

# 外部排序

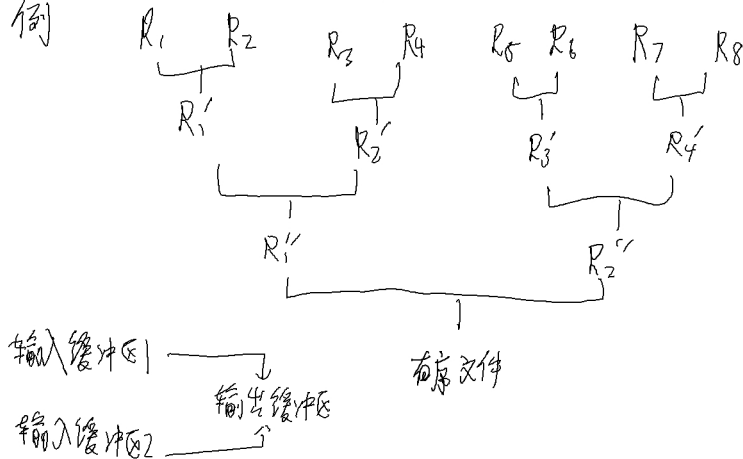
## 概念

待排记录过多，存储在外存中，排序时将数据逐部分调入内存排序，过程中需要多次内存和外存的交换

外部排序主要时间在访问磁盘 (I/O 次数)

- 归并排序法
- ① 根据缓冲区大小，划分文件成长度  $L$  的子文件，依次读入内存排序，排好后重新写回外存，排好后的有序子文件称为归并段
  - ② 对归并段逐趟归并，使归并段由小到大直至整个文件有序

例



外部排序总时间 = 内部排序时间 + 外存读写时间 + 内部归并时间

增大归并路可以减少归并趟数, 进而减少 I/O 次数

$r$  个初始归并段, 做  $k$  路归并, 可以用严格  $k$  叉树表示

归并趟数  $S = \text{树高} - 1 = \lceil \log_k r \rceil$  (2有度为  $k$  和度为 0)

多路平衡归并和败者树

$S$  趟归并所需的比较次数为

$$S(n-1)(k-1) = \lceil \log_k r \rceil (n-1)(k-1) = \lceil \log_2 r \rceil (n-1) \frac{(k-1)}{\lceil \log_2 k \rceil}$$

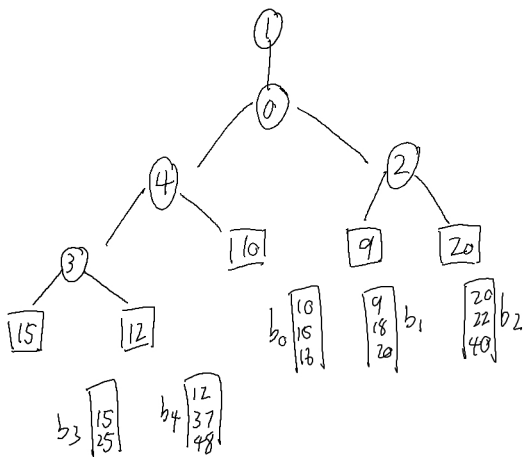
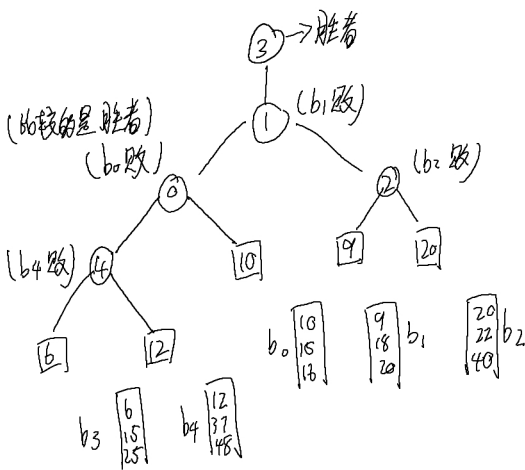
内部归并排序效率受  $k$  的影响, 因此引入败者树

败者树  $\Rightarrow$  一棵完全二叉树,  $k$  个叶子存放  $k$  个归并段  
当前比较的记录, 内部结点记录左右子  
树的胜败者, 让胜者向上继续比较, 直至根。  
若大的失败, 小的胜利, 则根指向最小数

$k$  路归并败者树深度  $\lceil \log_2 k \rceil$

因此  $S$  趟归并比较次数  $S(n-1) \lceil \log_2 k \rceil = \lceil \log_2 r \rceil (n-1)$

归并路数  $k$  不是越大越好,  $k$  增大, 缓冲区个数也增多



5路合并败者树

# 置换-选择排序 (生成初始归并段)

归并段长度  $l$ , 个数  $r = \lceil n/l \rceil$

置换-选择算法可以产生更长的初始归并段, 以减少个数  $r$

- ① 从输入文件  $FI$  中输出  $w$  个记录到工作区  $WA$  中
  - ② 从  $WA$  中选出最小的记录, 记为  $MINMAX$
  - ③  $MINMAX$  输出到输出文件  $FO$  中
  - ④ 若  $FI$  不空,  $FI$  输出下一个记录到  $WA$  中
  - ⑤  $WA$  中找出大于  $MINMAX$  的最小值作为新  $MINMAX$
  - ⑥ 重复 ③-⑤ 直至选不出新的  $MINMAX$ , 输出一个归并段结束符
- 使用败者树从  $WA$  中选  $MINMAX$

$FO$

$WA (w=3)$

$FI$

17, 21, 05, 44, 10, 12, 56, 32, 29

17, 21, 05

44, 10, 12, 56, 32, 29

05

17, 21, 44

10, 12, 56, 32, 29

05, 17

10, 21, 44

12, 56, 32, 29

05, 17, 21

10, 12, 44

56, 32, 29

05, 17, 21, 44

10, 12, 56

32, 29

05, 17, 21, 44, 56

10, 12, 32

29

10

29, 12, 32

10, 12

29, 32

10, 12, 29

32

10, 12, 29, 32

## 最佳归并树

置换-选择排序 得到的归并段长度不等

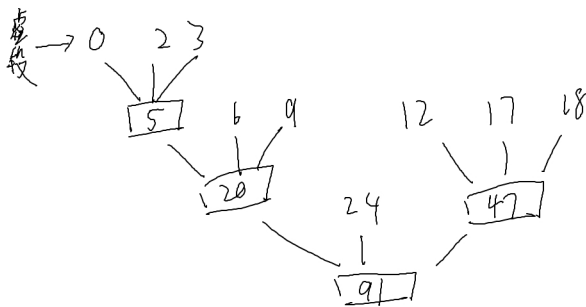
归并段长度不等的归并段 使 I/O 次数最少, 建立最佳归并树

$I/O \text{ 次数} = 2 \times WPL$  (带权路径长度) (权值为归并段长度)

最佳归并树  $\rightarrow$  哈夫曼树推广到  $m$  叉情况

若初始段个数不足以构成严格  $m$  叉树, 则需要补充长度为 0 的虚段

$\{2, 3, 6, 9, 24, 12, 17, 18\}$  不足以构成严格 3 叉树



如何判定要加入几个虚段

设度为 0 的结点  $n_0$  (初始  $= n$ ) 个, 度为  $k$  的结点  $n_k$  个

$(n_0 - 1) \% (k - 1) = 0$  不用加虚段

$(n_0 - 1) \% (k - 1) = u \neq 0$  加  $k - u - 1$  个虚段

上述例子  $(8 - 1) \% (3 - 1) = 7 \% 2 = 1$

加  $k - u - 1 = 3 - 1 - 1 = 1$  个虚段