

## 归并排序及其他

- 归并排序
- 排序算法时间下界
- 外排序方法简介

数据之法  
结构之美  
算法之道

# 归并排序算法提出者



**John von Neumann**

冯·诺依曼

美国科学院院士

普林斯顿大学教授

现代计算机之父

**If people do not believe that mathematics is simple, it is only because they do not realize how complicated life is.**

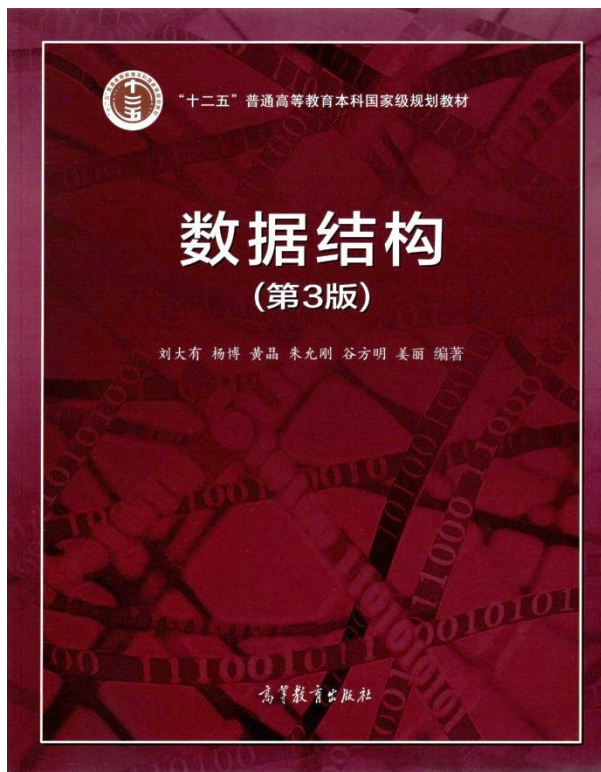
如果有人不相信数学是简单的，  
那是因为他们没有意识到生活  
有多复杂。





# 归并排序及其他

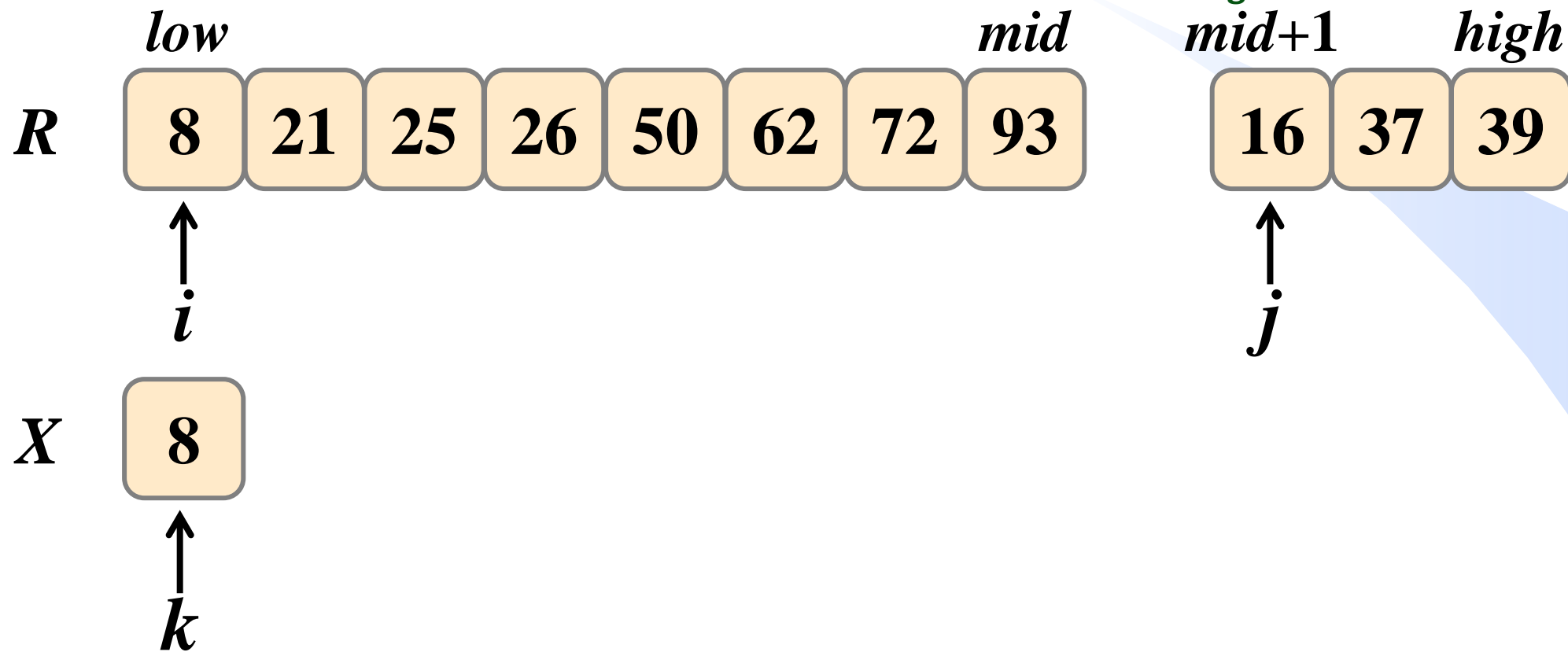
- 归并排序
- 排序算法时间下界
- 外排序方法简介



数据之法  
结构之美  
算法之道

## 两个有序子数组合并成一个大的有序子数组

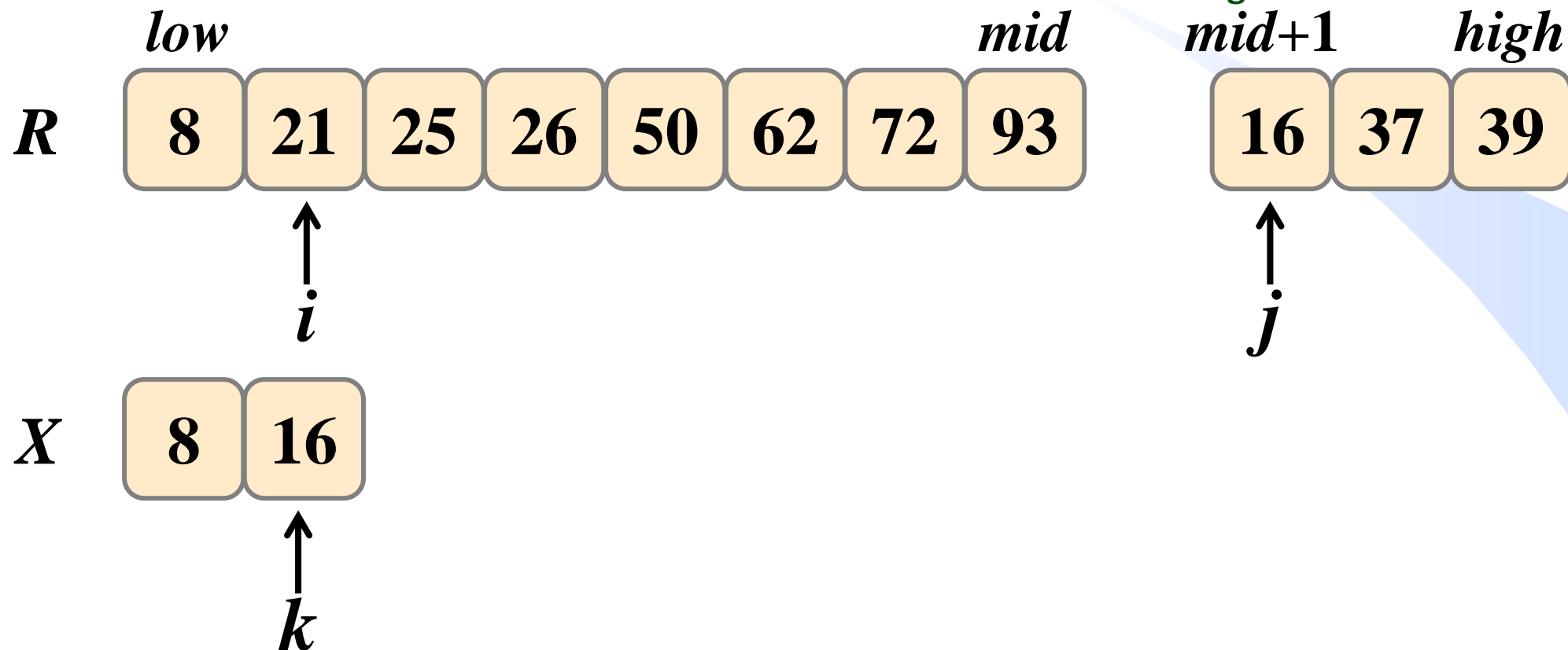
```
void Merge(int R[], int low, int mid, int high){  
    //将两个相邻的有序数组( $R_{low}, \dots, R_{mid}$ )和( $R_{mid+1}, \dots, R_{high}$ )合并成一个有序数组
```



- 通过 *i* 和 *j* 扫描两个子数组，将  $R[i]$  与  $R[j]$  的关键词较小者放入  $X[k]$  中
- 当一个子数组扫描完毕后，将另一个子数组的剩余部分直接放入  $X$  中

## 两个有序子数组合并成一个大的有序子数组

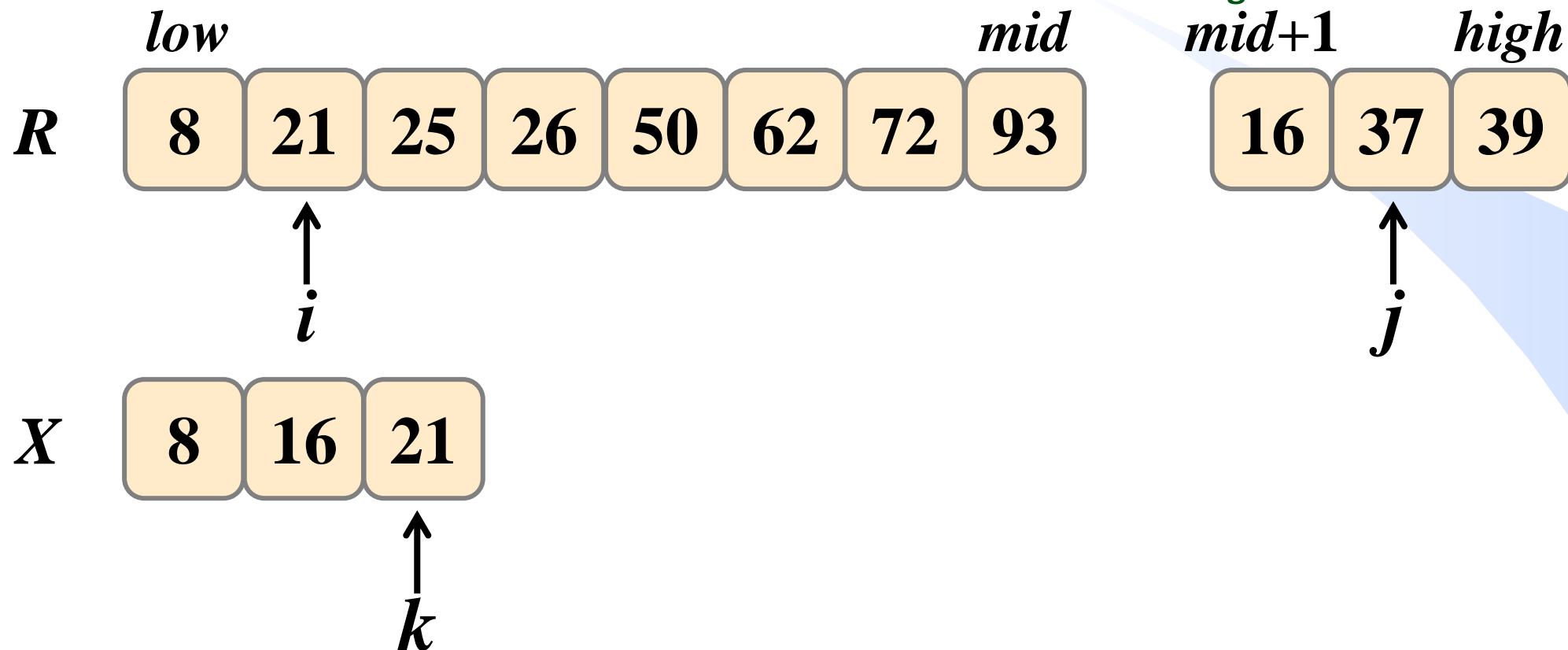
```
void Merge(int R[], int low, int mid, int high){  
    //将两个相邻的有序数组( $R_{low}, \dots, R_{mid}$ )和( $R_{mid+1}, \dots, R_{high}$ )合并成一个有序数组
```



- 通过  $i$  和  $j$  扫描两个子数组，将  $R[i]$  与  $R[j]$  的关键词较小者放入  $X[k]$  中
- 当一个子数组扫描完毕后，将另一个子数组的剩余部分直接放入  $X$  中

## 两个有序子数组合并成一个大的有序子数组

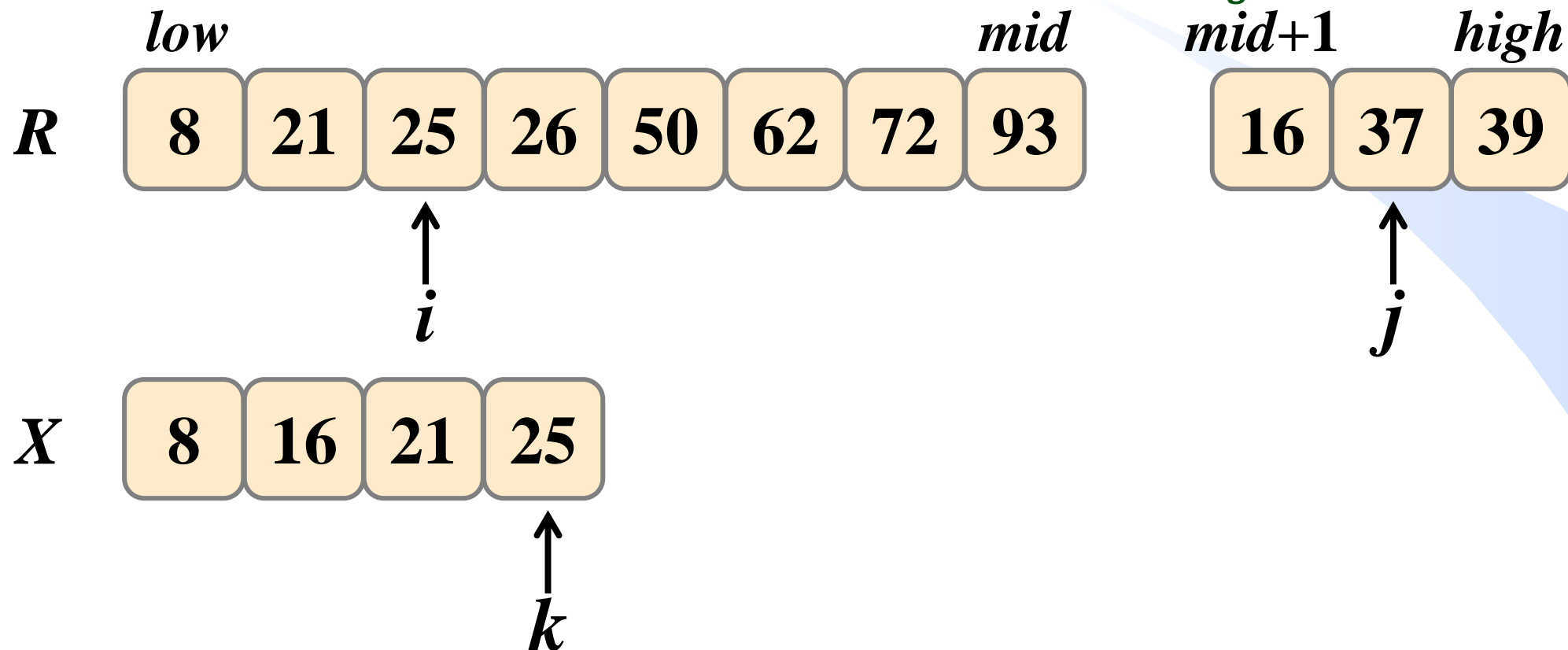
```
void Merge(int R[], int low, int mid, int high){  
    //将两个相邻的有序数组( $R_{low}, \dots, R_{mid}$ )和( $R_{mid+1}, \dots, R_{high}$ )合并成一个有序数组
```



- 通过  $i$  和  $j$  扫描两个子数组，将  $R[i]$  与  $R[j]$  的关键词较小者放入  $X[k]$  中
- 当一个子数组扫描完毕后，将另一个子数组的剩余部分直接放入  $X$  中

## 两个有序子数组合并成一个大的有序子数组

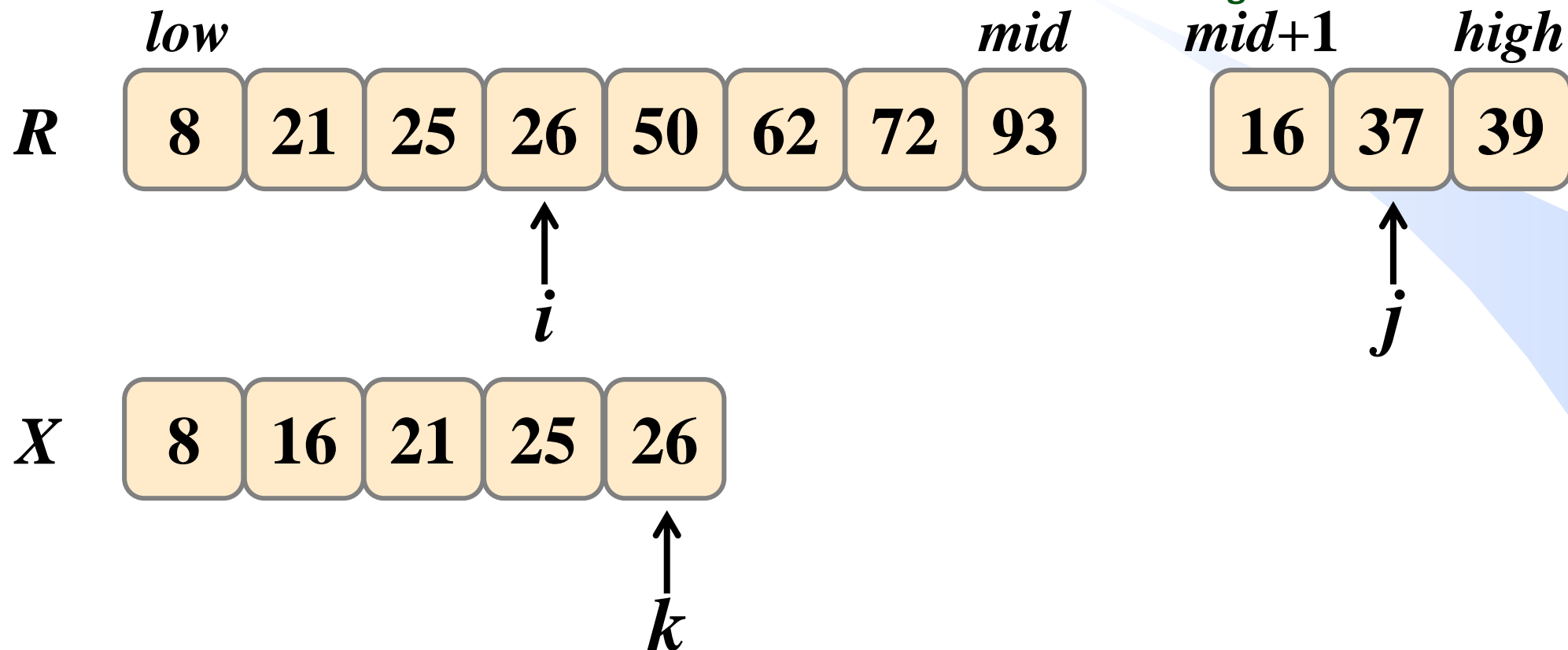
```
void Merge(int R[], int low, int mid, int high){  
    //将两个相邻的有序数组( $R_{low}, \dots, R_{mid}$ )和( $R_{mid+1}, \dots, R_{high}$ )合并成一个有序数组
```



- 通过  $i$  和  $j$  扫描两个子数组，将  $R[i]$  与  $R[j]$  的关键词较小者放入  $X[k]$  中
- 当一个子数组扫描完毕后，将另一个子数组的剩余部分直接放入  $X$  中

## 两个有序子数组合并成一个大的有序子数组

```
void Merge(int R[], int low, int mid, int high){  
    //将两个相邻的有序数组( $R_{low}, \dots, R_{mid}$ )和( $R_{mid+1}, \dots, R_{high}$ )合并成一个有序数组
```

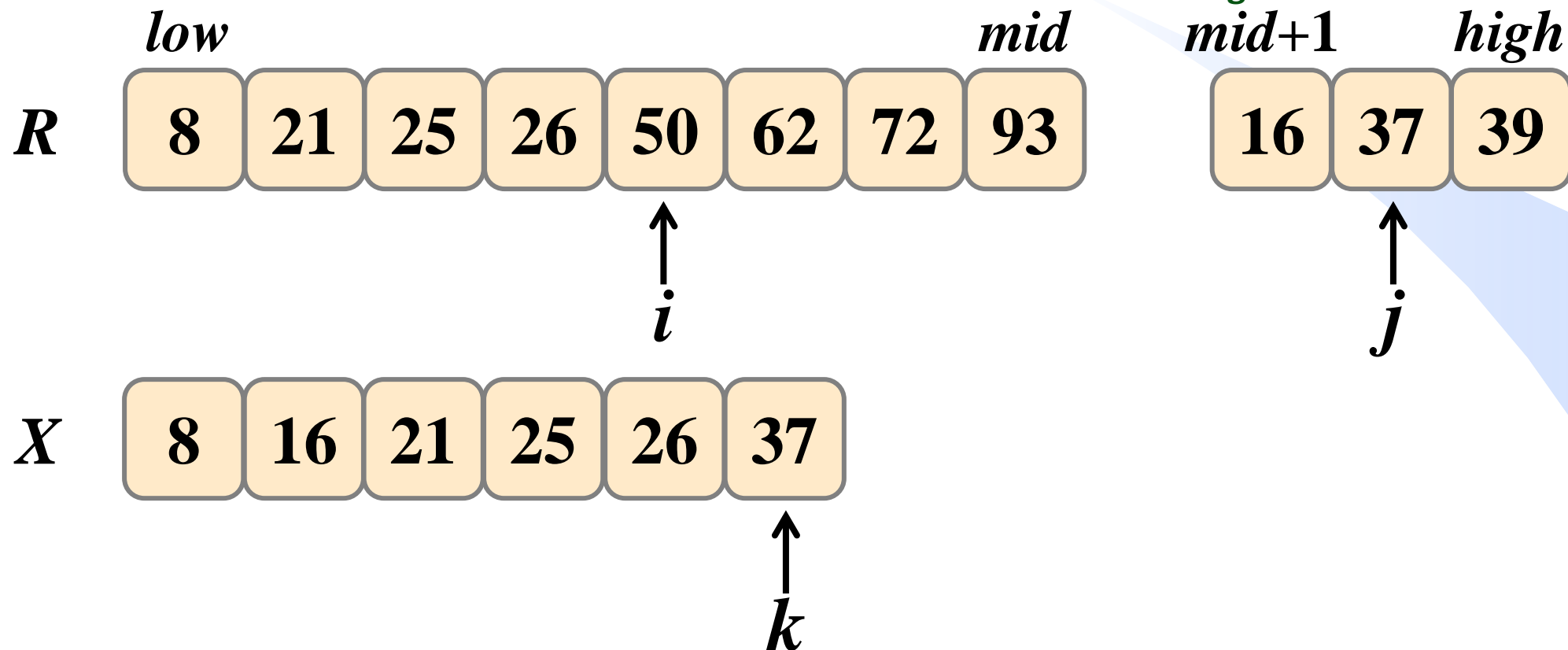


- 通过  $i$  和  $j$  扫描两个子数组，将  $R[i]$  与  $R[j]$  的关键词较小者放入  $X[k]$  中
- 当一个子数组扫描完毕后，将另一个子数组的剩余部分直接放入  $X$  中



## 两个有序子数组合并成一个大的有序子数组

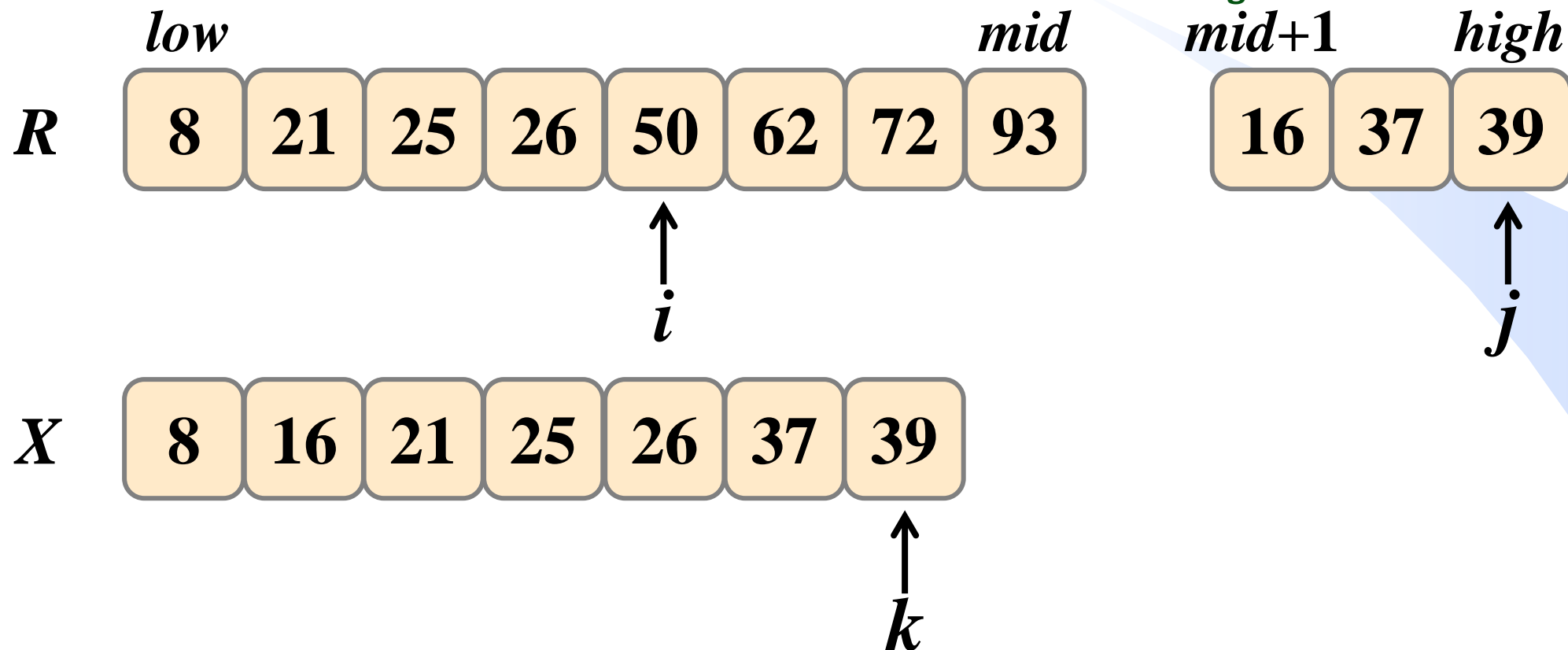
```
void Merge(int R[], int low, int mid, int high){  
    //将两个相邻的有序数组( $R_{low}, \dots, R_{mid}$ )和( $R_{mid+1}, \dots, R_{high}$ )合并成一个有序数组
```



- 通过  $i$  和  $j$  扫描两个子数组，将  $R[i]$  与  $R[j]$  的关键词较小者放入  $X[k]$  中
- 当一个子数组扫描完毕后，将另一个子数组的剩余部分直接放入  $X$  中

## 两个有序子数组合并成一个大的有序子数组

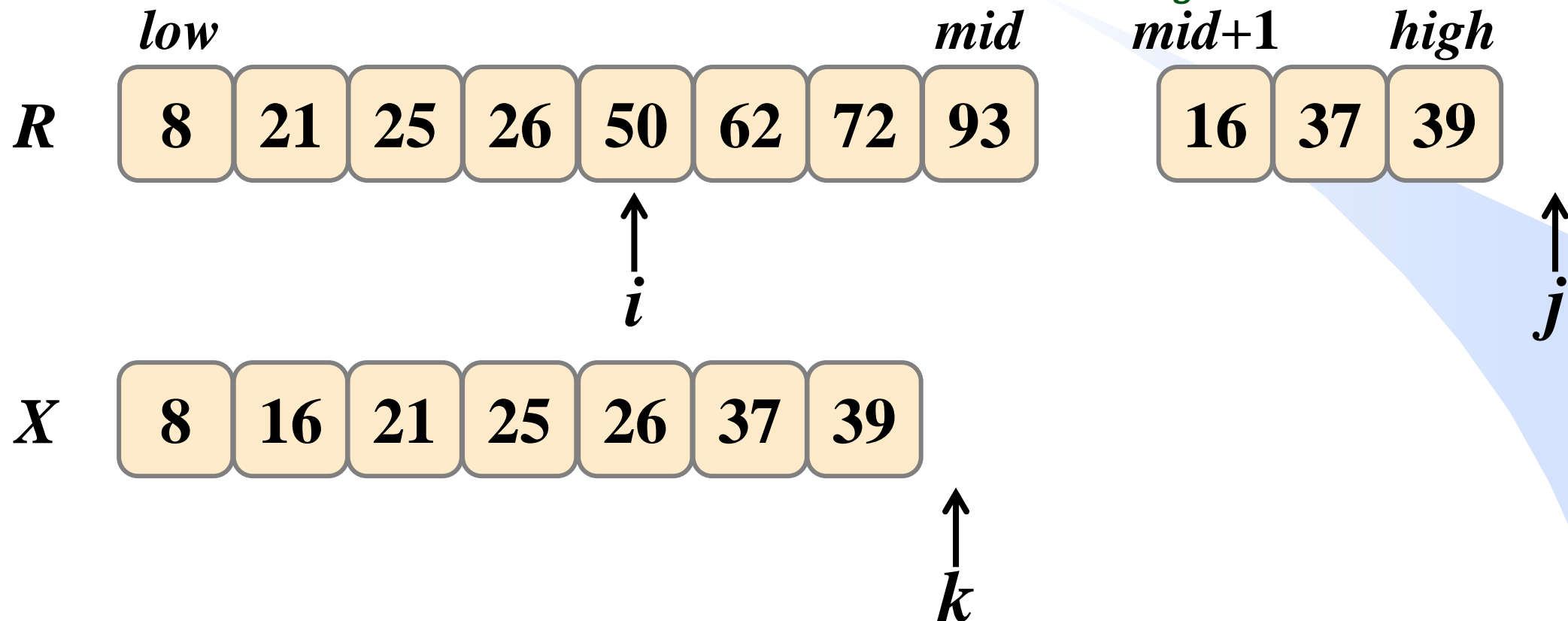
```
void Merge(int R[], int low, int mid, int high){  
    //将两个相邻的有序数组( $R_{low}, \dots, R_{mid}$ )和( $R_{mid+1}, \dots, R_{high}$ )合并成一个有序数组
```



- 通过  $i$  和  $j$  扫描两个子数组，将  $R[i]$  与  $R[j]$  的关键词较小者放入  $X[k]$  中
- 当一个子数组扫描完毕后，将另一个子数组的剩余部分直接放入  $X$  中

## 两个有序子数组合并成一个大的有序子数组

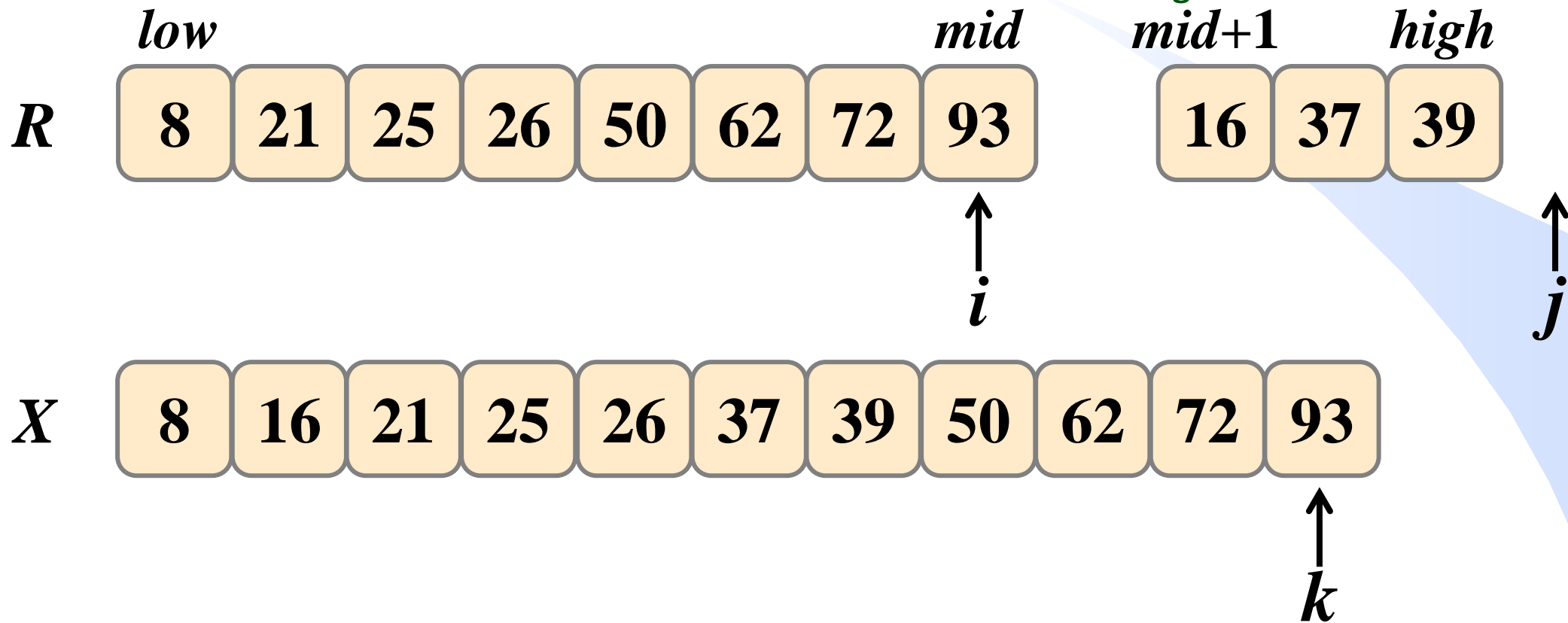
```
void Merge(int R[], int low, int mid, int high){  
    //将两个相邻的有序数组( $R_{low}, \dots, R_{mid}$ )和( $R_{mid+1}, \dots, R_{high}$ )合并成一个有序数组
```



- 通过  $i$  和  $j$  扫描两个子数组，将  $R[i]$  与  $R[j]$  的关键词较小者放入  $X[k]$  中
- 当一个子数组扫描完毕后，将另一个子数组的剩余部分直接放入  $X$  中

## 两个有序子数组合并成一个大的有序子数组

```
void Merge(int R[], int low, int mid, int high){  
    //将两个相邻的有序数组( $R_{low}, \dots, R_{mid}$ )和( $R_{mid+1}, \dots, R_{high}$ )合并成一个有序数组
```



- 通过  $i$  和  $j$  扫描两个子数组，将  $R[i]$  与  $R[j]$  的关键词较小者放入  $X[k]$  中
- 当一个子数组扫描完毕后，将另一个子数组的剩余部分直接放入  $X$  中

# 两个有序子数组合并成一个大的有序子数组

```

void Merge(int R[], int low, int mid, int high){
//将两个相邻的有序数组( $R_{low}, \dots, R_{mid}$ )和( $R_{mid+1}, \dots, R_{high}$ )合并成一个有序数组
    int i=low, j=mid+1, k=0;
    int *X=new int[high-low+1];
    while(i<=mid && j<=high)
        if(R[i]<=R[j]) X[k++]=R[i++];
        else X[k++]=R[j++];
    while(i<=mid) X[k++]=R[i++]; //复制余留记录
    while(j<=high) X[k++]=R[j++];
    for(i=0; i<=high-low; i++) //将X拷贝回R
        R[low+i]=X[i];
    delete []X;
}

```

时间复杂度  $O(n)$   
空间复杂度  $O(n)$

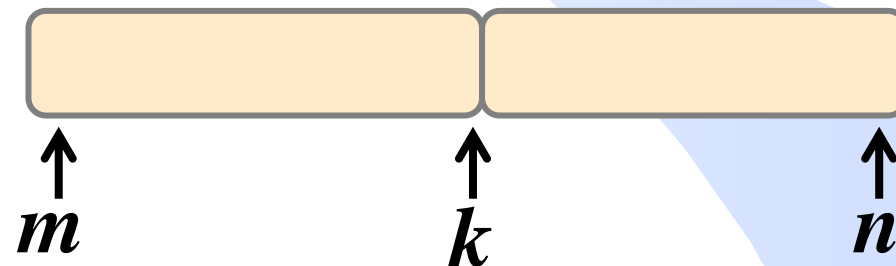
$X_0$	$X_1$	$X_2$	$\dots$	$X_i$	$\dots$	$X_{high-low}$
-------	-------	-------	---------	-------	---------	----------------

$R_{low}$	$R_{low+1}$	$R_{low+2}$	$\dots$	$R_{low+i}$	$\dots$	$R_{high}$
-----------	-------------	-------------	---------	-------------	---------	------------



# 归并排序（递归形式）

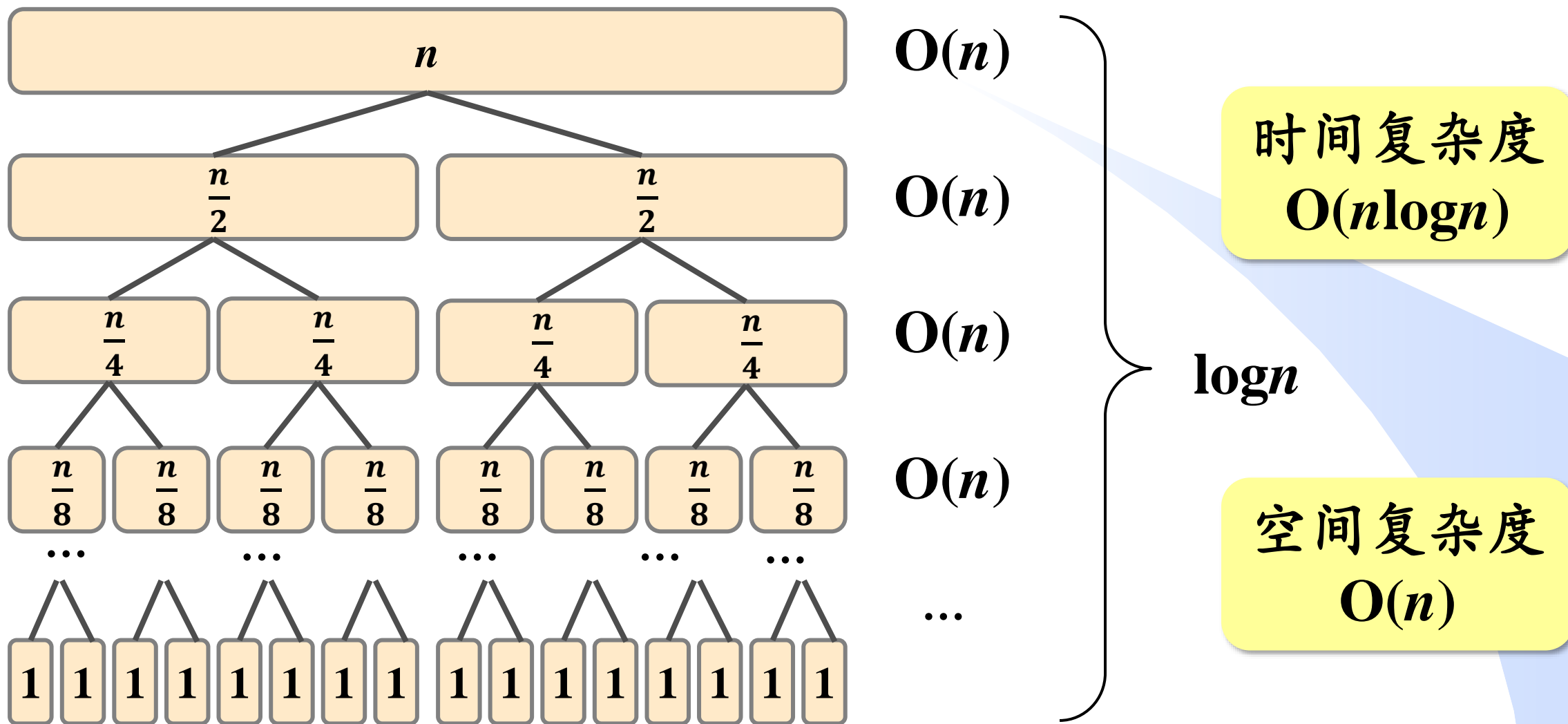
```
void MergeSort(int R[], int m, int n){  
    if(m < n){  
        int k = (m+n)/2; //将待排序序列等分为两部分  
        MergeSort(R, m, k);  
        MergeSort(R, k+1, n);  
        Merge(R, m, k, n);  
    }  
}
```



分治法

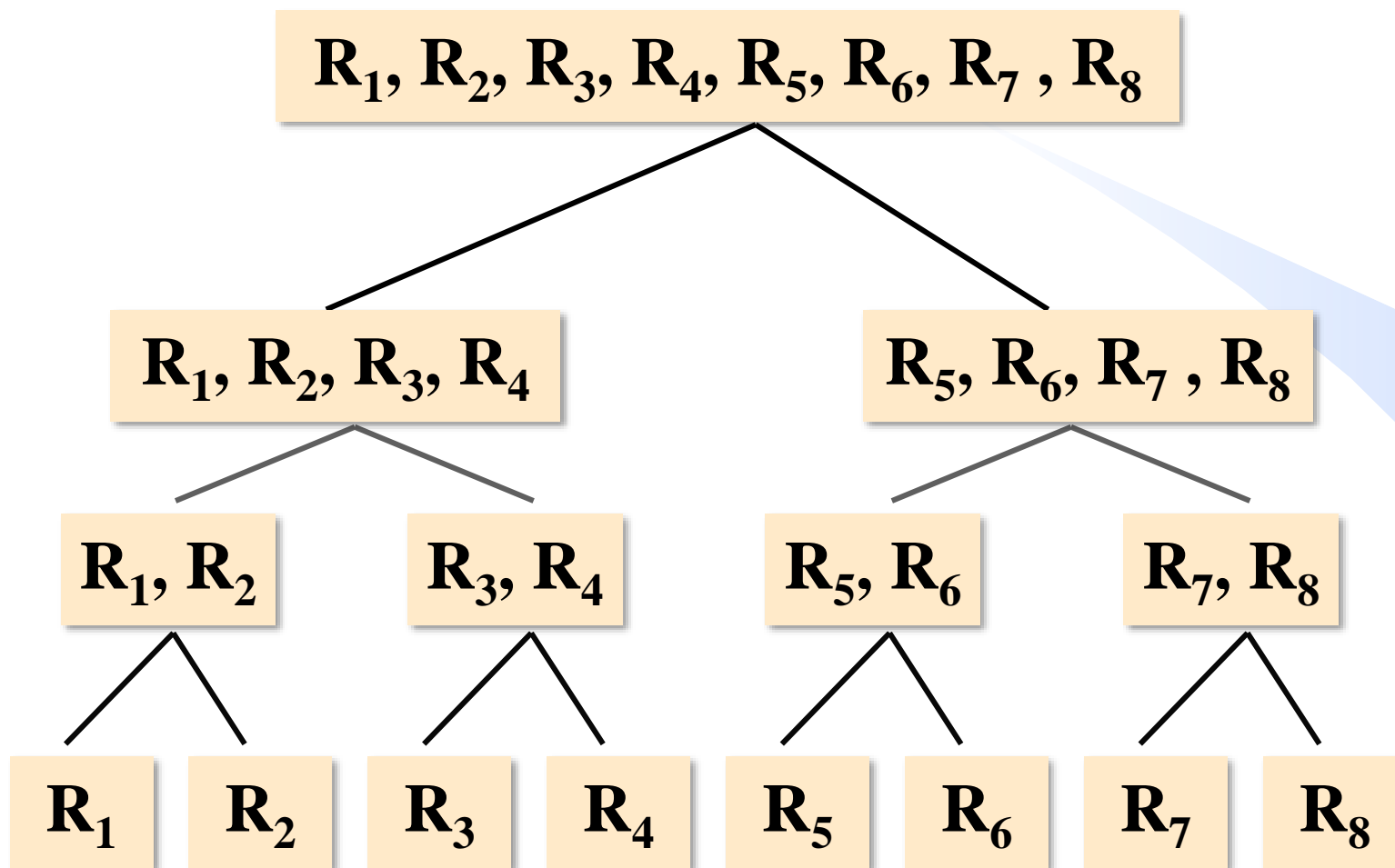
# 归并排序——时间复杂度

A

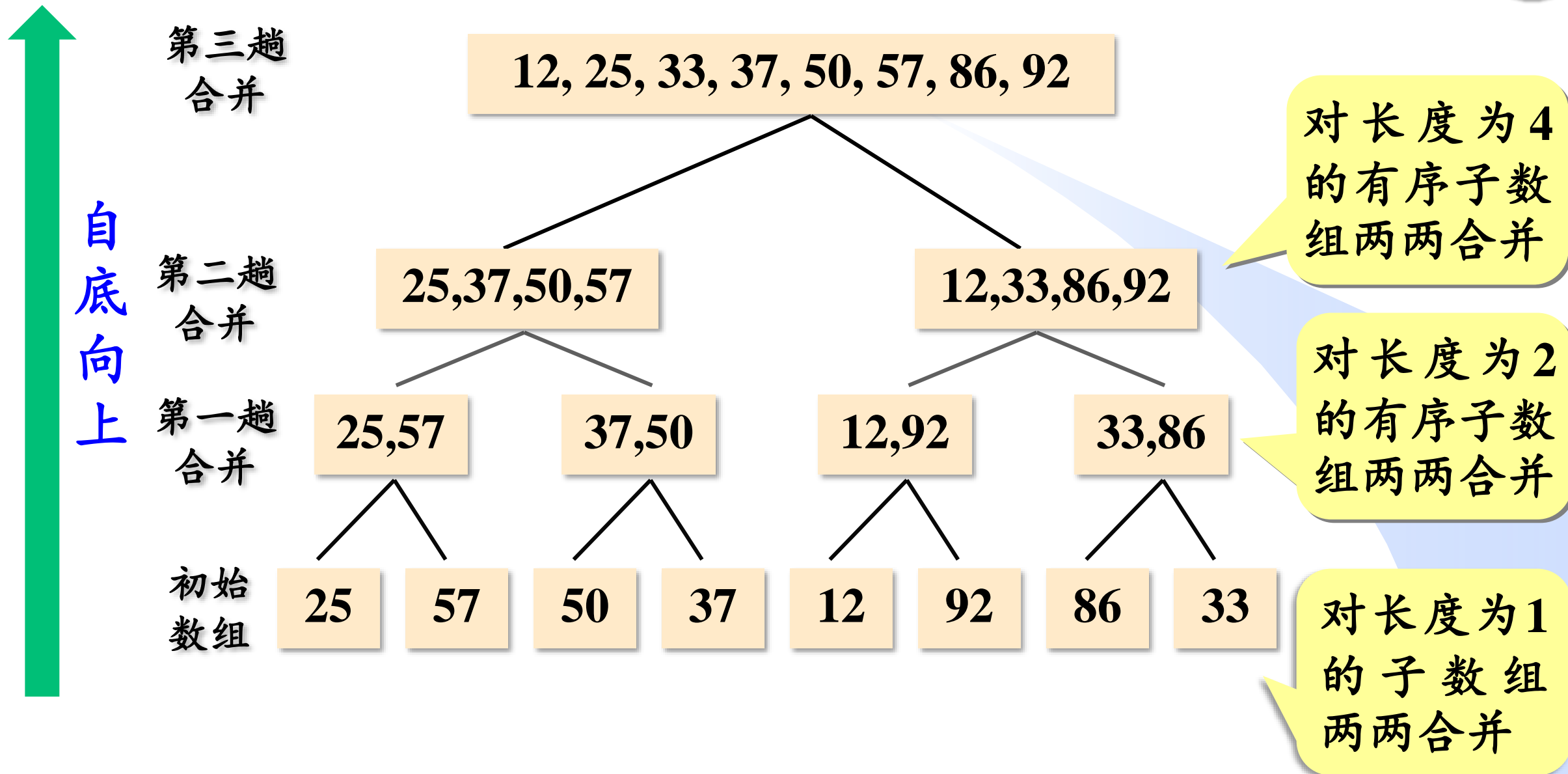


# 归并排序——递归过程

自顶向下



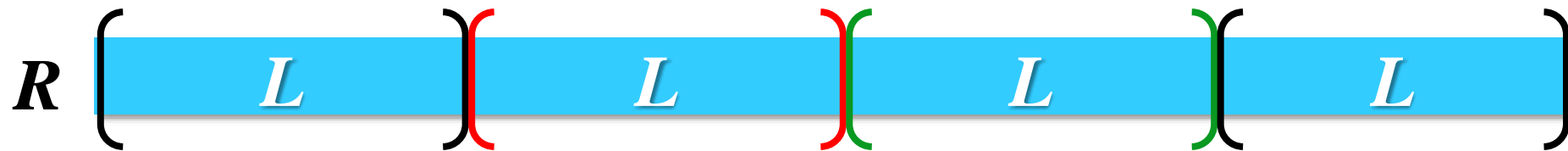
# 归并排序——非递归过程



# 归并排序——非递归形式

```
void MergeSort(int R[], int n){  
    for(int L=1; L<n; L*=2)  
        MergePass(R, n, L);  
}
```

执行一趟合并过程，将数组  $R$  中相邻的、长度为  $L$  的各有序子数组依次两两合并

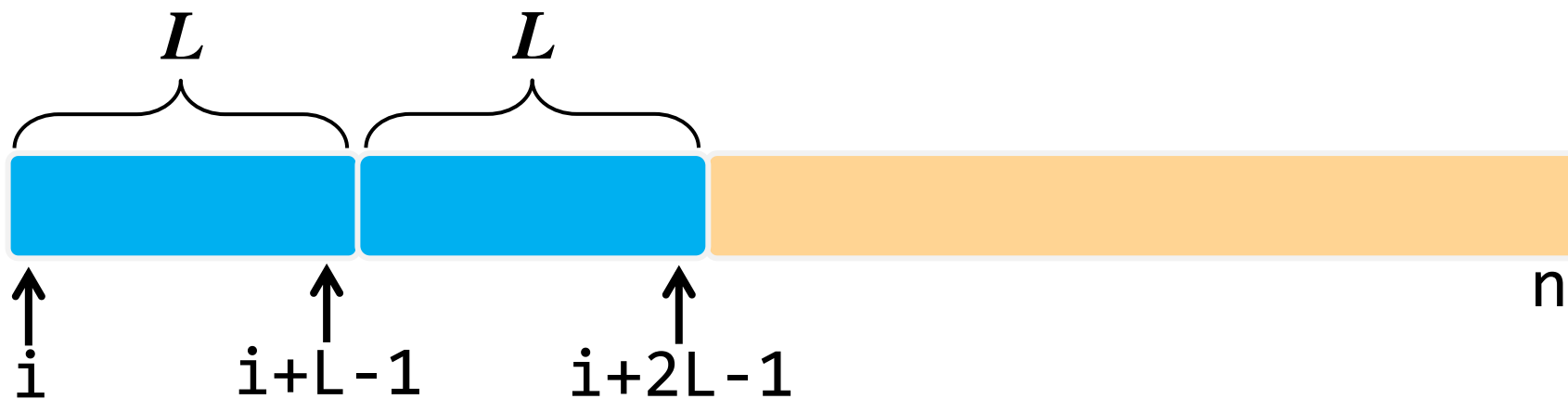




## 一趟合并

```
void MergePass(int R[], int n, int L){  
    int i;  
    for(i=1; i+2*L-1<=n; i+=2*L)  
        Merge(R, i, i+L-1, i+2*L-1);  
}
```

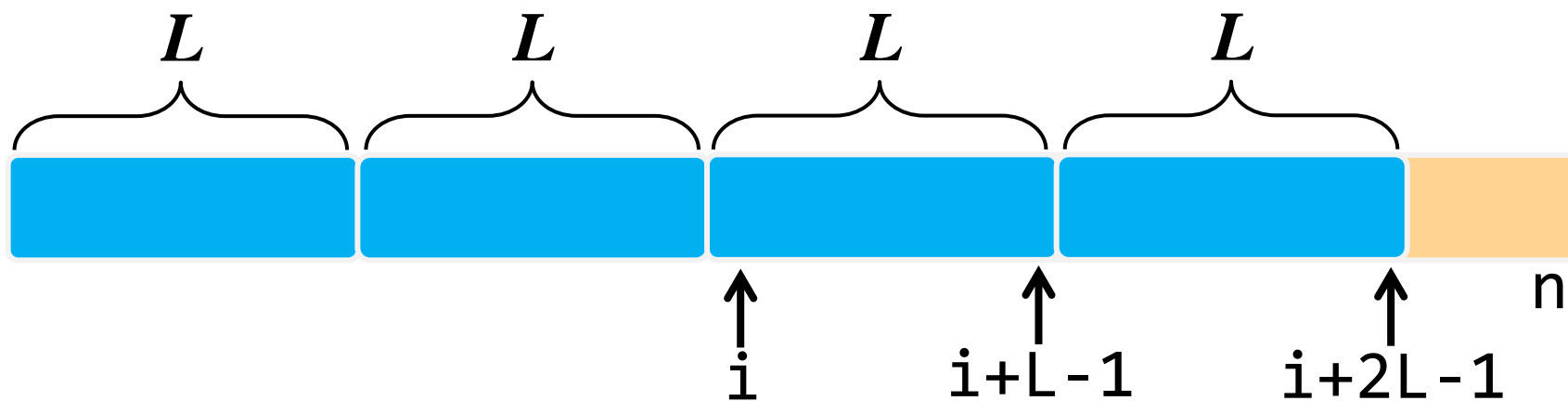
执行一趟合并过程，将数组 $R$ 中长度为 $L$ 的相邻有序子数组两两合并



## 一趟合并

```
void MergePass(int R[], int n, int L){  
    int i;  
    for(i=1; i+2*L-1<=n; i+=2*L)  
        Merge(R, i, i+L-1, i+2*L-1);  
}
```

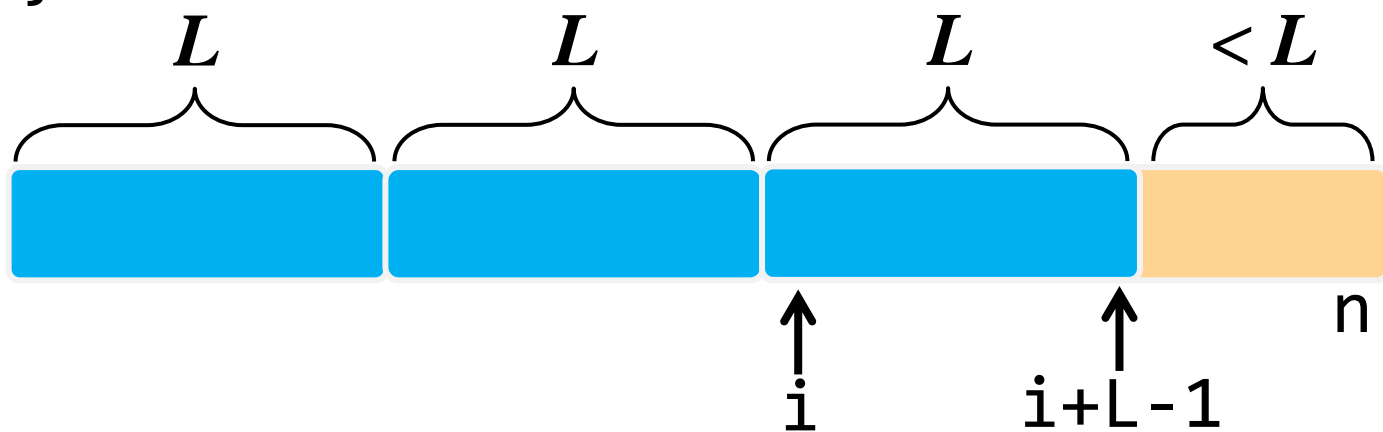
执行一趟合并过程，将数组 $R$ 中长度为 $L$ 的相邻有序子数组两两合并



# 一趟合并

```
void MergePass(int R[], int n, int L){
    int i;
    for(i=1; i+2*L-1<=n; i+=2*L)
        Merge(R, i, i+L-1, i+2*L-1);
    //处理余留的长度小于2*L的子数组
    if(i+L-1<n)
        Merge(R, i, i+L-1, n);    //L<剩余部分长度<2L
}
```

若剩余子数组长度不够  $2L$ ，则该for循环结束



情况1

$L < \text{剩余部分} < 2L$

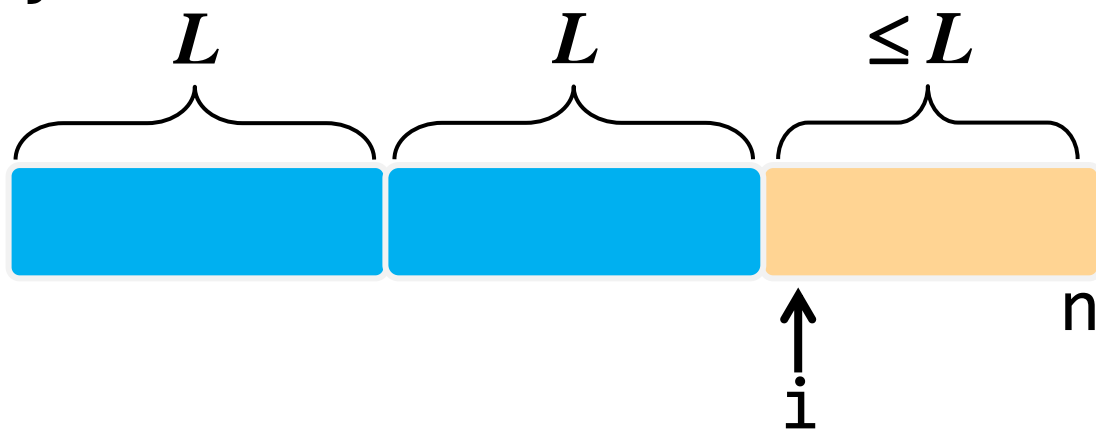
方案：将剩余的  
长度为  $L$  和长度  $< L$  的  
两个子数组合并

## 一趟合并

```

void MergePass(int R[], int n, int L){
    int i;
    for(i=1; i+2*L-1<=n; i+=2*L)
        Merge(R, i, i+L-1, i+2*L-1);
    //处理余留的长度小于2*L的子数组
    if(i+L-1<n)
        Merge(R, i, i+L-1, n);    //L<剩余部分长度<2L
}

```

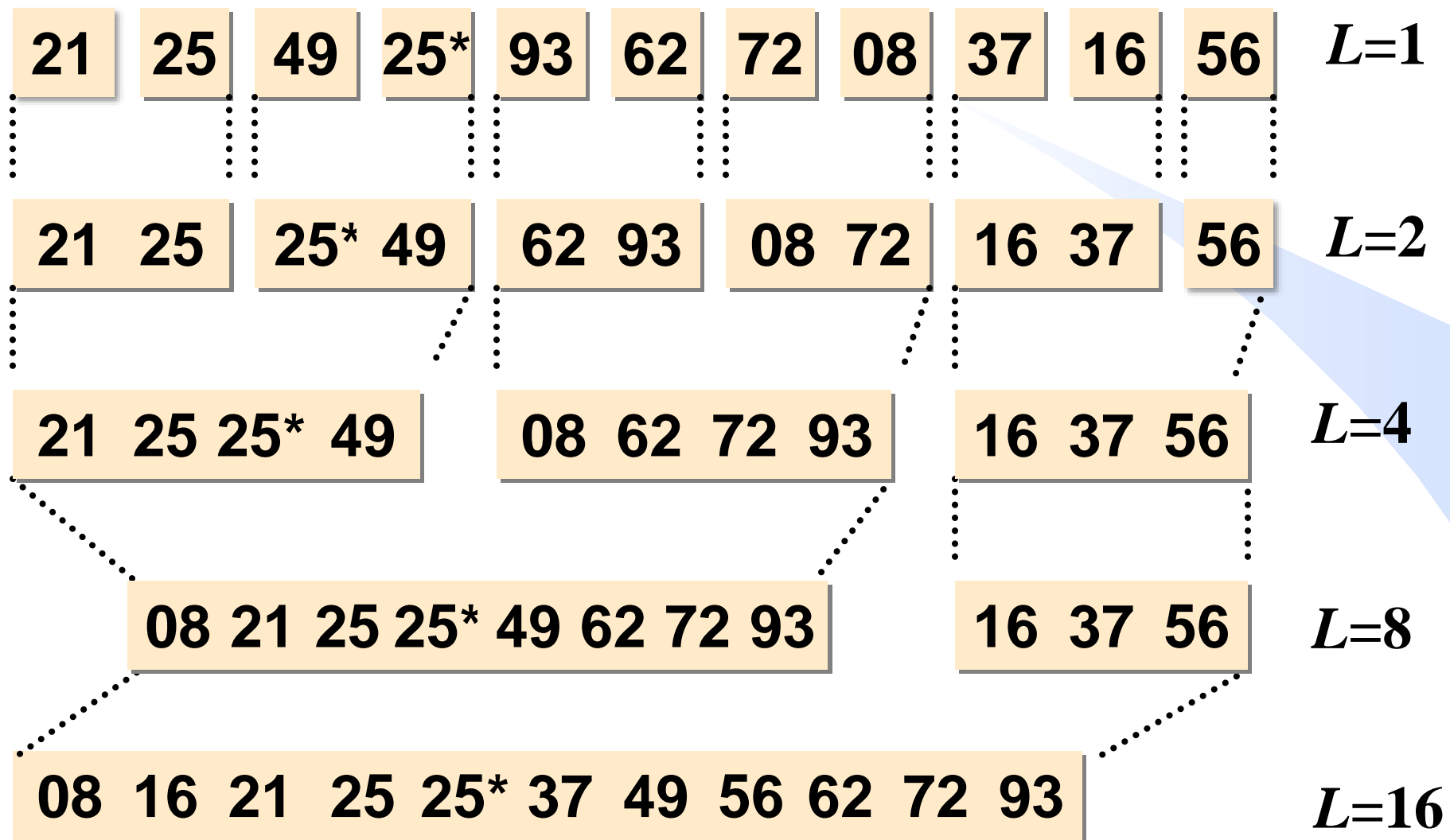


时间复杂度  $O(n)$

情况2  
剩余部分  $\leq L$

方案：剩余部分已有序，故不做处理

## 归并排序——示例





# 归并排序——非递归形式

```
void MergeSort(int R[], int n){  
    for(int L=1; L<n; L*=2)  
        MergePass(R, n, L);  
}
```

时间复杂度  
 $O(n\log n)$

## 归并排序的稳定性

```
void Merge(int R[],int low, int mid, int high){  
    //将两个相邻的有序数组( $R_{low}, \dots, R_{mid}$ )和( $R_{mid+1}, \dots, R_{high}$ )合并成一个有序数组  
    int i=low, j=mid+1, k=0;  
    int *X=new int[high-low+1];  
    while(i<=mid && j<=high)  
        if(R[i]<=R[j]) X[k++]=R[i++];  
        else X[k++]=R[j++];  
    while(i<=mid) X[k++]=R[i++];  
    while(j<=high) X[k++]=R[j++];  
    for(i=0; i<high-low+1; i++)  
        R[low+i]=X[i];  
    delete []X;  
}
```

Diagram illustrating the merge step of Merge Sort. Two input arrays  $R$  are shown:  $[3, 5, 6, 7]$  and  $[2, 6^*, 8]$ . The first array has an arrow labeled  $i$  pointing to the element 6. The second array has an arrow labeled  $j$  pointing to the element  $6^*$ . A blue shaded area represents the temporary array  $X$ , which contains the merged elements  $[2, 3, 5, 6]$ . An arrow labeled  $k$  points to the element 6 in array  $X$ .

✓若关键词相等先把 $i$ 指向的元素放入 $X$   
✓使两个关键词相等的元素合并后左面的元素还在左面

- ✓若关键词相等先把 $i$ 指向的元素放入 $X$
- ✓使两个关键词相等的元素合并后左面的元素还在左面

# 归并排序算法总结

排序算法	时间复杂度			空间复杂度	稳定性
	最好	平均	最坏		
归并排序	$O(n\log n)$	$O(n\log n)$	$O(n\log n)$	$O(n)$	稳定

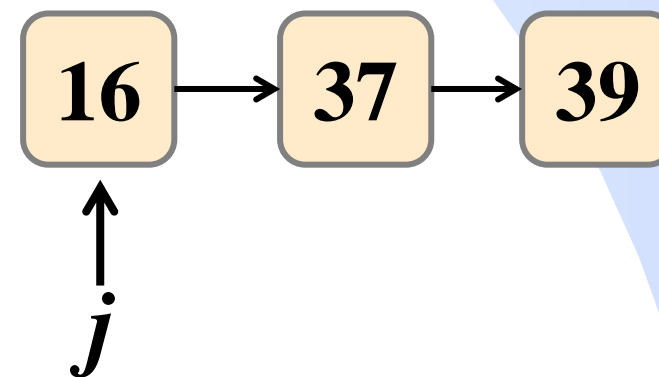
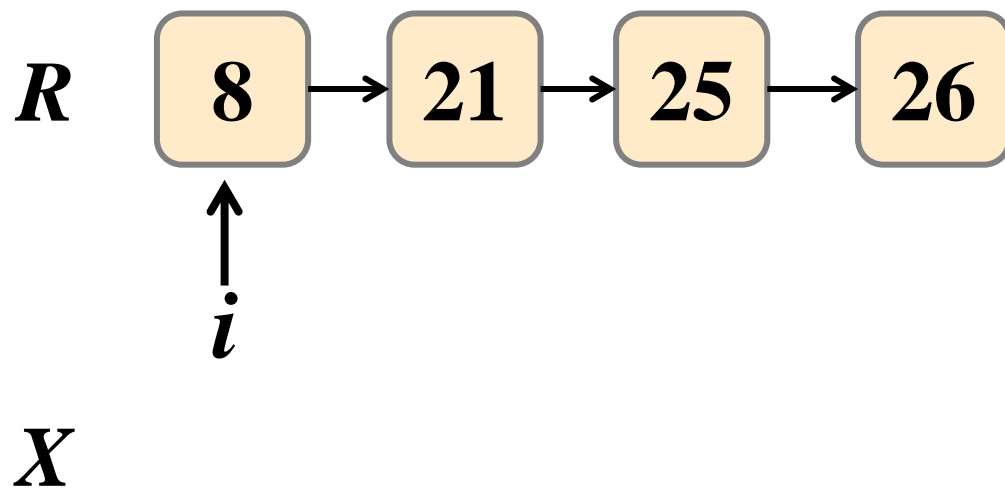
最快的稳定性排序算法

## 归并排序优化策略

- **问题**：当数据量非常小时，若仍然采用分治策略，效率不高。
- **优化策略**：对于非常小的数据集，以及前几次合并动作，调用直接插入排序算法。
- **问题**：Merge操作基于元素移动，当元素比较大时，赋值操作会比较费时。
- **优化策略**：将数组存储改为链表存储，这样记录移动就变为指针移动了。

## 课下思考

- 对单链表进行排序，哪种排序算法最适合？【腾讯、华为、阿里、字节跳动、百度、快手、美团、谷歌、微软、苹果面试题】
- 给定一个包含哨位结点的单链表，请设计一个时空效率尽可能高效的算法，对该链表进行递增排序， $n$ 为链表长度。  
【吉林大学21级期末考试题，15分】





## 归并排序的经典应用——求逆序对数目

求数组中逆序对的个数【华为、阿里、字节跳动、百度、B站、美团、谷歌、京东、滴滴面试题 [LeetCode LCR170](#)】

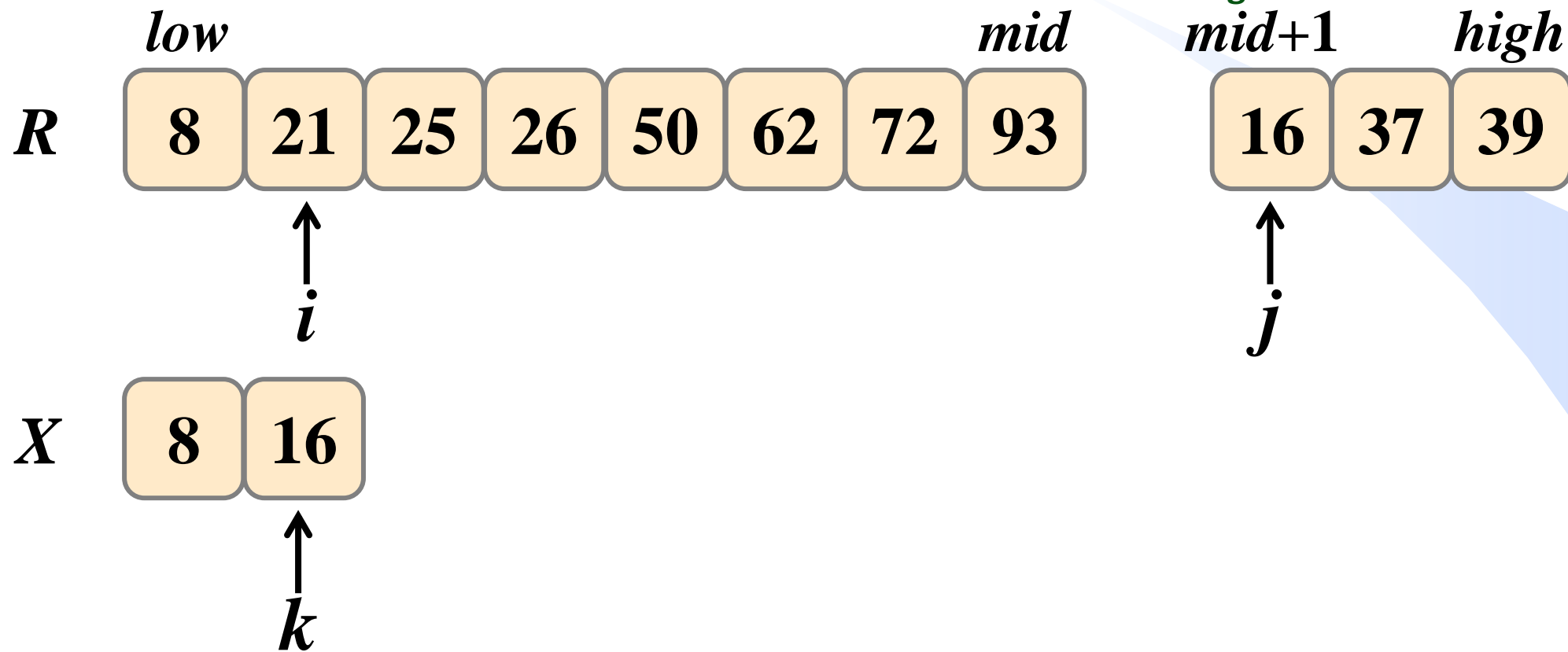


```
int cnt=0;
for(int i=1;i<=n;i++)
    for(int j=i+1;j<=n;j++)
        if(R[i]>R[j]) cnt++;
```

时间复杂度  
 $O(n^2)$

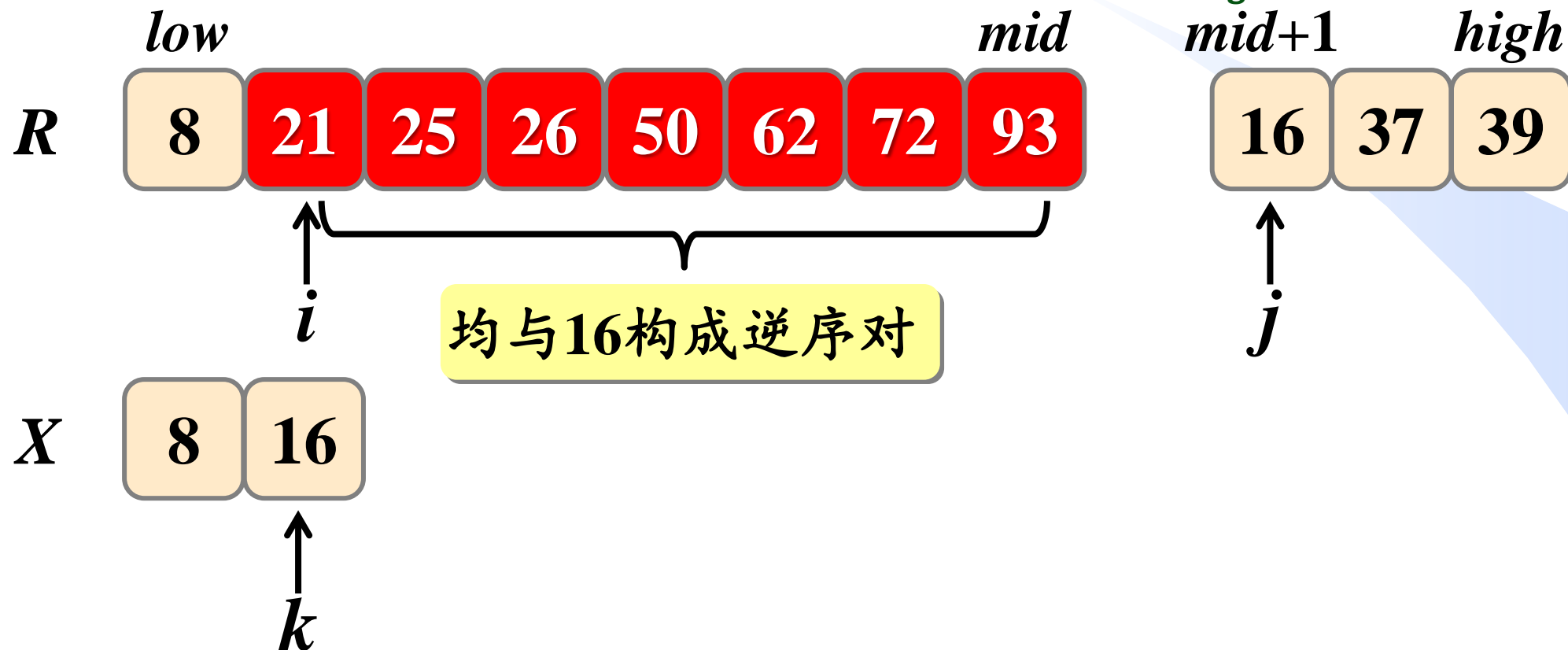
# 计算逆序对个数

```
void Merge(int R[], int low, int mid, int high){  
    //将两个相邻的有序数组( $R_{low}, \dots, R_{mid}$ )和( $R_{mid+1}, \dots, R_{high}$ )合并成一个有序数组
```



# 计算逆序对个数

`void Merge(int R[], int low, int mid, int high){`  
 //将两个相邻的有序数组( $R_{low}, \dots, R_{mid}$ )和( $R_{mid+1}, \dots, R_{high}$ )合并成一个有序数组



$R[i] > R[j] \Leftrightarrow R[i] \dots R[mid] > R[j] \Leftrightarrow$  逆序对增加  $mid - i + 1$  个

# 计算逆序对个数

```
int cnt=0;
void Merge(int R[],int low, int mid, int high){
    int i=low, j=mid+1, k=0;
    int *X=new int[high-low+1];
    while(i<=mid && j<=high)
        if(R[i]<=R[j]) X[k++]=R[i++];
        else { X[k++]=R[j++]; cnt += mid-i+1;}
    while(i<=mid) X[k++]=R[i++];
    while(j<=high) X[k++]=R[j++];
    for(i=0; i<high-low+1; i++)
        R[low+i]=X[i];
    delete []X;
}
```

归并排序时，把调用 Merge 时算出来的 *cnt* 值累加起来

时间复杂度  
 $O(n\log n)$

# 基于分治法 (*Divide and Conquer*) 的排序算法

A

- 分治法思想：将一个输入规模为 $n$ 的问题分解为多个规模较小的子问题，这些子问题互相独立且与原问题相同，然后递归地求解这些子问题，最后用适当的方法将各子问题的解合并，以获得原问题的解。
- 分治排序包括三个步骤：
  - ✓ “分”：将数组划分为若干子数组，一般为二分；
  - ✓ “治”：对子数组递归排序；
  - ✓ “合”，将子数组排序后的结果整合在一起。

# 基于分治法 (*Divide and Conquer*) 的排序算法

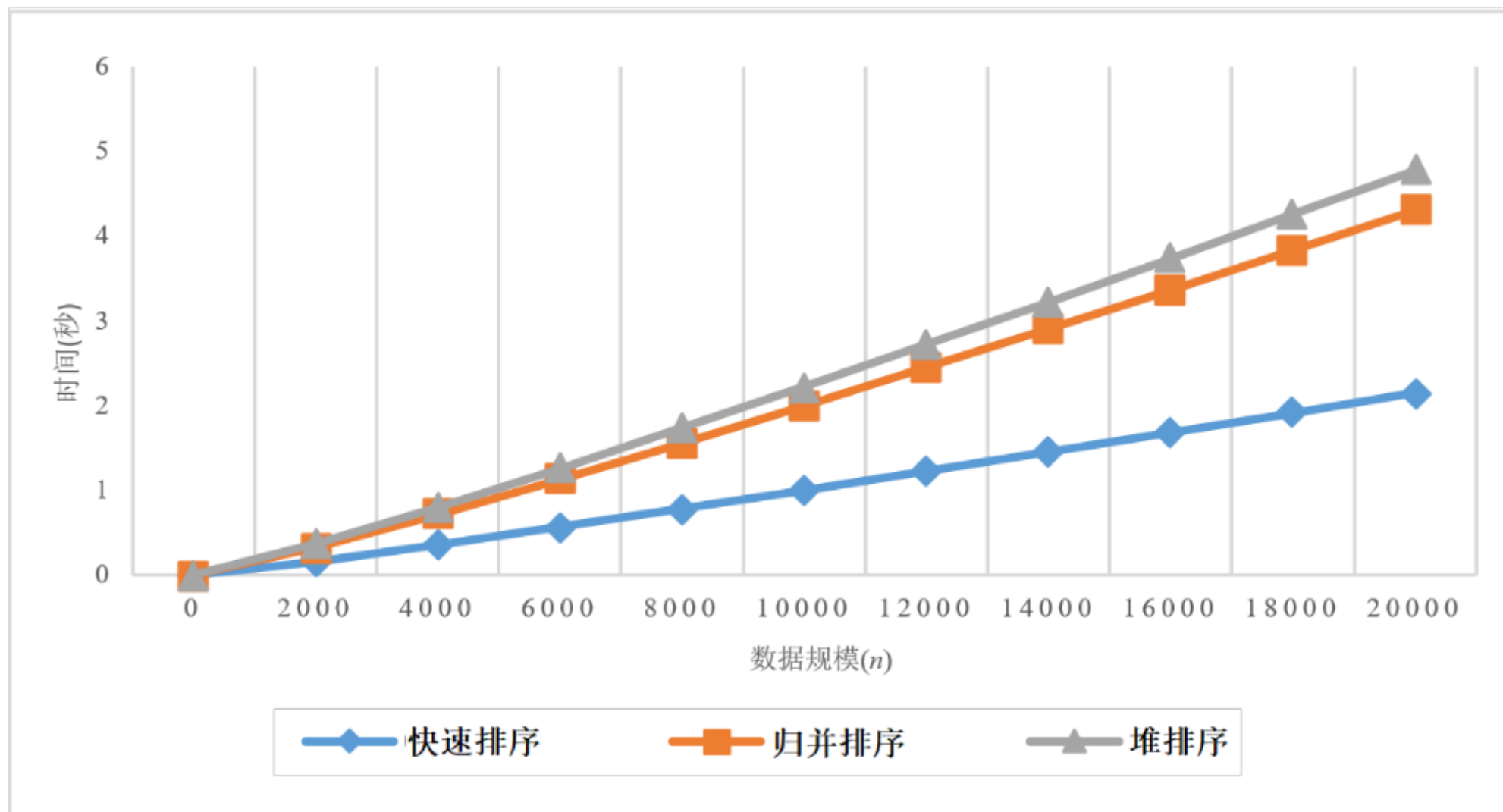
A

- 由于“治”是递归调用过程，分治排序的重点在“分”和“合”两步，具体的算法往往侧重其中的某一步：
  - ✓ 快速排序：侧重在“分”（调用Partition算法），在“分”的过程中调整了元素位置，而“合”过程无需任何操作。
  - ✓ 归并排序：侧重在“合”（调用Merge算法），分的过程就是简单的对半二分。

# 基于关键词比较的排序算法对比

排序算法	时间复杂度			空间复杂度	稳定性
	最好	平均	最坏		
直接插入排序	$O(n)$	$O(n^2)$	$O(n^2)$	$O(1)$	稳定
冒泡排序	$O(n)$	$O(n^2)$	$O(n^2)$	$O(1)$	稳定
直接选择排序	$O(n^2)$	$O(n^2)$	$O(n^2)$	$O(1)$	不稳定
堆排序	$O(n \log n)$	$O(n \log n)$	$O(n \log n)$	$O(1)$	不稳定
快速排序	$O(n \log n)$	$O(n \log n)$	$O(n^2)$	$O(\log n) - O(n)$	不稳定
归并排序	$O(n \log n)$	$O(n \log n)$	$O(n \log n)$	$O(n)$	稳定

# $O(n\log n)$ 排序算法实验对比





## 练习

现有 $n$ 条词以及对应的拼音串，对其排序，排序规则：首先按拼音串的字母序排序，如果拼音串相同，则按当前词所在的顺序排序，下列排序算法中\_\_\_\_\_符合条件。【搜狗校园招聘笔试题】

A. 直接插入排序

B. 快速排序

C. 堆排序

D. 冒泡排序

长春	<i>changchun</i>
北京	<i>beijing</i>
吉大	<i>jida</i>
数据	<i>shuju</i>
背景	<i>beijing</i>



北京	<i>beijing</i>
背景	<i>beijing</i>
长春	<i>changchun</i>
吉大	<i>jida</i>
数据	<i>shuju</i>

## 课下思考

下列排序算法中，不稳定的是\_\_\_\_\_。【2023年考研题全国卷】

I. 希尔排序      II. 合并排序      III. 快速排序

IV. 堆排序      V. 冒泡排序

A. 仅 I 和 II

B. 仅 II 和 V

C. 仅 I、III 和 IV

D. 仅 III、IV 和 V

## 课下思考

下列排序方法中，每趟排序结束都至少能确定一个元素最终位置的方法是\_\_\_\_\_。【考研题全国卷】

I. 直接选择排序

II. 希尔排序

III. 快速排序

IV. 堆排序

V. 合并排序

**A. I、III、IV**

B. I、III、V

C. II、III、IV

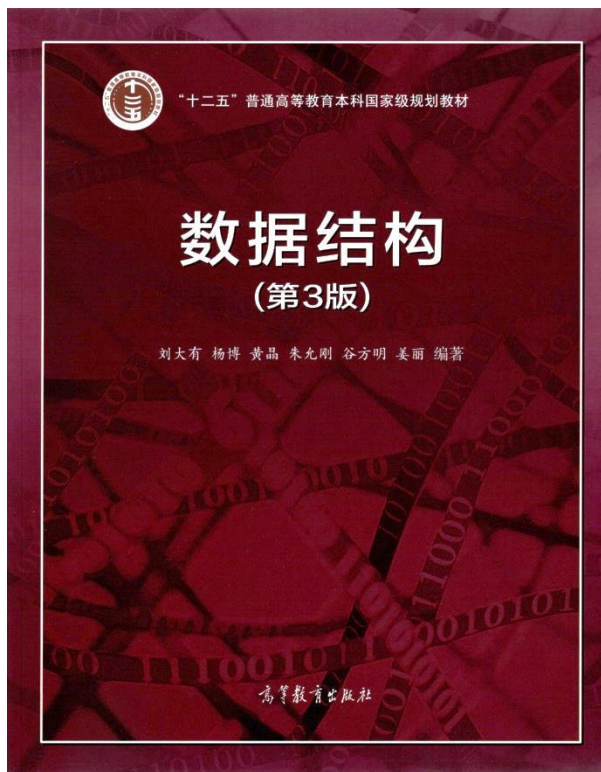
D. III、IV、V

## 课下思考

下列哪个算法可能出现下列情况：在最后一趟开始之前，所有的元素都不在其最终的位置上。

- A. 堆排序      B. 冒泡排序      C. 插入排序      D. 快速排序

2 3 4 5 6 ← 1



# 归并排序及其他

- 归并排序
- **排序算法时间下界**
- 外排序方法简介

数据之法  
结构之美  
算法之道

# 基于关键词比较的排序算法时间下界

**下界** 对于输入规模为 $n$ 的问题，若不存在解决该问题的算法，其时间复杂性小于 $T(n)$ ，则称针对该问题的算法的时间复杂性下界为 $T(n)$ 。

$K_1, K_2, K_3$  有6种可能排序结果。

$$K_1 < K_2 < K_3$$

$$K_3 < K_2 < K_1$$

$$K_1 < K_3 < K_2$$

$$K_3 < K_1 < K_2$$

$$K_2 < K_1 < K_3$$

$$K_2 < K_3 < K_1$$

# 基于关键词比较的排序算法时间下界

$K_1 : K_2$

$K_1 < K_2 < K_3$

$K_3 < K_2 < K_1$

$K_1 < K_3 < K_2$

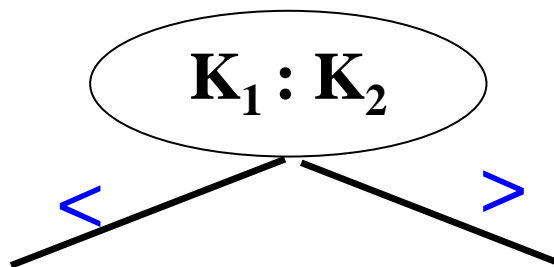
$K_3 < K_1 < K_2$

$K_2 < K_1 < K_3$

$K_2 < K_3 < K_1$

# 基于关键词比较的排序算法时间下界

**B**



$K_1 < K_2 < K_3$

$K_3 < K_2 < K_1$

$K_1 < K_3 < K_2$

$K_3 < K_1 < K_2$

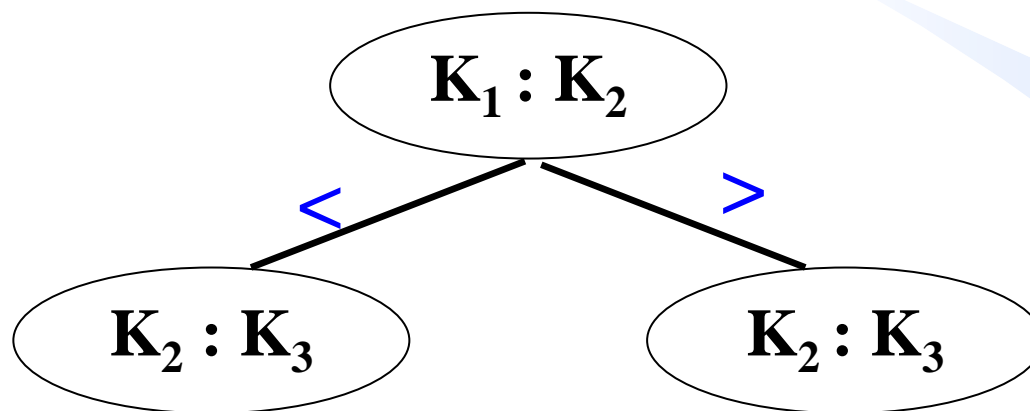
$K_2 < K_1 < K_3$

$K_2 < K_3 < K_1$



# 基于关键词比较的排序算法时间下界

**B**



$K_1 < K_2 < K_3$

$K_3 < K_2 < K_1$

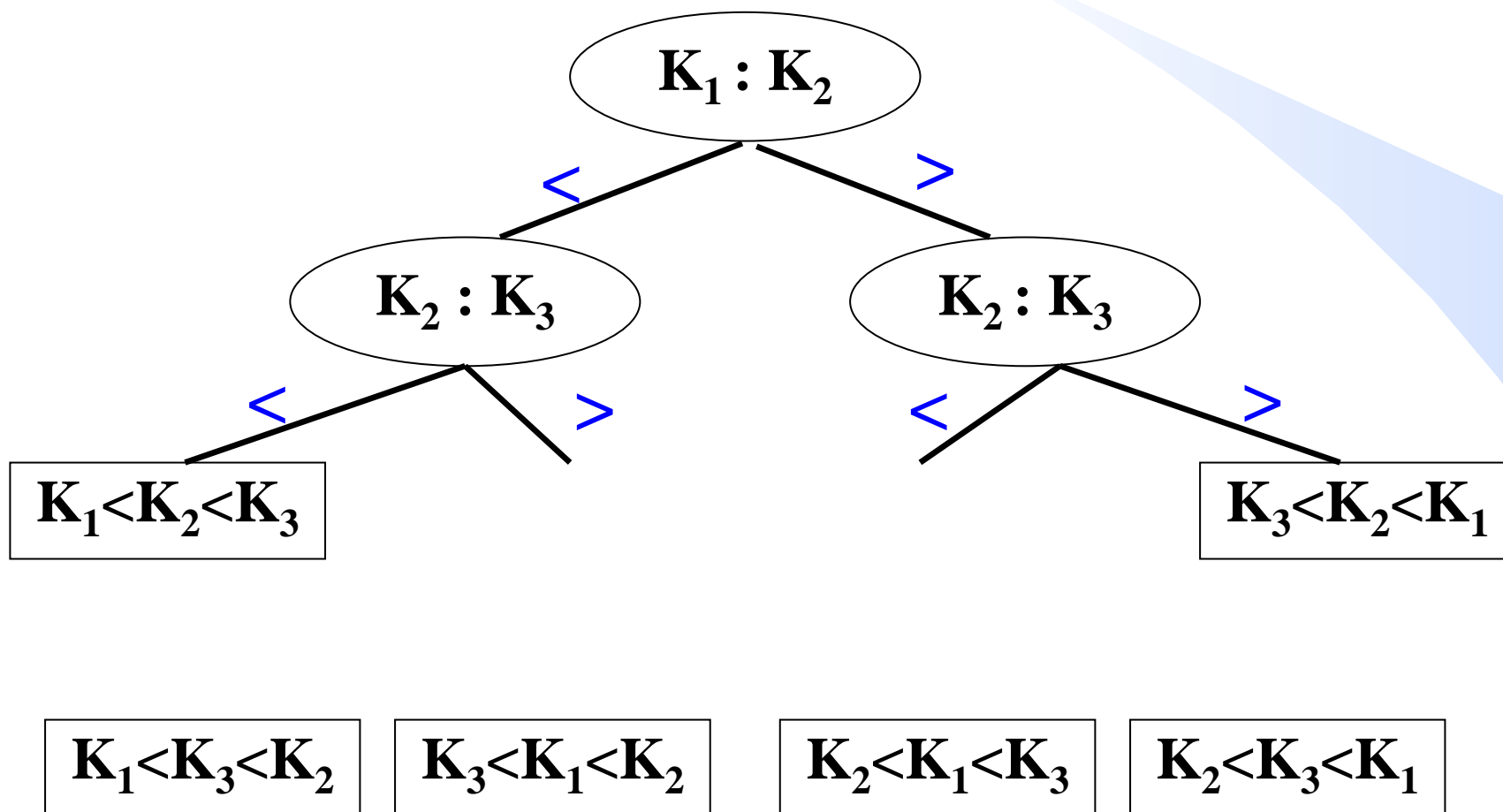
$K_1 < K_3 < K_2$

$K_3 < K_1 < K_2$

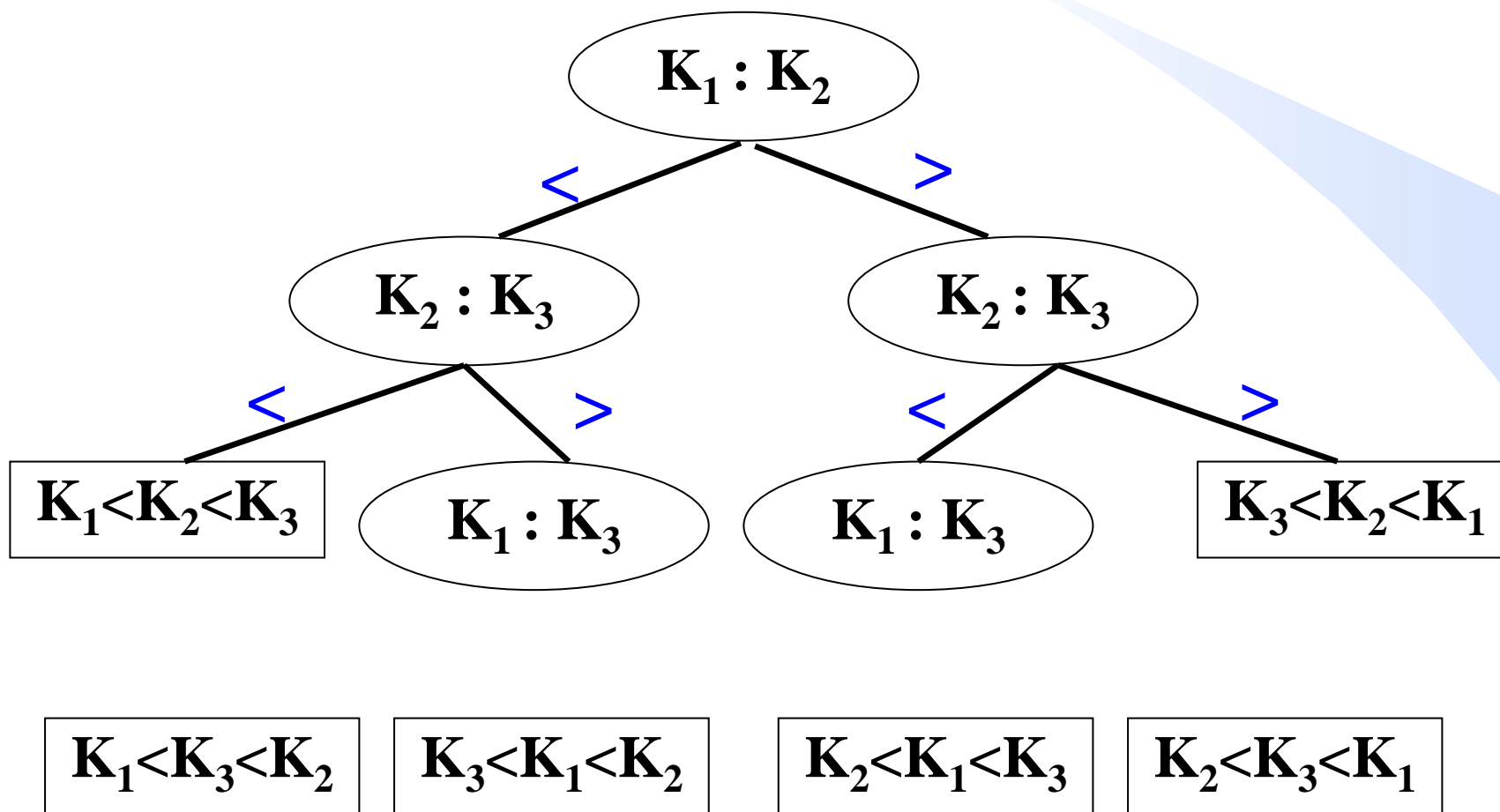
$K_2 < K_1 < K_3$

$K_2 < K_3 < K_1$

# 基于关键词比较的排序算法时间下界

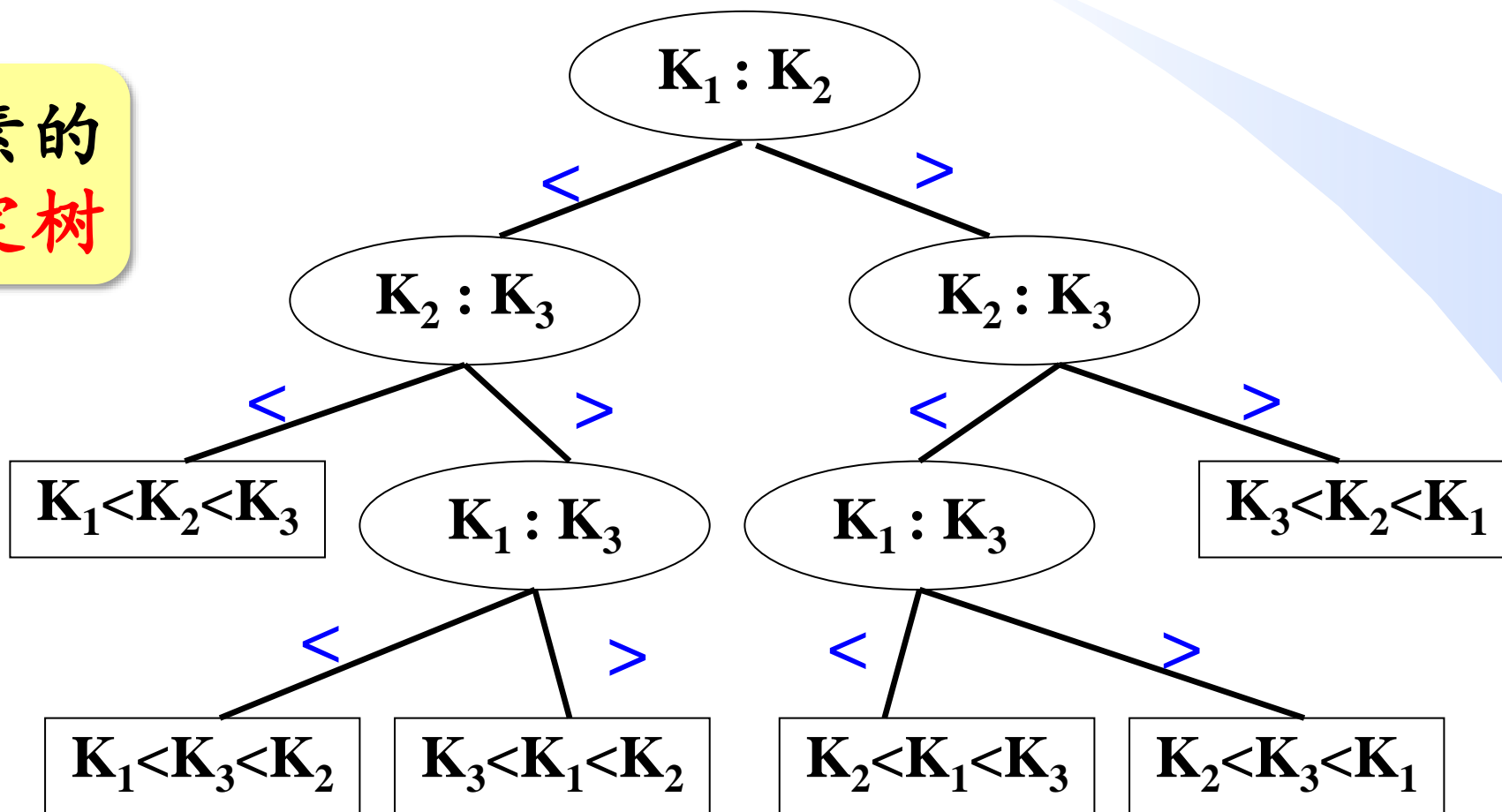


# 基于关键词比较的排序算法时间下界



# 基于关键词比较的排序算法时间下界

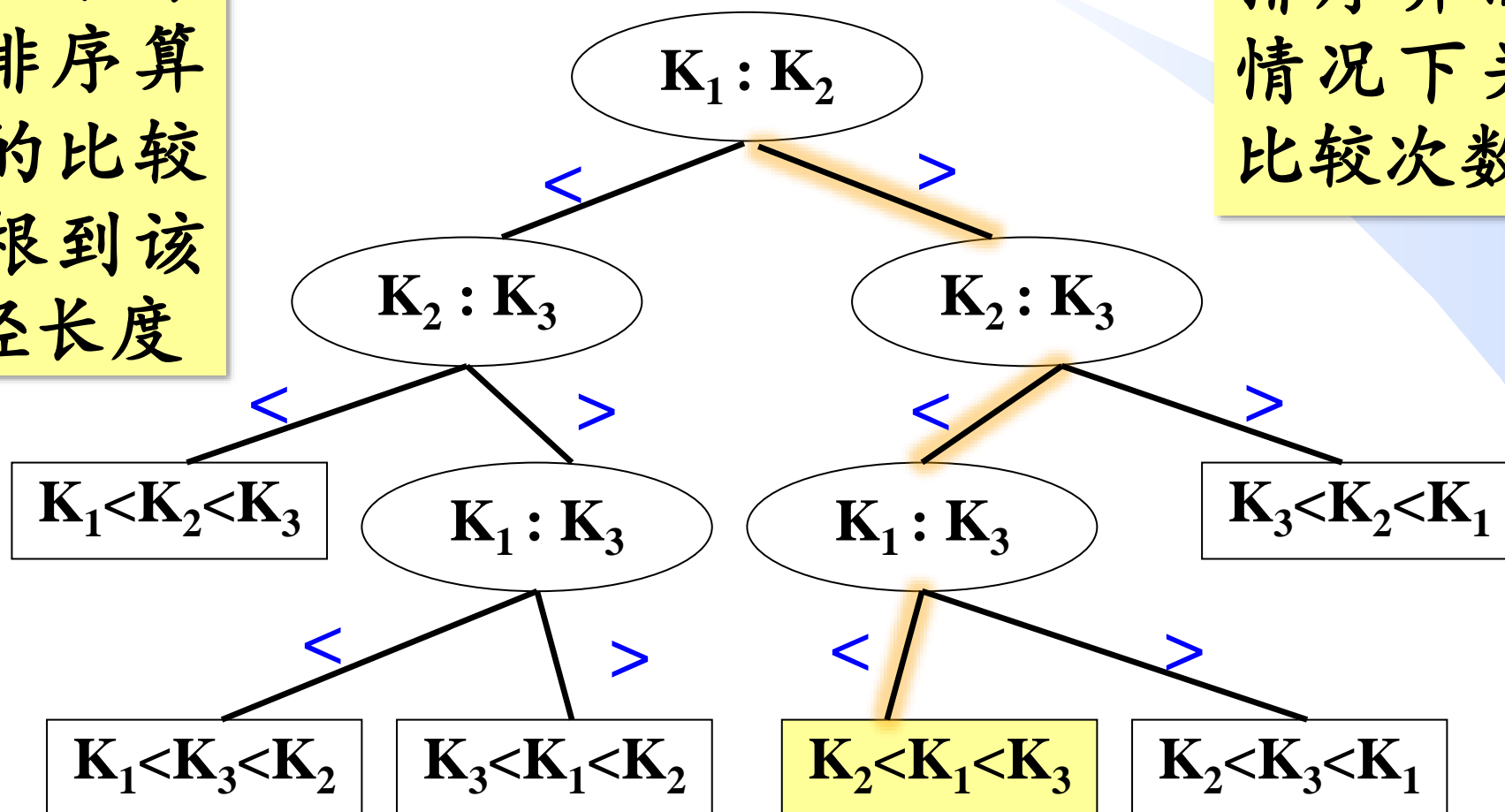
三个元素的  
排序判定树



# 基于关键词比较的排序算法时间下界

每种排列看成一种输入，对于每种输入，排序算法所执行的比较次数为从根到该排列的路径长度

判定树的高度是排序算法在最坏情况下关键词的比较次数。



# 基于关键词比较的排序算法时间下界

- $n$ 个元素的排列数为 $n!$ ，故 $n$ 个元素的判定树有 $n!$ 个叶结点；
- 假定判定树高度为 $h$ ，因高度为 $h$ 的二叉树最多有 $2^h$ 个叶结点，即叶结点个数小于等于 $2^h$ ，又当前实际的叶结点个数为 $n!$ ，故有： $n! \leq 2^h$

即  $h \geq \log_2(n!)$ ,

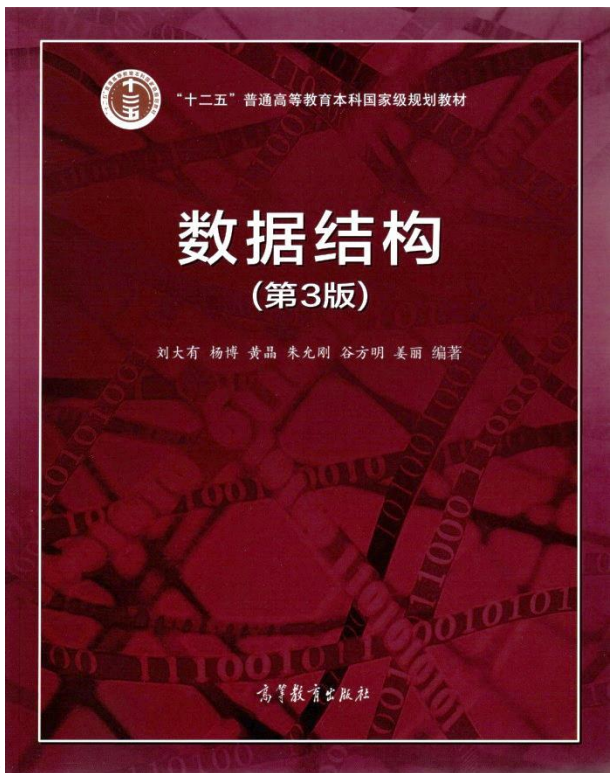
$$\log(n!) = \log n + \log(n-1) + \log(n-2) + \cdots + \log 2 + \log 1$$

$$\geq \log n + \log(n-1) + \log(n-2) + \cdots + \log(n/2)$$

$$\geq \frac{n}{2} \log \frac{n}{2}$$

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

- ✓ 基于比较的排序算法的最坏时间复杂度下界是 $\Omega(n \log n)$ 。
- ✓ 亦可证明基于比较的排序算法的平均时间复杂度下界是 $\Omega(n \log n)$ 。



# 归并排序及其他

- 归并排序
- 排序算法时间下界
- 外排序方法简介

数据之法  
结构之美  
算法之道

# 外排序原理简介

32G

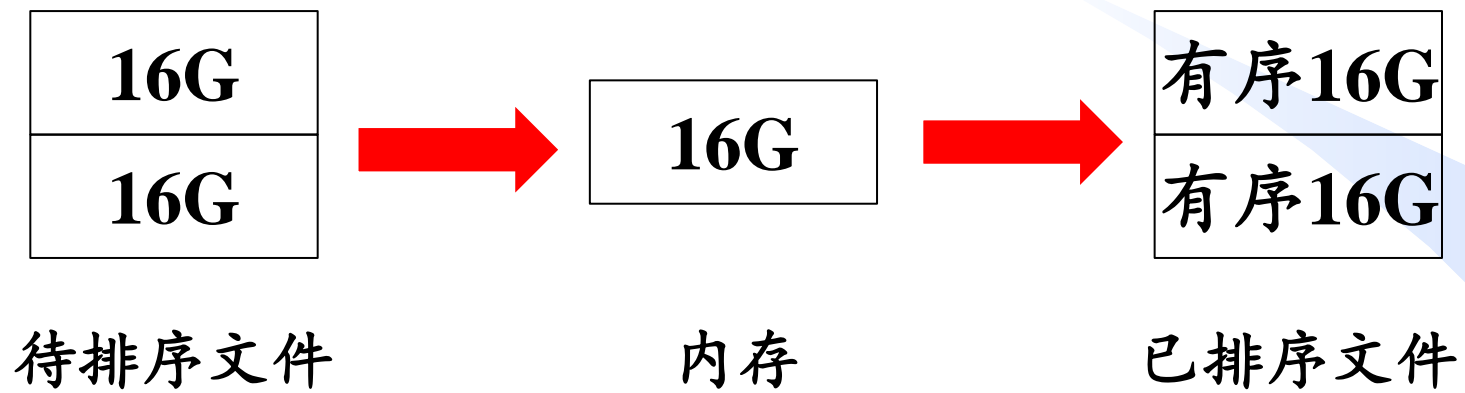
待排序文件

16G

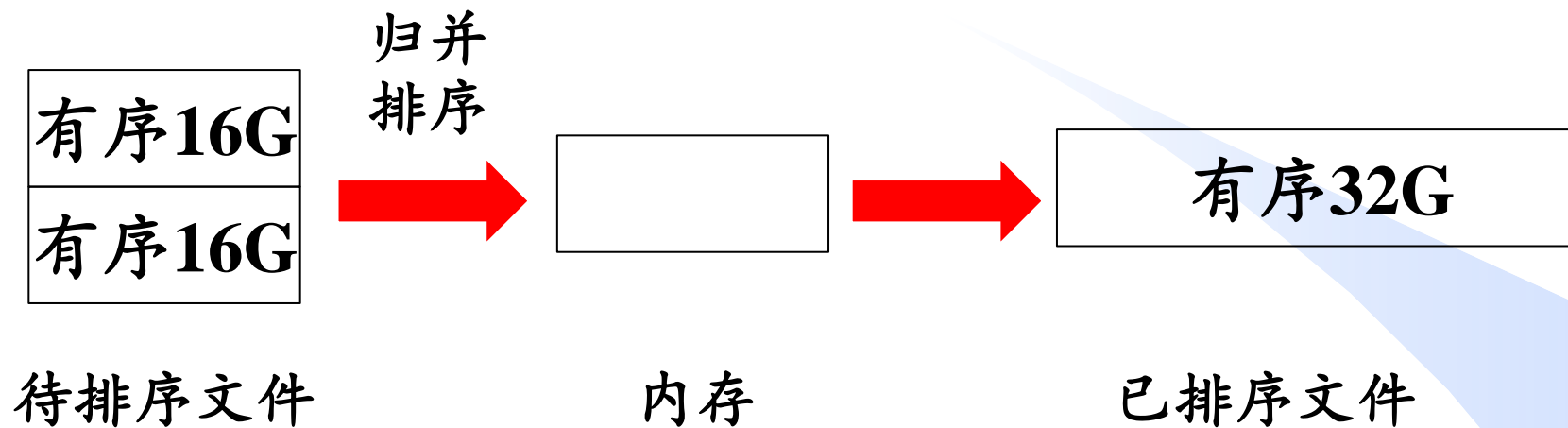
内存



# 外排序原理简介



# 外排序原理简介



## 练习

对10TB的数据文件进行排序，应使用的方法是\_\_\_\_\_。【考研题全国卷】

A.希尔排序

B.堆排序

C.快速排序

D.归并排序