

# Sorting atomic items

## Chapter 5

Distribution based sorting paradigms

# The distribution-based sorting

QuickSort (S, i, j)

1. **If** (i < j) {
2.     r = pick the position of a "good pivot"
3.     Swap S[r] with S[i];
4.     p = Partition (S, i, j);
5.     QuickSort (S, i, p-1);
6.     QuickSort (S, p+1, j)
7. }

Based on Divide&Conquer. **Combine** step is not present.

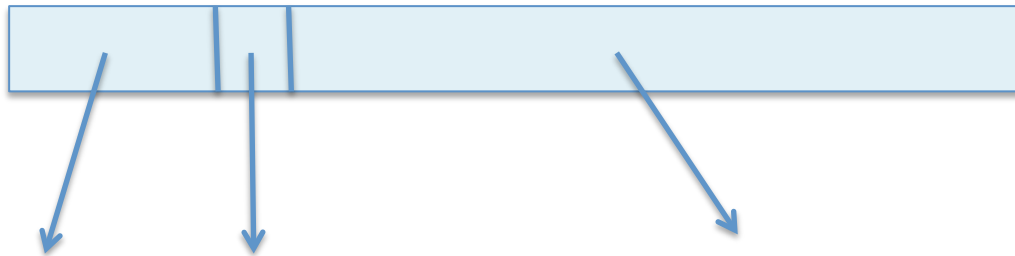
**Divide** step : Procedure Partition

QuickSort is **in place** alg.

# The distribution-based sorting

Partition divides the array in 3 parts:

$S(i, p-1)$   $S(p)$   $S(p+1, j)$



Items  $\leq$  pivot

pivot

Items  $\geq$  pivot

Partition takes  $O(n)$

If the two sub-arrays are balanced at each level of the recursion

$T(n) = 2T(n/2) + O(n) = O(n \log n)$  as MergeSort

To study the **worst case**, we look at the position of  $q$  that maximize the time

$T(n) = \max_{0 \leq q \leq n-1} (T(q) + T(n-q-1)) + O(n)$  where  $q$  range from 0 to  $n-1$

# The distribution-based sorting

Guess:  $T(n) \leq cn^2$

$$T(n) = \max_{0 \leq q \leq n-1} (cq^2 + c(n-q-1)^2) + O(n) = c \max_{0 \leq q \leq n-1} (q^2 + (n-q-1)^2) + O(n)$$

Gives the maximum when  $q=0$  or  $q=n-1$ :

$$(q^2 + (n-q-1)^2) \leq (n-1)^2 = n^2 - 2n - 1$$

$$T(n) \leq c(n^2 - 2n - 1) + O(n) \leq cn^2 \quad \text{worst case}$$

# QuickSort

## Expected running time

- Sequence  $S(1,n)$ ; Rank  $Z(1, n)$  :  $Z_i$  is the  $i$ -th smallest element;
- $p_{i,j}$  is the probability that a comparison  $Z_i : Z_j$  occurs during an execution of QuickSort;
- The expected total number is:  $E = \sum_{i=1}^n \sum_{j>i} p_{i,j}$

### Remarks:

1. In Partition two items are compared if one of them is a pivot.
2. If two items go in different sub-arrays they will never be compared in the future.

# Expected running time

If  $j=i+1$  the elements are compared for sure: there not exist an element that, being the pivot, can put them in separate sub-arrays as pivot.  $p_{i,i+1}=1$

If  $j>i+1$  consider the set of elements  $A = \{Z_i, Z_{i+1}, \dots, Z_j\}$  if as pivot is selected an element not in  $A$  all elements remain in the same partition and  $Z_i$  and  $Z_j$  are not compared.

if  $Z_i$  or  $Z_j$  are selected as pivot,  $Z_i$  and  $Z_j$  are compared

If  $Z_k$  is selected with  $k \neq i, j$   $A$  is split into 2 sub-arrays and  $Z_i$  and  $Z_j$  are not compared.

So  $p_{i,j} = 2/(j-i+1)$  (when  $j=i+1$   $p_{i,j}=1$ )

# QuickSort

## Expected running time

The expected total number is:

$$E = \sum_{i=1}^n \sum_{j>i} p_{i,j} = \sum_{i=1}^n \sum_{j>i} p_{i,j} = 2/(j-i+1) = \sum_{i=1}^n \sum_{k=1}^{n-i} 2/(k+1) \leq$$

$$2 \sum_{i=1}^n \sum_{k=1}^n 1/k$$

since  $\sum_{k=1}^n 1/k = \ln n + O(1)$  hence

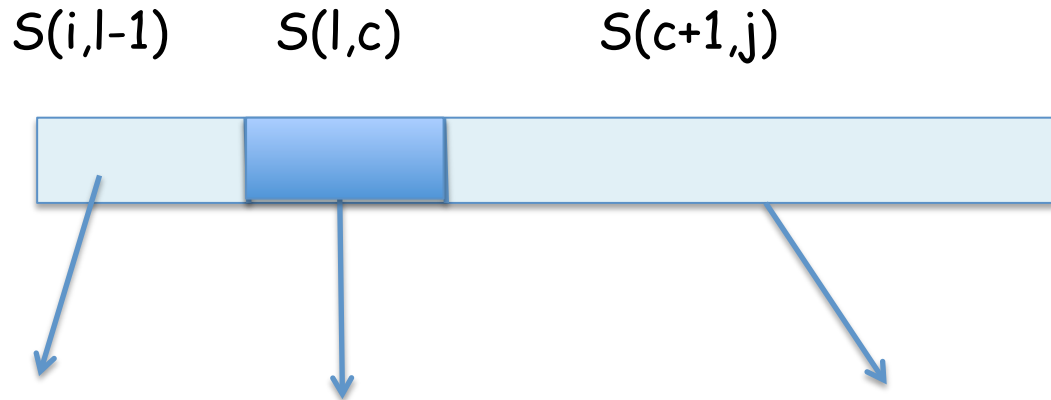
$$E = O(n \log n)$$

In average quicksort takes at most  $1.45 n \log n$  operations

# 3-ways Partition

In procedure Partition of QuickSort elements equal to the pivot are arbitrarily distributed among the 2 partitions.

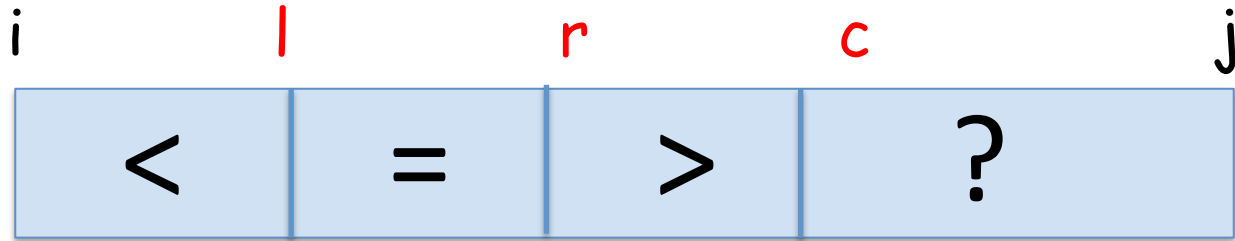
In 3-ways Partition we have:



3-ways Partition takes  $O(n)$ . The central part can be discarded in the Recursion.



# 3-ways partition



Variable  $l$  indicates the last item  $<$  than the pivot

Variable  $r$  indicates the first item  $>$  to the pivot

Variable  $c$  indicates the next item to be considered.

If  $S[c] > \text{pivot}$   $c = c + 1$

If  $S[c] = \text{pivot}$  exchange  $S[c]$  and  $S[r]$ ;  $c = c + 1$ ;  $r = r + 1$ ;

If  $S[c] < \text{pivot}$   $l = l + 1$ ; exchange  $S[c]$  and  $S[l]$ ;  $S[c]$  and  $S[r]$ ;  $r = r + 1$ ;  $c = c + 1$

# 3-ways partition

l r c  
5 2 9 12 12 12 20 18 13 15 17 19 12 8

l r c  
5 2 9 12 12 12 20 18 13 15 17 19 12 8

no exchange

l r c  
5 2 9 12 12 12 12 18 13 15 17 18 20 8

1 exchange 12:20

5 2 9 8 12 12 12 18 13 15 17 18 20 12

2 exchanges 8:12 and

5 2 9 8 12 12 12 18 13 15 17 19 20 12

and 12:18

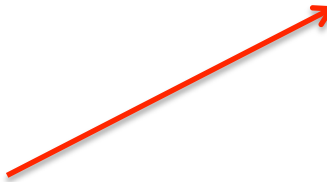
5 2 9 8 12 12 12 12 13 15 17 18 20 18  
l r c

# 3-ways Partition( $S, i, j$ )

---

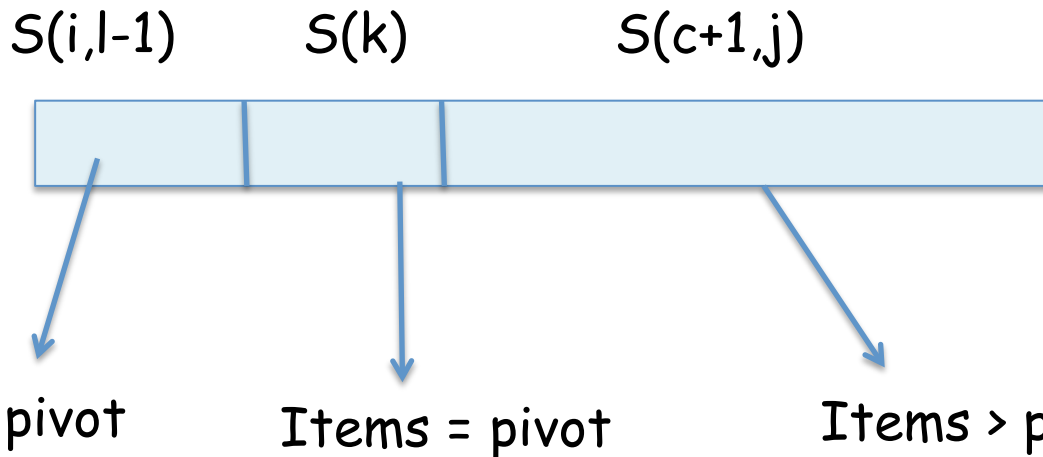
```
1:  $P = S[i]; l = i; r = i + 1;$ 
2: for ( $c = r; c \leq j; c++$ ) do
3:     if ( $S[c] == P$ ) then
4:         swap  $S[c]$  with  $S[r];$ 
5:          $r++;$ 
6:     else if ( $S[c] < P$ ) then
7:         swap  $S[c]$  with  $S[l];$ 
8:         swap  $S[c]$  with  $S[r];$ 
9:          $r++; l++;$ 
10:    end if
11: endfor
12: return  $\langle l, r-1 \rangle$ 
```

Recursion on  $S(i, l-1)$   
and  $S(r, j)$



# Modify QuickSort to select the k-th item

**Idea:** select an item at random  $S[r]$  and call Partition.  
Let  $k$  the position of the pivot.



If  $k$  is in the range of the items equal to the pivot return :  $S[r]$  is the  $k$ -th item.

If  $k$  is in the range the items less than the pivot: Recurse on  $S(i, l-1)$  and  $k$ .

If  $k$  is in the range the items greater than the pivot: Recurse on  $S(c+1, j)$  and  $k-c$ .

# RandSelect

---

**Algorithm 5.5** Selecting the  $k$ -th ranked item: **RANDSELECT**( $S, k$ )

---

```
1:  $r =$  pick a random item from  $S$ ;  
2:  $S_{<} =$  items of  $S$  which are smaller than  $S[r]$ ;  
3:  $S_{>} =$  items of  $S$  which are larger than  $S[r]$ ;  
4:  $n_{<} = |S_{<}|$ ;  
5:  $n_{=} = |S| - (|S_{<}| + |S_{>}|)$ ;  
6: if ( $k \leq n_{<}$ ) then  
7:     return RANDSELECT( $S_{<}, k$ );  
8: else if ( $k \leq (n_{<} + n_{=})$ ) then  
9:     return  $S[r]$ ;  
10: else  
11:     return RANDSELECT( $S_{>}, k - n_{<} - n_{=}$ );  
12: end if
```

---

# Expected running time

$$T(n) = T(n-1) + O(n) = O(n^2) \quad \text{Worst case time}$$
$$O(n) \quad \text{Average time} \quad \text{RAM model}$$
$$O(n/B) \quad \text{I/O's for the disk model}$$

"good selection" a partition where  $n_1$  and  $n_2$  are not larger than  $2/3n$ . Positions of the pivot for a good selection: the blue



Probability to have a good selection is  $1/3$ . Let  $T_a$  the average time:

$$T_a(n) \leq O(n) + 1/3 T_a((2/3)n) + 2/3 T_a(n) \quad \text{subtract } T_a(n)$$

$$1/3 T_a(n) \leq O(n) + 1/3 T_a((2/3)n) \quad \text{multiply by 3}$$

$$T_a(n) \leq O(n) + T_a((2/3)n)$$

# Expected running time

$$T_a(n) \leq O(n) + T_a(2/3n)$$

It can be computed with Master Th. ( or by substitution)

$$T_a(n) \leq O(n)$$

RandSelect is very efficient in average!

2-level model:

$$T_a(n) \leq O(n/B) + T_a(2/3n) = O(n/B)$$

Since the procedure partition can be executed in the 2-level model with a single pass over the input items.

# Use RandSelect to improve QuickSort

- Instead of 1 pivot, select at random  $2s+1$  pivots.
- Select the median pivot among the  $2s+1$
- $s=1$  select 3 pivot and with 2 comparisons select the median.
- $s>1$  : sort the items and select the median  $O(s \log s)$
- select the median ( $k = s/2$ ) by RandSelect  $O(s)$  average.
- Select as pivot the median item of the **whole** array  $k=n/2$
- Select a pivot that generates 2 a balanced partition, the 2 parts are fractions of  $n$ :  $\alpha n$  and  $(1-\alpha)n$  with  $\alpha < 0.5$ . Apparently meaningless, is **good for parallel CPU**.