

哈工大计算机考研全套视频和资料，真题、考点、典型题、命题规律独家视频讲解！
详见：网学天地（www.e-studysky.com）；咨询QQ：2696670126

数据结构与算法 第8章 外部分类 Slide. 8-1

第8章 外部分类

8.1 磁盘文件的归并分类

8.2 磁带文件的归并分类

■归并外部分类重点研究的问题是

- (1) 如何进行多路归并以减少文件的归并遍数；
- (2) 如何巧妙地运用内存的缓冲区使I/O和CPU尽可能并行工作；
- (3) 根据外存的特点选择较好的产生初始归并段的方法。

■磁盘和磁带归并分类

■选择树算法

磁盘是随机存储设备；磁带是顺序存储设备

HIT CST 哈尔滨工业大学 计算机科学与技术学院 张岩

数据结构与算法 第8章 外部分类 Slide. 8-2

外部分类的概念：

是指在分类的过程中，数据的主要部分存放在外存储器上，借助内存（作为工作单元），来调整外存储器上数据的位置。

归并方法的一般过程：分两个阶段。

第一阶段：首先，将文件中的数据分段输入到内存，在内存中采用内部分类方法对其进行分类（分类完的文件段，称为归并段run），然后将有序段写回外存。

整个文件经过在内存逐段分类又逐段写回外存，这样在外存中形成多个初始的归并段。

第二阶段：对这些初始归并段采用某种归并分类方法，进行多遍归并，最后形成整个文件的单一归并段（整个文件有

HIT CST 哈尔滨工业大学 计算机科学与技术学院 张岩

数据结构与算法 第8章 外部分类 Slide. 8-3

8.1 磁盘文件的归并分类

□ 示例：设有一个包含4500个记录的输入文件。现用一台其内存至多可容纳750个记录的计算机对该文件进行分类。输入文件放在磁盘上，磁盘每个页块可容纳250个记录，这样全部记录可存储在 $4500 / 250 = 18$ 个页块中。输出文件也放在磁盘上，用以存放归并结果。

□ 由于内存中可用于分类的存储区域能容纳750个记录，因此内存中恰好能存3个页块的记录。

□ 在外分类一开始，把18块记录，每3块一组，读入内存。利用某种内分类方法进行内分类，形成初始归并段，再写回外存。总共可得到6个初始归并段。然后一趟一趟进行归并分类。

初始归并段

R1	R2	R3	R4	R5	R6
1~750	751~1500	1501~2250	2251~3000	3001~3750	3751~4500

HIT CST 哈尔滨工业大学 计算机科学与技术学院 张岩

数据结构与算法 第8章 外部分类 Slide. 8-4

8.1 磁盘文件的归并分类

□ 若把内存区域等份地分为3个缓冲区。其中的两个为输入缓冲区，一个为输出缓冲区，可以在内存中利用简单2路归并函数 merge() 实现2路归并。

□ 首先，从参加归并排序的两个输入归并段 R_1 和 R_2 中分别读入一块，放在输入缓冲区1和输入缓冲区2中。然后在内存中进行2路归并，归并结果顺序存放到输出缓冲区中。

□ 当输出缓冲区装满250个记录时，就输出到磁盘。

□ 如果归并期间某个输入缓冲区空了，就立即向该缓冲区继续装入所对应归并段的一块记录信息，使之于另一个输入缓冲区的剩余记录归并，直到 R_1 和 R_2 归并为 R_{12} 、 R_3 和 R_4 归并为 R_{34} 、 R_5 和 R_6 归并为 R_{56} 为止。

□ 再把 R_{12} 和 R_{34} 归并为 R_{1234} ，最后把 R_{1234} 和 R_{56} 归并为 R_{123456} （如下页图示）

HIT CST 哈尔滨工业大学 计算机科学与技术学院 张岩

哈工大计算机考研全套视频和资料，真题、考点、典型题、命题规律独家视频讲解！
 详见：网学天地（www.e-studysky.com）；咨询QQ：2696670126

数据结构与算法 第8章 外部分类 Slide. 8-5

8.1 磁盘文件的归并分类

输入缓冲区 1
输入缓冲区 2
输出缓冲区

初始归并段
第一趟归并结果
第二趟归并结果
第三趟归并结果

HIT CST 哈尔滨工业大学 计算机科学与技术学院 张岩

数据结构与算法 第8章 外部分类 Slide. 8-6

8.1 磁盘文件的归并分类

2路归并
层数 $\lceil \log_2 m \rceil + 1$
遍数 $\lceil \log_2 m \rceil$
m个归并段的归并过程

讨论问题
 (1) 多路归并——减少归并遍数
 (2) 并行操作的缓冲区处理——使输入、输出和CPU处理尽可能重叠
 (3) 初始归并段的生成

HIT CST 哈尔滨工业大学 计算机科学与技术学院 张岩

数据结构与算法 第8章 外部分类 Slide. 8-7

8.1.1 多路归并——减少归并遍数

m个初始段进行2路归并，需要 $\lceil \log_2 m \rceil$ 遍归并；
 一般地，m个初始段，采用K路归并，需要 $\lceil \log_K m \rceil$ 遍归并。
 显然，K越大，归并遍数越少，可提高归并的效率。

在K路归并时，从K个关键字中选择最小记录时，要比较K-1次。若记录总数为n，每遍要比较 $n \cdot (K-1)$ 次， $\lceil \log_K m \rceil$ 遍要比较的次数为：

$$n \cdot (K-1) \lceil \log_K m \rceil = n \cdot (K-1) \lceil \log_2 m / \log_2 K \rceil$$

$$= n \cdot (K-1) \lceil \log_2 \lceil n / K \rceil / \log_2 K \rceil$$

可以看出，随着K增大， $(K-1)/\log_2 K$ 也增大，当归并路数多时，CPU处理的时间也随之增多。当K值增大到一定程度时，可能使CPU处理时间大于因K值增大而减少归并遍数所节省的时间。

为此可以 (1) 选择好的分类方法，以减少分类中比较次数；
 (2) 选择好的初始归并段形成方法，增大归并段长度；
 提高分类的效率。

HIT CST 哈尔滨工业大学 计算机科学与技术学院 张岩

数据结构与算法 第8章 外部分类 Slide. 8-8

K路平衡归并与败者树

选择树 (Selection tree) 或败者树 (tree of loser)

分析：
 第一次建立选择树的比较所花时间为：
 $O(K-1) = O(K)$
 而后每次重新建造选择树所需时间为：
 $O(\log_2 K)$
 n个记录处理时间为初始建立选择树的时间加上n-1次重建选择树的时间：
 $O((n-1) \cdot \log_2 K) + O(K) = O(n \cdot \log_2 K)$
 这就是K路归并一遍所需的CPU处理时间。归并遍数为 $\log_K m$ ，总时间为：
 $O(n \cdot \log_2 K \cdot \log_K m) = O(n \cdot \log_2 m)$
 (K路归并 CPU 时间与 K 无关——选择树太好了)

败者树选最小对象举例

HIT CST 哈尔滨工业大学 计算机科学与技术学院 张岩

哈工大计算机考研全套视频和资料，真题、考点、典型题、命题规律独家视频讲解！

详见：网学天地（www.e-studysky.com）；咨询QQ：2696670126

数据结构与算法 第8章 外部分类 Slide: 8-9

8.1.2 并行操作的缓冲区处理——使输入、输出和 CPU 处理尽可能重叠

对 k 个归并段进行 k 路归并至少需要 k 个输入和 1 个输出缓冲区。要使输入、输出和归并同时进行， $k+1$ 个缓冲区是不够的，需要 $2(k+1)$ 个缓冲区实现并行操作。

(a)

-	-
ou(1)	ou(2)
1	2
iu(1)	iu(2)

 (b)

1	-
2	-
-	-
-	-

 (c)

-	3
-	4
-	-
-	-

归并到 ou(1) 输出到 in(3) 归并到 ou(2) 输出到 in(4)

(d)

5	-
6	-
-	-
-	-

 (e)

-	7
-	8
-	-
-	-

 (f)

9	-
-	-
-	-
-	-

输出到 ou(2) 归并到 ou(1) 输出到 in(1) 输出到 ou(1) 归并到 ou(2) 输出到 in(2) 输出到 ou(2) 归并到 ou(1) 输出到 in(3)

HIT CST 哈尔滨工业大学 计算机科学与技术学院 张岩

数据结构与算法 第8章 外部分类 Slide: 8-10

8.1.3 初始归并段的生成

- 任何内部分类算法都可作为生成初始归并段的算法
- 初始归并段的长度 \geq 缓冲区的长度？！
- 选择树法

假设初始待分类文件为输入文件 **FI**，初始归并段文件为输出文件 **FO**。

内存缓冲区为 **W**，可容纳 **P** 个记录。**FO**、**W** 初始为空，则置换-选择如下：

- 从 **FI** 输入 **P** 个记录到缓冲区 **W**；
- 从 **W** 中选择出关键字最小的记录 **MIN**；
- 将 **MIN** 记录输出到 **FO** 中去；
- 若 **FI** 不空，则从 **FI** 输入下一个记录到 **W**；
- 从 **W** 中所有关键字比 **MIN** 关键字大的记录中选出最小关键字记录，作为新的 **MIN**；
- 重复 (3)~(5)，直到在 **W** 中选不出新的 **MIN** 为止。得到一个初始归并段，输出归并段结束标志到 **FO** 中
- 重复 (2)~(6)，直到 **W** 为空，由此得到全部初始归并段。

HIT CST 哈尔滨工业大学 计算机科学与技术学院 张岩

数据结构与算法 第8章 外部分类 Slide: 8-11

8.1.3 初始归并段的生成

例如：缓冲区的长度 $P=4$ ，输入序列为：

15 19 04 83 12 27 11 25 16 34 26 07 10 90 06 ...

注意：如果新输入记录的关键字小于最后输出记录的关键字，则新输入记录不能成为当前归并短的一部分；他要等待生成下一个归并段时供选择。

步	1	2	3	4	5	6	7	8	9	10	11	12	13	...
缓冲区内容	15	15	15	(11)	(11)	(11)	(11)	(11)	(11)	11	11	(06)
输出结果	04	12	27	27	27	27	34	(26)	(26)	26	26	26
归并段	83	83	83	83	83	83	83	83	(07)	10	90	90

采用选择树法生成初始归并段的平均长度是缓冲区长度的两倍。

HIT CST 哈尔滨工业大学 计算机科学与技术学院 张岩

数据结构与算法 第8章 外部分类 Slide: 8-12

8.2 磁带文件的归并分类

与磁盘不同，磁带是顺序存储设备，读取信息块的时间与信息块的位置有关。研究磁带分类，需要了解信息块的分布。

K路平衡归并分类

磁带机数量： $2K$

输入： T_1, T_2, \dots, T_k 输出： $T_{k+1}, T_{k+2}, \dots, T_{2k}$

磁带机	T_1	T_2	...	T_k
归并段	R_1	R_2	...	R_k
	R_{k+1}	R_{k+2}	...	R_{2k}

	R_{mk+1}

$T_1: R_1(1000), R_3(1000), R_5(1000)$
 $T_2: R_2(1000), R_4(1000), R_6(1000)$
 $T_3: \emptyset$
 $T_4: \emptyset$
 $T_1: \emptyset$
 $T_2: \emptyset$
 $T_3: R_7(2000), R_9(2000)$
 $T_4: R_8(2000)$
 $T_1: R_1(4000)$
 $T_2: R_2(2000)$
 $T_3: \emptyset$
 $T_4: \emptyset$
 $T_1: \emptyset$
 $T_2: \emptyset$
 $T_3: R_1(6000)$
 $T_4: \emptyset$

HIT CST 哈尔滨工业大学 计算机科学与技术学院 张岩

哈工大计算机考研全套视频和资料，真题、考点、典型题、命题规律独家视频讲解！
 详见：网学天地（www.e-studysky.com）；咨询QQ：2696670126

数据结构与算法 第8章 外部分类 Slide. 8-13

多阶段归并分类
 K+1台磁带机，实现 k 路归并

i 遍后	t ₁	t ₂	t ₃
开始	13(1L)	21(1L)	空
1	空	8(1L)	13(2L)
2	8(3L)	空	5(2L)
3	3(3L)	5(5L)	空
4	空	2(5L)	3(8L)
5	2(13L)	空	1(8L)
6	1(13L)	1(21L)	空
7	空	空	1(34L)

步	t ₁	t ₂	t ₃	总段数
n	0	0	1	1
n-1	1	1	0	2
n-2	2	0	1	3
n-3	0	2	3	5
n-4	3	5	0	8
n-5	8	0	5	13
n-6	0	8	13	21
n-7	13	21	0	34

HIT CST 哈尔滨工业大学 计算机科学与技术学院 张岩

数据结构与算法 第8章 外部分类 Slide. 8-14

小结

- 内部分类过程中不设计数据德内、外存交换，待分类德记录全部存放在内存中；
- 若待非类的文件很大，就无法将整个文件的所有记录同时调入内存进行分类；
- 外部分类的实现，主要是依靠数据德内、外存交换和内部归并；
- 外部分类基本上包括相对独立的两个阶段：初始归并段的形成；多路归并。

外部分类主要研究的技术问题是：

- (1) 如何进行多路归并以减少文件的归并遍数；
- (2) 如何巧妙地运用内存的缓冲区使I/O和CPU尽可能并行工作；
- (3) 根据外存的特点选择较好的产生初始归并段的方法。

HIT CST 哈尔滨工业大学 计算机科学与技术学院 张岩

