

空间存储和索引



- ●目标和基本思想
- ●物理存储介质
- ●缓冲区管理
- ●存储组织
- ●存取路径:索引结构

目标和基本思想







●概念和逻辑层次

- O空间概念模型
- O空间查询语言

●物理存储与数据结构

- O物理存储器分级管理
- O存储组织
- O聚类和空间聚类
- O索引和空间索引

目标和基本思想





●存储器分级管理

- O解决速度快、容量小、易失和速度慢、容量大、 持久之间的矛盾
- O考虑I/O代价,采用缓存机制

●存储组织

- O适合外存操作的数据结构
- O记录和记录的多种组织形式

目标和基本思想







●聚类和空间聚类

- O降低大查询的磁盘寻道时间和等待时间
- O空间上相邻的对象在外存中也相邻

●索引和空间索引

- O提供多种存取路径,提高查询效率
- O高效的排序树结构及其空间扩展

数据库存储管理





●如何存储和管理海量的数据?

O目标数据、元数据、索引、日志等

●采用何种表示方式和数据结构对数据处理 提供最佳、最有效的支持?

物理存储介质:基本存储



- 寄存器 (register)
 - OCPU的一部分,用于暂存运算中间结果
 - 〇与运算部件直接连接,速度最快,极少(几十个)
- 高速缓冲存储器 (cache memory)
 - OCPU的一部分,用于缓存主存储器
 - O在CPU中,速度极快,容量小(几十K~几M)
 - O操作系统底层管理

物理存储介质:基本存储



- ●主存储器 (main memory)
 - O通过总线与CPU相连,存储运算所需的数据和指 。
 - O速度很快(纳秒级),一般容量在几百M~几G
 - O随机访问: 访问任何存储单元时间相同
 - O易失性: 断电丢失
 - 〇操作系统提供机制,应用程序管理

物理存储介质: 在线存储



- ●快闪存储器 (flash memory)
 - O通过外设接口与总线相连,存储永久保留的数据
 - O速度受到存储介质和接口限制
 - O随机访问, 非易失性, 断电不丢失
 - O文件系统管理,可以通过操作系统在线访问
- ●磁盘存储器 (disk memory)
 - 〇同上, 但是机械装置, 速度更慢

物理存储介质:脱机存储



- ●光盘存储器(CDROM/CDR/CDRW/DVD)
 - O脱机存储,保存备份或者历史档案
 - O分为光物理盘, 光化学盘和光磁盘
 - O只读, 可写一次, 可重复读写
 - 〇机械装置,随机访问,速度更低
 - 〇有标准数据记录格式,操作系统提供文件系统接 口访问

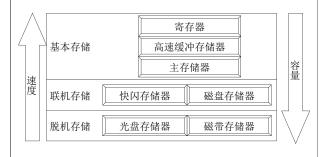
物理存储介质:脱机存储



- ●磁带存储器 (tape)
 - O脱机存储,保存备份或者历史档案
 - O电磁记录原理
 - O机械装置,顺序访问
 - O速度最低,容量价格比最高(至几百G)

物理存储介质层次





磁盘存储器

- ●磁盘物理特性
- ●磁盘性能度量
- ●磁盘块存取的优化
- ●磁盘阵列技术RAID

磁盘存储物理特性



- 盘片 ○磁道track
 - 〇扇区sector 〇柱面cylinder
- 磁头
- 驱动臂

磁盘存储器性能度量



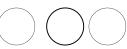
- 容量:磁盘所能容纳的字节数
- ◆ 存取时间:从发出读写请求到数据开始传输之间的时间
 - 〇寻道时间:移动磁盘臂,定位到正确磁道所需的时间。 (平均4-10毫秒)
 - ○旋转等待时间:等待被存取的扇区出现在读写头下所需的时间。(平均2 5毫秒)
- 数据传输率:从磁盘获得数据或者向磁盘存储数据的速率
- 平均故障时间(可靠性):预期系统无故障连续 运行的平均时间MTBF

磁盘块存取的优化



- ●在主存储器中对块进行缓冲以减少块的读 写次数
- ●按柱面组织数据
- ●磁盘臂调度 -- 电梯算法
- ●利用非易失性RAM作为写缓冲
- ●预读和双缓冲

磁盘阵列技术RAID



- RAID: Redundant Array of Inexpensive(Independent) Disks
 - O通过冗余提高可靠性
 - O通过并发访问提高性能
- ●冗余方案
 - O镜像: 复制写到多个磁盘
 - O校验位: 写入数据及其校验和
 - O损失有效容量

磁盘阵列技术RAID





●并发方案

- O位级拆分:将每个字节按位分开,存储到多个磁盘上
- 〇块级拆分:对块进行逻辑编号,对于n个磁盘的阵列,将逻辑上的第i块存储到第(i mod n)+1个磁盘上

磁盘阵列技术RAID



● 标准RAID级别

ORAID 0级: 块级拆分,无冗余

ORAID 1级: 带块级拆分的磁盘镜像

ORAID 2级:内存风格的纠错码组织结构ORAID 3级:位交叉的奇偶校验组织结构ORAID 4级:块交叉的奇偶校验组织结构

ORAID 5级: 块交叉的分布奇偶校验位的组织结构

ORAID 6级: P+Q冗余方案

磁盘阵列技术RAID



●RAID级别的选择

- O所需的额外磁盘存储带来的开销
- O在I/O操作数量方面的性能需求
- O磁盘故障时的性能
- O数据重建过程中的性能

缓冲区管理





●缓冲区buffer

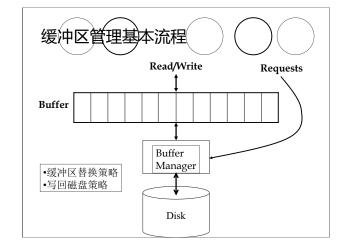
〇主存储器中用于存储磁盘块的拷贝的部分,由固 定数目的缓冲块构成

●缓冲区管理器

O负责缓冲区空间分配调度的子系统

●缓冲区管理的目的

O减少磁盘和主存储器之间传输的块的数目



缓冲区替换策略





● 最近最少使用(LRU)

O替换出最长时间没有读或写过的块

● 先进先出 (FIFO)

O替换出被同一个块占用时间最长的缓冲块

● "时钟" 算法

OLRU的一个常见的、有效的近似

● 系统控制

O查询优化器或者其它的DBMS部件可以给缓冲区管理器提供建议来避免象LRU,FIFO,或者时钟这样的严格的策略可能引起的问题

最近最少使用(LRU)



●基本方法

- 〇缓冲区管理器保持一个表,登记每个缓冲块被访问的最后一次时间
- O每个数据库访问在这个表中生成一个表项。

●LRU是一个有效的策略

- O直觉上,长时间没有使用的缓冲区比那些最近访问过的缓冲区有更小的最近访问的可能性
- O但也有例外

先进先出 (FIFO)







- 〇缓冲区管理器保持一个表,登记当前占用缓冲区的各个块 装入缓冲区的时间
- 〇当块从磁盘读入内存的时候, 生成一个表项
- O 当块被访问时,不需要修改这个表
- 与LRU相比, FIFO需要较少的维护
 - O但它存在更多的问题
 - 〇例如被重复使用的,B-树索引的根块将最终变成一个缓冲区中最旧的块它将被写回到磁盘上,很快又被重新读入另一个缓冲区。

"时钟"





● 基本方法

- 〇 将缓冲区看作一个环。一 个指针指向一个缓冲块
- 〇每一个缓冲块有一个"标志", 0或1
- 〇 带有**0**标志的是很可能被替 换出去的缓冲块
- O 刚读入和访问过的块标记 为1
- O 指针旋转过程中不停将1变 为0



1



1

缓冲区管理系统控制



● 查询优化器或者其它的DBMS部件可以给缓 冲区管理器提供建议

- O例如,将某些块定义为"固定的"来强迫它们保持 在内存中,如B-树的根、数据字典中的块。
- O又如,对于象一遍散列连接那样的算法,查询处理器可以"固定"较小的关系的块,使得确保在全部时间内它都将留在内存中。

块写回控制策略





●被钉住的块

〇为了使数据库系统能够从崩溃中恢复,必须限制 一个块写回磁盘的时间。不允许写回磁盘的块称 为被钉住的块

1

- ●块的强制写出
 - 〇某些情况下,尽管不需要一个块所占用的缓冲区空间,仍必须把这个块写回磁盘。这样的写操作称为块的强制写出。

数据存储组织



●数据库中的数据存储在操作系统管理的文 件中

〇域->记录->(块)->文件

●域根据类型占据不同大小空间

- O定长域类型
- O变长域类型
- O二进制大对象类型(BLOB)
 - ●常用于空间复杂对象的存储,至少可以提供存储管理 和事物支持

数据存储组织





● 记录由域顺序排列组成

- O定长记录
- O变长记录: 含有变长域的记录
- 定长记录文件
 - ○文件中所有的记录都具有相同的长度,从而一个块中所有 的记录都是等长的

● 变长记录文件

〇文件中的记录可以有不同的长度,从而一个块中的各个记录可以具有不同的长度

数据存储组织



● 定长记录文件的顺序 存储

〇块的大小应为记录大小 记录0 的整倍数

●对齐

O删除记录代价高

●删除标记,或

●有效指针链接

记录2 记录3 记录4 记录5 记录6 记录7 记录8

记录1

Perryridge	A-102	400
Rouond Hill	A-305	350
Mianus	A-215	700
Downtown	A-101	500
Redwood	A-222	700
Perryridge	A-201	900
Brighton	A-217	750
Downtown	A-110	600
Perryridge	A-218	700

数据存储组织



●变长记录文件

- O多种记录类型在一个文件中存储
- O记录类型允许一个或多个字段是变长的
- O记录类型允许可重复的字段

● 变长记录文件的定长表示法

- 〇 预留空间:使用长度为最大记录长度的定长记录。对较短记录未使用的空间用特殊的空值或记录终结符号来填充。
- O使用指针:变长记录用一系列通过指针链接起来的定长记 录来表示。

:预留空间法

记录0	Perryridge	A-102	400	A-201	900	A-218	700
记录1	Round Hill	A-305	350	Т	Т	Т	Т
记录2	Mianus	A-215	700	Т	4	4	上
记录3	Downtown	A-101	500	A-110	600	긕	上
记录4	Redwood	A-222	700		1	\dashv	上
记录5	Brighton	A-217	750	1	Т	Τ	上

变长记录文件:使用指针



变长记录文件 :锚块和溢出块

锚块

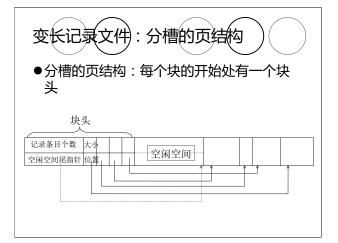
Perryridge	A-102	400	\Box
Round Hill	A-305	350	
Mianus	A-215	700	
Downtown	A-101	500	_
Redwood	A-222	700	
Brighton	A-217	750	

溢出块 A-210 900 700 A-218 A-110 600

变长记录文件: 字节流表示法

●字节流表示:把每个记录作为一个连续的 字节流存储,每个记录的末尾附加一个特 殊的记录终止符号

0	Perryridge	A-102	400	A-201	900	A-218	700	1
1	Round Hill	A-305	350	Т				
2	Mianus	A-215	700	Т				
3	Downtown	A-101	500	A-110	600	Τ		
4	Redwood	A-222	700	Т				
5	Brighton	A-217	750	Т				



数据存储组织



●文件中记录的组织

- O关系中的各个记录存放在文件中的什么位置
- 〇堆文件组织:记录没有顺序,一条记录可以放在 文件中的任何地方。
- 〇散列文件组织: 散列函数的计算结果确定记录应存储到文件的哪个块中。
- O顺序文件组织:记录根据搜索码的值顺序存储

数据字典的存储



- ●数据字典:数据库的描述信息
- O关系模式信息:逻辑结构
 - O关系存储信息: 物理结构
 - O用户信息:安全控制
 - O统计信息:数量/容量统计
 - O索引信息......
- ●RDBMS中,数据字典和普通关系同样存储



参考文献



- [TP311.13/261]空间数据库 = Spatial databases a tour (美) Shashi Shekhar, Sanjay Chawla著 谢 昆青 ... 等译 北京 机械工业出版社 2004
- 北京大学计算机系数据库教研室
 - 〇数据库原理与技术讲义(杨冬青)
- 数据库系统概论
 - 〇岳丽华, 丁卫群 编著
- 〇科学出版社,2000年