

# Selection problem

Given a set of “n” numbers we can say that,

- ✓ **Mean:** Average of the “n” numbers
- ✓ **Median:** Having sorted the “n” numbers, the value which lies in the middle of the list such that half the numbers are higher than it and half the numbers are lower than it

(if n is even, can declare  $(n-1)/2$  or  $(n+1)/2$  lowest number as median)

**k-Selection problem:** given a list of n numbers, find the  $k^{\text{th}}$  smallest number

Median is when  $k = n/2$

To simplify things, we assume throughout that numbers are **distinct**

# Selection problem

**Exercise:** Design algorithms that solve the k-selection problem in time:

- a)  $O(kn)$
- b)  $O(n \log n)$
- c)  $O(n + k \log n)$

**A:**

- a) **Do k linear passes**, each time removing the smallest element; return the last one you removed
- b) **Sort** the list and look at the kth position
- c) **Build a Heap** with the numbers (time  $O(n)$ ), remove the minimum element k times ( $O(\log n)$  per removal), return the last removed item

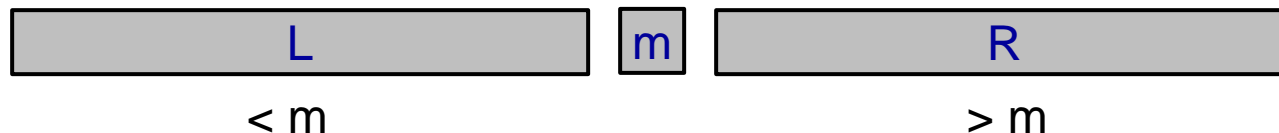
We will see how to do this in time  $O(n)$

# Selection in linear time

We will first cheat: we assume we know we can solve the **median** problem **for free**

ExotericSelect(A,k)

- [Base case:  $|A| = 1$ ]
- **Ask the oracle for the median  $m$  of  $A$**
- Partition numbers around the median so that values less than  $m$  are in set  $L$  and values greater than  $m$  are in set  $R$
- If  $|L|=k-1$  then
- If  $|L|>k-1$  then
- Else



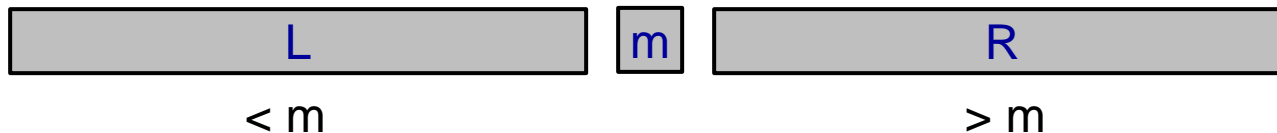
**Obs:** We did not use so far that  $m$  is the median, we will only need this in the analysis of running time

# Selection in linear time

We will first cheat: we assume we know we can solve the **median** problem **for free**

ExotericSelect(A,k)

- [Base case:  $|A| = 1$ ]
- **Ask the oracle for the median  $m$  of  $A$**
- Partition numbers around the median so that values less than  $m$  are in set L and values greater than  $m$  are in set R
- If  $|L|=k-1$  then return  $m$
- If  $|L|>k-1$  then ExotericSelect(L,k)
- Else ExotericSelect(R,  $k-|L|-1$ )



**Obs:** We did not use so far that  $m$  is the median, we will only need this in the analysis of running time

# Selection in linear time: ExotericSelect Analysis

The recurrence equation of the algorithm is

$$\begin{aligned} T(1) &\leq cst \\ T(n) &\leq cst*n + T(n/2) \quad \text{if } n > 1 \end{aligned}$$

Q: What is the solution to this recurrence?

$T(n)$	$cst*n$	
$T(n/2)$	$cst*n/2$	Total: $cst*(n + n/2 + n/4 + \dots + 1)$
$T(n/4)$	$cst*n/4$	$= cst*n*(1 + \frac{1}{2} + \frac{1}{4} + \dots + 1/n) \leq O(n)$
...	...	
$T(1)$	$cst$	

↑  
Geometric series

So the ExotericSelect algorithm is  $O(n)$

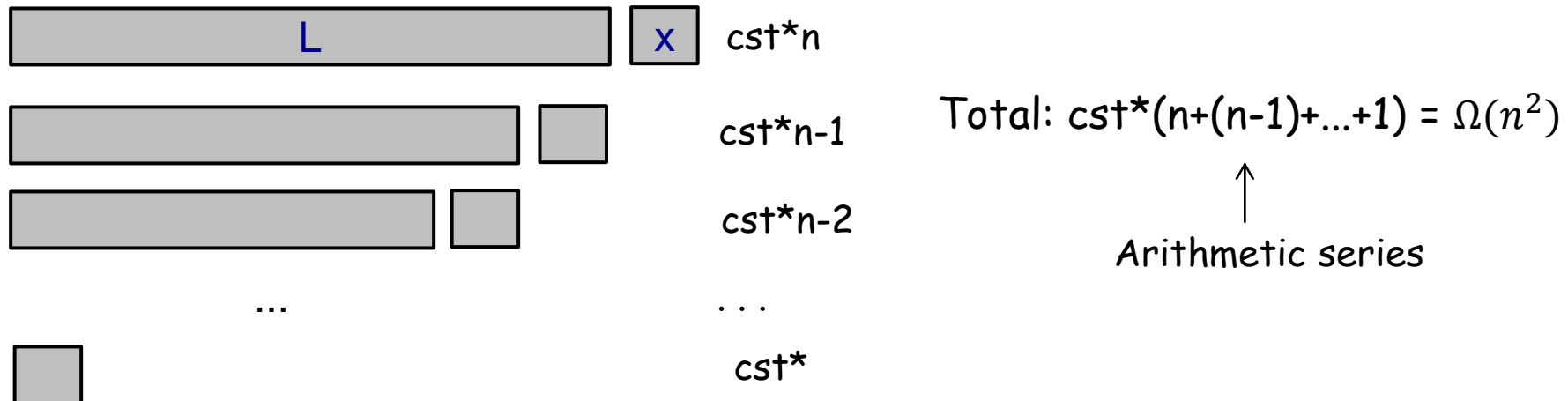
# Selection in linear time

But ExotericSelect cheated by assuming we have the median for free

We cannot yet compute median, need to use something simpler

**Q:** What is the problem of using **any** number  $x$  to split our list into  $< x$  and  $> x$ ?

**A:** Running time can be  $\Omega(n^2)$  if  $x$  is always the largest element (and  $k=1$ )



**Idea:** Use a relaxed version of median splitting around a value  $x$  such that there are **at least  $3n/10$**  numbers smaller and at least  $3n/10$  larger

# Selection in linear time

Median of medians algorithm (Blum, Floyd, Pratt, Rivest, Tarjan '73)

## MOM\_algo(A,k)

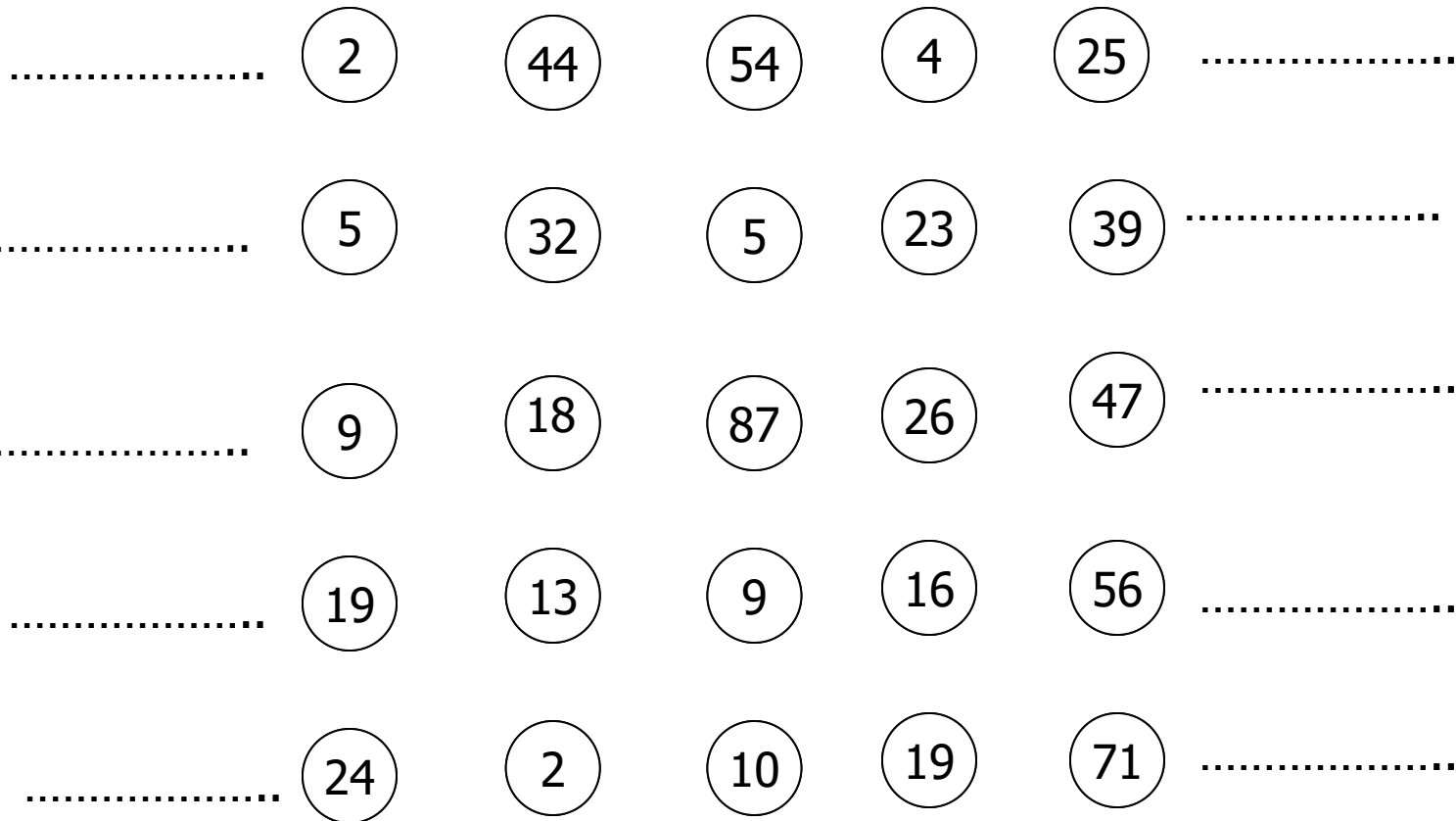
Replaces median(A)

1. Group the numbers into sets of 5
2. Sort individual groups and find the median of each group; put these medians in a set M
3. Find **median m'** of set M using MOM\_algo(M,|M|/2)
4. Partition original data around **m'** such that values less than it are in set L and values greater than it are in set R
5. If  $|L| = k-1$ , then return **m'**  
If  $|L| > k-1$ , then return MOM\_algo(L,k)  
If  $|L| < k-1$  then return MOM\_algo(R,k-|L|-1)

} oracle

Ex: (2,5,9,19,24,54,5,87,9,10,44,32,18,13,2,4,23,26,16,17,25,39,47, 56,71)

**Step1:** Group numbers in sets of 5 (Vertically)





**Step2:** Sort each group, find median of each group

..... (2) (2) (5) (4) (25) .....

..... (5) (13) (9) (16) (39) .....

..... (9) (18) (10) (19) (47) .....

..... (19) (32) (54) (23) (56) .....

..... (24) (44) (87) (26) (71) .....

M = Medians of each group

### Step3: Find the median of medians

..... (2) (2) (5) (4) (25) .....

..... (5) (13) (9) (16) (39) .....

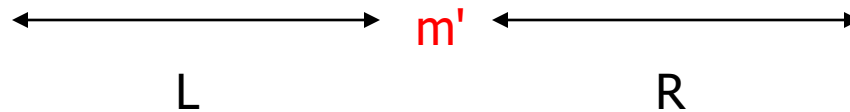
..... (9) (18) (10) (19) (47) .....

..... (19) (32) (54) (23) (56) .....

..... (24) (44) (87) (26) (71) .....

$m'$  = the median of medians

**Step4:** Partition original data around the median-of-medians



**Step5:** If  $|L| = k-1$ , then return  $m'$  else  
If  $|L| > k-1$ , then return  $\text{kth\_smallest}(L, k)$  else  
If  $|L| < k-1$  then return  $\text{kth\_smallest}(R, k-(|L|+1))$

# Selection in linear time

Median of medians algorithm (Blum, Floyd, Pratt, Rivest, Tarjan '73)

## MOM\_algo(A,k)

Replaces median(A)

1. Group the numbers into sets of 5
2. Sort individual groups and find the median of each group; put these medians in a set M
3. Find median  $m'$  of set M using MOM\_algo(M,|M|/2)
4. Partition original data around  $m'$  such that values less than it are in set L and values greater than it are in set R
5. If  $|L| = k-1$ , then return  $m'$   
If  $|L| > k-1$ , then return MOM\_algo(L,k)  
If  $|L| < k-1$  then return MOM\_algo(R,k-|L|-1)

} oracle

This is now ok because we want median of a smaller set of numbers, can use recursion

# Selection in linear time: Analysis

We show that the median-of-medians  $m'$  gives a  $\sim$ -balanced partition

**Lemma:** The set L has at least  $3n/10$  elements. Similarly  $|R| \geq 3n/10$

*	*	*	*	*	*	*
*	*	*	*	*	*	*
1	4	3	7	5	6	2
*	*	*	*	*	*	*
*	*	*	*	*	*	*

# Selection in linear time: Analysis

We show that the median-of-medians  $m'$  gives a  $\sim$ -balanced partition

**Lemma:** The set L has at least  $3n/10$  elements. Similarly  $|R| \geq 3n/10$

Proof:

- groups have median smaller than  $m'$ , so   groups

*	*	*	*	*	*	*
*	*	*	*	*	*	*
1	4	3	7	5	6	2
*	*	*	*	*	*	*
*	*	*	*	*	*	*

# Selection in linear time: Analysis

We show that the median-of-medians  $m'$  gives a  $\sim$ -balanced partition

**Lemma:** The set L has at least  $3n/10$  elements. Similarly  $|R| \geq 3n/10$

Proof:

- Half of the groups have median smaller than  $m'$ , so  $(n/5)/2 = n/10$  groups
- Each such group has at least  $\square$  elements smaller than  $m'$

*	*	*	*	*	*	*
*	*	*	*	*	*	*
1	4	3	7	5	6	2
*	*	*	*	*	*	*
*	*	*	*	*	*	*

# Selection in linear time: Analysis

We show that the median-of-medians  $m'$  gives a  $\sim$ -balanced partition

**Lemma:** The set  $L$  has at least  $3n/10$  elements. Similarly  $|R| \geq 3n/10$

Proof:

- Half of the groups have median smaller than  $m'$ , so  $(n/5)/2 = n/10$  groups
- Each such group has at least three elements smaller than  $m'$
- So there is a total of  $(n/10)*3 = 3n/10$  elements smaller than  $m'$   
 $\Rightarrow L$  has at least  $3n/10$  elements

*	*	*	*	*	*	*
*	*	*	*	*	*	*
1	4	3	7	5	6	2
*	*	*	*	*	*	*
*	*	*	*	*	*	*

**Corollary:**  $|L| \leq 7n/10$ . And by the same argument  $|R| \leq 7n/10$



# Selection in linear time: Complexity

Complexity:

<b>Step</b>	<b>Task</b>	<b>Complexity</b>
1	Group into sets of 5	$O(n)$
2	Find Median of each group	$O(n)$
3	Find med-of-med $m'$	$T(n/5)$
4	Partition around $m'$	$O(n)$
5	Recurse on L or R	$\leq T(7n/10)$

Recurrence equation:

$$T(1) \leq cst$$

$$T(n) \leq cst*n + T(n/5) + T(7n/10) \quad \text{for } n > 1$$

# Selection in linear time: Complexity

Solving the recurrence equation:

$$T(1) \leq cst$$

$$T(n) \leq cst*n + T(n/5) + T(7n/10) \quad \text{for } n > 1$$

So the algorithm is  $O(n)$

# Selection in linear time

## Natural Questions

- Can we split the list into groups of 3 elements instead of 5?
- Can we split the list into groups of 7 elements instead of 5?

# Selection in linear time

## Natural Questions

- Can we split the list into groups of 3 elements instead of 5?
  - $T(n) = cn + T(n/3) + T(2n/3)$ . This recurrence does **not** have a linear solution
- Can we split the list into groups of 7 elements instead of 5?
  - $T(n) = cn + T(n/7) + T(5n/7)$ . **Linear solution**, with a different constant

# QuickSelect

Linear time algorithm has a large constant

- How to proceed in practice?

Just like in quicksort, all we want is an element (pivot) that partitions input list in balanced way

QuickSelect(A,k)

- **Select a Random element p from the list**
- Partition original data around the pivot p such that values less than it are in set L and values greater than it are in set R
- If  $|L|=k-1$  then return p
- If  $|L|>k-1$  then QuickSelect(L,k)
- Else QuickSelect(R, k-|L|-1)

# Exercise

**Exercise:** Use the linear time algorithm for finding the median to design a version of Quicksort that runs in time  $O(n \log n)$  (in the worst case)