

# Recursion Tree Method and Selection

# Divide & Conquer: Quicksort

**Quicksort( $A$ ):**

If  $|A| < 3$  : Sort( $A$ ) directly

Else: choose a pivot element  $p \leftarrow A$

$A_{<p}, A_{>p} \leftarrow$  Partition around  $p$

Quicksort( $A_{<p}$ )

Quicksort( $A_{>p}$ )

- Running time depends on the rank (position in sorted order) of the pivot

# Quick Sort Analysis

- Partition takes  $O(n)$  time
- Size of the subproblems depends pivot; let  $r$  be the rank of the pivot, then:
- $T(n) = T(r - 1) + T(n - r) + O(n)$ ,  $T(1) = 1$
- Let us analyze some cases for  $r$ 
  - **Best case:**  $r$  is the median:  $r = \lfloor n/2 \rfloor$  (we will learn how to compute the median in  $O(n)$  time)
  - **Worst case:**  $r = 1$  or  $r = n$
  - **In between:** say  $n/10 \leq r \leq 9n/10$
- Note in the worst-case analysis, we only consider the worst case for  $r$ . We are looking at the difference cases, just to get a sense for it.

# Quick Sort Cases

- Suppose  $r = n/2$  (pivot is the median element), then
  - $T(n) = 2T(n/2) + O(n)$ ,  $T(1) = 1$
  - We have already solved this recurrence
  - $T(n) = O(n \log n)$
- Suppose  $r = 1$  or  $r = n - 1$ , then
  - $T(n) = T(n - 1) + T(1) + 1$
  - What running time would this recurrence lead to?
  - $T(n) = \Theta(n^2)$  (notice: this is tight!)

# Quick Sort Cases

- Suppose  $r = n/10$  (that is, you get a one-tenth, nine-tenths split)
- $T(n) = T(n/10) + T(9n/10) + O(n)$
- Let's look at the recursion tree for this recurrence
- We get  $T(n) = O(n \log n)$ , in fact, we get  $\Theta(n \log n)$
- In general, the following holds (we'll show it later):
- $T(n) = T(\alpha n) + T(\beta n) + O(n)$ 
  - If  $\alpha + \beta < 1 : T(n) = O(n)$
  - If  $\alpha + \beta = 1, T(n) = O(n \log n)$

# Quick Sort: Theory and Practice

- We can find the **median element** in  $\Theta(n)$  time
  - Using divide and conquer! we'll learn how in next lecture
- In practice, the constants hidden in the Oh notation for median finding are too large to use for sorting
- Common heuristic
  - Median of three (pick elements from the start, middle and end and take their median)
- If the pivot is chosen **uniformly at random**
  - quick sort runs in time  $O(n \log n)$  in expectation and *with high probability*
  - We will prove this in the second half of the class

# Challenge Recurrence

- Solve the following recurrence:

$$T(n) = \sqrt{n}T(\sqrt{n}) + n$$

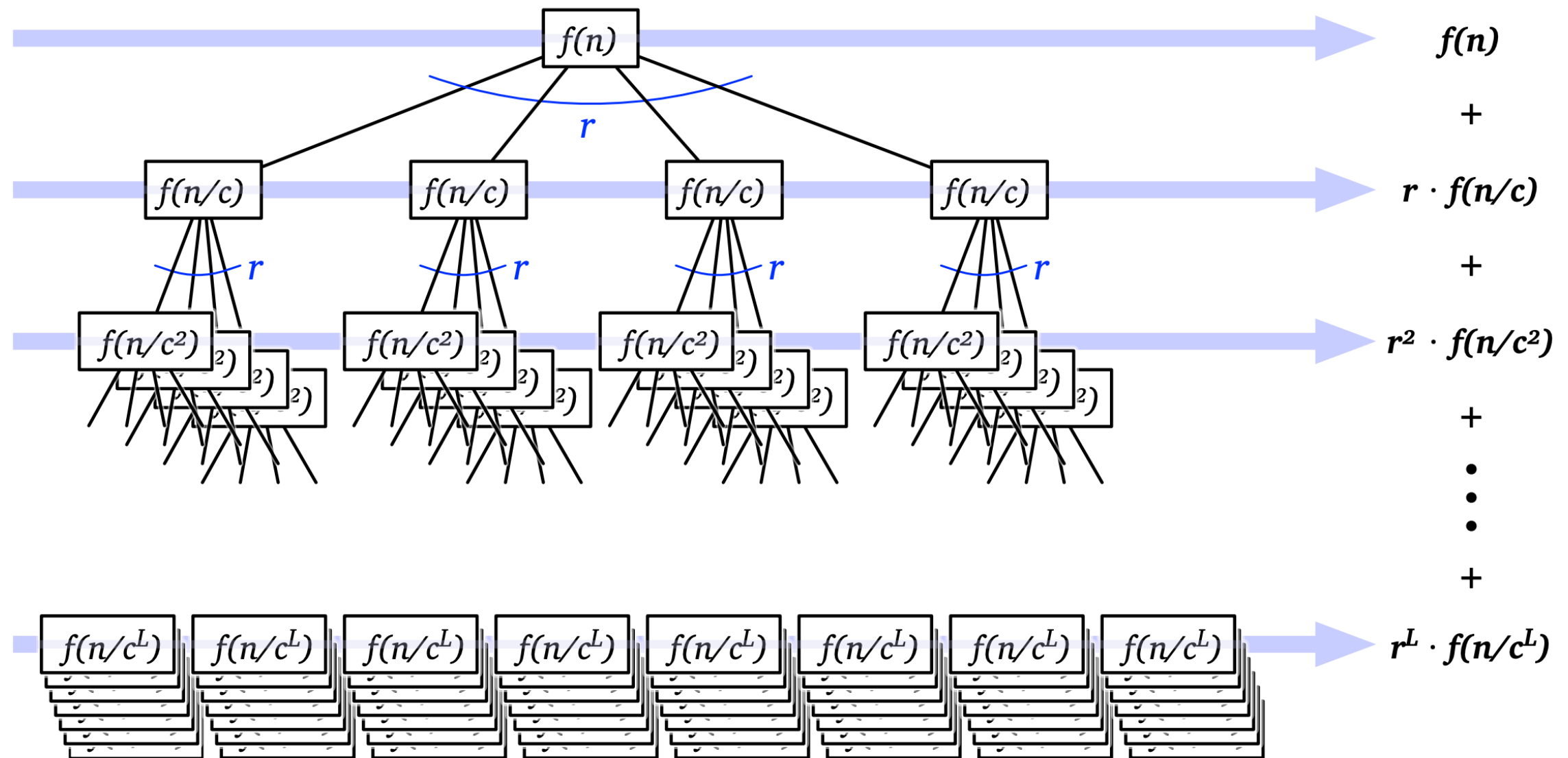
- **Hint.** Try some change of variables

# General Recursion Trees

- Consider a divide and conquer algorithm that
  - spends  $O(f(n))$  time on non-recursive work and makes  $r$  recursive calls, each on a problem of size  $n/c$
- Up to constant factors (which we hide in  $O()$ ), the running time of the algorithm is given by what **recurrence**?
  - $T(n) = rT(n/c) + f(n)$
- Because we care about asymptotic bounds, we can assume base case is a small constant, say  $T(n) = 1$



# General Recursion Trees



A recursion tree for the recurrence  $T(n) = rT(n/c) + f(n)$

- For each  $i$ , the  $i$ th level of tree has exactly  $r^i$  nodes
- Each node at level  $i$ , has cost  $f(n/c^i)$

# General Recursion Trees

- Running time  $T(n)$  of a recursive algorithm is the sum of all the values (sum of work at all nodes at each level) in the recursion tree
- For each  $i$ , the  $i$ th level of tree has exactly  $r^i$  nodes
- Each node at level  $i$ , has cost  $f(n/c^i)$
- Thus, 
$$T(n) = \sum_{i=0}^L r^i \cdot f(n/c^i)$$
- Here  $L = \log_c n$  is the depth of the tree
- Number of leaves in the tree:  $r^L = n^{\log_c r}$
- Cost at leaves:  $O(n^{\log_c r} f(1))$

# General Recursion Trees

- Running time  $T(n)$  of a recursive algorithm is the sum of all the values (sum of work at all nodes at each level) in the recursion tree
- For each  $i$ , the  $i$ th level of tree has exactly  $r^i$  nodes
- Each node at level  $i$ , has cost  $f(n/c^i)$
- Thus,  $T(n) = \sum_{i=0}^L r^i \cdot f(n/c^i)$
- Here  $L = \log_c n$  is the depth of the tree
- Number of leaves in the tree:  $r^L = n^{\log_c r}$  (why?)
- Cost at leaves:  $O(n^{\log_c r} f(1))$

$$r^L = r^{\log_c n} = (2^{\log_2 r})^{\log_c n} = (2^{\log_c n})^{\log_2 r} = (2^{\log_2 n})^{\frac{\log_2 r}{\log_2 c}} = n^{\log_c r}$$

# Common Cases

$$T(n) = \sum_{i=0}^L r^i \cdot f(n/c^i)$$

Don't forget:  $\sum_{i=0}^L a^i = \frac{a^{L+1} - 1}{a - 1}$

- **Decreasing series.** If the series decays exponentially (every term is a constant factor smaller than previous), cost at root dominates:

$$T(n) = O(f(n))$$

- **Equal.** If all terms in the series are equal:

$$T(n) = O(f(n) \cdot L) = O(f(n) \log n)$$

- **Increasing series.** If the series grows exponentially (every term is constant factor larger), then the cost at leaves dominates:

$$T(n) = O(n^{\log_c r})$$

# Recurrences

So far we saw divide and conquer algorithms, where we split the problem in more than one subproblem.

**Question.** Can you think of some examples (that you have likely seen before) where we split the problem into **one** smaller subproblem?

# D&C: One Smaller Subproblem

- Binary search
  - $T(n) = T(n/2) + 1$
- Binary search tree
  - $T(n) = T(n/2) + 1$
- Fast exponentiation (you may not have seen this)
  - Compute  $a^n$ , how many multiplications?
  - Naive way:  $a \cdot a \cdot \dots \cdot a$  ( $n$  times)
  - Faster way:  $a^n = (a^{n/2})^2$  (suppose  $n$  is even)
  - $T(n) = T(n/2) + 1$
  - What does this solve to?
  - Think at home: What if  $n$  is odd?

# In Class Exercises

- Take a few minutes to draw recursions trees for each of the following recurrences
- Then break into small groups (~size 3) and discuss which of the three cases each of them fall into

- $T(n) = 2(Tn/2) + n^2$

- $T(n) = 3T(n/2) + n$

# Master Theorem (optional)

Set of rules to solve some common recurrences automatically

**(Master Theorem)** Let  $a \geq 1$ ,  $b > 1$  and  $f(n) \geq 0$ . Let  $T(n)$  be defined by the recurrence  $T(n) = aT(n/b) + f(n)$  and  $T(1) = O(1)$ .

Then  $T(n)$  can be bounded asymptotically as follows.

- If  $f(n) = n^{\log_b a - \epsilon}$  for some constant  $\epsilon > 0$ , then  $T(n) = \Theta(n^{\log_b a})$
- If  $f(n) = \Theta(n^{\log_b a})$ , then  $T(n) = \Theta(n^{\log_b a} \log n)$
- If  $f(n) = \Omega(n^{\log_b a + \epsilon})$ , for some constant  $\epsilon > 0$ , and if  $af(n/b) \leq c_0 f(n)$  for some constant  $c_0 < 1$  and all sufficiently large  $n$ , then  
$$T(n) = \Theta(f(n))$$



# Master Theorem

- It exists; it can make things easier. You don't need to know it
- OK to use in this class, but I don't encourage (nor discourage) it
  - Recursion trees promote a better understanding of the recurrence—and they can be simpler
- Master Theorem only applies to some recurrences (generalizations do exist)

# Selection

# Selection: Problem Statement

Given an array  $A[1, \dots, n]$  of size  $n$ , find the  $k$ th smallest element for any  $1 \leq k \leq n$

- Special cases: min  $k = 1$ , max  $k = n$ :
  - Linear time,  $O(n)$
- What about **median**  $k = \lfloor n + 1 \rfloor / 2$ ?
  - Sorting:  $O(n \log n)$  compares
  - Binary heap:  $O(n \log k)$  compares

**Question.** Can we do it in  $O(n)$  compares?

- **Surprisingly yes.**
- Selection is easier than sorting.

# Selection: Problem Statement

Example. Take this array of size 10:

$A = 12|2|4|5|3|1|10|7|9|8$

Suppose we want to find 4th smallest element

- First, take any pivot  $p$  from  $A[1, \dots, n]$
- If  $p$  is the 4th smallest element, return it
- Else, we partition  $A$  around  $p$  and recurse

# Selection Algorithm: Idea

Select ( $A, k$ ):

If  $|A| = 1$ : return  $A[1]$

Else:

- Choose a pivot  $p \leftarrow A[1, \dots, n]$ ; let  $r$  be the rank of  $p$
- $r, A_{<p}, A_{>p} \leftarrow \text{Partition}(A, p)$
- If  $k == r$ , return  $p$
- Else:
  - If  $k < r$ : Select ( $A_{<p}, k$ )
  - Else: Select ( $A_{>p}, k - r$ )

# Selection: Problem Statement

Example. Take this array of size 10:

$A = 12|2|4|5|3|1|10|7|9|8$

Suppose we want to find 4th smallest element

- Choose pivot 8
- What is its rank?
  - Rank 7
- So let's find all of the smaller elements of  $A$ :
  - $A' = 2|4|5|3|1|7$
- Want to find the element of rank 4 in this new array

# Selection: Problem Statement

Example. Take this array of size 10:

$A = 12|2|4|5|3|1|10|7|9|8$

Suppose we want to find 4th smallest element

- Choose as pivot 3
- What is its rank?
  - Rank 3
- So let's find all of the **larger** elements of  $A$ :
  - $A' = 12|4|5|10|7|9|8$
- Want to find the element of rank  $4 - 3 = 1$  in this new array

# When is this method good?

- If we guess the pivot right! (but we can't always do that)
- If we partition the array pretty evenly (the pivot is close to the middle)
  - Let's say our pivot is not in the first or last 3/10ths of the array
  - What is our recurrence?
  - $T(n) \leq T(7n/10) + O(n)$
  - $T(n) = O(n)$



# Our high-level goal

- Find a pivot that's close to the median—has a rank between  $3n/10$  and  $7n/10$ , in time  $O(n)$
- But the array is unsorted? How do we do that?
- Want to *always* be successful

# Finding an Approximate Median

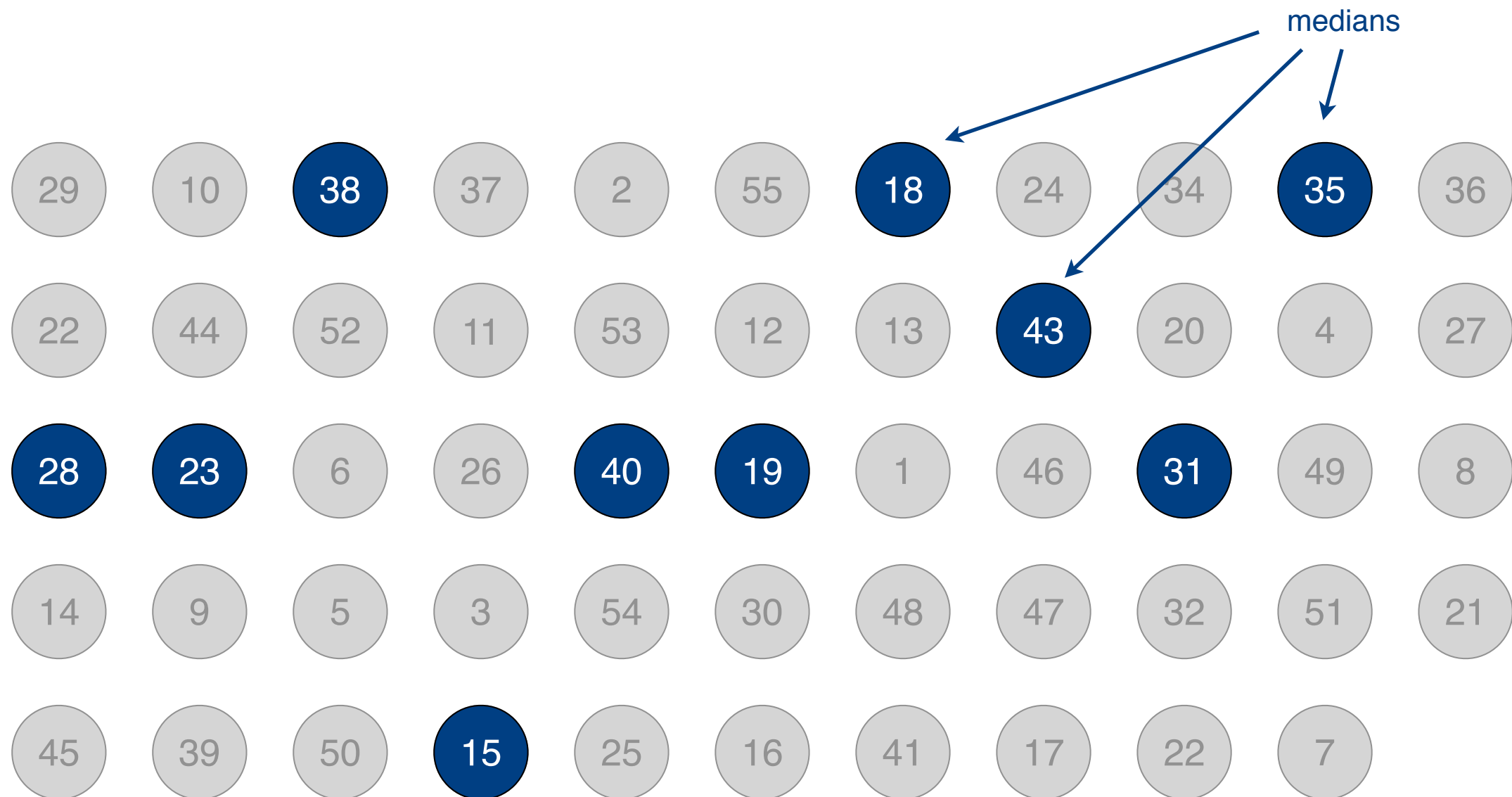
- Divide the array of size  $n$  into  $\lceil n/5 \rceil$  groups of 5 elements (ignore leftovers)
- Find median of each group

29	10	38	37	2	55	18	24	34	35	36
22	44	52	11	53	12	13	43	20	4	27
28	23	6	26	40	19	1	46	31	49	8
14	9	5	3	54	30	48	47	32	51	21
45	39	50	15	25	16	41	17	22	7	

$n = 54$

# Finding an Approximate Median

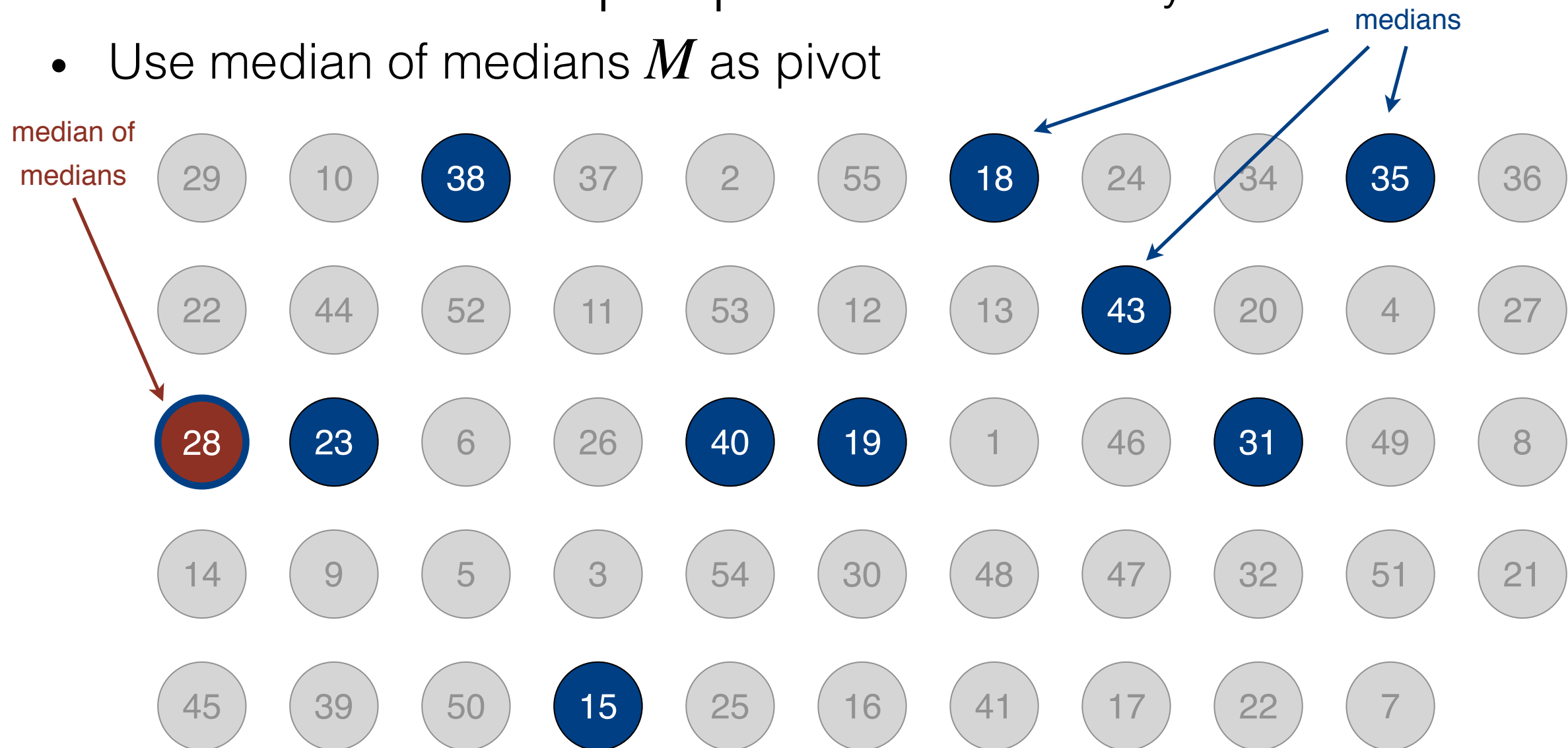
- Divide the array of size  $n$  into  $\lceil n/5 \rceil$  groups of 5 elements (ignore leftovers)
- Find median of each group



$n = 54$

# Finding an Approximate Median

- Divide the array of size  $n$  into  $\lceil n/5 \rceil$  groups of 5 elements (ignore leftovers)
- Find median of each group
- Find  $M \leftarrow$  median of  $\lceil n/5 \rceil$  medians recursively
- Use median of medians  $M$  as pivot

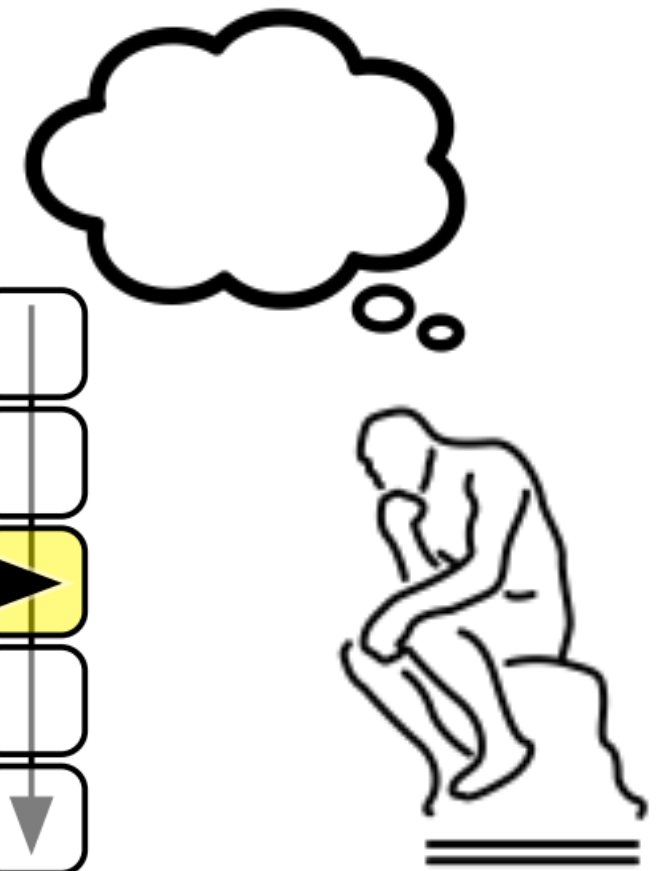
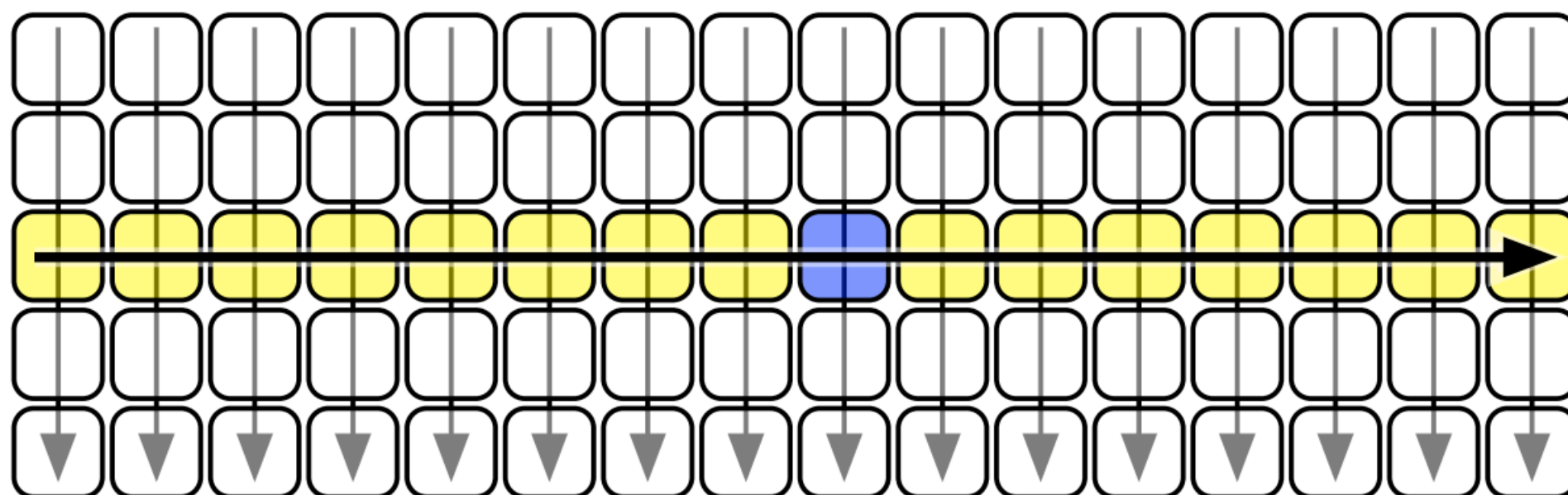


# What did we gain?

- How can I show that the median of medians is “close to the center” of the array?
- What elements can I say, for sure, are  $\leq$  the median of medians?
  - The smaller half of the medians
  - $n/10$  elements
- Any other elements?
  - Another 2 elements in each median’s list

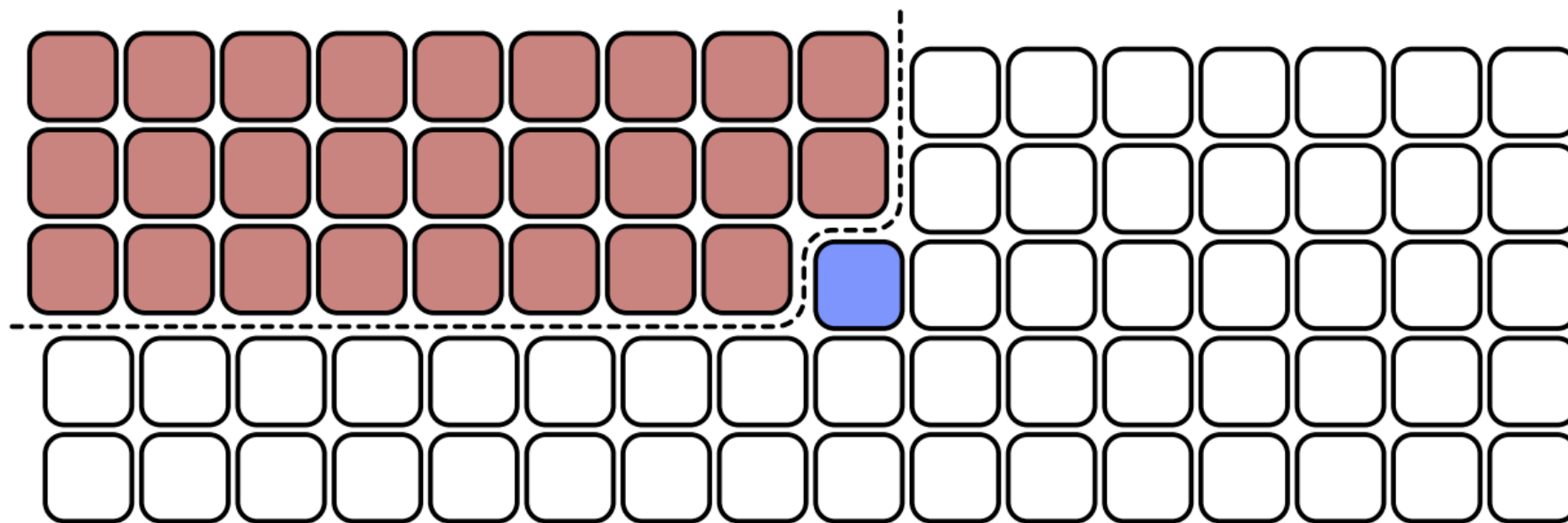
# Visualizing MoM

- In the  $5 \times n/5$  grid, each column represents five consecutive elements
- **Imagine** each column is sorted top down
- **Imagine** the columns as a whole are sorted left-right
  - We don't actually do this!
- MoM is the element closest to center of grid



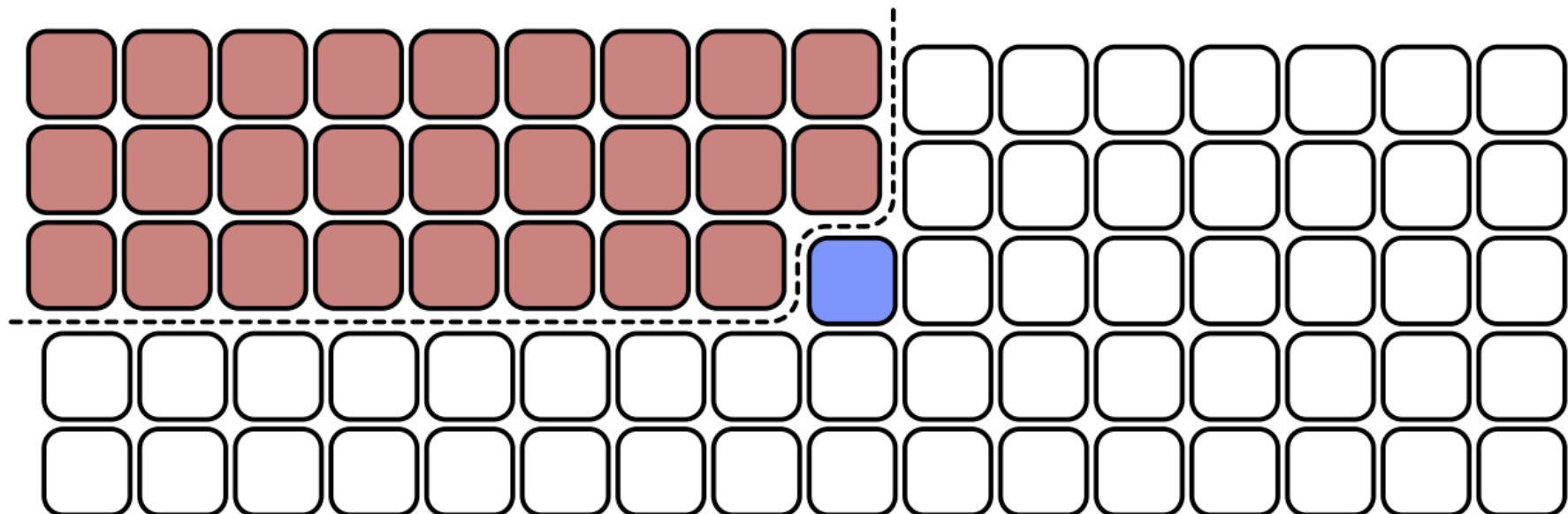
# Visualizing MoM

- Red cells (at least  $3n/10$ ) are smaller than  $M$



# Visualizing MoM

- Red cells (at least  $3n/10$ ) in size are smaller than  $M$
- If we are looking for an element larger than  $M$ , we can throw these out, before recursing
- Symmetrically, we can throw out  $3n/10$  elements smaller than  $M$  if looking for a smaller element
- Thus, the recursive problem size is at most  $7n/10$





# How Good is Median of Medians

**Claim.** Median of medians  $M$  is a good pivot, that is, at least  $3/10$ th of the elements are  $\geq M$  and at least  $3/10$ th of the elements are  $\leq M$ .

**Proof.**

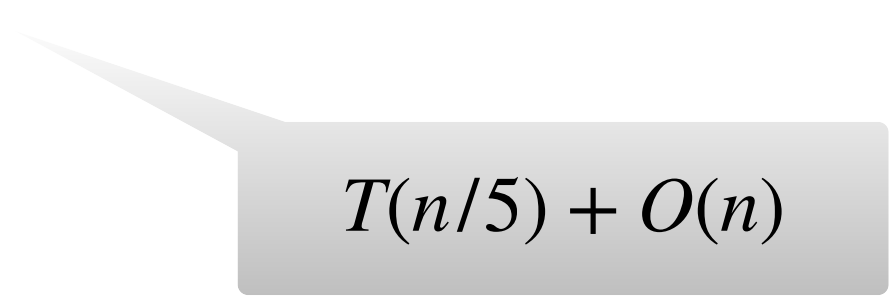
- Let  $g = \lceil n/5 \rceil$  be the size of each group.
- $M$  is the median of  $g$  medians
  - So  $M \geq g/2$  of the group medians
  - Each median is greater than 2 elements in its group
  - Thus  $M \geq 3g/2 = 3n/10$  elements
- Symmetrically,  $M \leq 3n/10$  elements. ■

# How to Use the MoM?

- There are  $3n/10$  elements smaller than the MoM
- By the same argument:  $3n/10$  elements larger than the MoM
- So we can throw out  $3n/10$  elements, adjust the value of  $k$  we are looking for, and recurse!
- Don't forget: we *also* recursed to find the MoM!

# Median of Medians Subroutine

- MoM( $A, n$ ):
  - If  $n = 1$ : return  $A[1]$
  - Else:
    - Divide  $A$  into  $\lceil n/5 \rceil$  groups
    - Compute median of each group
    - $A' \leftarrow$  group medians
    - MoM( $A', \lceil n/5 \rceil$ )


$$T(n/5) + O(n)$$

# Linear time Selection

Select ( $A, k$ ):

If  $|A| = 1$ : return  $A[1]$ ; else:

- Call median of medians to find a good pivot

$$p \leftarrow \text{MoM}(A, n); \quad n = |A|$$

- $r, A_{<p}, A_{>p} \leftarrow \text{Partition}(A, p)$

- If  $k == r$ , return  $p$

- Else:

- If  $k < r$ : Select ( $A_{<p}, k$ )

- Else: Select ( $A_{>p}, k - r$ )

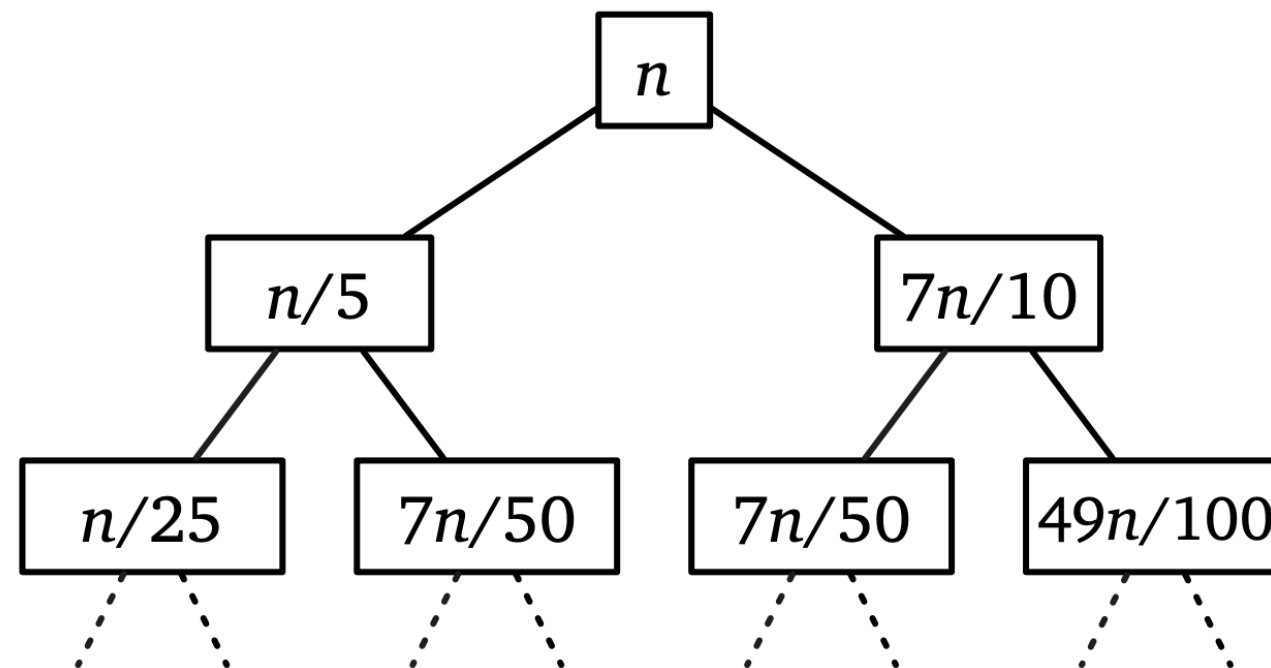
$$T(n/5) + O(n)$$

Larger subproblem  
has size  $\leq 7n/10$

$$\text{Overall: } T(n) = T(n/5) + T(7n/10) + O(n)$$

# Selection Recurrence

- Okay, so we have a good pivot
- We are still doing two recursive calls
  - $T(n) \leq T(n/5) + T(7n/10) + O(n)$
- Key: total work at each level still goes down!
- Decaying series gives us :  $T(n) = O(n)$



# Why the Magic Number 5?

- What was so special about 5 in our algorithm?
- It is the smallest odd number that works!
  - (Even numbers are problematic for medians)
- Let us analyze the recurrence with groups of size 3
  - $T(n) \leq T(n/3) + T(2n/3) + O(n)$
  - Work is equal at each level of the tree!
  - $T(n) = \Theta(n \log n)$

# Theory vs Practice

- $O(n)$ -time selection by [\[Blum–Floyd–Pratt–Rivest–Tarjan 1973\]](#)
  - Does  $\leq 5.4305n$  compares
- Upper bound:
  - [Dor–Zwick 1995]  $\leq 2.95n$  compares
- Lower bound:
  - [Dor–Zwick 1999]  $\geq (2 + 2^{-80})n$  compares.
- Constants are still too large for practice
- Random pivot works well in most cases!
  - We will analyze this when we do randomized algorithms

# Recall Challenge Recurrence

- Recall the challenge recurrence

$$T(n) = \sqrt{n}T(\sqrt{n}) + O(n)$$

- How much work at each level?  $O(n)$
- Analyzing how quickly the problem size goes down
- $n \rightarrow n^{1/2} \rightarrow n^{1/4} \rightarrow \dots \rightarrow n^{1/2^L}$
- What is  $L$  for this to be a small constant?
- $L = \log \log n$  (number of levels)
- $T(n) = \Theta(n \log \log n)$ ,



# Floors and Ceilings

- Why doesn't floors and ceilings matter?
- Suppose  $T(n) = T(\lfloor n/2 \rfloor) + T(\lceil n/2 \rceil) + O(n)$
- First, for upper bound, we can safely overestimate
  - $T(n) \leq 2T(\lceil n/2 \rceil) + n \leq 2T(n/2 + 1) + n$
- Second, we can define a function  $S(n) = T(n + \alpha)$ , so that  $S(n)$  satisfies  $S(n) \leq S(n/2) + O(n)$

$$\begin{aligned} S(n) &= T(n + \alpha) \leq 2T(n/2 + \alpha/2 + 1) + n + \alpha \\ &= 2T(n/2 + \alpha - \alpha/2 + 1) + n + \alpha \\ &= 2S(n/2 - \alpha/2 + 1) + n + \alpha \\ &\leq 2S(n/2) + n + 2, \text{ for } \alpha = 2 \end{aligned}$$

# Floors & Ceilings Don't Matter

- Why doesn't floors and ceilings matter?
- Suppose  $T(n) = T(\lfloor n/2 \rfloor) + T(\lceil n/2 \rceil) + O(n)$
- First, for upper bound, we can safely overestimate
  - $T(n) \leq 2T(\lceil n/2 \rceil) + n \leq 2T(n/2 + 1) + n$
- Second, we can define a function  $S(n) = T(n + \alpha)$ , so that  $S(n)$  satisfies  $S(n) \leq S(n/2) + O(n)$ 
  - Setting  $\alpha = 2$  works
- Finally, we know  $S(n) = O(n \log n) = T(n + 2)$
- $T(n) = O((n - 2)\log(n - 2)) = O(n \log n)$

# Can Assume Powers of 2

- Why doesn't taking powers of 2 matter?
- Running time  $T(n)$  is monotonically increasing
- Suppose  $n$  is not a power of 2, let  $n' = 2^\ell$  be such that  $n \leq n' \leq 2n$ ; then
- We can upper bound our asymptotic using  $n'$  and lower bound using  $n'/2$
- In particular, let  $T(n) \leq T(n')$
- And  $T(n) \geq T(n'/2)$
- That is,  $T(n) = \Theta(T(n'))$

# Guess & Verify Recurrences

- **Method 3.** Requires some practice and creativity
- Verification by induction may run into issues
  - Example,  $T(n) = 2T(n/2) + 1$
  - Guess?
    - $T(n) \leq cn$
  - Check  $T(n) \leq cn + 1 \not\leq cn$  for any  $c > 0$
  - Is the guess wrong? Not asymptotically, can fix it up by adding lower-order terms
  - New guess  $T(n) \leq cn - d$  (why minus?)
    - $T(n) \leq cn - 2d + 1 \leq cn - d$  for any  $d \geq 1$
  - $c$  must be chosen large enough to satisfy boundary conditions

# Extra: Verify by Induction

- Suppose I want to prove that the recurrence  $T(n) = 2T(n/2) + 4n, T(1) = 8$  evaluates to  $T(n) = O(n \log n)$
- I need to show that for all sufficiently large  $n$ , I can find a constant  $c$ , such that  $T(n) \leq c \cdot n \log n$
- Base case?
  - $T(1) = 8 \not\leq c \log 1 = 0$  (doesn't work yet, let us fix it up later)
- Assume holds for all  $< n$
- $$\begin{aligned} T(n) &\leq 2(c(n/2)\log(n/2)) + n \\ &= cn \log(n/2) + n \\ &= cn \log n - cn \log 2 + n \leq cn \log n \text{ if } c \geq 1 \end{aligned}$$

# Extra: Verify by Induction

- What about the base case?
- As long as  $n \geq 4$ , our recurrence does not depend on  $T(1)$ ;
- We can just use  $T(2)$  as the base case our induction!  
 $T(2) = 2T(1) + 8 = 24 \leq c \log 2$  for  $c > 24$
- Thus our induction holds for all  $n \geq 2$  and  $c > 24$
- This is how we usually verify our recurrences and prove they are correct: by induction.

# Acknowledgments

- Some of the material in these slides are taken from
  - Kleinberg Tardos Slides by Kevin Wayne (<https://www.cs.princeton.edu/~wayne/kleinberg-tardos/pdf/04GreedyAlgorithmsI.pdf>)
  - Jeff Erickson's Algorithms Book (<http://jeffe.cs.illinois.edu/teaching/algorithms/book/Algorithms-JeffE.pdf>)
  - CLRS Algorithms book