

All of the sorting algorithms we have considered so far have been **comparison sorts**. That is, we only use comparisons to determine the relative order of elements.

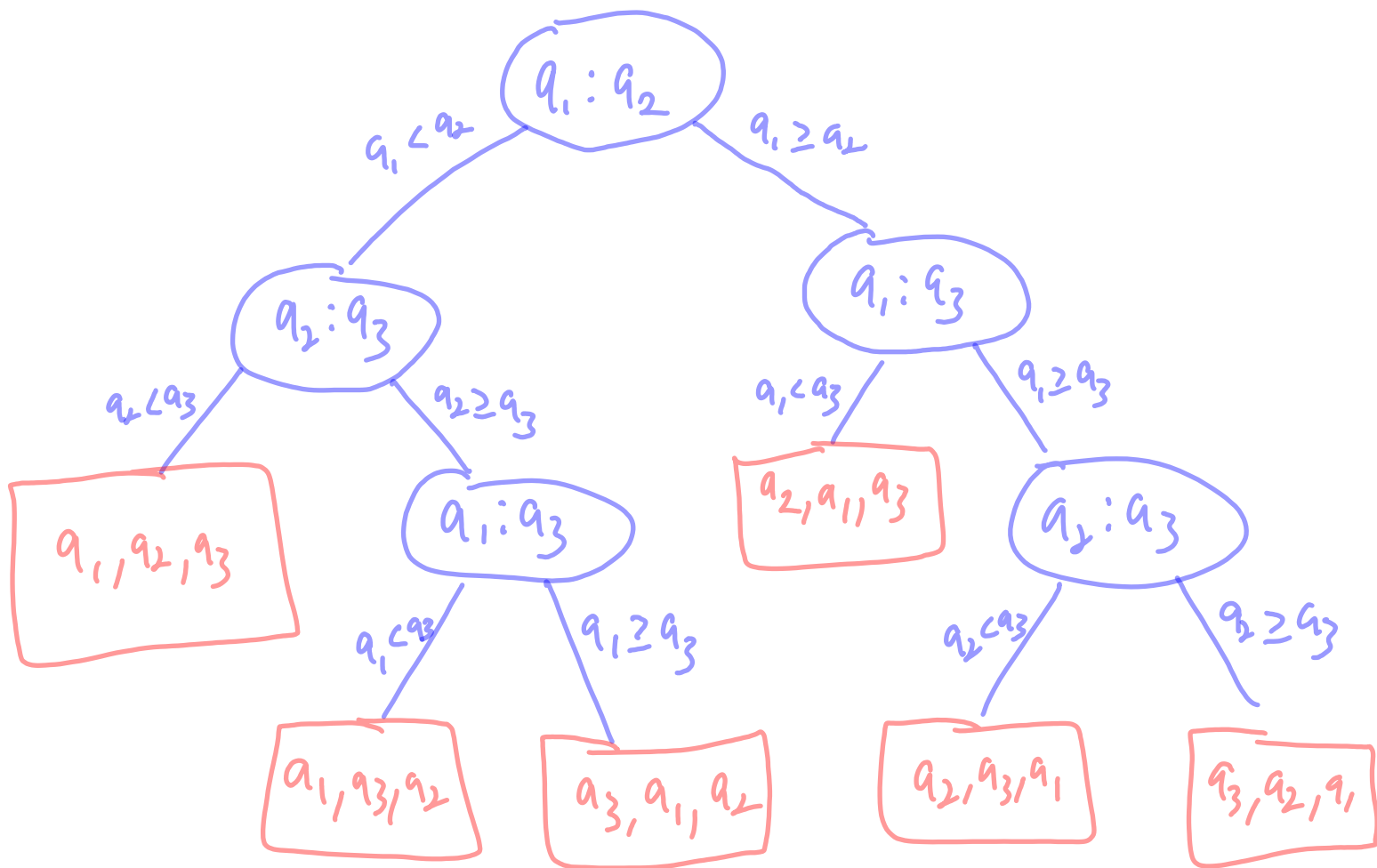
The best algorithm we have seen in terms of worst-case running time is  $O(n \log n)$ .

Can we do better? That is, is it possible to design a comparison-based sorting algorithm which has worst-case running time  $o(n \log n)$  (e.g.  $O(n)$  or  $O(n \log \log n)$ )?

We will use a **decision tree** to help us answer this question.

# Insertion Sort Decision Tree

Suppose  $n=3$ , we want to sort  $[a_1, a_2, a_3]$ .



We can construct a decision tree which models the execution of any comparison sorting algorithm.

- One tree per input size  $n$ .
- The tree will contain *all* possible comparisons that could be executed for *any* input of size  $n$ .
- For one input, only one path to a leaf is executed.
- The running time is equal to the length of the path taken.
- Worst-case running time is the length of the longest path in the tree (the height of the tree).

**Theorem:** Any decision tree of a comparison sorting algorithm on an input of size  $n$  must have height  $\Omega(n \log n)$ .

Proof: There must be at least  $n!$  leaves.

The tree is a binary tree, because comparing is a binary operation. A binary tree of height  $h$  has at most  $2^h$  leaves.

$$\Rightarrow 2^h \geq n! \Rightarrow h \geq \log n!$$

$$n! = \underbrace{n(n-1)(n-2) \cdots \left(\frac{n}{2}\right)}_{\geq \frac{n}{2}} \underbrace{\left(\frac{n}{2}-1\right) \cdots 1}_{\geq 1}$$
$$\geq \left(\frac{n}{2}\right)^{\left(\frac{n}{2}\right)}$$

$$h \geq \log n! \geq \log \left(\left(\frac{n}{2}\right)^{\left(\frac{n}{2}\right)}\right) = \frac{n}{2} \cdot \log \frac{n}{2} = \Theta(n \log n)$$

Recall we wanted to answer the following question.

- The best algorithm we have seen in terms of worst-case running time is  $O(n \log n)$ . Can we do better? That is, is it possible to design a comparison-based sorting algorithm which has worst-case running time  $o(n \log n)$  (e.g.  $O(n)$  or  $O(n \log \log n)$ )?

The theorem we proved implies that the answer to this question is no.

**Corollary:** Merge Sort is an asymptotically optimal comparison sorting algorithm.

not

Can we do better if we do use comparisons in our algorithm?

### Counting Sort:

- Input:  $A[1 \dots n]$ , where each  $A[j] \in \{1, 2, \dots, k\}$ .
- Output:  $B[1 \dots n]$ , sorted.
- Auxiliary storage:  $C[1 \dots k]$ .

# Counting Sort

1. for  $i=1$  to  $k$   
 $C[i] = 0$

}  $\Theta(k)$

2. for  $j=1$  to  $n$

$$C[A[j]] = C[A[j]] + 1$$

}  $\Theta(n)$

3. for  $i=2$  to  $k$

$$C[i] = C[i] + C[i-1]$$

}  $\Theta(k)$

4. for  $j=n$  to  $1$

$$B[C[A[j]]] = A[j]$$

$$C[A[j]]--$$

}  $\Theta(n)$

A = 

4	1	3	4	3
---	---	---	---	---

C: 

1	0	2	2
---	---	---	---

B: 

1	3	3	4	4
---	---	---	---	---

  
1 2 3 4 5

C': 

<del>1</del>	1	<del>3</del>	<del>8</del>
--------------	---	--------------	--------------

  
0                  2 4  
                     1 3

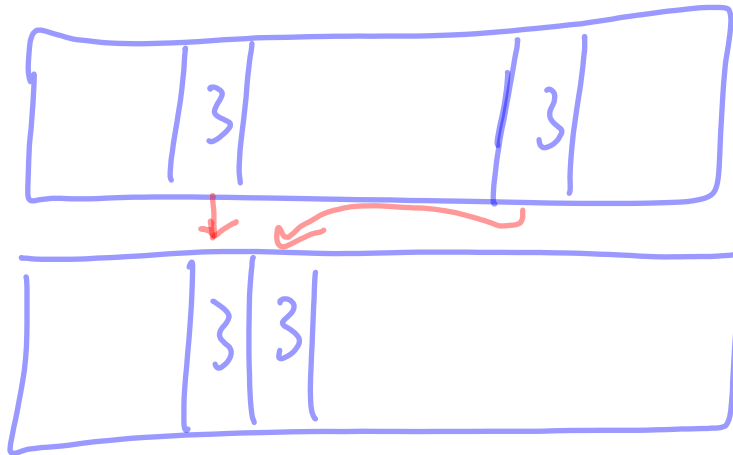


Note that we did not ever directly compare any two elements of  $A$ .

Also note that if  $k = O(n)$ , then counting sort takes  $\Theta(n)$  time (which is better than the  $\Omega(n \log n)$  lower bound of comparison sorting algorithms).

In the worst case,  $k$  can be arbitrarily large, and therefore the running time of the algorithm can be arbitrarily large (with respect to  $n$ ).

Counting sort is a **stable** sort. That is, if the same value appears multiple times, then their input order is preserved in the output array.



## Radix Sort:

- Digit-by-digit sort.
- Idea: sort the input on the *least-significant digit* first (i.e. digits from right to left) using a stable sorting algorithm (like counting sort).

329	720	720	329
457	355	329	355
657	436	436	436
839	457	839	457
436	657	355	657
720	329	457	720
355	839	657	839
↑	↑	↑	

Note that we can convert a decimal number to any other base (i.e. binary or hexadecimal) and sort these numbers instead (sorted order of binary/hex will be the same as sorted order for decimal).

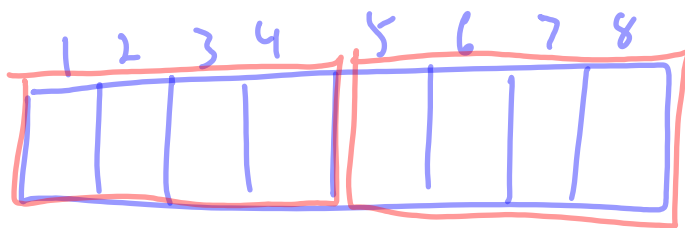
Note that if we convert the number to binary, then the number of digits for radix sort increases, but the range of values for counting sort decreases ( $k = 2$ ).

Note that if we convert the number to hexadecimal, then the number of digits for radix sort decreases, but the range of values for counting sort increases ( $k = 16$ ).

What base should we use to minimize the running time of the algorithm?

We can view the problem as sorting  $n$  binary numbers of  $b$  bits each.

Consider some positive integer  $r \leq b$ . We can view each number as having  $b/r$  digits in the range  $[0, 2^r - 1]$ . Then the algorithm would make  $b/r$  calls to counting sort with  $k = 2^r$ .



$$b=8$$
$$r=4$$

What choice of  $r$  should we make to minimize the running time of the algorithm.

Counting sort takes  $\Theta(n + k)$  time, and our  $k = 2^r$  which implies each call to counting sort takes time  $\Theta(n + 2^r)$ .

Since there are  $b/r$  calls to counting sort, we have the overall running time is  $\Theta(\frac{b}{r}(n + 2^r))$ . What choice of  $r$  minimizes this function?

$$\begin{aligned} 2^r &= n \\ \log_2 2^r &= \log_2 n \\ \parallel \\ r &= \log_2 n \end{aligned}$$

So with  $r = \log n$ , we get that the running time of radix sort when the input is  $n$  numbers of  $b$  bits each is  $\Theta(\frac{bn}{\log n})$ .

Therefore, if the range of values is a constant (and therefore  $b$  is a constant), then the algorithm runs in  $o(n)$  time.

If the range of values is from 0 to  $n^d - 1$ , then  $b = d \log n$ . The algorithm would have running time  $\Theta(dn)$ .

If the range of values is arbitrarily large, then the running time of radix sort is arbitrarily large.