# SORTING

Michael Tsai

2017/3/28

# Sorting

- Definition:
- Input: $\langle a_1, a_2, \ldots, a_n \rangle$ a sequence of $n$ numbers
- Output: $\langle a_1', a_2', \ldots, a_n' \rangle$ is a permutation (reordering) of the original sequence, such that $a_1' \leq a_2' \leq \cdots \leq a_n'$

- In reality, $a_i$ is the key of a record (of multiple fields) (e.g., student ID)
- In a record, the data fields other than the key is called **satellite data**
- If satellite data is large in size, we will only sort the pointers pointing to the records. (avoiding moving the data)

# Applications of Sorting

- Example 1: Looking for an item in a list

- Q: How do we look for an item in an **unsorted** list?
- A: We likely can only linearly traverse the list from the beginning.
- →$O(n)$

- Q: What if it is sorted?
- A: We can do binary search →$O(\log n)$

- But, how much time do we need for sorting? (pre-processing)

# Applications of Sorting

- Example 2:
- Compare to see if two lists are identical (list all different items)
- The two lists are $n$ and $m$ in length

- Q: What if they are **unsorted**?
- Compare the $1^{st}$ item in list 1 with (m-1) items in list 2
- Compare the $2^{nd}$ item in list 1 with (m-1) items in list 2
- ...
- Compare the n-th item in list 1 with (m-1) items in list 2
- $O(nm)$ time is needed

- Q: What if they are sorted?
- A: $O(n + m)$

- Again, do not forget we also need time for sorting. But, how much?
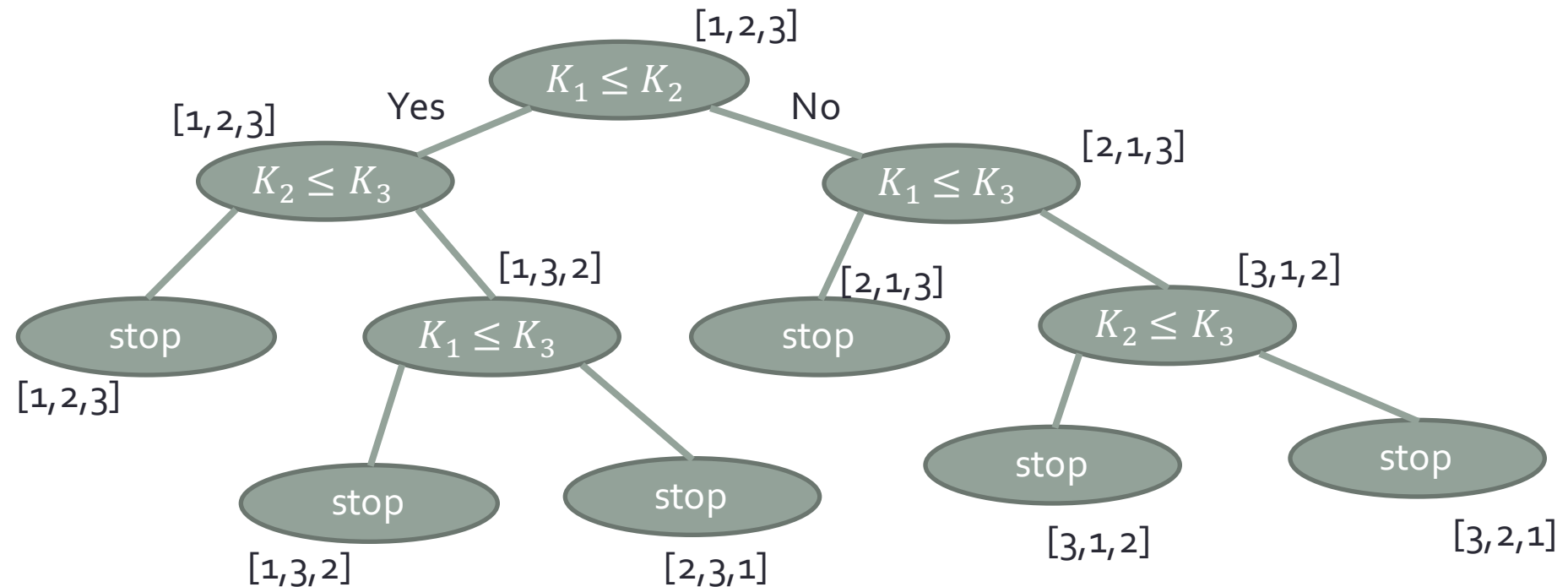
# Categories of Sorting Algo.

- Internal Sort:
  - Place all data in the memory

- External Sort:
  - The data is too large to fit it entirely in the memory.
  - Some need to be temporarily placed onto other (slower) storage, e.g., hard drive, flash disk, network storage, etc.

- In this lecture, we will only discuss **internal sort**.
- Storage is **cheap** nowadays. In most cases, only internal sort is needed.
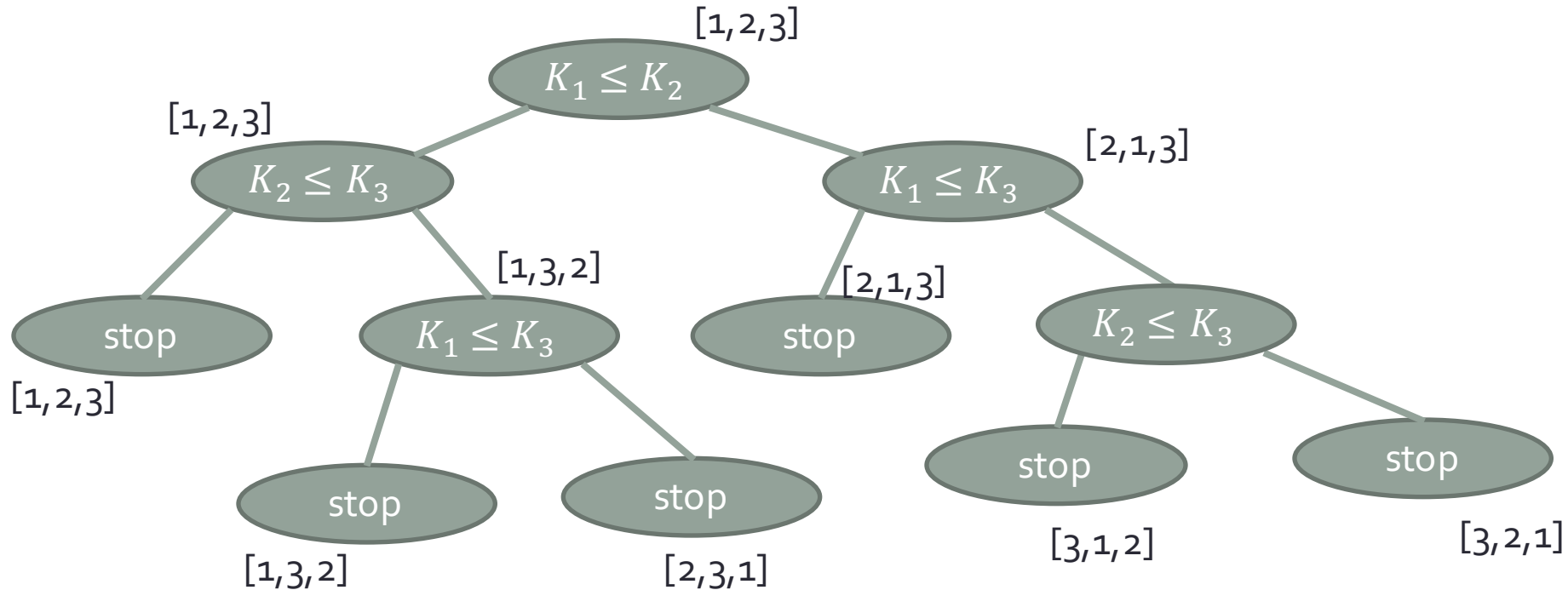
# Some terms related to sorting

- **Stability:**

- If $a_i = a_j$ (equal key value)，then they maintain the same order before and after sorting.

- **In-place:**

- Directly sort the keys at **their current memory locations**. Therefore, only O(1) additional space is needed for sorting.

- **Adaptability:**

- If **part of the sequence is sorted**, then the time complexity of the sorting algorithm reduces.

# How fast can we sort?

- Assumption: compare and swap
- Compare: compare two items in the list
- Swap: Swap the locations of these two items
- How much time do we need in the worst case?

# Decision tree for sorting

$[1,2,3]$

$K_1 \leq K_2$

$[1,2,3]$

$K_2 \leq K_3$

$[2,1,3]$

$K_1 \leq K_3$

stop

$[1,2,3]$

$[1,3,2]$

$K_1 \leq K_3$

$[2,1,3]$

stop

$K_2 \leq K_3$

stop

$[1,3,2]$
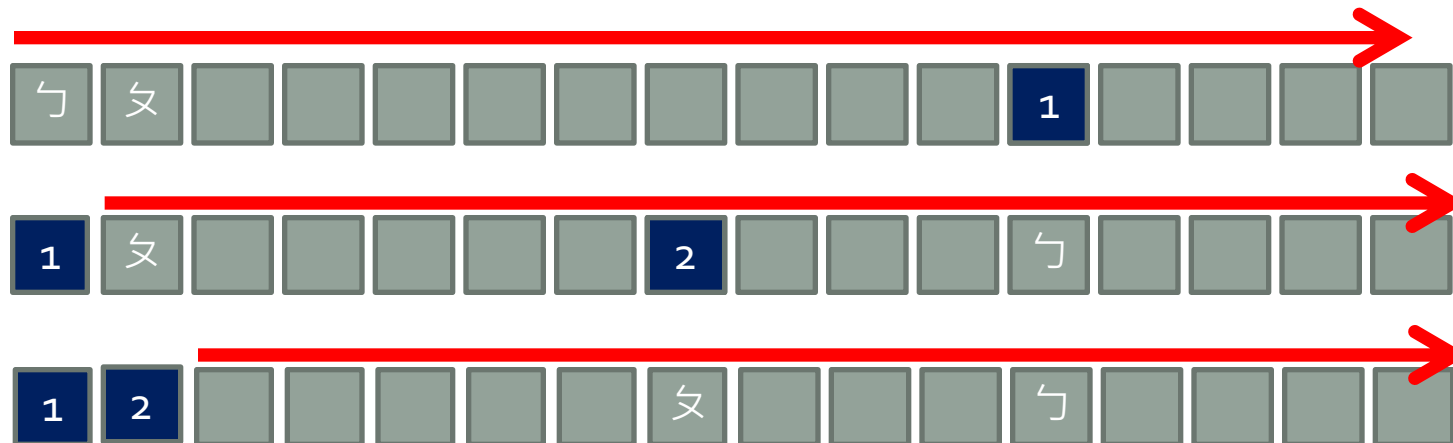
stop

$[2,3,1]$

stop

$[3,1,2]$

stop

$[3,2,1]$

- Every node represents a comparison & swap
- Sorting is completed when reaching the leaf
- How many leaves?
- $n!$, since there are that many possible permutations

# How fast can we sort?

- 所以, worst case所需要花的時間, 為此binary tree的height.
- 如果decision tree height為h, 有l個leaves
- $l \geq n!$, we have a least n! outcomes (leaves)
- $l \leq 2^h$, a binary tree (decision tree) of height $h$ has at most $2^{h-1}$ leaves
- $2^h \geq l \geq n!$
- $h \geq \log_2 n!$

- $n! = n(n-1)(n-2) \ldots 3 \cdot 2 \cdot 1 \geq \left(\frac{n}{2}\right)^{\frac{n}{2}}$

- $\log_2 n! \geq \log_2 \left(\frac{n}{2}\right)^{\frac{n}{2}} = \frac{n}{2} \log_2 \frac{n}{2} = \Omega(n \log n)$

- Summary: Any "comparison-based" sorting algorithm has worst-case time complexity of $\Omega(n \log n)$.

# Review: Selection Sort

- Select the smallest, move it to the first position.
- Select the second smallest, move it to the second position.
- ....
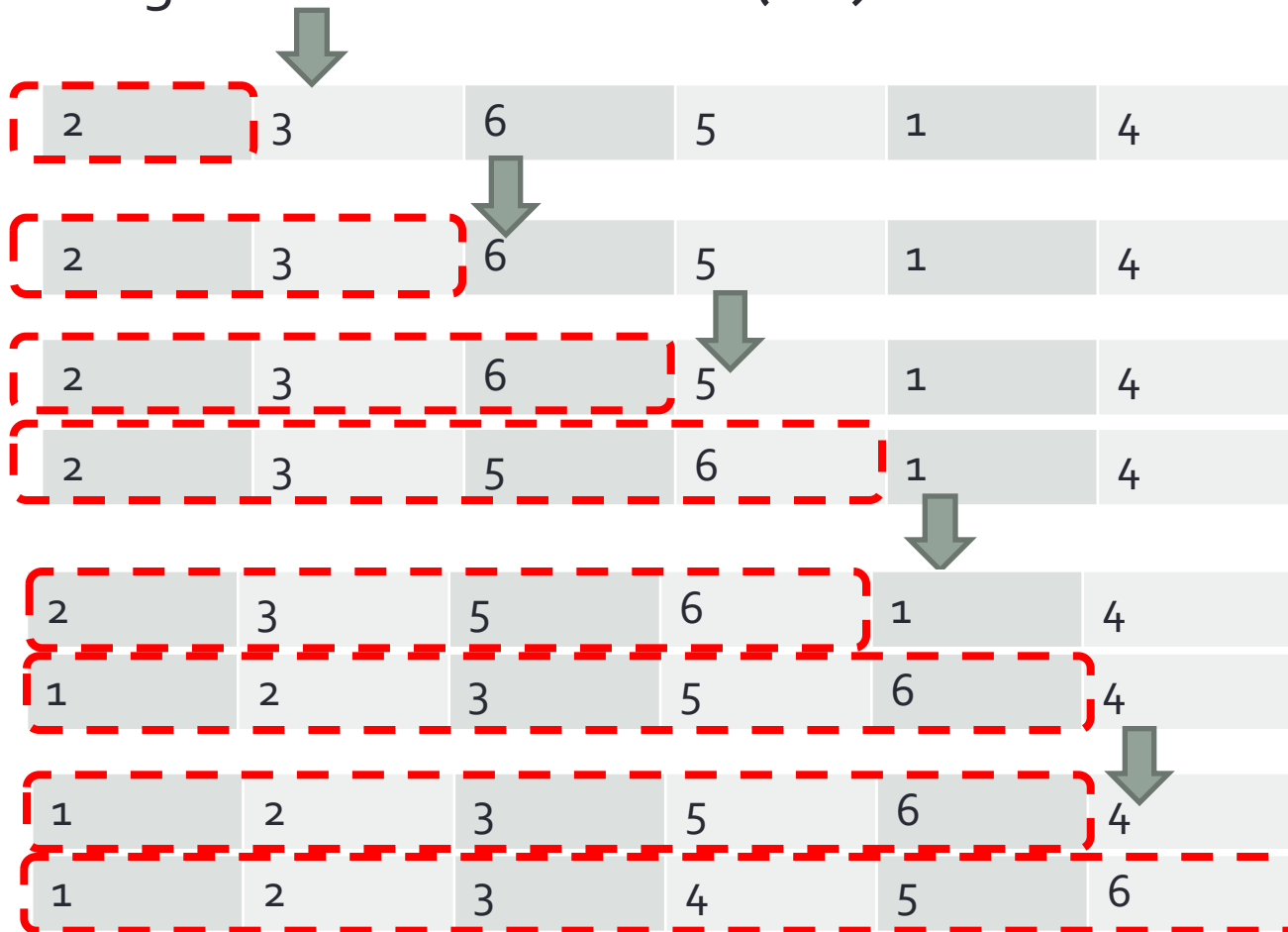- The last item will automatically be placed at the last position.

# Review: Selection Sort

- Selection sort does not change the execution of the algorithm due to the current conditions.
- Always going through the entire array in each iteration.
- Therefore, its best-case, worst-case, average-case running time are all $O(n^2)$

- **Not adaptive!**
- **In-place**

# Insertion Sort

- In each iteration, add one item to **a sorted list of *i* item**.
- Turning it into **a sorted list of *(i+1)* item**

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

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

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

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

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

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

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

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

# Pseudo code

INSERTION-SORT($A$)

1  **for** $j = 2$ **to** $A.length$
2         $key = A[j]$
3         // Insert $A[j]$ into the sorted
   sequence $A[1 \ldots j - 1]$.
4         $i = j - 1$
5         **while** $i > 0$ and $A[i] > key$
6                 $A[i + 1] = A[i]$
7                 $i = i - 1$
8         $A[i + 1] = key$

# Insertion Sort

- Q: How much time is needed?
- A: In the worst case, the item needs to be placed at the beginning for each and every iteration.
- (Spending time linear to the size of sorted part)
- $\sum_1^{n-1} i = \frac{n(n-1)}{2} = O(n^2)$
- Average-case complexity: $O(n^2)$. (Why?)

- Possible variation: (do those improve the time complexity?)
- 1. Use binary search to look for the location to insert.
- 2. Use linked list to store the items. Then moving takes only $O(1)$!

# What's good about insertion sort

- Simple (small constant in time complexity representation)
  - Good choice when sorting a small list
- Stable
- In-place
- Adaptive
  - Example: In ⟨1,2,5,3,4⟩, only two inversions <5,3>, <5,4>.
  - The running time for insertion sort: O(n+d), d is the number of inversions
    Best case: O(n) (No inversion, **sorted**)
- Online:
  No need to know all the numbers to be sorted. Possible to sort and take input at the same time.

# Merge Sort

- Use **Divide-and-Conquer** strategy
- Divide-and-Conquer:
  - Divide: Split the big problem into small problems
  - Conquer: Solve the small problems
  - Combine: Combine the solutions to the small problems into the solution of the big problems.

- Merge sort:
  - Divide: Split the *n* numbers into two sub-sequences of *n/2* numbers
  - Conquer: Sort the two sub-sequences (use recursive calls to delegate to the clones)
  - Combine: Combine the two **sorted sub-sequences** into the one **sorted** sequence

# Merge Sort

MERGE-SORT$(A, p, r)$

1   **if** $p < r$
2        $q = \lfloor(p + r)/2\rfloor$
3        MERGE-SORT$(A, p, q)$
4        MERGE-SORT$(A, q + 1, r)$
5        MERGE$(A, p, q, r)$

Divide

Conquer x2

Combine

# Merge Sort

A[]: the array to be sorted
temp: temporarily storage
left,right: the left & right indices of the range to be sorted.

```
void Mergesort(int A[], int temp, int left, int right) {
        int mid; .
        if (right > left) {
                mid=(right+left)/2;
                Mergesort(A,temp,left,mid);
                Mergesort(A,temp,mid+1,right);
                Merge(A,temp,left,mid+1,right);
        }
}
```
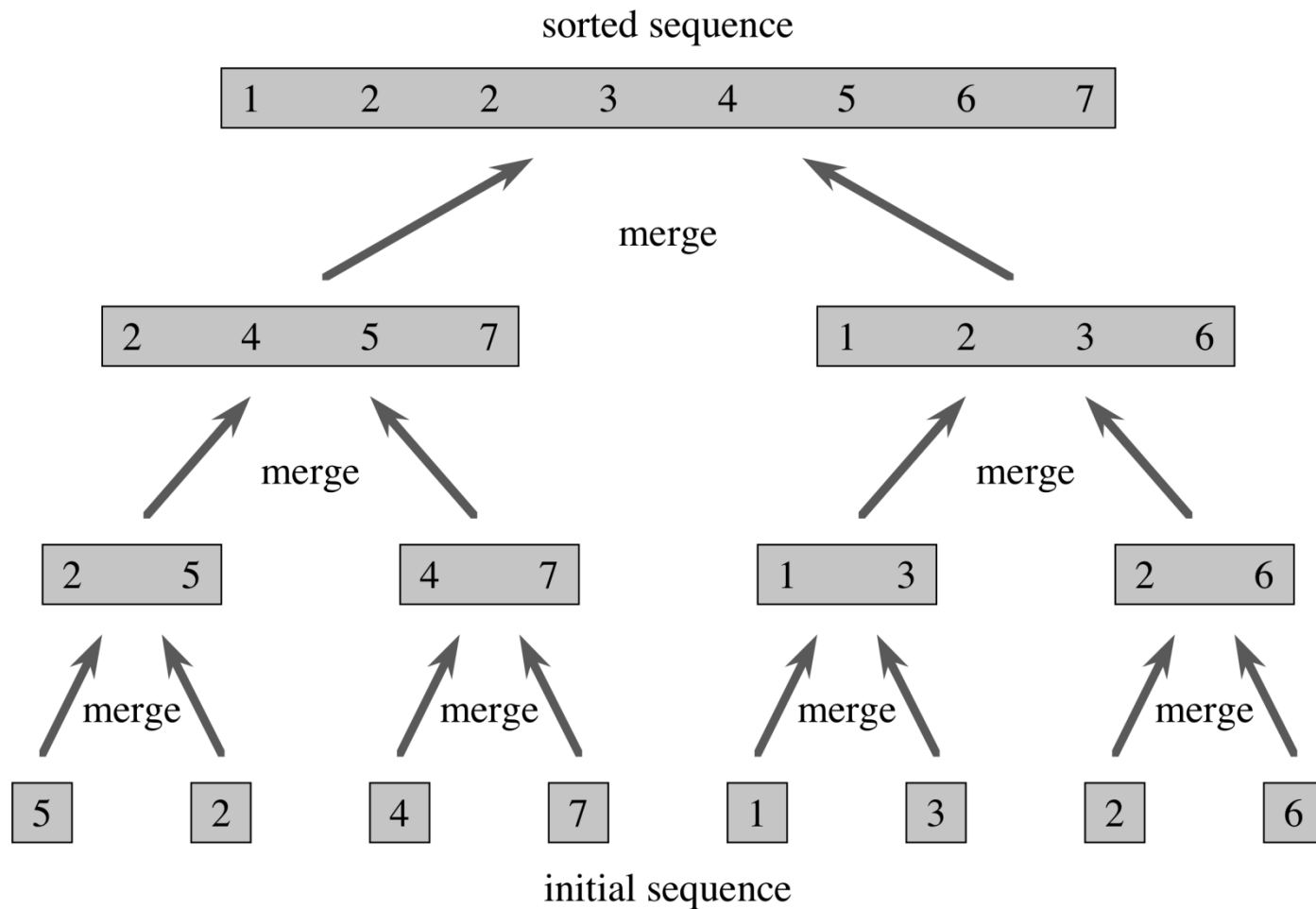
Divide

Conquer

Combine

# Merge Sort: Example



sorted sequence

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

merge

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

merge      merge

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

merge   merge   merge   merge

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

initial sequence

# How to combine (merge)?

Temporary storage

i

j

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

Original array

- Running time: $O(n_1 + n_2) = O(n), n_1$ 和 $n_2$ are the lengths of the two sub-sequences.
- A temporary storage of size O(n) is needed during the merge process

# Implementation: Merge

MERGE$(A, p, q, r)$

```
 1   n₁ = q − p + 1
 2   n₂ = r − q
 3   let L[1 .. n₁ + 1] and R[1 .. n₂ + 1] be new arrays
 4   for i = 1 to n₁
 5         L[i] = A[p + i − 1]
 6   for j = 1 to n₂
 7         R[j] = A[q + j]
 8   L[n₁ + 1] = ∞
 9   R[n₂ + 1] = ∞
10   i = 1
11   j = 1
12   for k = p to r
13       if L[i] ≤ R[j]
14             A[k] = L[i]
15             i = i + 1
16       else A[k] = R[j]
17             j = j + 1
```

$1 \quad n_1 = q - p + 1$
$2 \quad n_2 = r - q$
$3 \quad \text{let } L[1 \mathinner{.\,.} n_1 + 1] \text{ and } R[1 \mathinner{.\,.} n_2 + 1] \text{ be new arrays}$
$4 \quad \textbf{for } i = 1 \textbf{ to } n_1$
$5 \qquad L[i] = A[p + i - 1]$
$6 \quad \textbf{for } j = 1 \textbf{ to } n_2$
$7 \qquad R[j] = A[q + j]$
$8 \quad L[n_1 + 1] = \infty$
$9 \quad R[n_2 + 1] = \infty$
$10 \quad i = 1$
$11 \quad j = 1$
$12 \quad \textbf{for } k = p \textbf{ to } r$
$13 \qquad \textbf{if } L[i] \leq R[j]$
$14 \qquad\quad A[k] = L[i]$
$15 \qquad\quad i = i + 1$
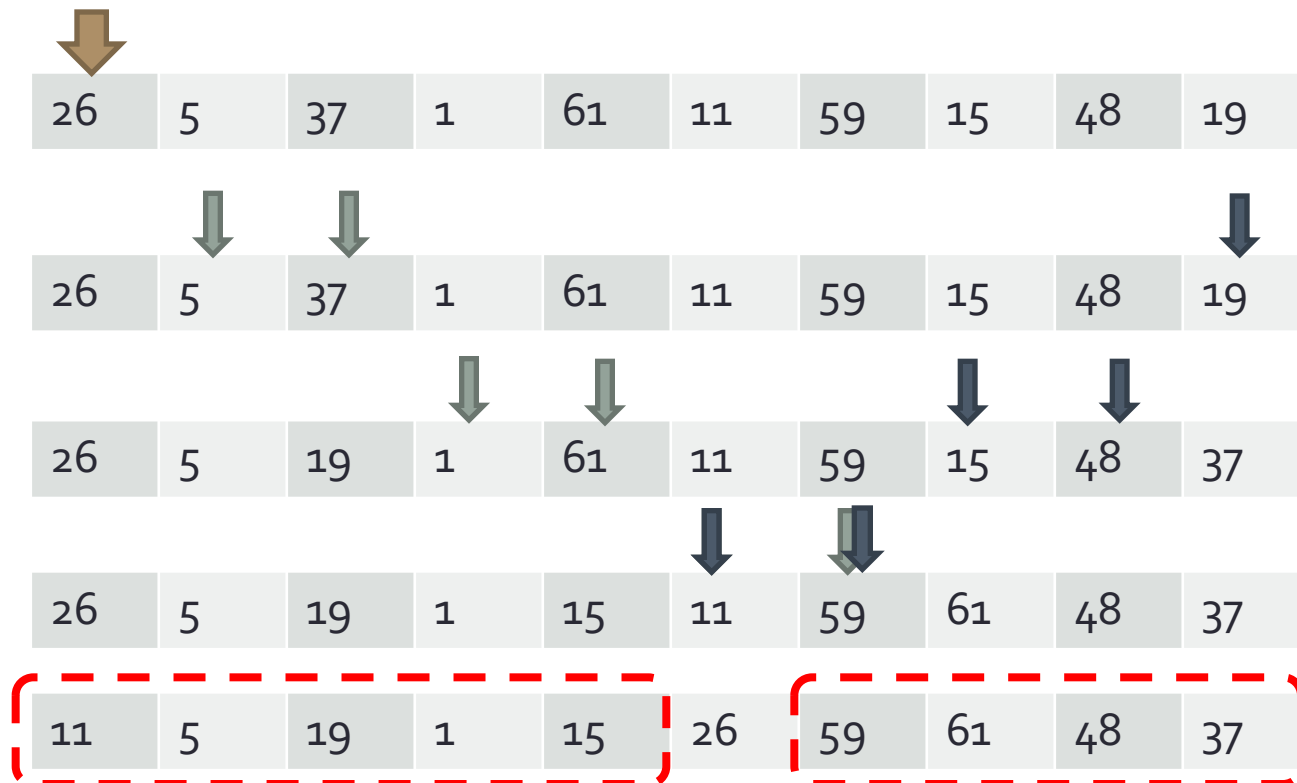$16 \qquad \textbf{else } A[k] = R[j]$
$17 \qquad\quad j = j + 1$

# Merge sort

- Every item to be sorted is processed once per "pass"$\rightarrow O(n)$
- How many passes is needed?
- The length of the sub-sequence **doubles** every pass, and finally it becomes the large sequence of *n* numbers
- Therefore, $\lceil \log_2 n \rceil$ passes.
- Total running time: $O(n \log_2 n) = O(n \log n)$
- Worst-case, best-case, average-case: $O(n \log n)$
  (**Not adaptive**)

- Not in-place: need additional storage for sorted sub-sequences
- Additional space: O(n)

# Quick Sort

- Find a pivot(支點), manipulate the locations of the items so that:
  - (1) all items to its left is smaller or equal (**unsorted**),
  - (2) all items to its right is larger
- Recursively call itself to sort the left and right sub-sequences.

| 26 | 5 | 37 | 1 | 61 | 11 | 59 | 15 | 48 | 19 |
|----|---|----|---|----|----|----|----|----|----|

| 26 | 5 | 37 | 1 | 61 | 11 | 59 | 15 | 48 | 19 |
|----|---|----|---|----|----|----|----|----|----|

| 26 | 5 | 19 | 1 | 61 | 11 | 59 | 15 | 48 | 37 |
|----|---|----|---|----|----|----|----|----|----|

| 26 | 5 | 19 | 1 | 15 | 11 | 59 | 61 | 48 | 37 |
|----|---|----|---|----|----|----|----|----|----|

| 11 | 5 | 19 | 1 | 15 | 26 | 59 | 61 | 48 | 37 |
|----|---|----|---|----|----|----|----|----|----|

# Pseudo Code

$\text{Quicksort}(A, p, r)$

1   **if** $p < r$
2        $q = \text{Partition}(A, p, r)$   Divide
3        $\text{Quicksort}(A, p, q-1)$
4        $\text{Quicksort}(A, q+1, r)$   Conquer x2

No Combine!

# Quick Sort

| 11 | 5 | 19 | 1 | 15 | 26 | 59 | 61 | 48 | 37 |
|----|---|----|---|----|----|----|----|----|----|
| 1 | 5 | 11 | 19 | 15 | 26 | 59 | 61 | 48 | 37 |
| 1 | 5 | 11 | 19 | 15 | 26 | 59 | 61 | 48 | 37 |
| 1 | 5 | 11 | 15 | 19 | 26 | 59 | 61 | 48 | 37 |
| 1 | 5 | 11 | 15 | 19 | 26 | 59 | 61 | 48 | 37 |
| 1 | 5 | 11 | 15 | 19 | 26 | 48 | 37 | 59 | 61 |
| 1 | 5 | 11 | 15 | 19 | 26 | 37 | 48 | 59 | 61 |
| 1 | 5 | 11 | 15 | 19 | 26 | 37 | 48 | 59 | 61 |

# Quick Sort: Worst & Best case

- But worst case running time is still $O(n^2)$
- Q: Give an example which produces worst-case running time for the quick sort algorithm.
- In this case: running time is $O(n^2)$

- Best case?
- Pivot can split the sequence into two sub-sequences of equal size.
- Therefore, T(n)=2T(n/2)+$O(n)$
- T(n)=$O(n \log n)$

# Randomized Quick Sort

- Avoid worst case to happen frequently

- Randomly select a pivot (not always the leftmost key)

- Reduce the probability of the worst case
- However, worst case running time is still $O(n^2)$

Randomly select a pivot

| 26 | 5 | 37 | 1 | 61 | 11 | 59 | 15 | 48 | 19 |

Swap in advance

# Average running time

- Better if the selection of pivot can evenly split the sequence into two sub-sequences of equal size

- Why the average running time is close to the best-case one?
- 假設很糟的一個狀況: 每次都分成1:9

Time needed for the "9/10 subsequence"    Time needed for partitioning
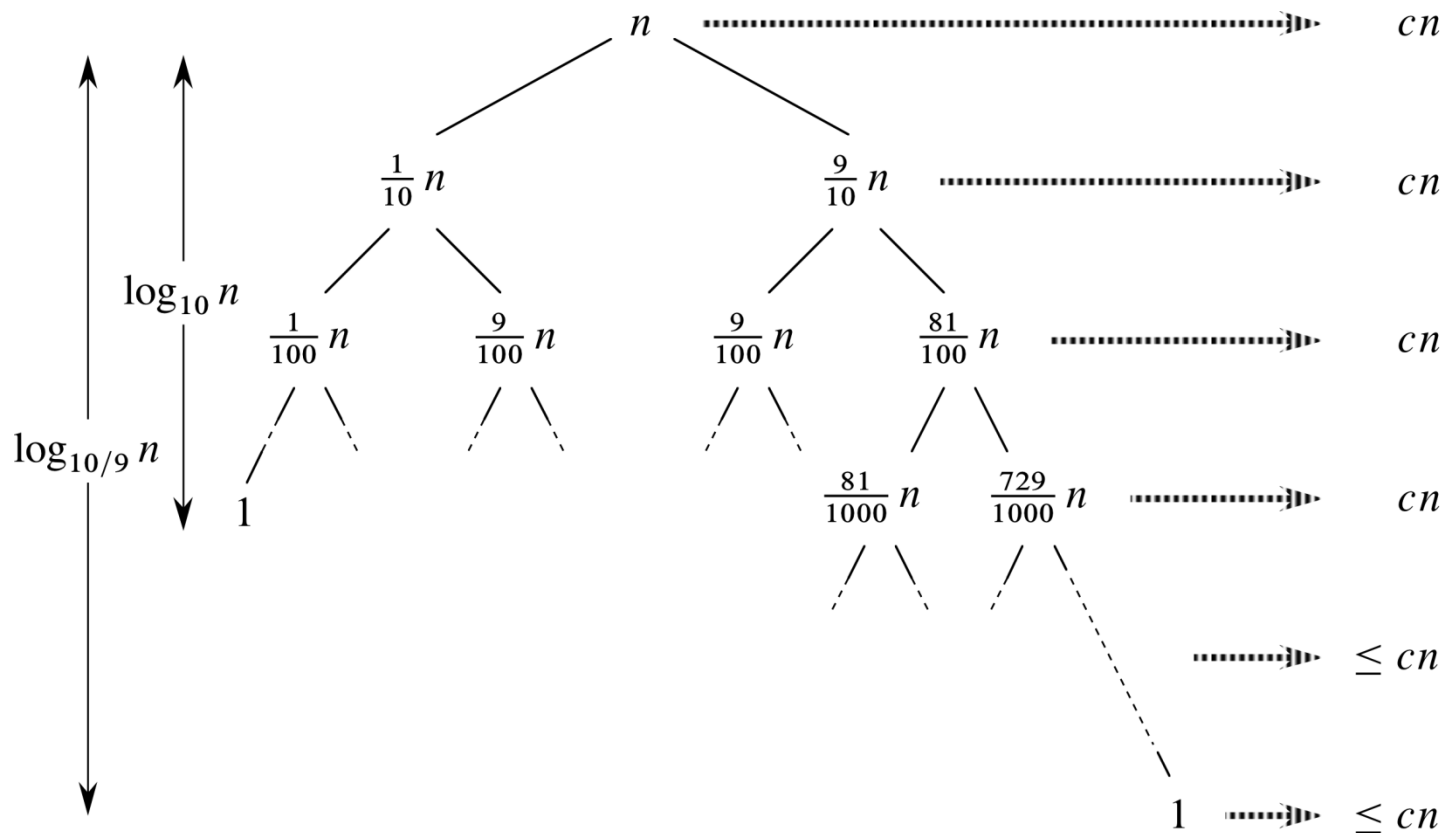
- $T(n) = T(9n/10) + T(n/10) + cn$

Time needed for the "1/10 subsequence"

- $= \left( T\left(\frac{81n}{100}\right) + T\left(\frac{9n}{100}\right) + \frac{9cn}{10} \right) + \left( T\left(\frac{9n}{100}\right) + T\left(\frac{1n}{100}\right) + \frac{cn}{10} \right) + cn$

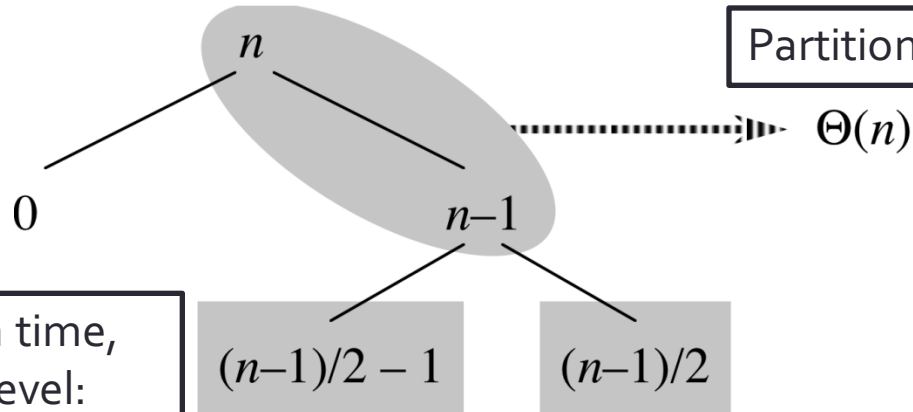- $= \cdots$

# Average running time



As long as the pivot can partition according to a particular ratio (even not close to 50%), we can still obtain $O(n \log n)$ running time!

# Average running time

Case 1:
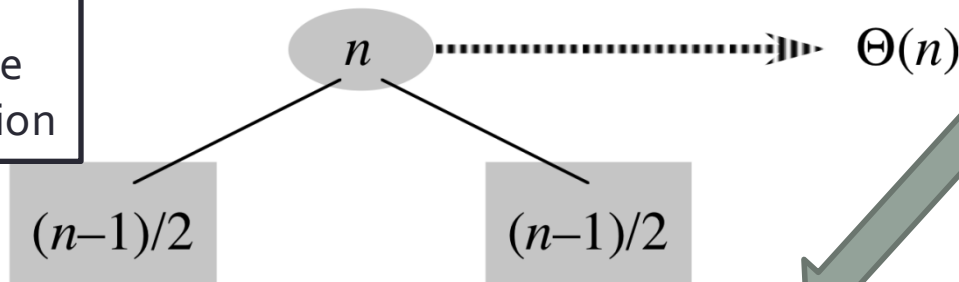Worst case for the first level partition, but best-case for second level.

Partition time for the first level: $\Theta(n)$

$n$

$\Theta(n)$

Partition time, second level:
$\Theta(n-1)$

0

$n-1$

$(n-1)/2 - 1$

$(n-1)/2$

$\Theta(n) + \Theta(n-1) = \Theta(n)$

Case 2:
Best case for the first level partition

$n$

$\Theta(n)$

Same!
(Case 1 has larger constant)
The better-partitioned level would "absorb" the extra running time for worse-partitioned level.

$(n-1)/2$

$(n-1)/2$

Partition time for the first level: $\Theta(n)$

# 比較四大金剛

| | Worst | Average | Additional Space? |
|---|---|---|---|
| Insertion sort | $O(n^2)$ | $O(n^2)$ | O(1) |
| Merge sort | $O(n \log n)$ | $O(n \log n)$ | O(n) |
| Quick sort | $O(n^2)$ | $O(n \log n)$ | O(1) |
| Heap sort | $O(n \log n)$ | – | O(1) |

**Not covered today!**

- Insertion sort: quick with small input size n. (small constant)
- Quick sort: Best average performance (fairly small constant)
- Merge sort: Best worst-case performance
- Heap sort: Good worst-case performance, no additional space needed.
- Real-world strategy: **a hybrid of insertion sort** + others. Use **input size** *n* to determine the algorithm to use.