# Unit #4: Sorting
## CPSC 221: Algorithms and Data Structures

Will Evans and Jan Manuch

2016W1

# Unit Outline

- ▶ Comparing Sorting Algorithms
- ▶ Heapsort
- ▶ Mergesort
- ▶ Quicksort
- ▶ More Comparisons
- ▶ Complexity of Sorting

# Learning Goals

- Describe, apply, and compare various sorting algorithms.
- Analyze the complexity of these sorting algorithms.
- Explain the difference between the complexity of a problem (sorting) and the complexity of a particular algorithm for solving that problem.

# How to Measure Sorting Algorithms
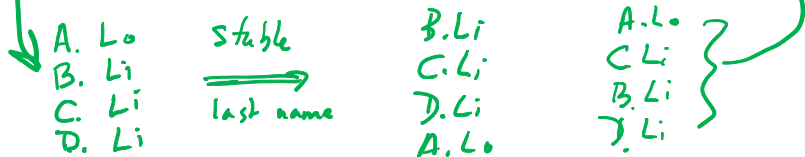
- Computational complexity (a.k.a. runtime)
  - Worst case
  - Average case
  - Best case
    How often is the input sorted, reverse sorted, or "almost" sorted ($k$ swaps from sorted where $k \ll n$)?

- Stability: What happens to elements with identical keys? ~~Why do we care?~~

- Memory Usage: How much extra memory is used?

A. L₀
B. Li
C. Li
D. Li

Stable
⟶
last name

B. Li
C. Li
D. Li
A. L₀

A. L₀
C Li
B. Li
D. Li

# Insertion Sort: Running Time

At the start of iteration $i$, the first $i$ elements in the array are sorted, and we insert the $(i+1)$st element into its proper place.

Worst case:

$$\sum_{i=1}^{n-1} i = \frac{n(n-1)}{2} = \Theta(n^2)$$

Best case:

$$\sum_{i=1}^{n-1} 1 = n-1 = \Theta(n)$$

Avg case?

$$\sum_{i=1}^{n-1} i/2 = \frac{1}{2}\left(\sum_{i=1}^{n-1} i\right) = \Theta(n^2)$$

# Insertion Sort: Stability & Memory

At the start of iteration $i$, the first $i$ elements in the array are sorted, and we insert the $(i+1)$st element into its proper place.

## Easily made stable:

"proper place" is **largest** $j$ such that $A[j-1] \leq$ new element.



## Memory:

Sorting is done **in-place**, meaning only a constant number of extra memory locations are used.

# Heapsort

$\Theta(n \lg n)$    1964 Williams    1964 Floyd
(heapify)

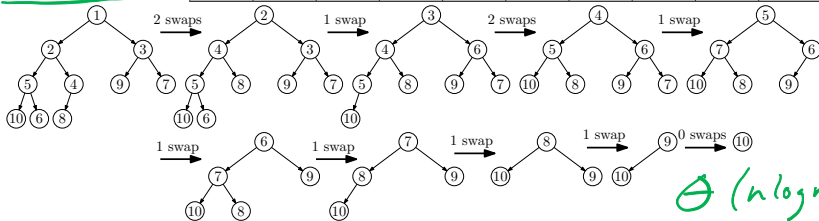1. Heapify input array.    $\Theta(n)$
2. Repeat *n* times: Perform `deleteMin`

time $\leq \lg n + \lg(n-1) + \lg(n-2) \ldots + \lg(2)$
first deleteMin    $= \lg(n!)$

Worst case:    $= \Theta(n \log n)$



Best case[1]:

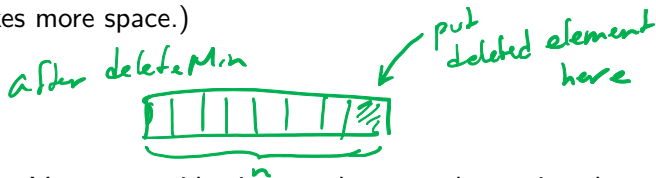| 1 | 2 | 3 | 5 | 4 | 9 | 7 | 10 | 6 | 8 |

$\Theta(n \log n)$

---
[1] Schaffer & Sedgewick, The Analysis of Heapsort, *J. Algorithms* **15** (1993) 76–100.

# Heapsort: Stability & Memory

1. Heapify input array.
2. Repeat $n$ times: Perform deleteMin

### Not stable:
Hack: Use index in input array to break comparison ties.
(but this takes more space.)

*after deleteMin*

*put deleted element here*

### Memory:

- **in-place**. You can avoid using another array by storing the result of the $i$th deleteMin in heap location $n - i$, except the array is then sorted in reverse order, so use a Max-Heap (and deleteMax).
- Far-apart array accesses ruin cache performance.

# Mergesort

Mergesort is a "divide and conquer" algorithm.

1. If the array has 0 or 1 elements, it's sorted. Stop. $\Theta(1)$
2. Split the array into two approximately equal-sized halves. $\Theta(1)$
3. Sort each half recursively (using Mergesort) $2T(n/2)$
4. Merge the sorted halves to produce one sorted result: $\Theta(n)$
   - Consider the two halves to be queues.
   - Repeatedly dequeue the smaller of the two front elements (or dequeue the only front element if one queue is empty) and add it to the result.
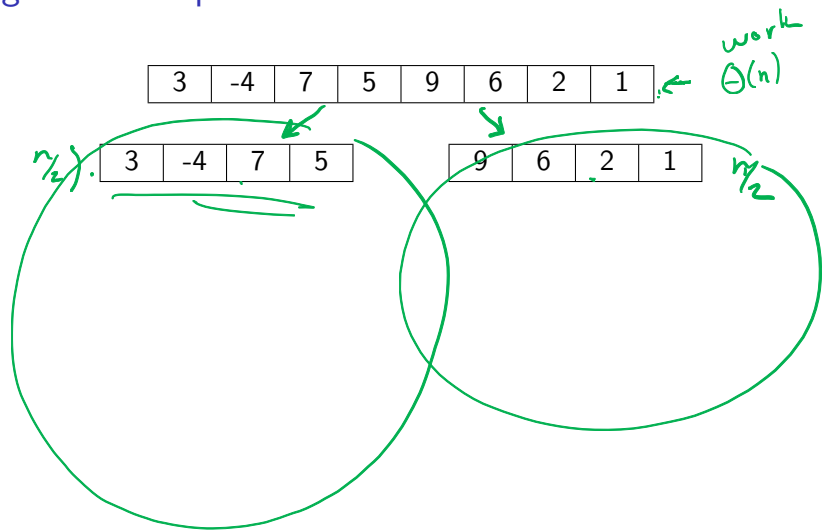
$T(n) = $ running time of Mergesort

recurrence relation .... )

$T(n) = \Theta(n \log n)$

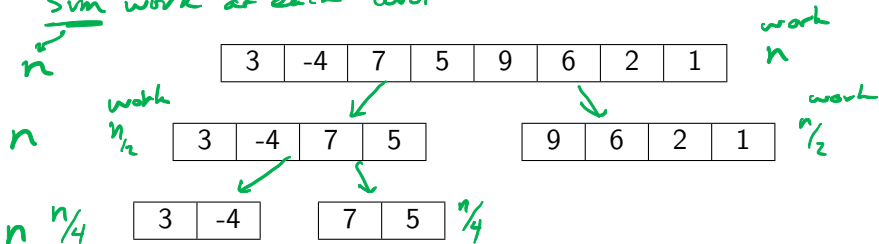# Mergesort Example

*recursion tree*

| 3 | -4 | 7 | 5 | 9 | 6 | 2 | 1 |
|---|----|---|---|---|---|---|---|

# Mergesort Example

# Mergesort Example

sum work at each level



n

n

n /4

n

n

n

n
_____

n lg n

work
n

work
n/2

work
n/4

# Mergesort Example

| 3 | -4 | 7 | 5 | 9 | 6 | 2 | 1 |
|---|----|---|---|---|---|---|---|

| 3 | -4 | 7 | 5 |
|---|----|---|---|

| 9 | 6 | 2 | 1 |
|---|---|---|---|

| 3 | -4 |
|---|----|

| 7 | 5 |
|---|---|

| 3 |
|---|

| -4 |
|----|

# Mergesort Example

recursion tree

recursive calls

| 3 | -4 | 7 | 5 | 9 | 6 | 2 | 1 |

| 3 | -4 | 7 | 5 |    | 9 | 6 | 2 | 1 |

| 3 | -4 |    | 7 | 5 |

| 3 |    | -4 |

| -4 | 3 | *

return from calls

# Mergesort Example

# Mergesort Example

| 3 | -4 | 7 | 5 | 9 | 6 | 2 | 1 |
|---|----|---|---|---|---|---|---|

| 3 | -4 | 7 | 5 |
|---|----|---|---|

| 9 | 6 | 2 | 1 |
|---|---|---|---|

| 3 | -4 |
|---|----|

| 7 | 5 |
|---|---|

| 3 |
|---|

| -4 |
|----|

| 7 |
|---|

| 5 |
|---|

| -4 | 3 | *
|----|---|

| 5 | 7 |
|---|---|

# Mergesort Example

| 3 | -4 | 7 | 5 | 9 | 6 | 2 | 1 |

| 3 | -4 | 7 | 5 |    | 9 | 6 | 2 | 1 |

| 3 | -4 |    | 7 | 5 |

| 3 |    | -4 |    | 7 |    | 5 |

| -4 | 3 |* | 5 | 7 |

| -4 | 3 | 5 | 7 |

# Mergesort Example

# Mergesort Code

*helper*

```
void msort(int x[], int lo, int hi, int tmp[]) {
  if (lo >= hi) return;
  int mid = (lo+hi)/2;
  msort(x, lo, mid, tmp);
  msort(x, mid+1, hi, tmp);
  merge(x, lo, mid, hi, tmp);
}

void mergesort(int x[], int n) {
  int *tmp = new int[n];
  msort(x, 0, n-1, tmp);
  delete[] tmp;
}
```
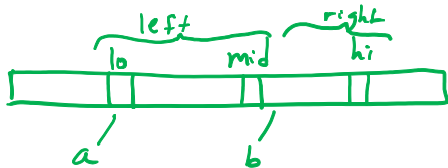
*allocate space*

# Merge Code



```
void merge(int x[],int lo,int mid,int hi,int tmp[]) {
  int a = lo, b = mid+1;
  for( int k = lo; k <= hi; k++ ) {
    if( a <= mid && (b > hi || x[a] < x[b]) )
      tmp[k] = x[a++];
    else tmp[k] = x[b++];
  }
  for( int k = lo; k <= hi; k++ )
    x[k] = tmp[k];
}
```

invariant: tmp[lo..k] contains smallest
k-lo+1 elements from x[lo..hi]
in sorted order.

# Sample Merge Steps

```
merge( x, 0, 0, 1, tmp ); // step *
```

|  x : | 3 | -4 | 7 | 5 | 9 | 6 | 2 | 1 |
|---|---|---|---|---|---|---|---|---|

| tmp : | -4 | 3 | ? | ? | ? | ? | ? | ? |
|---|---|---|---|---|---|---|---|---|

|  x : | -4 | 3 | 7 | 5 | 9 | 6 | 2 | 1 |
|---|---|---|---|---|---|---|---|---|

```
merge( x, 4, 5, 7, tmp ); // step **
```

|  x : | -4 | 3 | 5 | 7 | 6 | 9 | 1 | 2 |
|---|---|---|---|---|---|---|---|---|

| tmp : | ? | ? | ? | ? | 1 | 2 | 6 | 9 |
|---|---|---|---|---|---|---|---|---|

|  x : | -4 | 3 | 5 | 7 | 1 | 2 | 6 | 9 |
|---|---|---|---|---|---|---|---|---|

```
merge( x, 0, 3, 7, tmp ); // is the final step
```

# Mergesort: Stability & Memory

### Stable:
Dequeue from the left queue if the two front elements are equal.

### Memory:
Not easy to implement without using $\Omega(n)$ extra space, so it is not viewed as an in-place sort.

Even if you could avoid allocating tmp[], it still uses extra memory: the call stack $\Omega(\log n)$

# Quicksort (C.A.R. Hoare 1961)

In practice, one of the fastest sorting algorithms.

1. Pick a **pivot** _first element_

| 2 | -4 | 6 | 1 | 5 | -3 | 3 | 7 |

2. Reorder the array such that all elements $<$ pivot are to its left, and all elements $\geq$ pivot are to its right
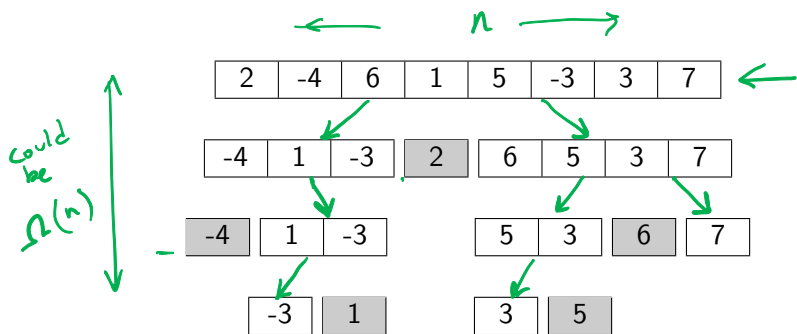
| -4 | 1 | -3 | 2 | 6 | 5 | 3 | 7 |

left partition   pivot   right partition

3. Recursively sort each partition.

base case?

# Quicksort Visually

# Quicksort by Jon Bentley

*Programming Pearls*

```
void qsort(int x[], int lo, int hi) {
  int i, p;
  if (lo >= hi) return;
  p = lo;
  for( i=lo+1; i <= hi; i++ )
    if( x[i] < x[lo] ) swap(x[++p], x[i]);
  swap(x[lo], x[p]);
  qsort(x, lo, p-1);
  qsort(x, p+1, hi);
}

void quicksort(int x[], int n) {
  qsort(x, 0, n-1);
}
```

*Lomuto*

*pivot index*

*loop invariant*

$(x[lo] = pivot)$

$x[lo+1 \ldots p] < pivot$

$x[p+1 \ldots i-1] \geq pivot$

# Quicksort Example (using Bentley's Algorithm)



```
if( x[i] < x[lo] ) swap(x[++p], x[i]);
```

$x[p+1] \quad x[i]$

| lo | | | | | | | hi |
|---|---|---|---|---|---|---|---|
| 2 | -4 | 6 | 1 | 5 | -3 | 3 | 7 |

p  i

| 2 | -4 | 6 | 1 | 5 | -3 | 3 | 7 |
|---|---|---|---|---|---|---|---|

p  i

| 2 | -4 | 6 | 1 | 5 | -3 | 3 | 7 |
|---|---|---|---|---|---|---|---|

p  i

| 2 | -4 | 1 | 6 | 5 | -3 | 3 | 7 |
|---|---|---|---|---|---|---|---|

p  i

```
if( x[i] < x[lo] ) swap(x[++p], x[i]);
```

# Quicksort Example (using Bentley's Algorithm)



```
     lo                          hi
   | 2 | -4 | 1 | -3 | 5 | 6 | 3 | 7 |
                     p            i
```

for-loop ends

```
swap(x[lo], x[p]);
```

```
     lo                          hi
   | -3 | -4 | 1 | 2 | 5 | 6 | 3 | 7 |
                     p            i
```

```
qsort(x, lo, p-1);
qsort(x, p+1, hi);
```

```
   | -4 | -3 | 1 | 2 | 3 | 5 | 6 | 7 |
```

# Quicksort: Running Time

Running time is proportional to number of comparisons so...
Let's count comparisons.

1. Pick a pivot.
   Zero comparisons

2. Reorder (partition) array around the pivot.
   Quicksort compares each element to the pivot.
   $n - 1$ comparisons

3. Recursively sort each partition.
   Depends on the size of the partitions.

   ▶ If the partitions have size $n/2$ (or any constant fraction of $n$),
     the runtime is $\Theta(n \log n)$ (like Mergesort).

   ▶ In the worst case, however, we might create partitions with
     sizes 0 and $n - 1$.

# Quicksort Visually: Worst case

recursion tree



$T(0) = 1$

$T(1) = 1$

$T(n) = n-1 + T(0) + T(n-1) = n + T(n-1)$
$$\underset{\|}{1}$$
$$= n + (n-1) + T(n-2)$$

# Quicksort: Worst Case

If this happens at every partition...
Quicksort makes $n - 1$ comparisons in the first partition and
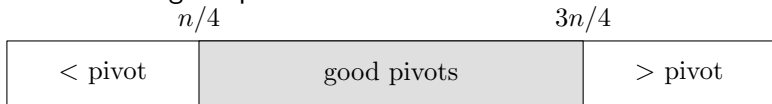recurses on a problem of size 0 and size $n - 1$:

$$T(n) = (n - 1) + T(0) + T(n - 1) = (n - 1) + T(n - 1)$$
$$= (n - 1) + (n - 2) + T(n - 2)$$
$$\vdots$$
$$= \sum_{i=1}^{n-1} i = (n - 1)(n - 2)/2$$

This is $\Theta(n^2)$ comparisons.

# Quicksort: Average Case (Intuition)

- On an average input (i.e., random order of $n$ items), our chosen pivot is equally likely to be the $i$th smallest for any $i = 1, 2, \ldots, n$.

- With probability $1/2$, our pivot will be from the middle $n/2$ elements – a good pivot.

| | | |
|---|---|---|
| $n/4$ | | $3n/4$ |
| < pivot | good pivots | > pivot |

- Any good pivot creates two partitions of size at most $3n/4$.
- We expect to pick one good pivot every two tries.
- Expected number of splits is at most $2 \log_{4/3} n \in O(\log n)$.
- $O(n \log n)$ total comparisons. True, but this intuition is not a proof.

$$C(n) = \frac{1}{n}\left(C(0) + C(n-1)\right) + \frac{1}{n}\left(C(1) + C(n-2)\right)$$

$$n-1+$$

$$= \frac{1}{n}\sum_{i=0}^{n-1}\left(C(i) + C(n-i-1)\right)$$

# Quicksort: Stability & Memory

### Stable:
Can be made stable, most easily by using more memory.

### Memory:
In-place sort.

# Compare: Running Times (100 samples)

| n | Insertion | | Heap | | Merge | | Quick | |
|---|---|---|---|---|---|---|---|---|
| | avg | max | avg | max | avg | max | avg | max |
| 100,000 | 11.20s | 16.37s | 0.04s | 0.08s | 0.03s | 0.04s | 0.02s | 0.04s |
| 200,000 | 36.97s | 60.01s | 0.08s | 0.16s | 0.06s | 0.11s | 0.06s | 0.16s |
| 400,000 | 172.36s | 505.38s | 0.56s | 1.74s | 0.54s | 0.91s | 0.46s | 0.69s |
| 800000 | | | 0.37s | 0.83s | 0.21s | 0.35s | 0.19s | 0.32s |
| 1600000 | | | 0.93s | 1.77s | 0.52s | 1.12s | 0.44s | 0.78s |
| 3200000 | | | 2.07s | 3.04s | 1.01s | 1.95s | 0.91s | 1.44s |
| 6400000 | | | 4.76s | 7.54s | 2.18s | 3.88s | 1.97s | 3.45s |
| 12800000 | | | 10.65s | 12.38s | 4.56s | 7.01s | 4.13s | 5.94s |

Code is from lecture notes and labs (not optimized).

# Compare: Quick, Merge, Heap, and Insert Sort

Running Time

|  | $\Theta(n)$ | $\Theta(n \log n)$ | $\Theta(n^2)$ |
|---|---|---|---|
| Best case: | Insert | Quick, Merge, Heap | |
| Average case: | | Quick, Merge, Heap | Insert |
| Worst case: | | Merge, Heap | Quick, Insert |
| "Real" data: | | Quick < Merge < Heap < Insert | |

Some Quick/Merge implementations use Insert on small arrays (base cases).

Some results depend on the implementation! For example, an initial check whether the last element of the left subarray is less than the first of the right can make Merge's best case linear.

# Compare: Quick, Merge, Heap, and Insert Sort

Stability
  Stable (easy):          Insert, Merge (prefer the left of the two
                          sorted subarrays on ties)
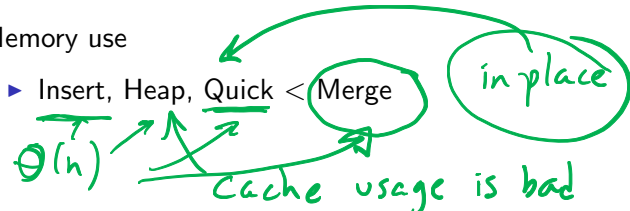  Stable (with effort):   Quick
  Unstable:               Heap

Memory use

- Insert, Heap, Quick $<$ Merge

$\Theta(n)$

in place

Cache usage is bad

# Complexity of the Sorting Problem

The **complexity** of a problem is the complexity of the best algorithm for that problem.
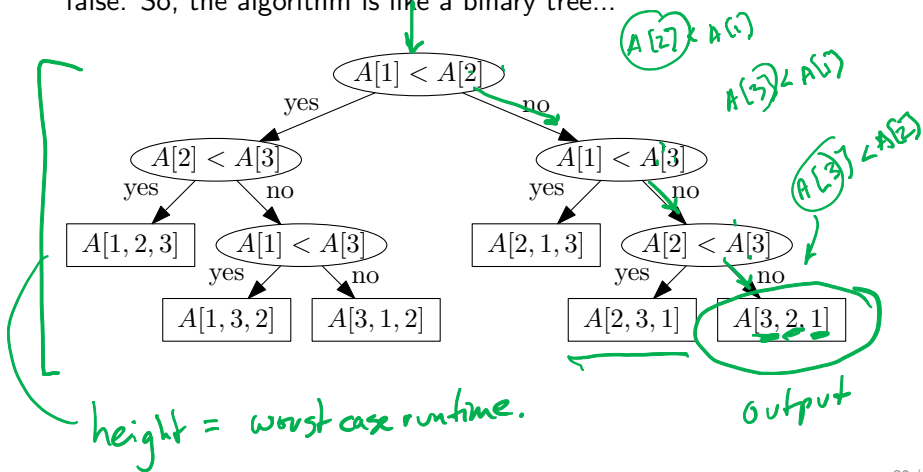
How powerful is our computer?

We'll only consider **comparison-based** algorithms. They can compare two array elements in constant time.
They cannot manipulate array elements in any other way.
For example, they cannot assume that the elements are numbers and perform arithmetic operations (like division) on them.

Insertion, Heap, Merge, and Quick sort are comparison-based.
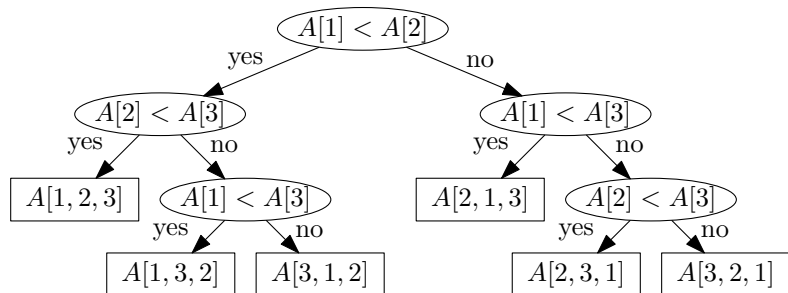Radix sort is not.

# Comparison-based algorithms using a Decision Tree model

Each comparison is a "choice point" in the algorithm: the algorithm can do one thing if the comparison is true and another if false. So, the algorithm is like a binary tree...



*(handwritten annotations:)* $A[2] < A[1]$, $A[3] < A[1]$, $A[3] < A[2]$, output, height = worst case runtime.

Decision tree:

$A[1] < A[2]$
- yes → $A[2] < A[3]$
  - yes → $A[1, 2, 3]$
  - no → $A[1] < A[3]$
    - yes → $A[1, 3, 2]$
    - no → $A[3, 1, 2]$
- no → $A[1] < A[3]$
  - yes → $A[2, 1, 3]$
  - no → $A[2] < A[3]$
    - yes → $A[2, 3, 1]$
    - no → $A[3, 2, 1]$

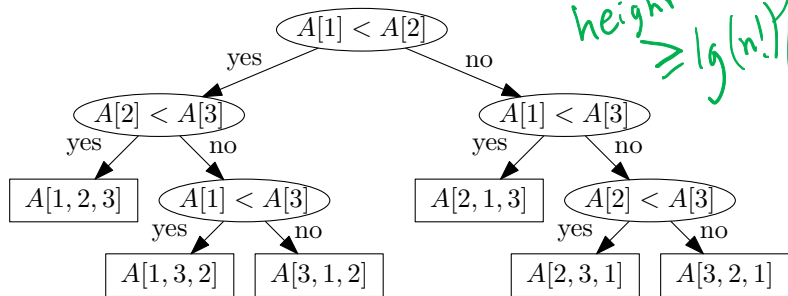# Complexity of the Sorting Problem

- This is the decision tree representation of Insertion Sort on inputs of size $n = 3$.
- Each leaf outputs the input array in some particular order. For example, $A[3, 1, 2]$ means output $A[3]$, $A[1]$, $A[2]$.

# Complexity of the Sorting Problem

- There are $n!$ possible output orderings of an input array of size $n$.

- There must be a leaf for each one, otherwise the algorithm fails to sort.

  - For example, if leaf $A[2, 3, 1]$ doesn't exist then the algorithm cannot sort [cat, ant, bee].



*height $\geq \lg(n!)$* (handwritten annotation)

Decision tree structure:
- Root: $A[1] < A[2]$
  - yes: $A[2] < A[3]$
    - yes: $A[1, 2, 3]$
    - no: $A[1] < A[3]$
      - yes: $A[1, 3, 2]$
      - no: $A[3, 1, 2]$
  - no: $A[1] < A[3]$
    - yes: $A[2, 1, 3]$
    - no: $A[2] < A[3]$
      - yes: $A[2, 3, 1]$
      - no: $A[3, 2, 1]$

# Complexity of the Sorting Problem

- ▶ The number of leaves is at least $n!$.
- ▶ The height of the decision tree is at least $\lceil \lg(n!) \rceil$.
- ▶ The number of comparisons made *in the worst case* is at least $\lceil \lg(n!) \rceil$.
- ▶ This is true for **any comparison-based sorting algorithm** so the complexity of the sorting problem is $\Omega(n \log n)$.