



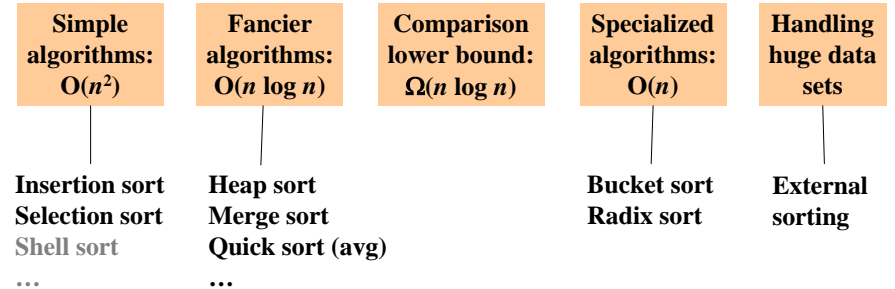
CSE332: Data Abstractions

Lecture 13: Comparison Sorting

Dan Grossman
Spring 2010

The Big Picture

Surprising amount of juicy computer science: 2-3 lectures...



Start with: How would “normal people (?)” sort?

Insertion Sort

- Idea: At the k^{th} step put the k^{th} element in the correct place among the first k elements
- “Loop invariant”: when loop index is i , first i elements are sorted
- Alternate way of saying this:
 - Sort first two elements
 - Now insert 3rd element in order
 - Now insert 4th element in order
 - ...
- Time?
Best-case _____ Worst-case _____ “Average” case _____

Insertion Sort

- Idea: At the k^{th} step put the k^{th} element in the correct place among the first k elements
- “Loop invariant”: when loop index is i , first i elements are sorted
- Alternate way of saying this:
 - Sort first two elements
 - Now insert 3rd element in order
 - Now insert 4th element in order
 - ...
- Time?
Best-case $O(n)$ start sorted Worst-case $O(n^2)$ start reverse sorted “Average” case $O(n^2)$ (see text)

Selection sort

- Idea: At the k^{th} step, find the smallest element among the not-yet-sorted elements and put it at position k
- “Loop invariant”: when loop index is i , first i elements are the i smallest elements in sorted order
- Alternate way of saying this:
 - Find smallest element, put it 1st
 - Find next smallest element, put it 2nd
 - Find next smallest element, put it 3rd
 - ...
- Time?
Best-case _____ Worst-case _____ “Average” case _____

Selection sort

- Idea: At the k^{th} step, find the smallest element among the not-yet-sorted elements and put it at position k
- “Loop invariant”: when loop index is i , first i elements are the i smallest elements in sorted order
- Alternate way of saying this:
 - Find smallest element, put it 1st
 - Find next smallest element, put it 2nd
 - Find next smallest element, put it 3rd
 - ...
- Time?
Best-case $O(n^2)$ Worst-case $O(n^2)$ “Average” case $O(n^2)$
Always $T(1) = 1$ and $T(n) = n + T(n-1)$

Mystery

This is one implementation of which sorting algorithm (for ints)?

```
void mystery(int[] arr) {
    for(int i = 1; i < arr.length; i++) {
        int tmp = arr[i];
        int j;
        for(j=i; j > 0 && tmp < arr[j-1]; j--)
            arr[j] = arr[j-1];
        arr[j] = tmp;
    }
}
```

Note: Like with heaps, “moving the hole” is faster than unnecessary swapping (constant factor)

Insertion vs. Selection

- They are different algorithms
- They solve the same problem
- They have the same worst-case and average-case asymptotic complexity
 - Insertion-sort has better best-case complexity; preferable when input is “mostly sorted”
- Other algorithms are more efficient *for non-small arrays that are not already almost sorted*

Aside

Why I'm not a fan of *bubble sort*

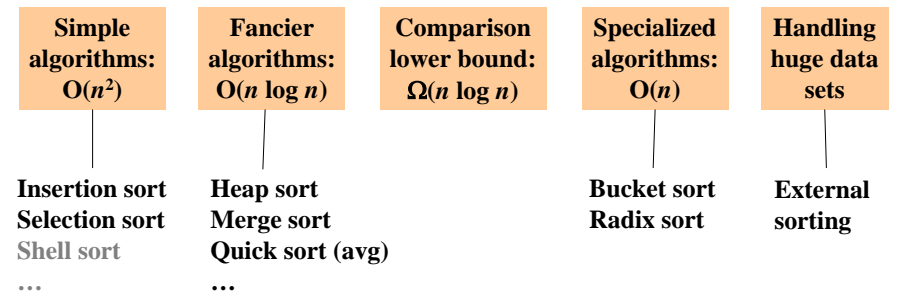
- It is not, in my opinion, what a "normal person" would think of
- It doesn't have good asymptotic complexity: $O(n^2)$
- It's not particularly efficient with respect to common factors

- Basically, almost everything it is good at some other algorithm is at least as good at

- So people seem to teach it just because someone taught it to them

The Big Picture

Surprising amount of juicy computer science: 2-3 lectures...



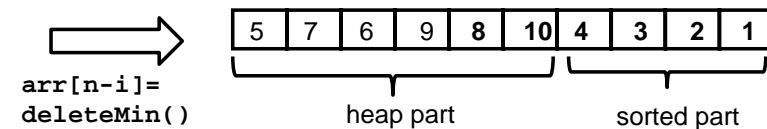
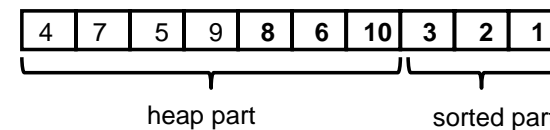
Heap sort

- As you saw on project 2, sorting with a heap is easy:
 - `insert` each `arr[i]`, better yet `buildHeap`
 - `for(i=0; i < arr.length; i++)`
 `arr[i] = deleteMin();`
- Worst-case running time: $O(n \log n)$
- We have the array-to-sort and the heap
 - So this is not an in-place sort
 - There's a trick to make it in-place...

In-place heap sort

But this reverse sorts – how would you fix that?

- Treat the initial array as a heap (via `buildHeap`)
- When you delete the i^{th} element, put it at `arr[n-i]`
 - It's not part of the heap anymore!



“AVL sort”

- We can also use a balanced tree to:
 - **insert** each element: total time $O(n \log n)$
 - Repeatedly **deleteMin**: total time $O(n \log n)$
- But this cannot be made in-place and has worse constant factors than heap sort
 - heap sort is better
 - both are $O(n \log n)$ in worst, best, and average case
 - neither parallelizes well
- Don't even think about trying to sort with a hash table

Divide and conquer

Very important technique in algorithm design

1. Divide problem into smaller parts
2. Solve the parts independently
 - Think recursion
 - Or potential parallelism
3. Combine solution of parts to produce overall solution

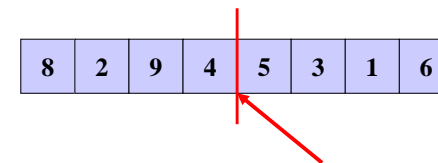
(The name “divide and conquer” is rather clever.)

Divide-and-conquer sorting

Two great sorting methods are fundamentally divide-and-conquer

1. Mergesort: Sort the left half of the elements (recursion)
Sort the right half of the elements (recursion)
Merge the two sorted halves into a sorted whole
2. Quicksort: Pick a “pivot” element
Divide elements into less-than pivot
and greater-than pivot
Sort the two divisions (recursion twice)
Answer is sorted-less-than then pivot then
sorted-greater-than

Mergesort



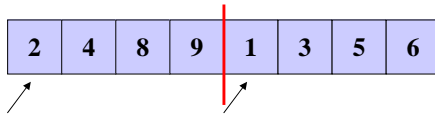
- To sort array from position lo to position hi :
 - If range is 1 element long, it's sorted! (Base case)
 - Else:
 - Sort from lo to $(hi+lo)/2$
 - Sort from $(hi+lo)/2$ to hi
 - Merge the two halves together
- Merging takes two sorted parts and sorts everything
 - $O(n)$ but requires auxiliary space...

Example, focus on merging

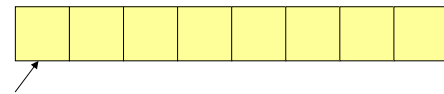
Start with:



After recursion:
(not magic ☺)



Merge:
Use 3 "fingers"
and 1 more array



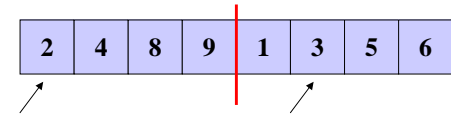
(After merge,
copy back to
original array)

Example, focus on merging

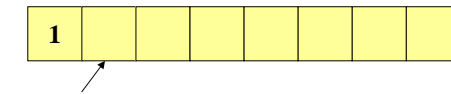
Start with:



After recursion:
(not magic ☺)



Merge:
Use 3 "fingers"
and 1 more array



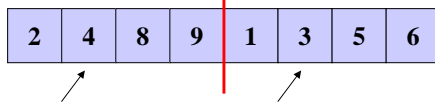
(After merge,
copy back to
original array)

Example, focus on merging

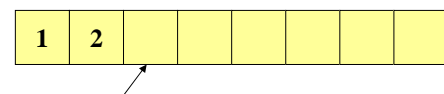
Start with:



After recursion:
(not magic ☺)



Merge:
Use 3 "fingers"
and 1 more array



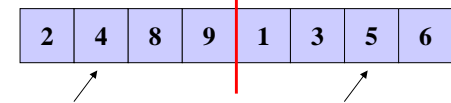
(After merge,
copy back to
original array)

Example, focus on merging

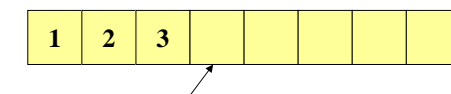
Start with:



After recursion:
(not magic ☺)



Merge:
Use 3 "fingers"
and 1 more array



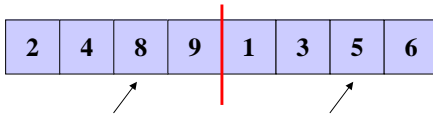
(After merge,
copy back to
original array)

Example, focus on merging

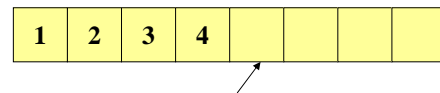
Start with:



After recursion:
(not magic ☺)



Merge:
Use 3 "fingers"
and 1 more array



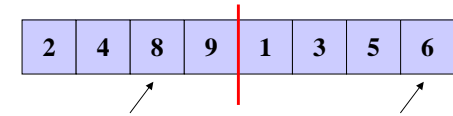
(After merge,
copy back to
original array)

Example, focus on merging

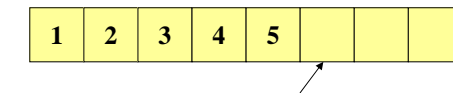
Start with:



After recursion:
(not magic ☺)



Merge:
Use 3 "fingers"
and 1 more array



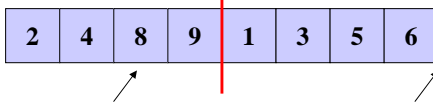
(After merge,
copy back to
original array)

Example, focus on merging

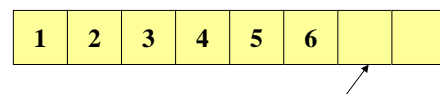
Start with:



After recursion:
(not magic ☺)



Merge:
Use 3 "fingers"
and 1 more array



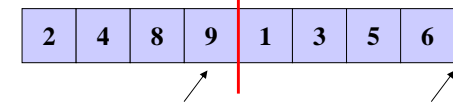
(After merge,
copy back to
original array)

Example, focus on merging

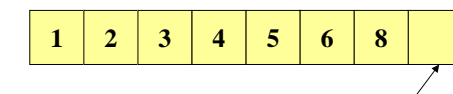
Start with:



After recursion:
(not magic ☺)



Merge:
Use 3 "fingers"
and 1 more array



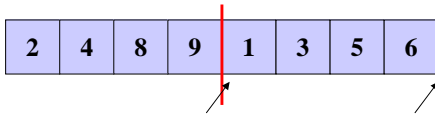
(After merge,
copy back to
original array)

Example, focus on merging

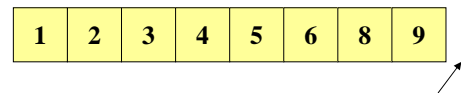
Start with:



After recursion:
(not magic ☺)



Merge:
Use 3 "fingers"
and 1 more array



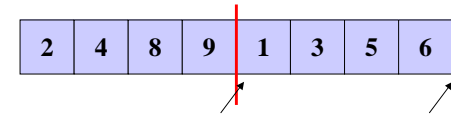
(After merge,
copy back to
original array)

Example, focus on merging

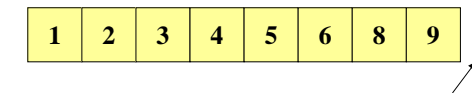
Start with:



After recursion:
(not magic ☺)

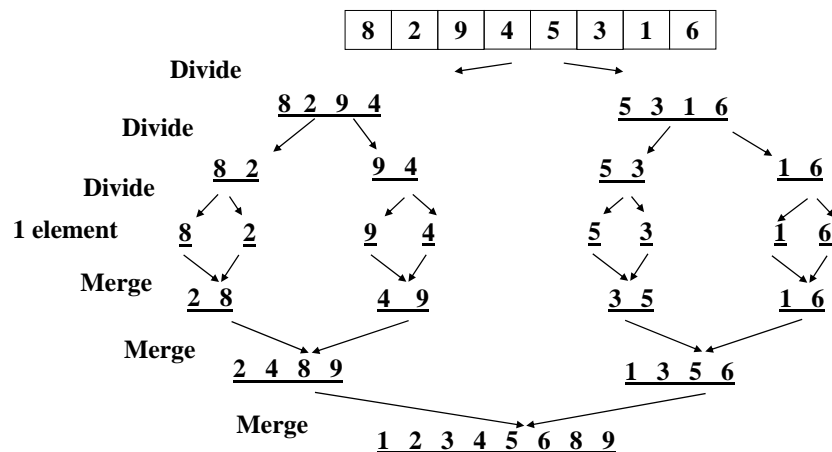


Merge:
Use 3 "fingers"
and 1 more array



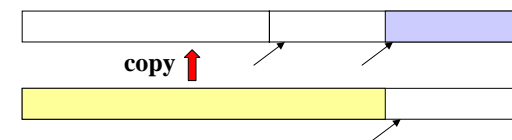
(After merge,
copy back to
original array)

Example, showing recursion

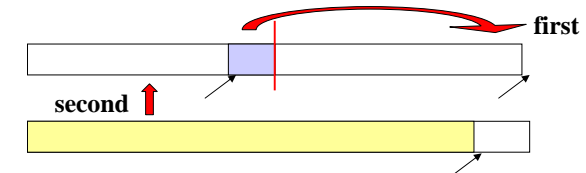


Some details: saving a little time

- In our example, we copied the "dregs" into the auxiliary array, but that's unnecessary right before copying back
 - If left-side finishes first, just stop the merge:



- If right-side finishes first, copy dregs directly into right side



Some details: saving space / copying

Simplest / worst approach:

Use a new auxiliary array of size $(hi-lo)$ for every merge

Better:

Use a new auxiliary array of size n for every merging stage

Better:

Reuse same auxiliary array of size n for every merging stage

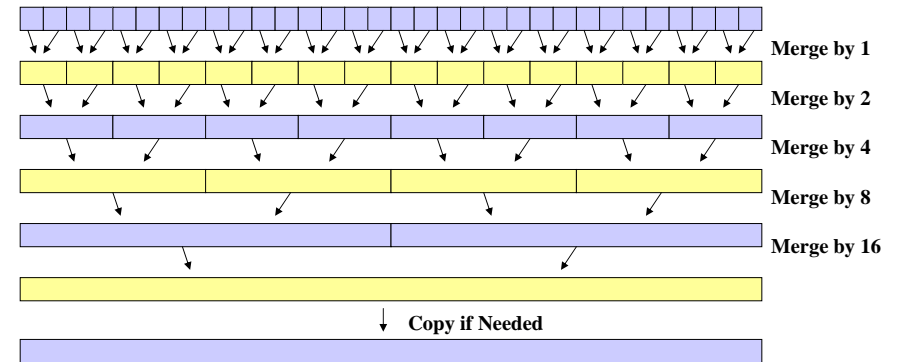
Best (but a little tricky):

Don't copy back – at 2nd, 4th, 6th, ... merging stages, use the original array as the auxiliary array and vice-versa

- Need one copy at end if number of stages is odd

Picture of the “best”

Arguably easier to code up without recursion at all



Linked lists and big data

We defined the sorting problem as over an array, but sometimes you want to sort linked lists

One approach:

- Convert to array: $O(n)$
- Sort: $O(n \log n)$
- Convert back to list: $O(n)$

Or: mergesort works very nicely on linked lists directly

- heapsort and quicksort do not
- insertion sort and selection sort do but they're slower

Mergesort is also the sort of choice for external sorting

- Linear merges minimize disk accesses

Analysis

Having defined an algorithm and argued it is correct, we should analyze its running time (and space):

To sort n elements, we:

- Return immediately if $n=1$
- Else do 2 subproblems of size $n/2$ and then an $O(n)$ merge

Recurrence relation:

$$T(1) = c_1$$

$$T(n) = 2T(n/2) + c_2n$$

One of the recurrence classics...

(For simplicity let constants be 1 – no effect on asymptotic answer)

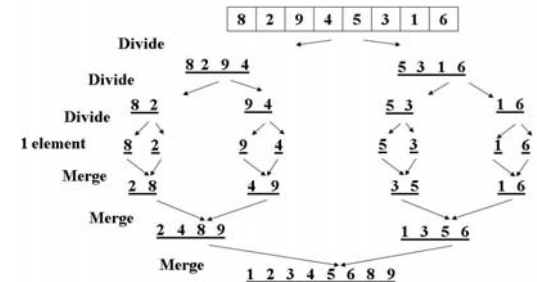
$$\begin{aligned}
 T(1) &= 1 && \text{So total is } 2^k T(n/2^k) + kn \text{ where} \\
 T(n) &= 2T(n/2) + n && \quad n/2^k = 1, \text{ i.e., } \log n = k \\
 &= 2(2T(n/4) + n/2) + n && \text{That is, } 2^{\log n} T(1) + n \log n \\
 &= 4T(n/4) + 2n && = n + n \log n \\
 &= 4(2T(n/8) + n/4) + 2n && = O(n \log n) \\
 &= 8T(n/8) + 3n \\
 &\dots \\
 &= 2^k T(n/2^k) + kn
 \end{aligned}$$

Or more intuitively...

This recurrence comes up often enough you should just “know” it’s $O(n \log n)$

Merge sort is relatively easy to intuit (best, worst, and average):

- The recursion “tree” will have $\log n$ height
- At each level we do a *total* amount of merging equal to n



Quicksort

- Also uses divide-and-conquer
- Does not need auxiliary space
- $O(n \log n)$ on average, but $O(n^2)$ worst-case
- Faster than mergesort in practice?
 - Often believed so
 - Does fewer copies and more comparisons, so it depends on the relative cost of these two operations!

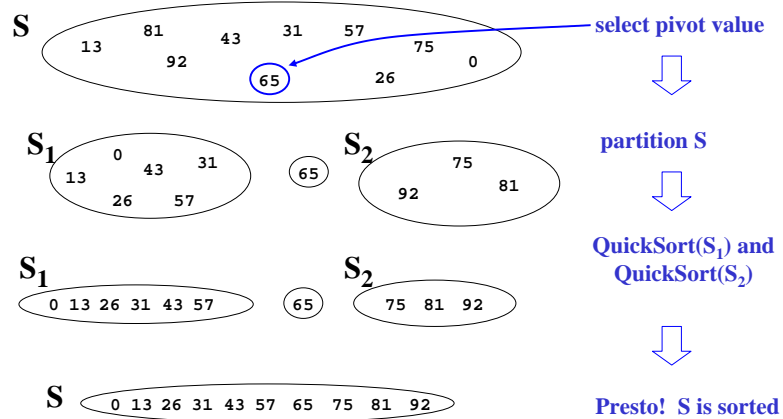
But we’re getting ahead of ourselves, how does it work...

Quicksort overview

1. Pick a pivot element
2. Partition all the data into:
 - A. The elements less than the pivot
 - B. The pivot
 - C. The elements greater than the pivot
3. Recursively sort A and C
4. The answer is, “as simple as A, B, C”

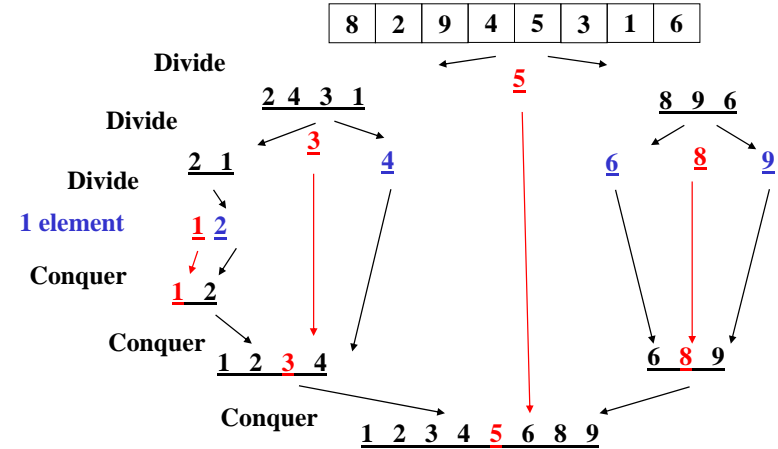
(Alas, there are some details lurking in this algorithm)

Think in terms of sets



[Weiss]

Example, showing recursion



Details

We haven't explained:

- How to pick the pivot element
 - Any choice is correct: data will end up sorted
 - But as analysis will show, want the two partitions to be about equal in size
- How to implement partitioning
 - In linear time
 - In place

Potential pivot rules

While sorting `arr` from `lo` (inclusive) to `hi` (exclusive)...

- Pick `arr[lo]` or `arr[hi-1]`
 - Fast, but worst-case is (mostly) sorted input
- Pick random element in the range
 - Does as well as any technique, but (pseudo)random number generation can be slow
 - (Still probably the most elegant approach)
- Median of 3, e.g., `arr[lo]`, `arr[hi-1]`, `arr[(hi+lo)/2]`
 - Common heuristic that tends to work well

Partitioning

- Conceptually simple, but hardest part to code up correctly
 - After picking pivot, need to partition in linear time in place
- One approach (there are slightly fancier ones):
 - Swap pivot with `arr[lo]`
 - Use two fingers `i` and `j`, starting at `lo+1` and `hi-1`
 - `while (i < j)`
 - `if (arr[j] > pivot) j--`
 - `else if (arr[i] < pivot) i++`
 - `else swap arr[i] with arr[j]`
 - Swap pivot with `arr[i]`

Example

- Step one: pick pivot as median of 3
 - `lo = 0, hi = 10`

0	1	2	3	4	5	6	7	8	9
8	1	4	9	0	3	5	2	7	6

- Step two: move pivot to the `lo` position

0	1	2	3	4	5	6	7	8	9
6	1	4	9	0	3	5	2	7	8

Example

Often have more than one swap during partition – this is a short example

Now partition in place

6	1	4	9	0	3	5	2	7	8
---	---	---	---	---	---	---	---	---	---

Move fingers

6	1	4	9	0	3	5	2	7	8
---	---	---	---	---	---	---	---	---	---

Swap

6	1	4	2	0	3	5	9	7	8
---	---	---	---	---	---	---	---	---	---

Move fingers

6	1	4	2	0	3	5	9	7	8
---	---	---	---	---	---	---	---	---	---

Move pivot

5	1	4	2	0	3	6	9	7	8
---	---	---	---	---	---	---	---	---	---

Analysis

- Best-case: Pivot is always the median
 - $T(0)=T(1)=1$
 - $T(n)=2T(n/2) + n$ -- linear-time partition
 - Same recurrence as mergesort: $O(n \log n)$
- Worst-case: Pivot is always smallest or largest element
 - $T(0)=T(1)=1$
 - $T(n) = 1T(n-1) + n$
 - Basically same recurrence as selection sort: $O(n^2)$
- Average-case (e.g., with random pivot)
 - $O(n \log n)$, not responsible for proof (in text)

Cutoffs

- For small n , all that recursion tends to cost more than doing a quadratic sort
 - Remember asymptotic complexity is for large n
- Common engineering technique: switch to a different algorithm for subproblems below a **cutoff**
 - Reasonable rule of thumb: use insertion sort for $n < 10$
- Notes:
 - Could also use a cutoff for merge sort
 - Cutoffs are also the norm with parallel algorithms
 - switch to sequential
 - None of this affects asymptotic complexity

Cutoff skeleton

```
void quicksort(int[] arr, int lo, int hi) {  
    if (hi - lo < CUTOFF)  
        insertionSort(arr, lo, hi);  
    else  
        ...  
}
```

Notice how this cuts out the vast majority of the recursive calls

- Think of the recursive calls to quicksort as a tree
- Trims out the bottom layers of the tree