

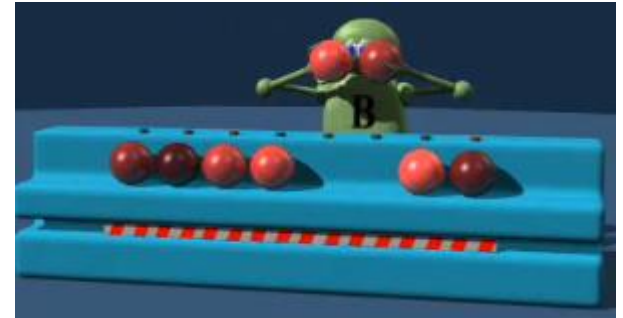
線性時間的sorting

# Sorting in Linear Time

## Chapter 8

Mei-Chen Yeh

# How fast can we sort?



- Comparison sorts
  - Bubble sort, selection sort, insertion sort, merge sort, quick sort 這些演算法的worst-case最好可以提升到  $O(n \log n)$
  - The **best worst-case** running time:  $O(n \log n)$

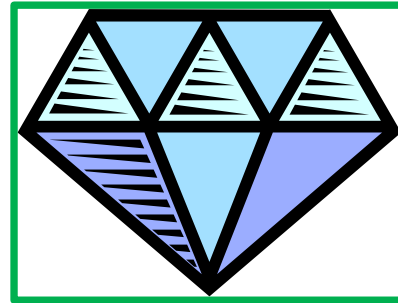
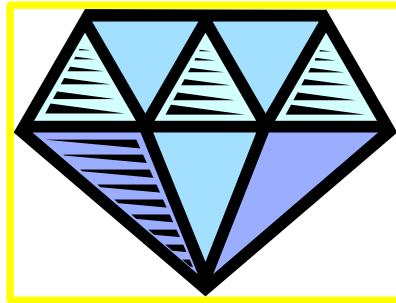
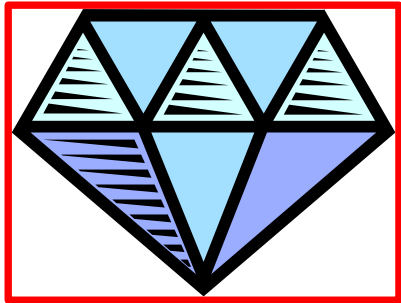
Q: Is  $O(n \log n)$  the best we can do?

A **decision tree** can help us answer this question!

決策樹

# Sorting 3 diamonds

把3個鑽石由輕排到重，只有天秤可以用，沒有磅秤



Minimal number of comparisons?  
Maximal number of comparisons?

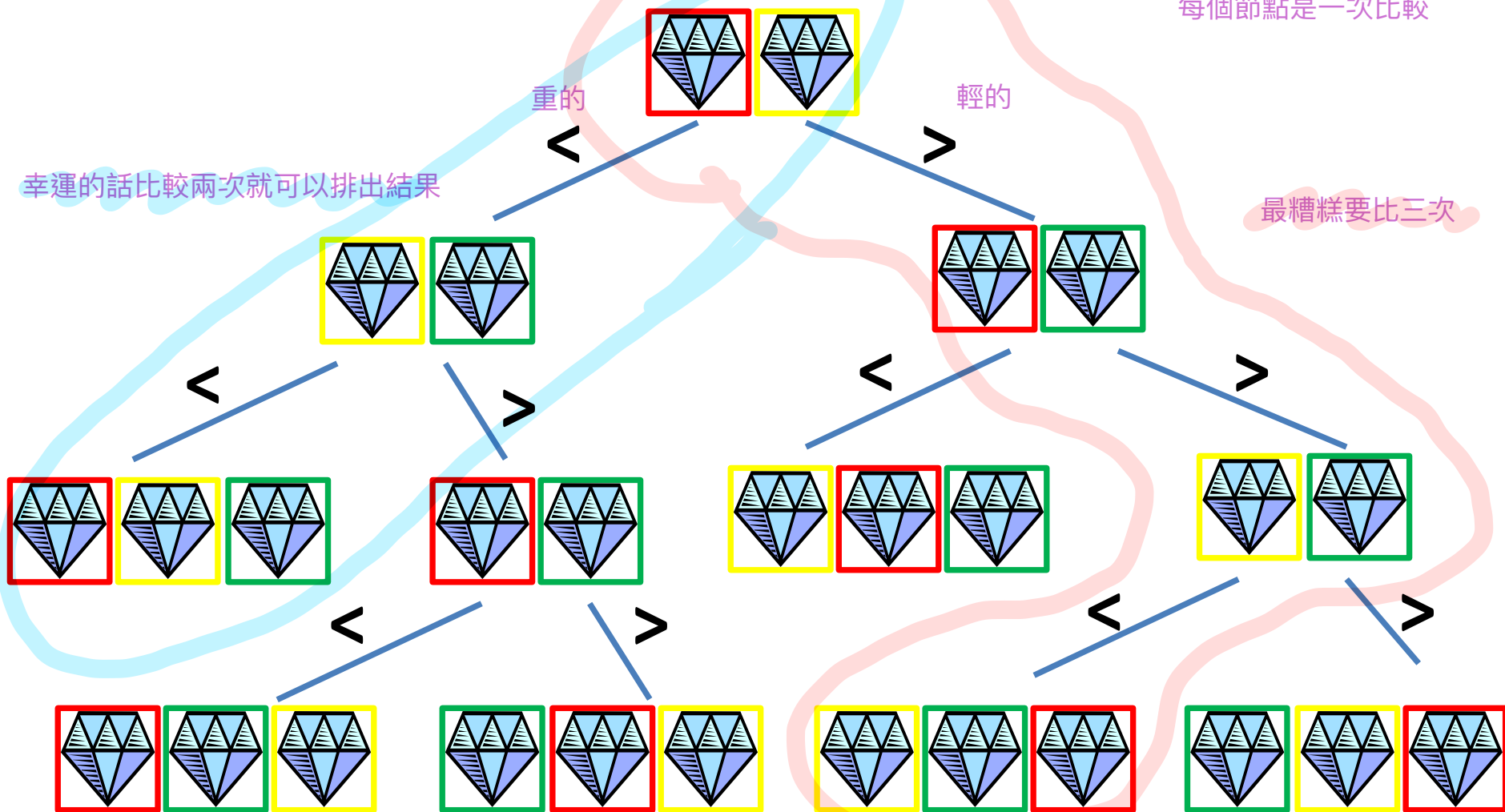
要排列3個鑽石，總共會有 $3!=6$ 種排列結果

# Decision tree

每個節點是一次比較

幸運的話比較兩次就可以排出結果

最糟糕要比三次



# Decision tree model

- Each leaf contains a permutation. 每個樹葉都是一種排序結果
- $n$  elements  $\Rightarrow n!$  permutations  $n$ 個鑽石要排，共會有 $n!$ 種結果
- # {comparisons} of the algorithm = the length of the path taken
- Worst-case running time = height of tree

# Lower bound for comparison sorts

如果發生兩兩比較就一定會是  $O(n \log n)$

## Theorem

Any comparison sort algorithm requires comparisons in  $\Omega(n \log n)$  in the worst case.

## Proof

Consider a decision tree of height  $h$  with  $L$  leaves.

$$n! \leq L \leq 2^h$$

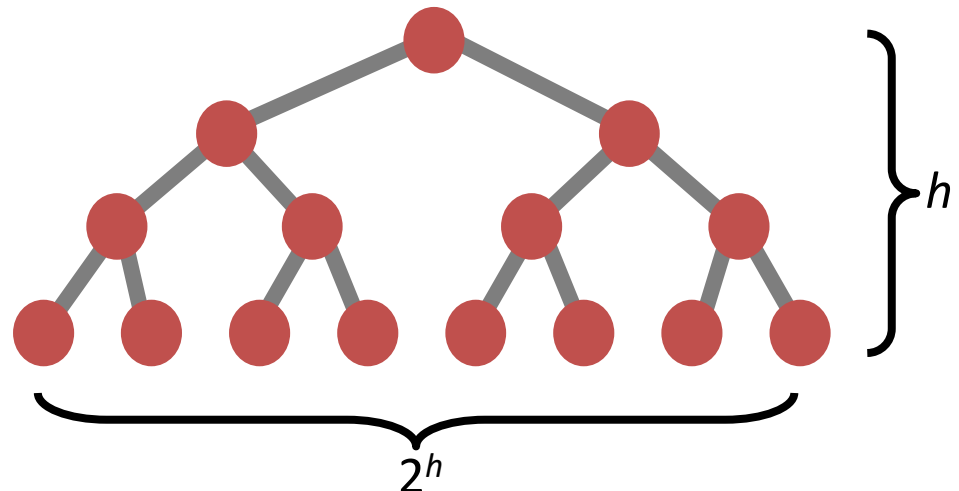
$$h \geq \log(n!)$$

$$= \Omega(n \log n)$$

Why?

$h$ 一定比 $\Omega(n \log n)$ 大

有 $L$ 個樹葉



$$\log(n!) = \Omega(n \log n)$$

$\log(n!)$ 可以寫成 $(n \log n)$ 起跳的一個式子

$\Omega(n \log n)$

1. Informally,

$$\log(n!) > \log((n/2)^{n/2}) = (n/2)\log(n/2) \rightarrow \Theta(n \log n)$$

$$n! = n(n-1)(n-2)\dots 1$$

$n/2$

2. Stirling's Approximation:  $n! > (n/e)^n$

$2.718$

$$\log(n!)$$

$$> \log((n/e)^n)$$

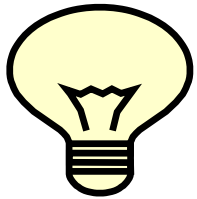
$$= n(\log n - \log e)$$

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

# Sorting in linear time

***No comparisons between elements!***

- Counting sort
- Radix sort



Counting sort : 不兩個兩個比了，直接用數的

[Example] Sort by age (young → old):

Count the number of people who is younger than me!





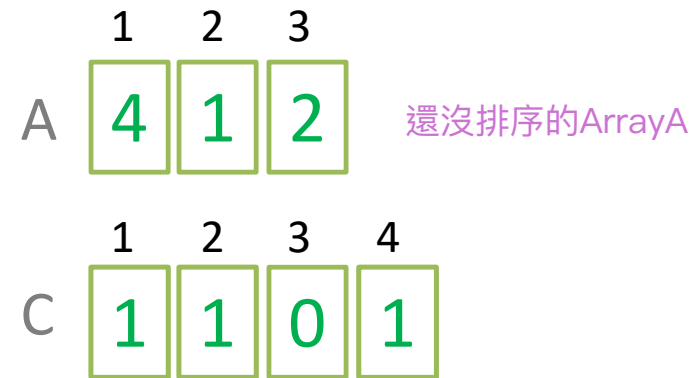
# Counting sort: the simplest setting

- Assumptions on the input
  - Distinct elements
  - Positive integers
  - Range is known ( $1..k$ )

# Algorithm Counting\_Sort

1. Allocate  $C[1..k]$   $O(1)$
2. **for**  $i = 1$  **to**  $k$   $O(k)$   
 $k$ 代表數值的上界  
(要排序的數字最大的那個)
3.      $C[i] = 0$
4. **for**  $i = 1$  **to**  $n$   $O(n)$
5.      $C[A[i]] = C[A[i]] + 1$
6. **for**  $i = k$  **downto**  $1$   $O(k)$
7.     **if**  $C[i] > 0$  **print**  $i$

$\max(O(n), O(k))$



>> 4 2 1

- ① 需要排列的數字有4、1、2
- ② 給他一個新的array來存 (array的index就代表要排列的數字，比方說index4的位置，因為要排的數字有4，所以就在index4的位置存1)
- ③ 最後把有存1這個數字的index印出來，就會呈現正確的排序結果。

# General setting

- Assumptions on the input
  - ~~– Distinct elements~~
  - Positive integers
  - ~~– Range is known ( $1..k$ )~~

0. Let  $k$  be the maximum of  $A[1]..A[n]$   $O(n)$

1. Allocate  $C[1..k]$

2. **for**  $i = 1$  **to**  $k$

3.      $C[i] = 0$

4. **for**  $i = 1$  **to**  $n$

5.      $C[A[i]] = C[A[i]] + 1$

6. **for**  $i = k$  **downto**  $1$

**for**  $j = 1$  **to**  $C[i]$

**print**  $i$

最慘的情況有可能 $O(kn)$

$O(kn)$



# Alternate algorithm for counting sort

counting sort的意思就是，數一下比自己小的有幾個，自己要排在他們後面。

1. **for**  $i = 1$  **to**  $k$
2.      $C[i] = 0$
3. **for**  $i = 1$  **to**  $n$  數數值為 $i$ 的，在 $k$ 裡面有幾個
4.      $C[A[i]] = C[A[i]] + 1$
5. **for**  $i = 2$  **to**  $k$  累加：把比 $i$ 小的都加起來
6.      $C[i] = C[i] + C[i-1]$
7. **for**  $i = n$  **downto**  $1$   $i$ 等於5到1
8.      $B[C[A[i]]] = A[i]$
9.      $C[A[i]] = C[A[i]] - 1$

	1	2	3	4	5
$A$ :	4	1	3	4	3

	1	2	3	4
$C$ :	1	0	2	2

$C'$ :	1	1	3	5
--------	---	---	---	---

小於等於3的3個

小於等於4的5個

為什麼counting sort可以在線性時間做完

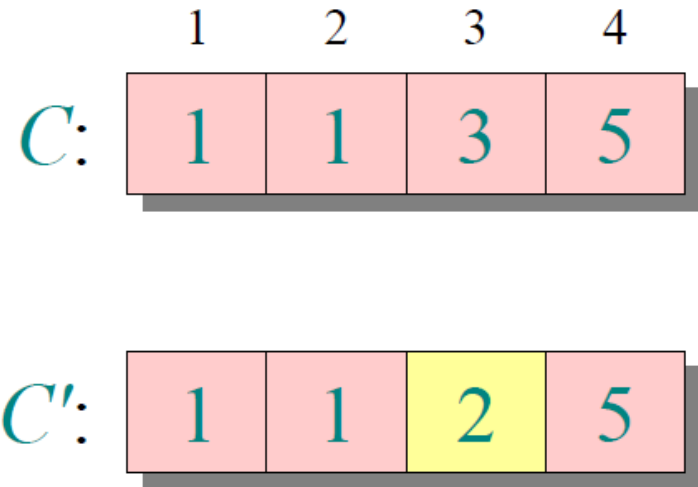
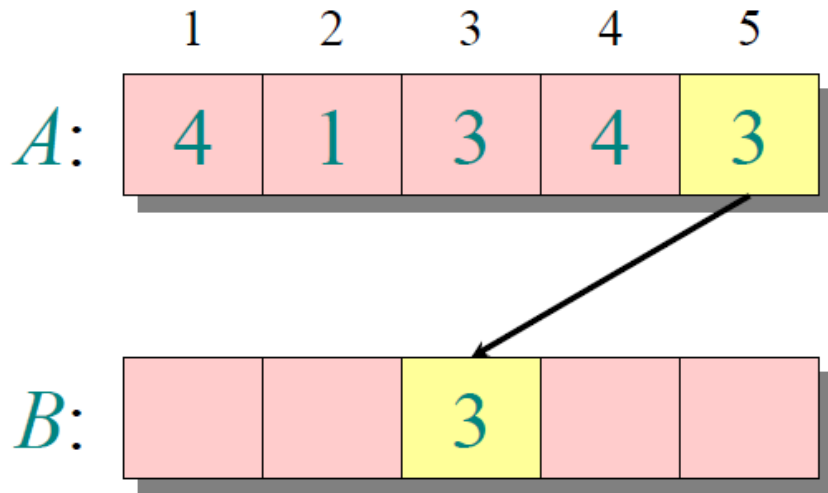
去掉了兩個兩個拿來比的程序，

他是多開一個array，還計算每個數出現了幾次

再用 $C'$ 來做累加，例如：小於等於1的幾個，小於等於2的幾個，

Results are stored in  $B[]$  把處理好的結果搬到Barray，最後輸出Barray，就會是排序好的結果

$$i = n (5)$$



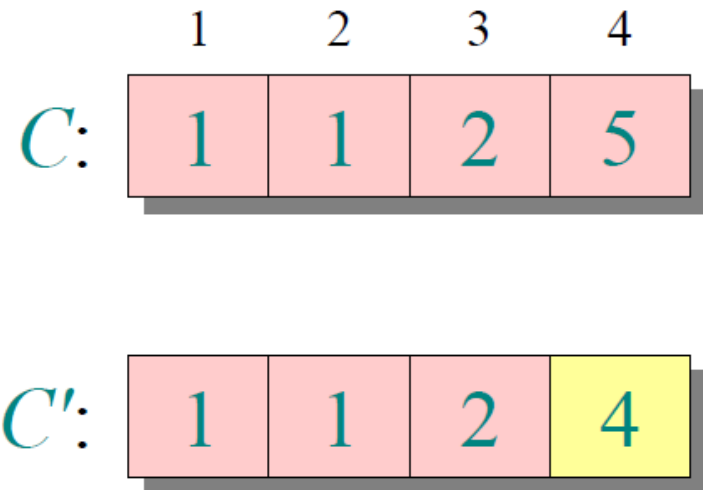
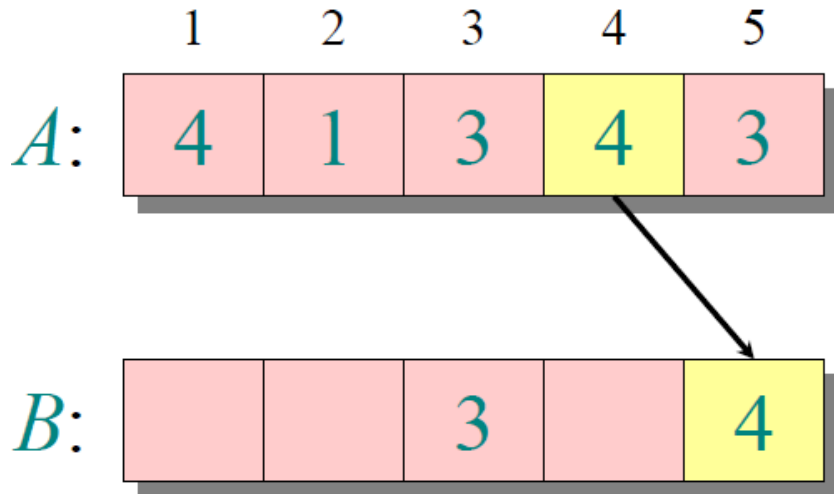
**for  $i = n$  downto 1**

$$B[C[A[i]]] = A[i]$$

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

$$i = n-1 \text{ (4)}$$

做完之後會把Aarray變成Barray，Barray是排序好的



**for  $i = n$  downto 1**

$$B[C[A[i]]] = A[i]$$

*C*[4] 是 5

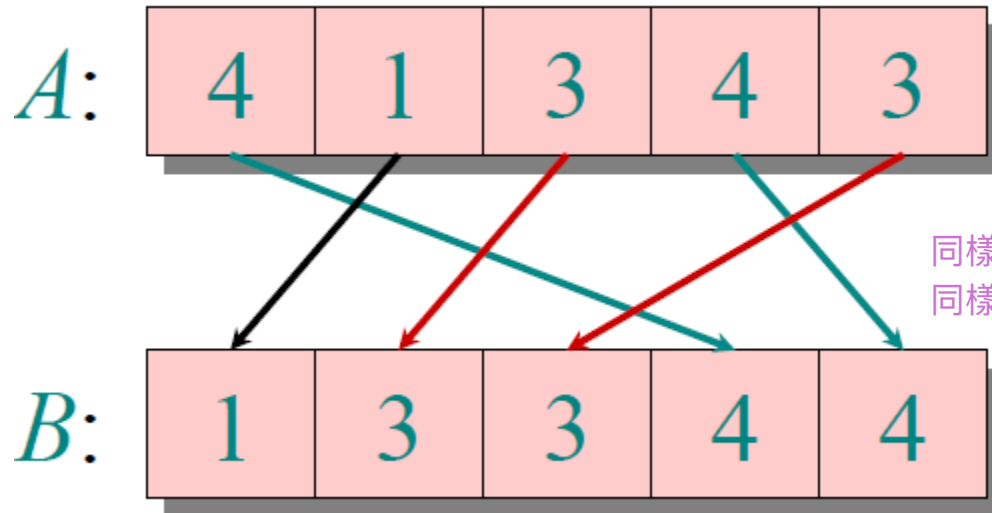
$$C[A[i]] = C[A[i]] - 1$$

... take-home practice

# Stable sorting

兩個一樣的數字，在原本的array中排左邊的，在新的array還是會排在左邊

- Preserves the input order among equal elements



Q: Counting sort is a stable sort?

Yes!

Merge sort? Quick sort?

3/1(5) end



作業手寫題要複習

期中考手寫，不上機

清明連假會放考古題到noodle

# Sorting in linear time

***No comparisons between elements!***

- Counting sort
- Radix sort

# Radix sort

- Idea: digit-by-digit sort
- Most significant digit first or less significant digit first?

要從個位數開始排 (i=1開始)

329

457

657

839

436

720

355

↑↑↑

# Radix-Sort ( $A, d$ )

$d$ 是常數，在這個例子中 $d=3$ （三位數）

如果不是stable，只是單純地讓排好的東西不動，這樣高位數的時候沒有辦法整坨移動。

1. for  $i = 1$  to  $d$  要從個位數開始排（ $i=1$ 開始）， $i$ 代表位數

2. use a *stable sort* to sort  $A$  on digit  $i$

這個排序一定要用stable sort來做（counting sort/merge sort）最後才是排百位數

0 (m)

329  
457  
657  
839  
436  
720  
355



720  
355  
436  
457  
657  
329  
839



720  
329  
436  
839  
355  
457  
657



329  
355  
436  
457  
657  
720  
839

根據個位數從小到大排，所以720上去

再來根據十位數排，720的2要在329的2上面

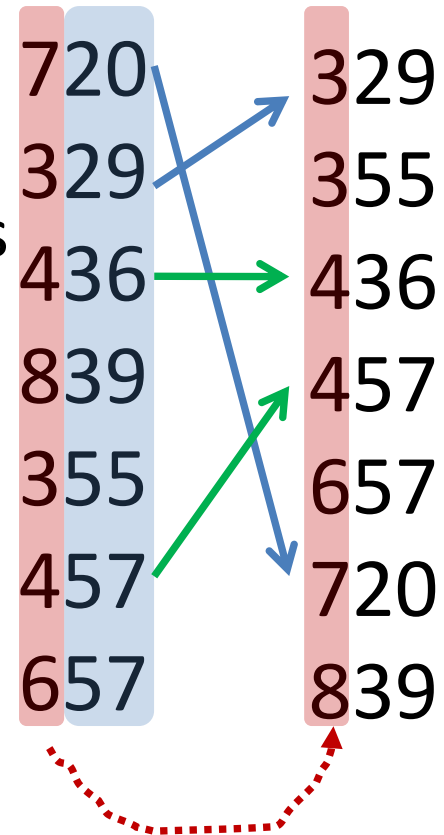
$O(d \cdot (n+k))$

# Correctness of Radix Sort

- Induction on digit position
  - Assume that the numbers are sorted by their low-order  $k-1$  digits
  - Sort on digit  $k$ 
    - Numbers that **equal** in digit  $k$
    - Numbers that **differ** in digit  $k$

為什麼從百位數開始排會錯，一定要從個位數開始排？

1. 有可能會遇到百位數的數字一樣
2. 有可能百位數排好了，結果到十位數結果被翻盤了



為什麼Radix Sort 一定要用stable sort排？

如果不stable的話，在低位數排好的東西，到高位數會亂掉。

# The Selection Problem

## Chapter 9

# Problem statement

給我一個沒有sort的array，然後跟我講要找第  $i$  大的數字

- Input: A set  $A$  of  $n$  distinct numbers and an integer  $i$ , with  $1 \leq i \leq n$ . 輸入會給我一個未排序的array加上一個1~n的整數
- Output: The element  $x$  that is larger than exactly  $i-1$  other elements of  $A$ . 輸出是從array抓一個數字出來，要剛好大於 $i-1$ 個數字，也就是排序後第 $i$ 個的數字

## The $i$ -th order statistic

$i = 1 \rightarrow$  find the *minimum*

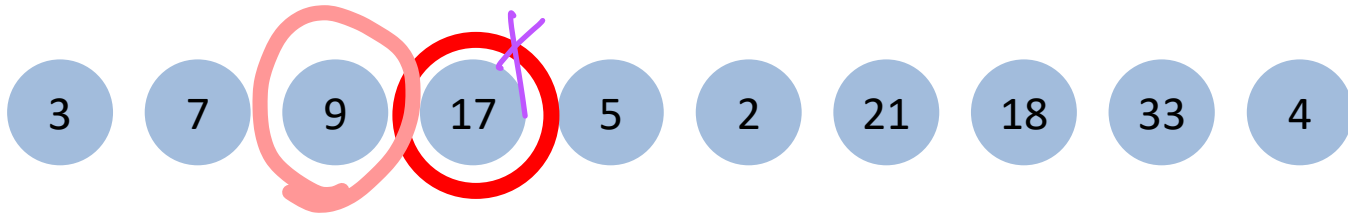
$i = n \rightarrow$  find the *maximum*

$i = n/2 \rightarrow$  find the *median*

# Example

- $i = 6$

index從1開始算



$i=1$   
↓  
2, 3, 4, 5, 7, 9, ...

$i=6$   
↓



# Naïve algorithm

什麼是Naive algorithm ?

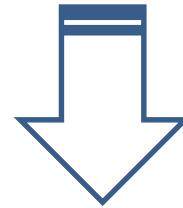
- Step 1: Sort  $A$
- Step 2: Return  $A[i]$

counting sort的時間複雜度

雖然counting sort比較有效率，  
但是有些情況不能用counting sort

$O(n \log n)$

quick sort的時間複雜度



$O(n)$

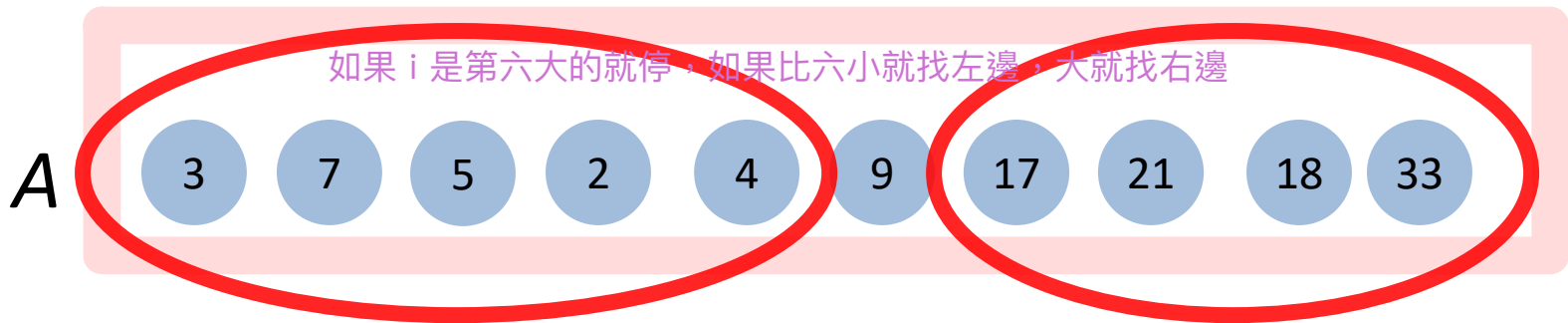
The number of comparison to solve the selection problem is between  **$1.5n$**  and  **$5.43n$** .

# Divide-and-Conquer Algorithm (1)

執行這個函式之後，會把元素跟pivot值比較後，拆成兩坨

pivot值是隨機抓的

$$q = \text{RANDOM-PARTITION}(A, p, r)$$



$q = 6$  (the index starts from 1) q是告訴我們pivot值落在這個array中的第幾個

The selection problem: select the  $i$ -th **order statistics**

What if  $i = q$ ,  $i < q$ , and  $i > q$ ?

# Divide-and-Conquer Algorithm (2)

會用到randomize-partition，所以叫這個名字

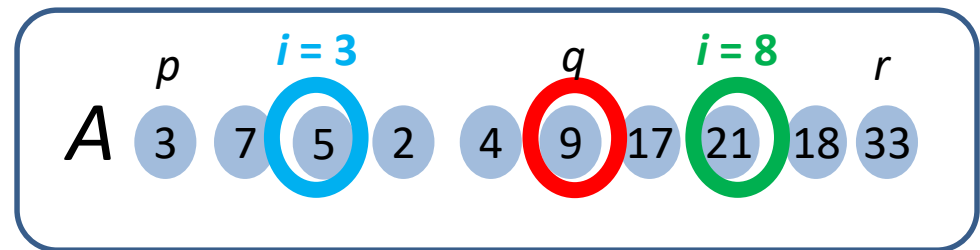
**RANDOMIZED-SELECT** ( $A, p, r, i$ )

1. **if**  $p == r$
2.     **return**  $A[p]$
3.  $q = \text{RANDOMIZED-PARTITION}$  ( $A, p, r$ )
4.  $k = q - p + 1$
5. **if**  $i == k$
6.     **return**  $A[q]$
7. **elseif**  $i < k$
8.     **return**  $\text{RANDOMIZED-SELECT}(A, p, q-1, i)$
9. **else return**  $\text{RANDOMIZED-SELECT}(A, q+1, r, i-k)$

k把絕對位置換成相對位置。

假設今天只看一段數列，k可以幫助我們知道在這一段pivot 值是在這一段的第幾個

假設pivot隨機選到9，他會給我一個結果是以9為分界，拆成兩坨



如果 $i$ 比 $k$ 大，就知道答案在右半邊，所以從 $q+1$ 開始

# RANDOMIZED-SELECT ( $A, p, r, i$ )

- The worst-case running time: ?

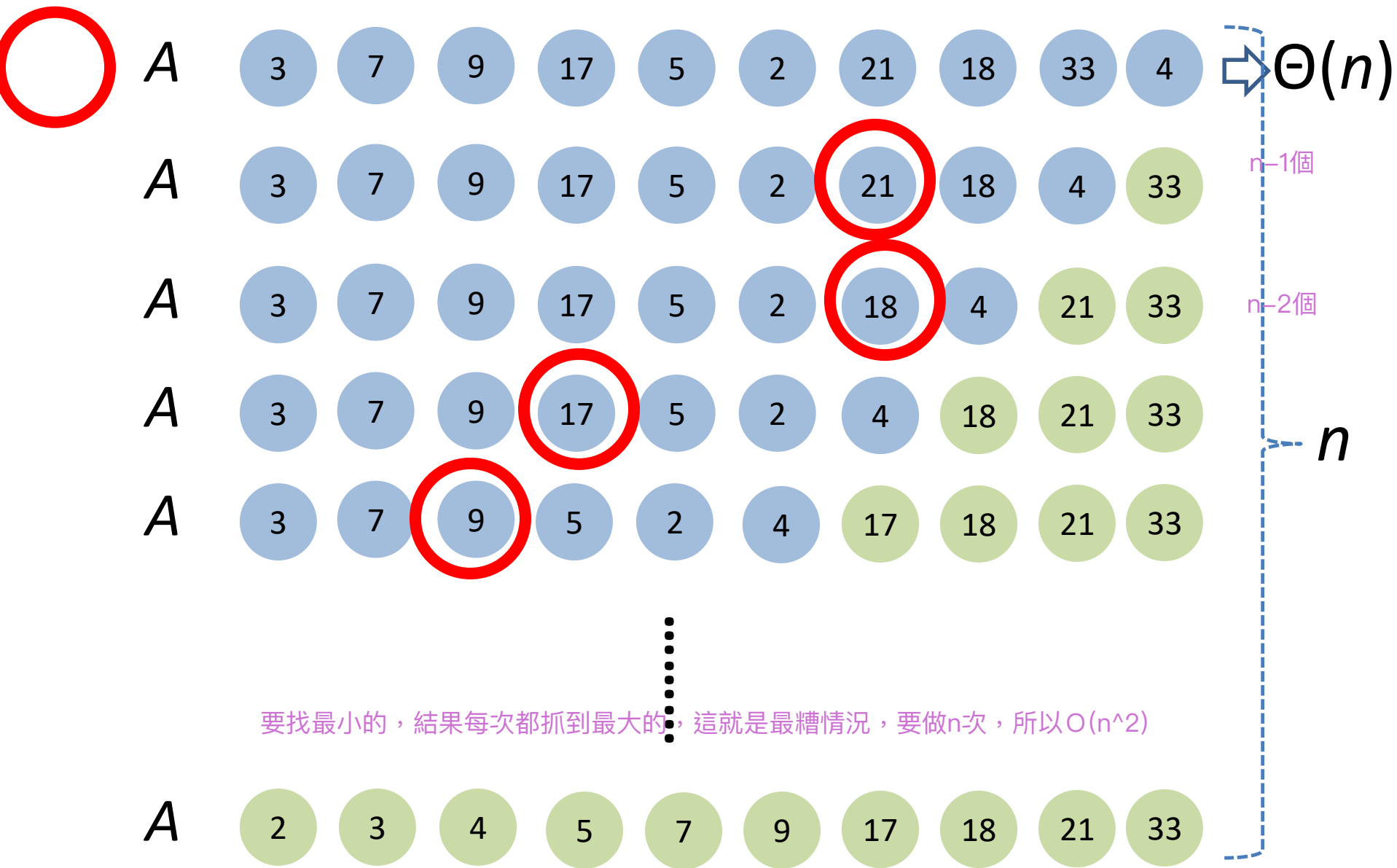
quick sort的worst case  $O(n^2)$   
因為無法掌控pivot值



$i = 1 \rightarrow$  find the minimum

*What's the selection of pivots that leads to the worst case?*

randomized-select的worst case的情況是：  
每次pivot一抓都是抓到端點，然後端點沒辦法幫我分兩坨，  
所以每次都得從 $n-1$ 做



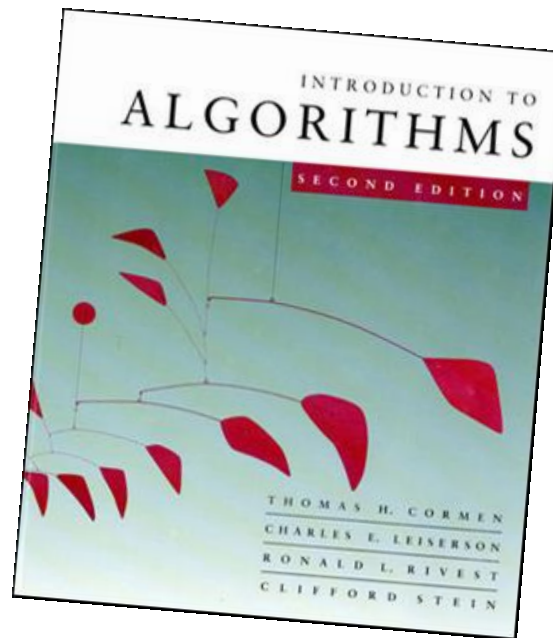
要找最小的，結果每次都抓到最大的，這就是最糟情況，要做 $n$ 次，所以 $O(n^2)$

$$T(n) = T(n-1) + \Theta(n) = \Theta(n^2) \quad \Theta(n^2)$$

# RANDOMIZED-SELECT ( $A, p, r, i$ )

- The average-case running time:  $\Theta(n)$

大家都會猜 $O(n \log n)$ ，但其實average-case出乎意料地是 $O(n)$



Pages 217-219

把randomized-selection的worst case變成線性時間

# Selection in worst-case linear time



Guarantee a good split upon partitioning the array!

SELECT ( $A, p, r, i$ )

1. if  $p == r$

2. return  $A[p]$

~~3.  $q = \text{RANDOMIZED-}$~~

4.  $k = q - p + 1$

5. if  $i == k$

6. return  $A[q]$

7. elseif  $i < k$

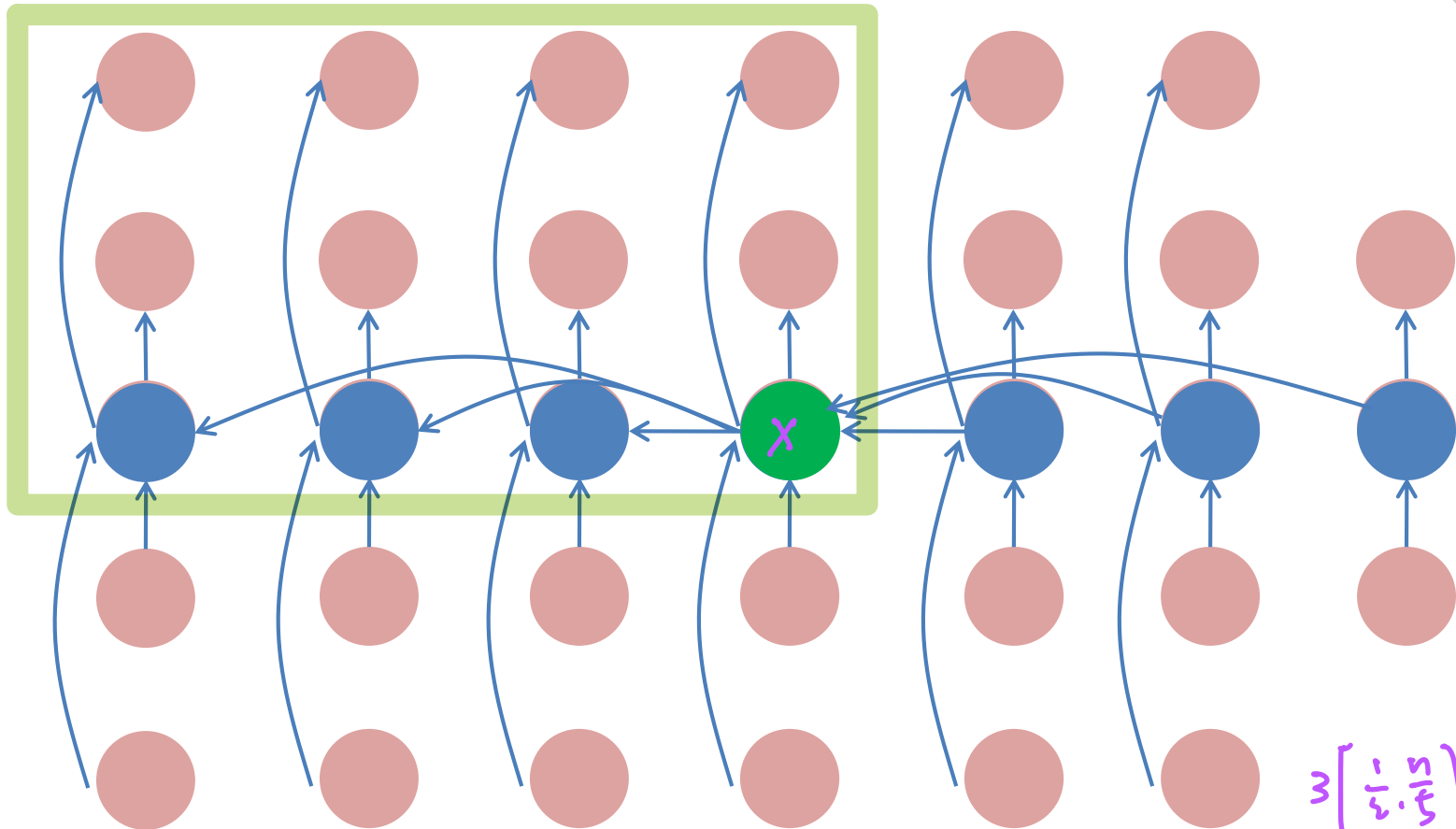
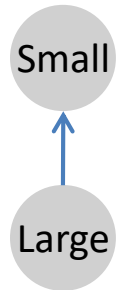
8. return SELECT ( $A, p, q-1, i$ )

9. else return SELECT ( $A, q+1, r, i-k$ )

1. Divide the  $n$  elements into groups of 5
2. Find the **median** of each 5-element group
3. Recursively select the **median** of those group medians as pivot
4. Partition the input using the pivot

Example:  $n = 33$  ( $A[0-32]$ )

5個一組  
選中位數  
在中位數中選中位數



$$3 \left( \frac{1}{5} \cdot \frac{n}{5} \right) = 3 \left( \frac{n}{10} \right) \leq n/3$$

Find medians in terms of *element values*, not locations!



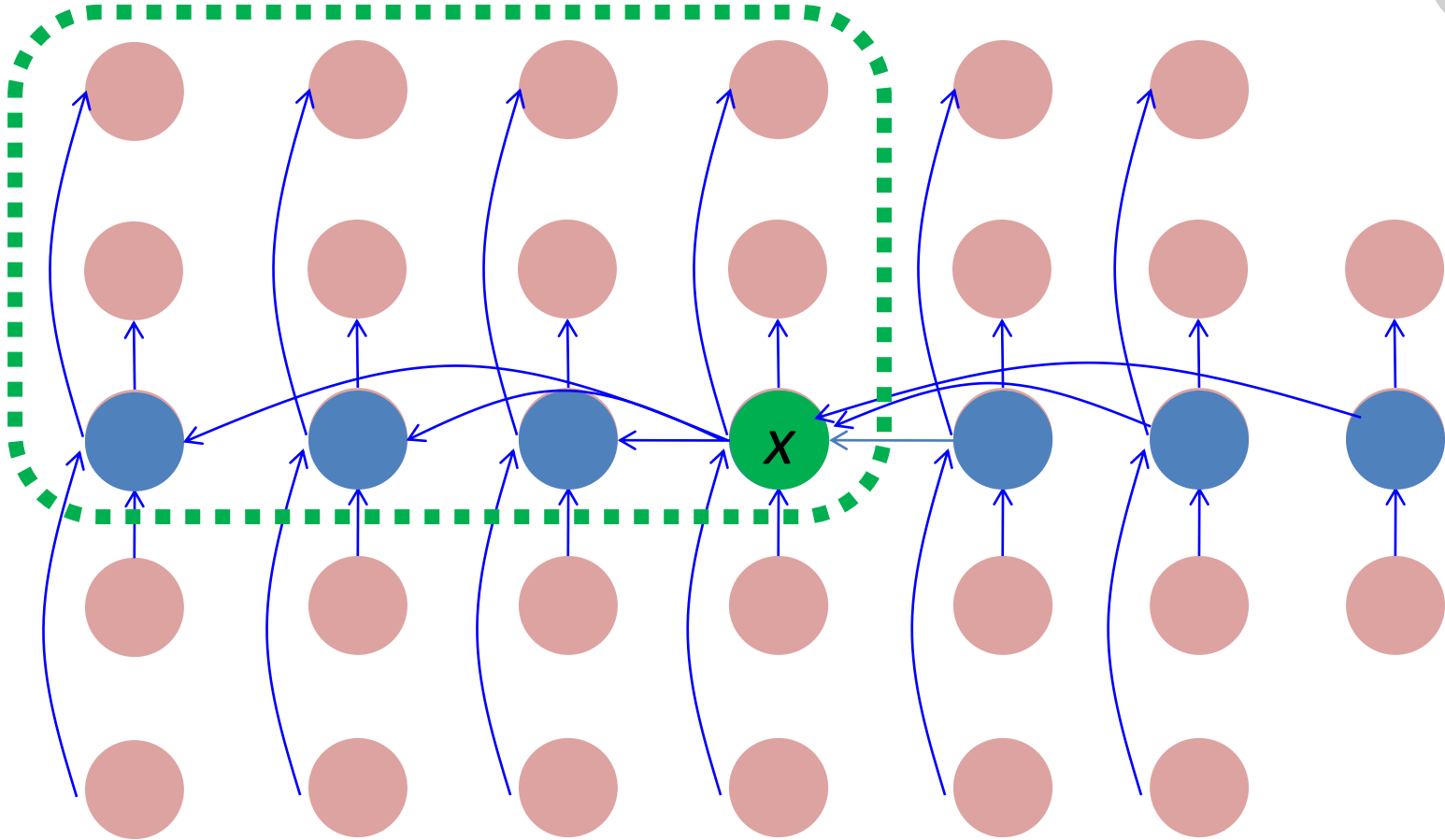
Linear time in the worst case!

# Analysis

Assume all elements are distinct

Small

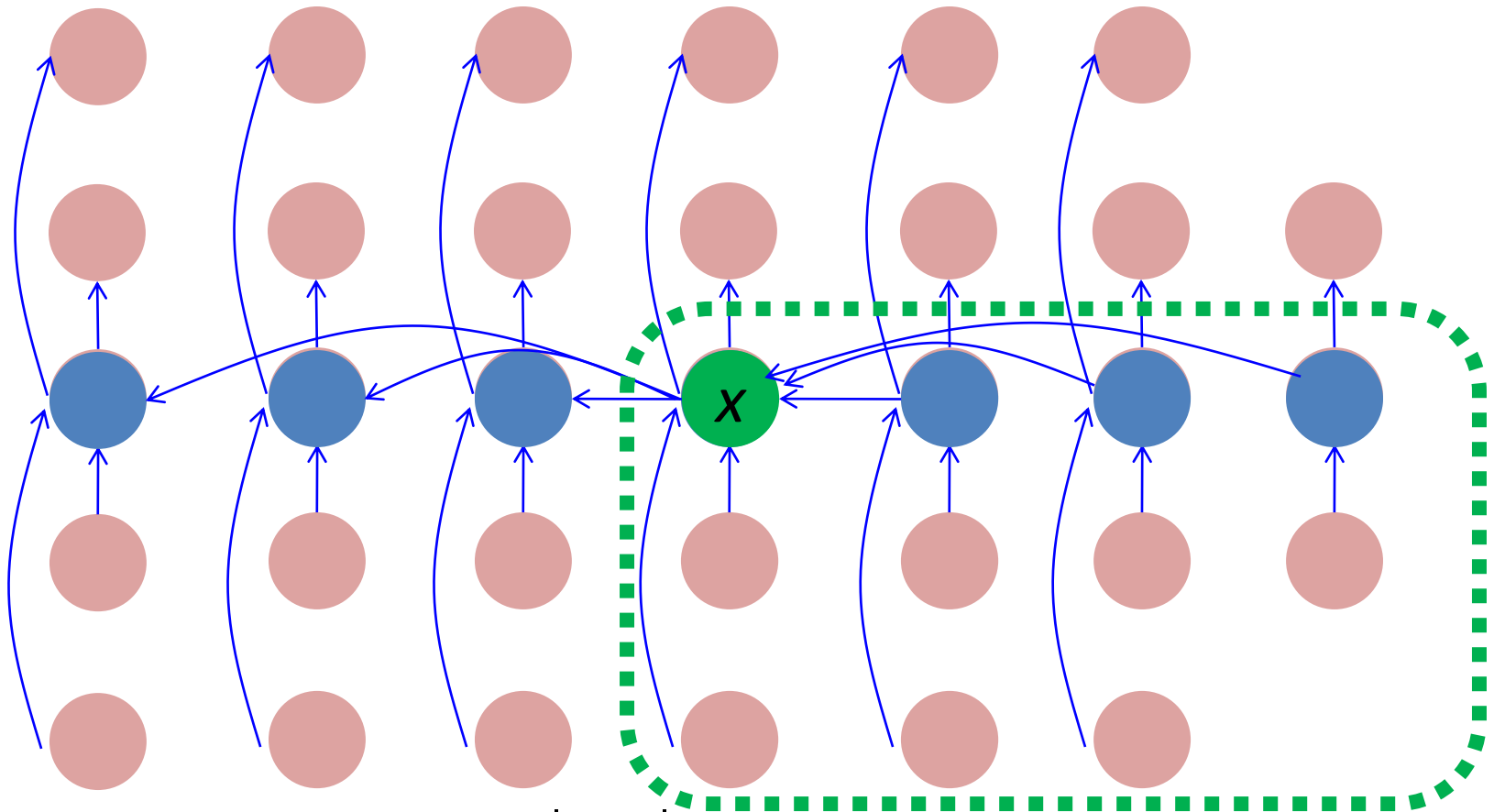
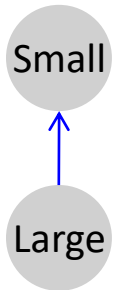
Large



How many elements are  $\leq x$ ? at least  $3 \left\lfloor \frac{1}{2} \cdot \frac{n}{5} \right\rfloor = 3 \left\lfloor \frac{n}{10} \right\rfloor$

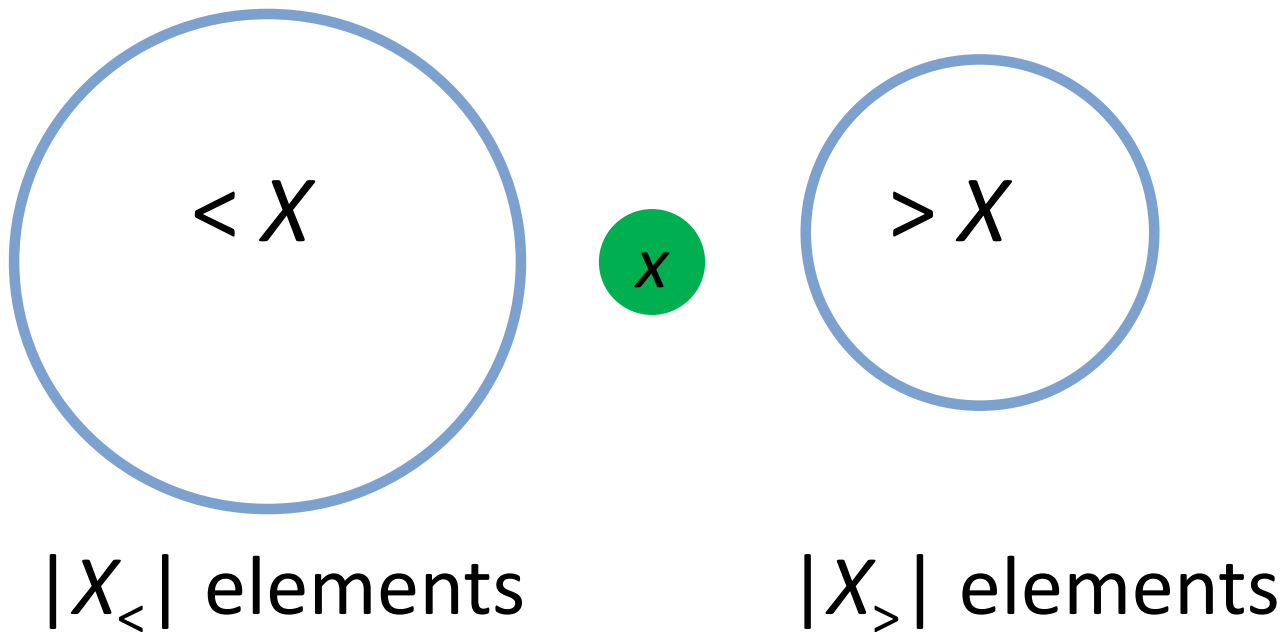
# Analysis

Assume all elements are distinct



Similarly, at least  $3 \left\lfloor \frac{n}{10} \right\rfloor$  elements  $\geq x$

- Worse case  $\rightarrow$  the  $i$ -th statistic order always falls in the **larger** group for each partition



通通加起來就是 $T(n)$

SELECT ( $A, p, r, i$ )  $T(n)$

就算worst case發生，還是可以保證 $3n/10$ 可以被丟棄

1. if  $p == r$

2. return  $A[p]$

- |          |   |
|----------|---|
| $O(n)$   | 1. Divide the $n$ elements into groups of 5                             |
| $T(n/5)$ | 2. Find the <b>median</b> of each 5-element group                       |
|          | 3. Recursively select the <b>median</b> of those group medians as pivot |

4.  $k = q - p + 1$

$O(n)$  4. Partition the input using the pivot

5. if  $i == k$

6. return  $A[q]$

7. elseif  $i < k$

8. return SELECT ( $A, p, q-1, i$ )

9. else return SELECT ( $A, q+1, r, i-k$ )

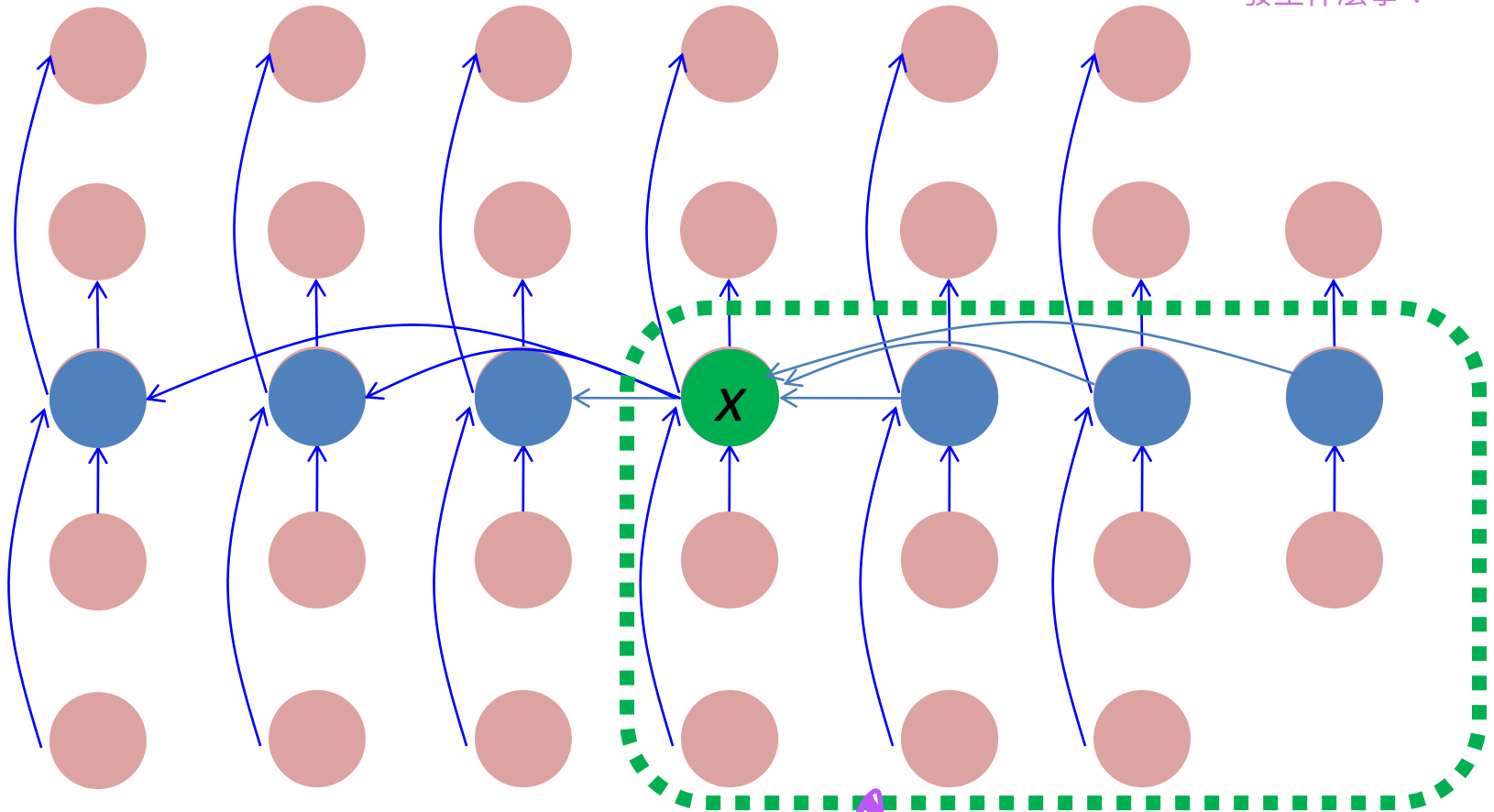
$$\max(T(|X_{<}|), T(|X_{>}|)) = O(7n/10)$$

# Deleting at least $3n/10$ elements!

Q.

為什麼一定要五個一組？

如果是三個、七個，會發生什麼事？

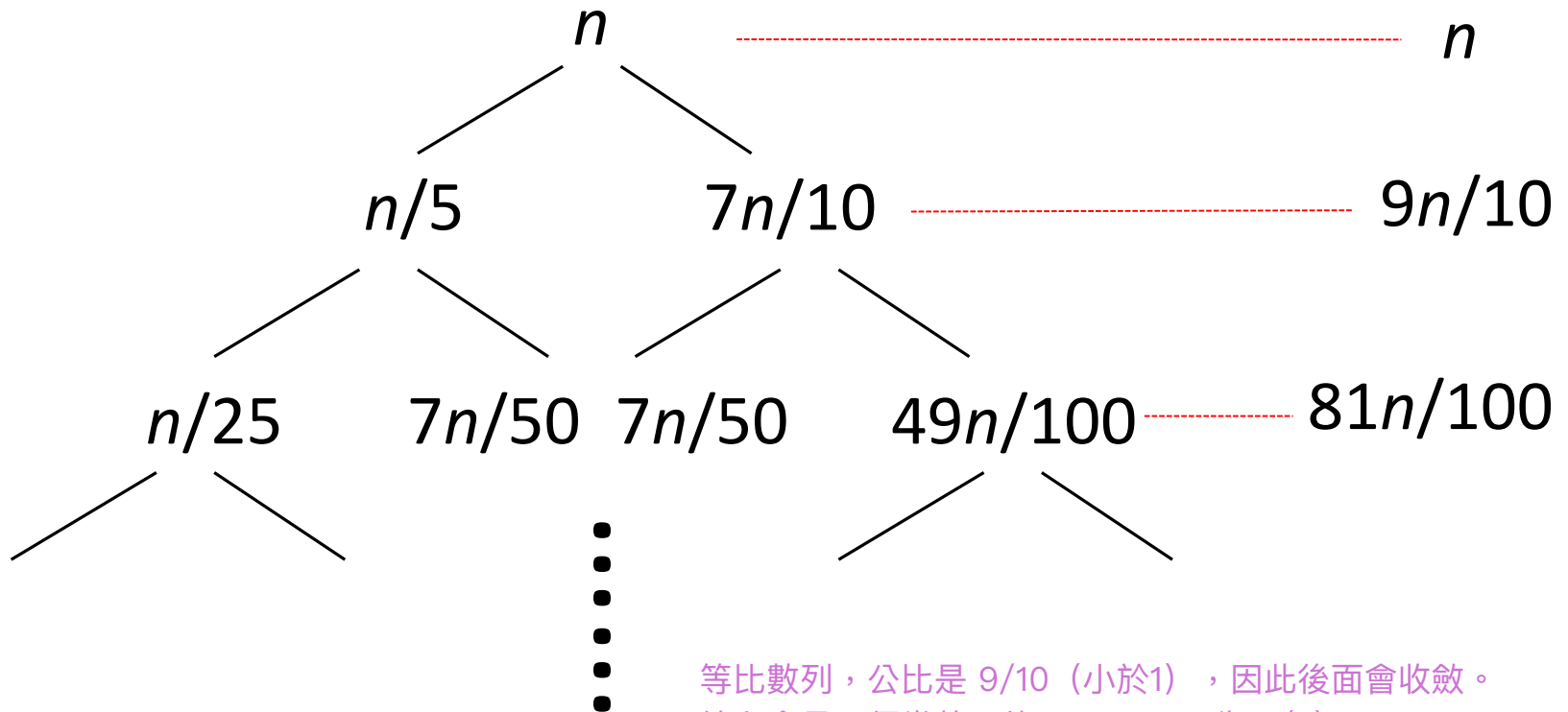


A.

如果三個一組：組數會太多，會跑出太多中位數，會花很多時間去sort。

如果是七個一組：每一組會太大，反而會花時間在組內排中位數。

$$T(n) = T\left(\frac{n}{5}\right) + T\left(\frac{7n}{10}\right) + O(n)$$



等比數列，公比是  $9/10$ （小於1），因此後面會收斂。  
總和會是一個常數，故worst-case為  $O(n)$ 。

$$T(n) = n(1 + 9/10 + 81/100 + \dots) = O(n)$$

