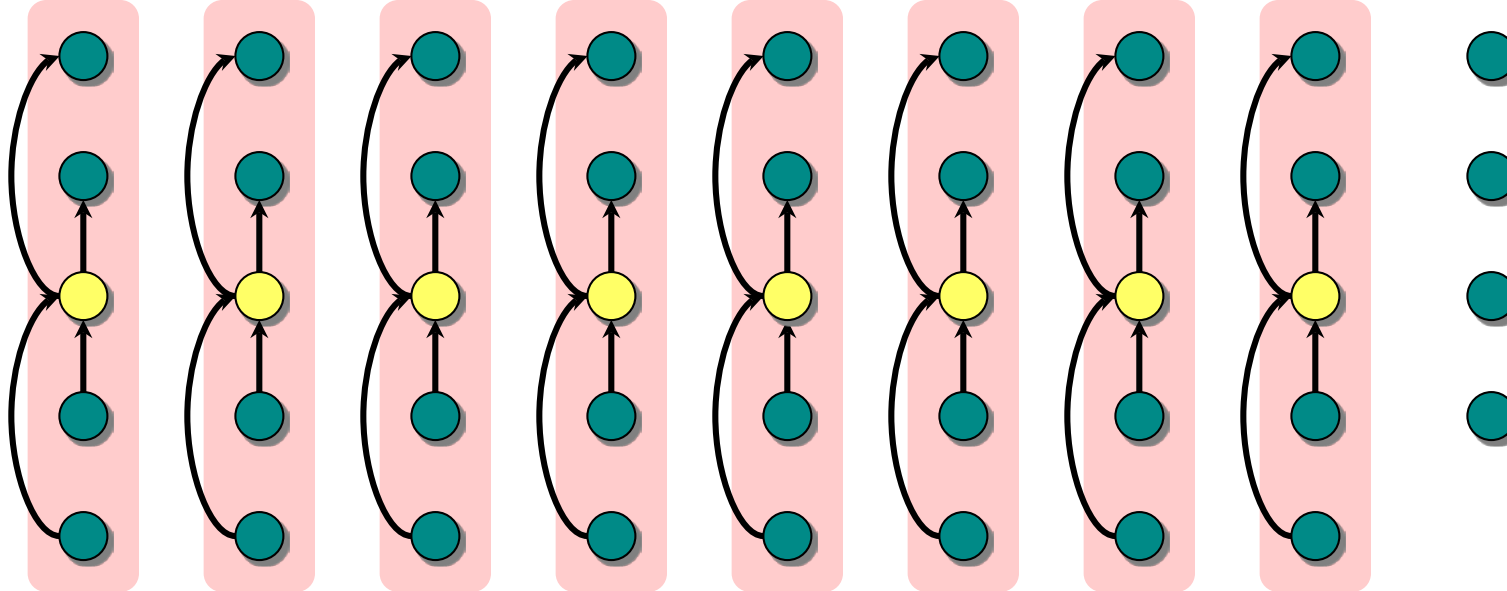
# CHOOSING THE PIVOT



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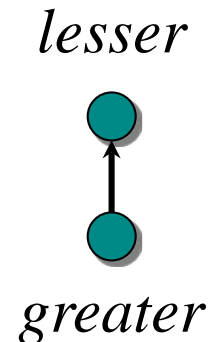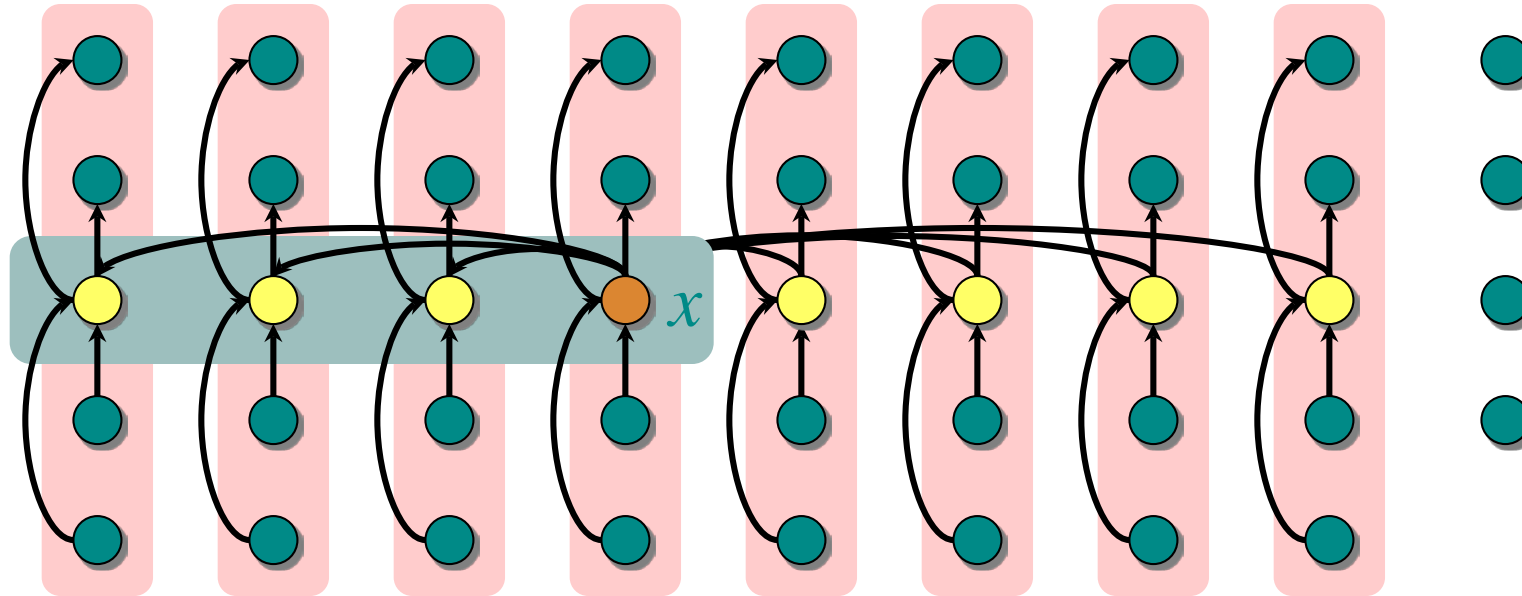
*lesser*

*greater*

# 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.
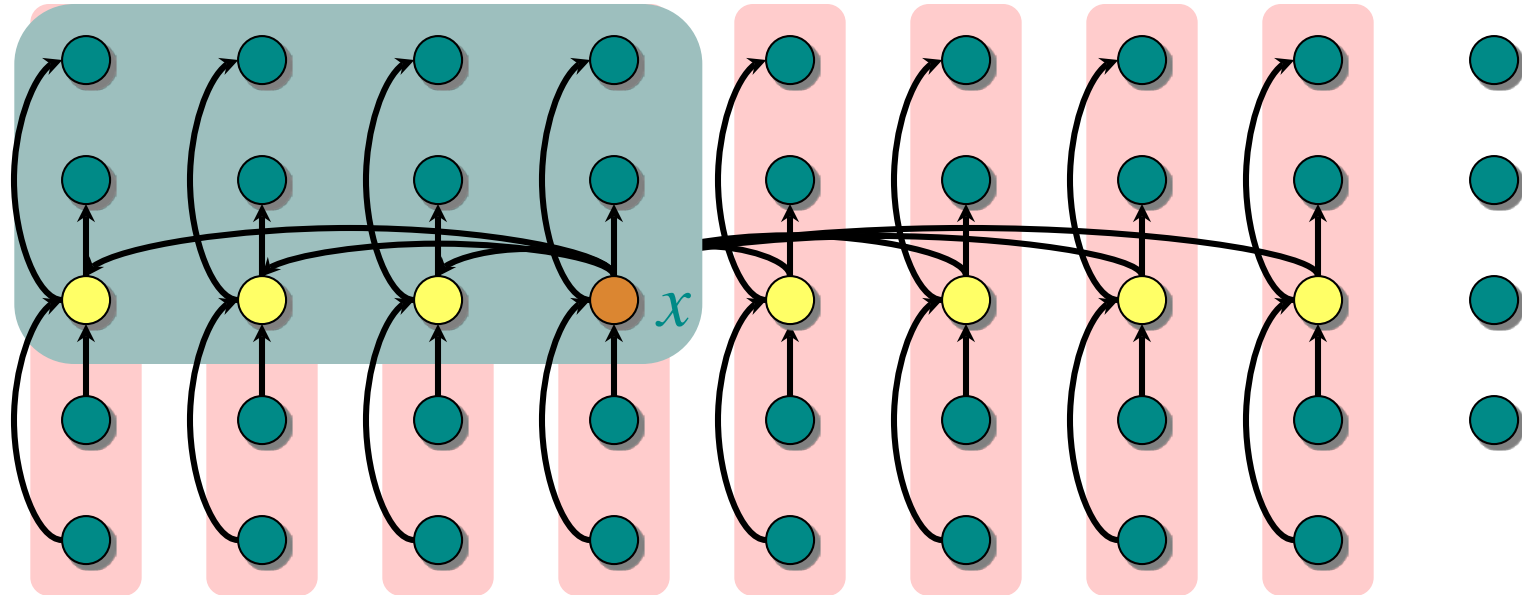
*lesser*

*greater*

# 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.
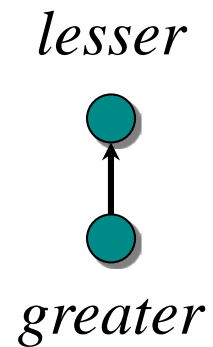
*lesser*

*greater*

# 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.

- Therefore, at least $3\lfloor n/10 \rfloor$ elements are $\leq x$.

*lesser*

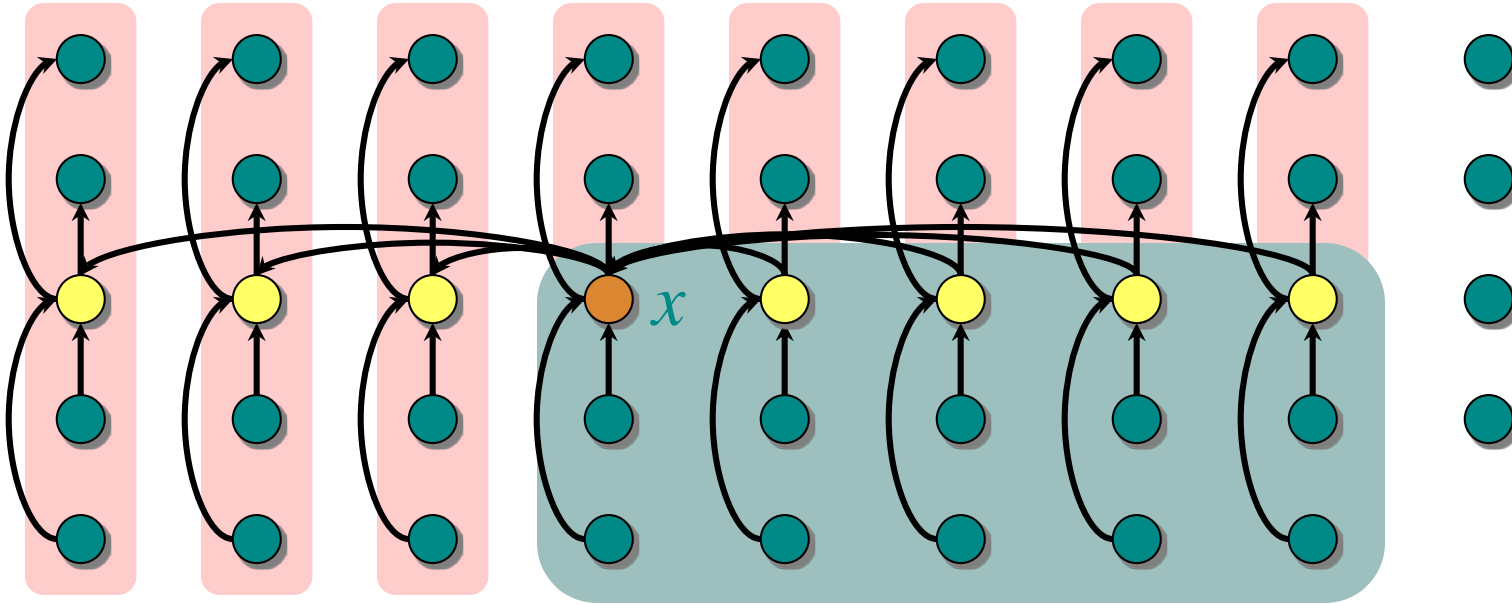*greater*

# ANALYSIS (Assume all elements are distinct.)



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.

- Therefore, at least $3\lfloor n/10 \rfloor$ elements are $\leq x$.
- Similarly, at least $3\lfloor n/10 \rfloor$ elements are $\geq x$.

*lesser*

*greater*

# DEVELOPING THE RECURRENCE

$T(n)$  Select$(i, n)$

$\Theta(n)$ $\begin{cases} \\ \\ \end{cases}$ 1. Divide the $n$ elements into groups of $5$. Find the median of each $5$-element group by rote.

$T(n/5)$ $\begin{cases} \\ \\ \end{cases}$ 2. Recursively Select the median $x$ of the $\lfloor n/5 \rfloor$ group medians to be the pivot.

$\Theta(n)$ 3. Partition around the pivot $x$. Let $k = \text{rank}(x)$.

$T(7n/10)$ $\begin{cases} \\ \\ \\ \\ \\ \\ \end{cases}$ 4. **if** $i = k$ **then return** $x$
**elseif** $i < k$
    **then** recursively Select the $i$th
       smallest element in the lower part
   **else** recursively Select the $(i–k)$th
      smallest element in the upper part

# SOLVING THE RECURRENCE

$$T(n) = T(n/5) + T(7n/10) + \Theta(n)$$

**Substitution:**

$T(n) \leq cn$

$$\begin{aligned}
T(n) &\leq cn/5 + 7cn/10 + \Theta(n) \\
&= 9cn/10 + \Theta(n) \\
&= cn - \left[ cn/10 - \Theta(n) \right] \\
&\leq cn
\end{aligned}$$

,

if $c$ is chosen large enough to handle both the $\Theta(n)$ and the initial conditions. $7n/10$

# CONCLUSIONS

- Since the work at each level of recursion is a constant fraction (9/10) 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.