

## 前言

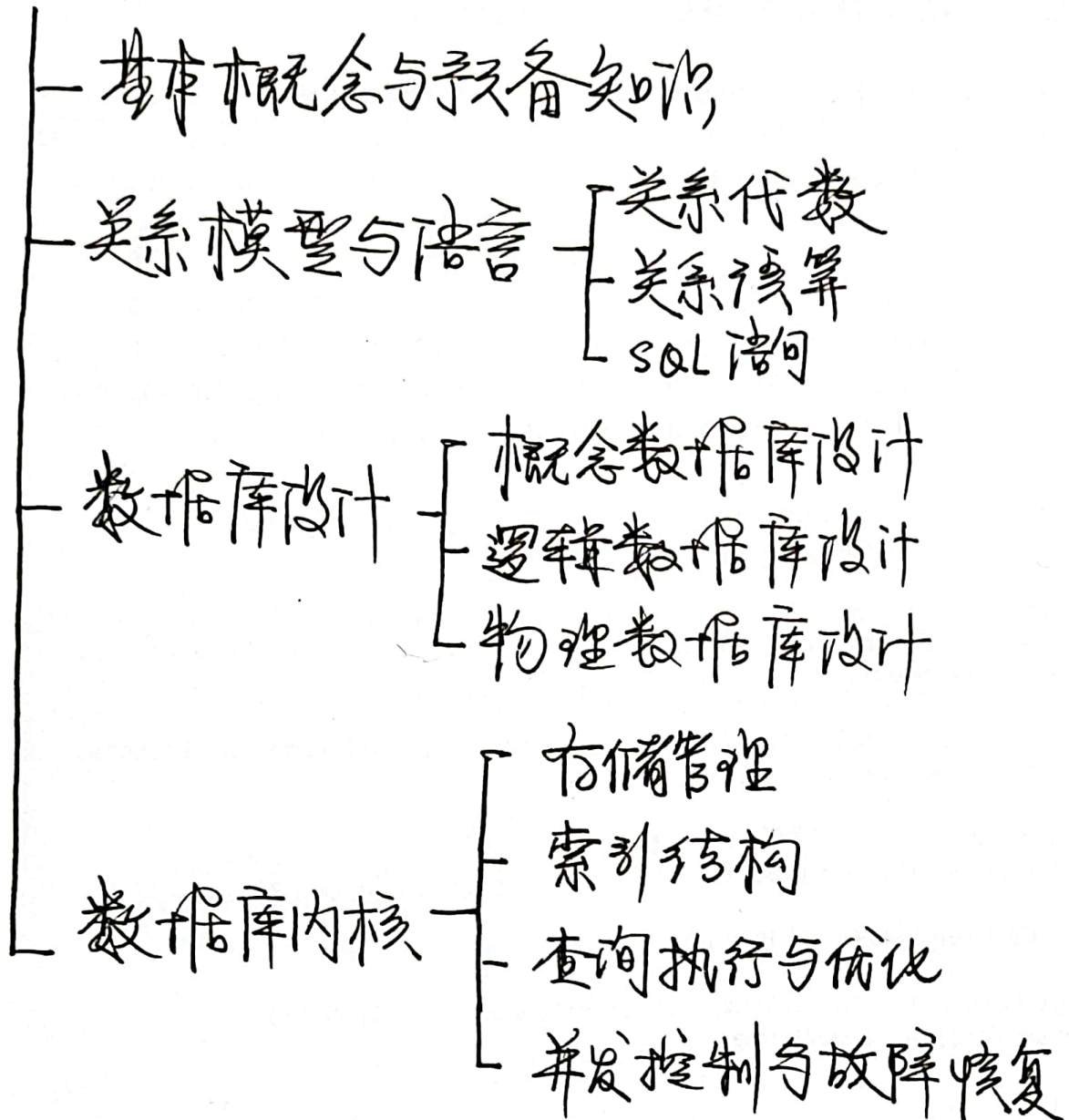
说实话，写这个复习指南（讲义），心里是十分忐忑的，因为《数据库系统》这门课程知识量非常丰富，而且有一定难度，总是担心会误导同学们。我也无法预测考试的内容和重点，所以只能尽可能详尽整理出一份复习指南（讲义），有一些地方用课件展示不是很方便，所以我的这份讲义是手写版，为了方便大家查阅，我把其进行了扫描。不幸的是，由于时间紧任务重，可能字迹比较潦草，请大家见谅。同时，这里的内容均是按照我个人的思路和理解整理的，未必适合所有人，也未必是科学的，大家各取所需即可。

在讲的过程中，由于时间限制，只选取了其中最重要的一部分讲解，余下的内容在本文件中均有体现。需要说明的是，这个复习指南是不能代替老师的课件的，大家还是要以老师的课件为主，看过课件以后如果觉得有必要可以再看一下这个文件，简要了解一下框架。

由于作者水平有限，其中难免会有错误，欢迎大家批评指正。

# 目录 (整体结构)

## 数据库系统



# 一. 基本概念与预备知识

1. 基于文件系统的数据库管理 VS 基于数据库管理系统的信息管理

2. DB VS DBMS VS DBS

3. 数据模型: 是进行数据抽象的工具。

将现实世界映射到计算机世界的过程。

(1) 数据模型的三要素

- 描述DB结构的概念
- 操纵数据结构的操作
- DB服从的一系列约束条件

(2) 分类

- 概念数据模型: 现实世界映射到信息世界
- 实现数据模型: 介于↓↑之间, 实现DBMS时使用
- 物理数据模型: DB在计算机中的存储细节

4. 数据库模式 VS 数据库实例 → 数据库在某一时点的存储数据

↓  
对数据库的结构、类型、约束的描述, 是DB的类型声明

## 5. 数据库的三层模式结构

┌ 内核式/存储模式 — 物理模型

├ 概念模式 — 实现模型

└ 外模式/视图 — 实现模型 (可以有多个)

└ 只有一个

(1) 模式映射用于完成请求转换和数据转换 → DBMS检查到的数据转换为外模式的组织形式

↓

应用程序在外模式上声明的数据请求 → 转换 → DBMS在内模式上的请求

(2) 模式映射的分类

- 外模式 — 概念模式映射
- 概念模式 — 内核式映射

(3) 数据独立性

- 逻辑数据独立性: 概念模式变化时, 修改外模式到概念模式的映射
- 物理数据独立性: 内核式发生变化时, 修改概念模式到内核式的映射

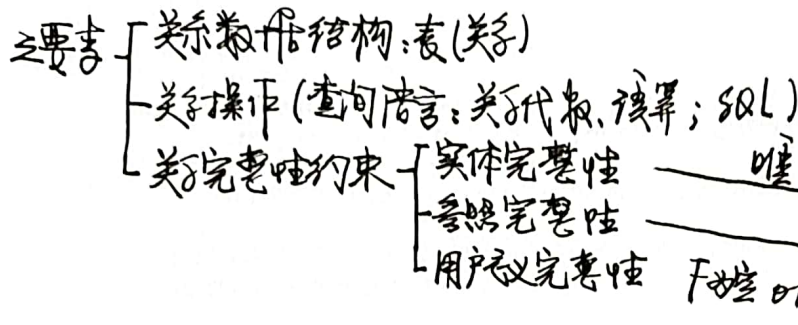
6. 数据库语言

- 数据定义语言 (DDL): Create
- 数据操纵语言 (DML): Select

## 二. 关系核型与语言

### (一) 基本概念

1. 关系数据库核型 (是一种实现数据库核型)



2. 键: 具有区分不同元组的属性的集合

候选键: 极小的键

↓ 人为指定一个

主键

是R的外键

外键: R(F) 参照 S(K)

参照关系

被参照关系

R 可以大于 S

### (二) 关系代数

1. 基本关系代数操作

(1) 选择  $\sigma$ :  $\sigma_b(R)$

(2) 投影  $\pi$ : 自动去重  $\pi_b(R)$

(3) 并  $\cup$ : 要求R和S必须具有相同表的属性且对应属性的值域相等  $R \cup S$

(4) 差  $-$ : 要求同并  $R - S$

(5) 笛卡尔积  $\times$ : 无条件连接 (拼接)  $R \times S$

(6) 重命名:  $\rho_{B \leftarrow A}(R)$ : 把关系R的属性A更名为B

$\rho_S(K)$ : 把关系R更名为S

$\rho_S(A_1, A_2, \dots, A_n)(R)$ : 把关系R更名为S, 并把R的全部属性更名为  $A_1, \dots, A_n$

2. 派生关系代数操作

(1) 交  $\cap$ : 求同并; 可以等价交换成  $R - (R - S)$

(2) 内连接

- $\theta$  连接  $R \bowtie_{\theta} S$  ( $\theta$  是连接条件, 属于同  $\sigma$  中的  $\theta$ )  
结果中包含R和S全部属性, 同名加前缀
- 自然连接  $R \bowtie S$  自带条件: 同名属性上的等值比较  
从连接结果中去掉重复的同名属性

(3) 外连接: 连接结果中保留R或S中的不满足连接条件的元组

左外连接  $R \bowtie_{(a)} S$ : R中不满足条件的也会在结果中, 对应S属性为空

右外连接  $R \bowtie_{(b)} S$ : S

R

全外连接  $R \bowtie_{(a,b)} S$ : 上面二者兼有

// 反连接  $R \bowtie_{\neg \theta} S$ : 不满足连接条件的元组进行连接

(4) 除:  $R \div S$ : 结果中只包含R中的属性, 不包含S中的属性  
 $R \div S$ 的结果是使得  $S \times T \subseteq R$  的最大的关系T  
 一个例子: 选修了所有课程的学生:  $All(sno, cno) \div All(cno)$

### 3. 扩展关系代数操作

(1) 分组操作  $\gamma_L: agg(R)$ 

- L: 分组属性列表, 用逗号分隔
- agg: 聚集函数表达式列表, 用逗号分隔
- 例:  $sum(score) \rightarrow Total\ score$  计算score的和, 结果命名为Total score

(2) 赋值操作: 避免过于冗长的查询表达式

### (三) 关系演算:

#### 1. 元组关系演算:

(1) 表达式形式:  $\{t \mid P(t)\}$   
 $t[A]$ : 元组t中属性A的值  
 $t \in R$ : t是关系R中的元组  
 结果元组      条件  
 结果是所有使谓词P为真的元组t的集合。

(2) 例子: ① 计算机系和数学系的学生  $\{t \mid t \in Student \wedge (t[Sdept]='CS' \vee t[Sdept]='MA')\}$   
 ② 全体学生的学号和姓名  $\{t \mid \exists s \in Student (t[Sno]=s[Sno] \wedge t[Shame]=s[Shame])\}$

#### 2. 域关系演算:

(1) 表达式形式:  $\{(x_1, x_2, \dots, x_n) \mid P(x_1, x_2, \dots, x_n)\}$   
 $(x_1, x_2, \dots, x_n)$  域变量  $x_1, x_2, \dots, x_n$  构成的元组  
 域变量      域关系演算公式      结果是所有使  $P(x_1, x_2, \dots, x_n)$  为真的元组  $(x_1, x_2, \dots, x_n)$  的集合

(2) 例子: ① 计算机系和数学系的学生  $\{(n, m, s, a, d) \mid (n, m, s, a, d) \in Student \wedge (d='CS' \vee d='MA')\}$   
 ② 全体学生的学号和姓名  $\{(n, m) \mid \exists s, a, d ((n, m, s, a, d) \in Student)\}$

### (四) SQL 语句

#### 1. 数据定义

##### (1) 基本数据类型

(2) 创建关系模式 Create table 表名 (属性名 数据类型, ..., 主键, 外键)

例: Create table sc (sno char(6), cno char(4), grade int,

(3) 声明用户定义完整性约束: primary key (sno, cno), foreign key (sno) references student (sno)

非空 NOT Null, 不重复 unique, 缺省值: default 缺省值, check (表达式)

检查是否满足      只解析不处理

(4) 删除关系: Drop table 关系名1, 关系名2

(5) 修改关系模式 (alter table)

① 修改关系名: alter table 旧名字 rename to 新名字

② 增加, 修改, 删除属性

- 增加属性: alter table 表名 add 列名 数据类型;

- 删除属性: alter table 表名 drop 列名;

- 增加表约束: alter table 表名 add constraint 约束名 约束;

- 删除表约束: alter table 表名 drop constraint 约束名;

- 修改属性名: alter table 表名 change 旧属性名 新属性名 属性定义;

- 修改属性定义: alter table 表名 modify 属性名 属性定义

(6) 视图

① 创建视图 create view 视图名 (属性名列表) AS 子查询;

② 修改视图 alter view 视图名 (属性名列表) AS 子查询;

③ 删除视图 drop view 视图名;

④ 查询视图 语句同表查询

2. SQL 数据更新

(1) 插入数据

① 直接插入 insert into 表名 (属性名列表) values (表值/值列表);

若没有给出属性名列表则插入一行完整元组; 没列出的属性值为空.

② 插入查询结果: insert into 表名 子查询;

(2) 修改数据

① 基于本关系的数据修改: update 表名 set 属性名 = 表达式1 ... where 修改条件;

仅涉及本关系

② 基于外部关系的数据修改: 只是修改语句中不只有本关系, 其新写上

可选

(3) 删除数据

① 基于本关系的数据删除: Delete from 表名 where 删除条件;

② 基于外部关系的数据删除

可选 (没有的时候删除所有数据)

(4) 数据完整性检查: 检查修改的过程中是否破坏了完整性约束

① 实体完整性约束检查:

检查时机: 插入元组或修改元组的主键属性值时

处理方法: 主键为空或不唯一都拒绝插入和修改

② 用户定义完整性约束检查

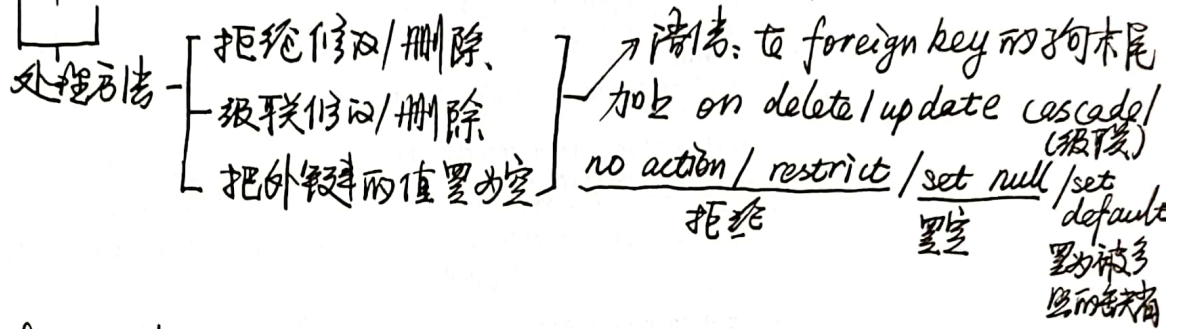
检查时机: 修改/插入元组

处理方法: 拒绝或拒绝

③ 参照完整性检查:

破坏时机: 插入/修改元组后外键的值在被参照关系中找不到。— 处理方法: 拒绝

— 修改/删除元组后外键的值在被参照关系中找不到。



3. 数据查询

核心 select — from — where —

1) 单关系查询

类型

① 投影: select (distinct) 属性名列表 from 关系名;

例子: select sname, (2022 - Sage) as bd from Student;

可以是表达式 重命名

② 选择查询: select (distinct) 表达式列表 from 关系名 where 选择子句;

选择条件

- 表达式比较: <>: 不等于
- 范围比较: 表达式1 (not) between 表达式2 and 表达式3
- 集合元素判断: 表达式1 (not) in (表达式2, ..., 表达式n)
- 字符串匹配: 字符串表达式 (not) like 模式 (escape 转义字符)
- 通配符: 匹配单个字符, 通配符 %: 匹配任意长度的字符串
- 字符串正则表达式匹配
- 空值判断:  $\neq$  Null  $\neq$  Null = Null
- 逻辑运算符

③ 集合操作: 求两个查询语句结果的并、交、差

- 并: 查询语句1 union (all) 查询语句2 (有all则不重复)
- 交: Intersect; 差: minus/except

说明: 查询语句1的结果属性名作为集合操作的结果属性名

(MySQL中不支持交和差: 交用  $R \cap S$ , 差用  $\pi(\sigma_{S=Null}(R \rightarrow S))$ )

④ 查询结果排序

升序(默认) 降序

在查询语句的后面加上 order by 列名 (ASC/DESC); ..., 列名n (ASC/DESC)

## ⑤ 聚集查询

SELECT 聚集函数 ((distinct) 表达式) from ... where ...;

※ 聚集函数不能出现在 where 子句中 (用嵌套查询来实现)

- count (\*): 所有元组的数量
- count (distinct 表达式): (不同的) 非空表达式数量
- max/min (distinct 表达式): 表达式的最大值/最小值
- sum (distinct 表达式): (不同) 表达式值的和
- avg (distinct 表达式): (不同) 表达式值的平均值

## ⑥ 分组查询

SELECT 分组属性列表, 聚集函数表达式列表 From ... where ... group by 分组属性列表

说明: 根据分组属性分组, 分组属性值相同的为一组; (Having 分组筛选条件) (可以用聚集函数)

对每个组中的非分组属性进行聚集.

比如: 查询 2 门以上课程得分超过 80 的学生的学号及这些课程的平均分

⑤ Select Sno, AVG(Grade) from SC Where Grade >= 80 Group By

③ Sno Having Count (\*) >= 2;

## 2. 连接查询

(1) 笛卡尔积

Select ... from 表名1 (cross join) 表名2 (...) ... , 表名n

(2) 内连接

<1> select ... from 表名1 (inner) join 表名2 on 连接条件  
可选

<2> 当内连接是等值连接且连接属性同名时, 可使用: (连接属性只保留一个副本)

select ... from 表名1 (inner) join 表名2 using (连接属性列表)

<3> 自然连接: 表名1 natural join 表名2

<4> 自连接: 别名重命名

(3) 外连接:

<1> 左/右/全外连接: 表名1 left/right/full (outer) join 表名2 on 连接条件

<2> 自然外连接: 表名1 natural left/right/full join 表名2

## 3. 嵌套查询

(1) 一些基本概念

查询块: 一个 SELECT-FROM-WHERE

嵌套查询: 一个查询块中嵌套另一个, 内层的称为子查询

子查询的类型: 不相关子查询 (子查询不依赖外层查询)

相关子查询 (子查询依赖外层查询)



(2) 嵌套查询的写法:

- ① 在集合判断条件中使用子查询: 使用 (not) in
- ② 在比较条件中使用子查询: 表达式 比较运算符 (all/any/some) (子查询)  
 ↓  
 当子查询结果只有一个值时可以用 not in
- ③ 在存在性测试条件中使用子查询: (not) exists (子查询)  
 ↳ 只需使用 select \*

从父查询中依次取出每组元组, 拿出在子查询中需要的属性值代入子查询, 若子查询结果 (不为空), 则把该条元组按需输出。

<1> 用 Exists 实现全称量词  $\forall$  功能:  $\forall x (P(x)) \Leftrightarrow \neg \exists x (\neg P(x))$

查询选修了全部课程的学生学号

思路: 设符合条件学生的元组,  $\forall c \in Course \{ \exists s \in SC, s[Sno]=t[Sno] \wedge s[Cno]=c[Cno] \}$   
 $\Leftrightarrow \neg \exists c \in Course (\neg \exists s \in SC, s[Sno]=t[Sno] \wedge s[Cno]=c[Cno])$

```
SELECT Sno FROM Student WHERE NOT EXISTS (
  SELECT * FROM Course WHERE NOT EXISTS (
    SELECT * FROM SC WHERE SC.Sno = Student.Sno AND SC.Cno
      = Course.Cno));
```

<2> 用 Exists 实现逻辑蕴含  $\rightarrow$  功能:  $x \rightarrow y = \neg x \vee y$  例: 查询至少选修过 CS-001 选修的课程的学生的学号

思路: 设 t ∈ Student 是符合条件的元组

$\forall s \in SC, ((s[Sno]='CS-001') \rightarrow (\exists x \in SC, (x[Sno]=t[Sno] \wedge x[Cno]=s[Cno])))$

$\Leftrightarrow \neg \exists s \in SC, ((s[Sno]='CS-001') \wedge (\neg \exists x \in SC, (x[Sno]=t[Sno] \wedge x[Cno]=s[Cno])))$

```
SELECT Sno FROM Student WHERE NOT EXISTS (
  SELECT * from SC as SC1 where SC1.sno='CS-001' AND not exists(
    SELECT * from SC as SC2 where SC2.sno=Student.sno and SC2.cno=SC1.cno));
```

④ 子查询结果作为聚合关系: FROM (子查询)

要求: 子查询必须是独立子查询, 派生表要重命名

例: 选修了 2 门以上课程的学生学号和选课程

```
SELECT Sno, T.Amt FROM (SELECT Sno, count(*) AS Amt FROM SC
  Group By Sno) AS T WHERE T.Amt >= 2;
```

⑤ 在 WITH 子句中使用了子查询

例: 问题同上

```
With T as (select sno, count(*) as Amt from SC Group By Sno) SELECT
  sno, Amt From T where Amt >= 2;
```

### 三. 数据库设计

#### (一) 概念数据库设计

##### 1. 数据库设计的过程

- (1) 概念设计: 设计数据库的(抽象)概念模型.
- (2) 逻辑设计: 设计数据库的概念模式.
- (3) 物理设计: 设计数据库的内核

##### 2. ER模型 → ER图

↳ 将现实世界抽象成实体及实体间的联系.

##### (1) 实体(普通实体)与属性

↓  
矩形 □

实体 vs 实体型 vs 实体集

↓                      ↓                      ↓  
对象                      类型                      对象集合

- 简单属性 ○
- 复合属性 ○-○-○
- 多值属性 ○
- 派生属性 (---)
- 键属性 (—)

(当前存储在数据库中的某实体型的实例的集合)

##### (2) 联系

##### ① 联系 VS 联系型 VS 联系集

↓                      ↓                      ↓  
定义                      具体类型                      DB中当前存储的联系型的实例的集合

##### ② 表示形式: 菱形 菱形

##### ③ 连线 线上的数 — 基数比: 实体型参与的联系型中的最大基数

(联系型的约束 有依赖约束(参与度约束): 参与联系型中的最小基数, 部分参与用单线, 全部参与用双线

##### ④ 联系型可以有属性, N:1 联系型的属性可被移到 N-1 的实体中.

##### ⑤ 多元联系: 一个 n 元联系和 n 个二元联系所表示的意义通常是不同的.

##### (3) 弱实体型: 一用弱实体型的部分键 + 标识实体型的主键来区分不同弱实体.

##### ① 弱实体型: 没有键属性的实体型 □

##### ② 标识实体型 / 弱实体型: 弱实体型依赖的用于区分的实体 } 通过标识联系型关联

##### ③ 部分键: 区分同一标识实体相关联的弱实体的属性集合. (---) 弱实体型全部参与



#### (二) 逻辑数据库设计

##### 1. 在逻辑数据库设计阶段, 需要将 ER 模型转换为关系数据库模式.

##### (1) 实体型的转换

- ① 无多值属性时:
  - 实体型的属性名 → 关系的属性名
  - 实体型的键属性 → 主键
  - 实体 → 元组

- ② 复合属性: 把最低值成员作为关系的属性
- ③ 多值属性: 单键属性, 其中包含实体型的主键并建立外键约束.

## (2) 弱实体型的转换

- ① 弱实体型的名  $\rightarrow$  关系名
- ② 弱属性  $\cup$  弱主键  $\rightarrow$  关系属性集
- ③ 弱主键  $\cup$  弱部分键  $\rightarrow$  关系键集
- ④ 弱实体型关系  $\cup$  弱实体型关系的外键约束。

## (3) 联系型的转换

- $M:N$ 型: (1) 名  $\rightarrow$  名
- $1:2$ 型: (1) 1的主键  $\cup$  2的主键  $\cup$  联系型属性集  $\rightarrow$  关系属性集
- $1:1$ 型: (1) 2的主键和联系型的属性并  $\cup$  1的属性集
- $1:1$ 型: (2) 建立1到2的外键约束
- $1:1$ 型: 同  $N:1$
- 二元联系型: 注意对重名属性重命名
- 相似联系型不转换

## 2. 关系数据库规范化理论 —— 评价一个关系数据库模式设计的“好坏”

### (1) 函数依赖

① 定义:  $R(U)$  为属性集  $U$  上的关系模式,  $X, Y \subseteq U$ . 若对于  $R(U)$  的任意关系实例中的任意两个元组  $t_1$  和  $t_2$ , 若由  $t_1[X] = t_2[X]$  能推出  $t_1[Y] = t_2[Y]$ , 则称  $X$  函数决定  $Y$ , 或  $Y$  函数依赖于  $X$ , 记作  $X \rightarrow Y$ .

### ② 函数依赖的类型

<1> 若  $Y \subseteq X$ , 则  $X \rightarrow Y$  是平凡函数依赖

<2> 完全函数依赖: 若  $X \rightarrow Y$ , 且对  $\forall X' \subseteq X$ , 都有  $X' \not\rightarrow Y$ , 则称  $Y$  完全函数依赖于  $X$ , 记作  $X \xrightarrow{f} Y$

部分函数依赖: 若  $X \rightarrow Y$ , 且  $\exists X' \subseteq X$ , 使得  $X' \rightarrow Y$ , 则称  $Y$  部分函数依赖于  $X$ , 记作  $X \twoheadrightarrow Y$

<3> 传递函数依赖: 若  $X \rightarrow Y$ ,  $Y \rightarrow Z$  且  $Y \not\subseteq X$ ,  $Y \not\rightarrow X$ , 则称  $Z$  传递函数依赖于  $X$ , 记作  $X \twoheadrightarrow Z$

### ③ 用公理系统 Armstrong 推导其他函数依赖

<1> 推理规则  $\rightarrow$  正确且完备的

- 自反律: 若  $Y \subseteq X \subseteq U$ , 则  $F \models X \rightarrow Y$
- 增广律: 若  $F \models X \rightarrow Y$ , 则对  $\forall Z \subseteq U$ , 有  $F \models XZ \rightarrow YZ$
- 传递律: 若  $F \models X \rightarrow Y$  且  $F \models Y \rightarrow Z$ , 则  $F \models X \rightarrow Z$
- 合并规则: 若  $F \models X \rightarrow Y$  且  $F \models X \rightarrow Z$ , 则  $F \models X \rightarrow YZ$
- 伪传递规则: 若  $F \models X \rightarrow Y$  且  $F \models WY \rightarrow Z$ , 则  $F \models XW \rightarrow Z$
- 分解规则: 若  $F \models X \rightarrow Y$ , 则对  $\forall Z \subseteq Y$ , 有  $F \models X \rightarrow Z$

④ 用属性集闭包推导函数依赖  
 根据F中的FD, 属性X决定的所有属性  
 <1> 定义: 设R(U, F)是一个关系模式, 集合  $\{A \mid A \in U, F \models X \rightarrow A\}$  称为X关于F的闭包, 记作  $X^+$   
 定理:  $F \models X \rightarrow Y \Leftrightarrow Y \subseteq X^+$  (把逻辑蕴含判断问题化成闭包计算问题)

<2> 函数依赖集的闭包: F逻辑蕴含的全部函数依赖的集合叫做F的闭包, 记作  $F^+$

<3> 函数依赖集的覆盖: 若  $G^+ \subseteq F^+$ , 则称F覆盖G

<4> 等价函数依赖集: 若  $F^+ = G^+$ , 则F与G等价.  $F^+ = G^+ \Leftrightarrow F^+ \subseteq G^+ \text{ 且 } G^+ \subseteq F^+$   
 判断等价的方法: 引例2

<4> 函数依赖集的最小覆盖

最小覆盖是指下列3个中的F的等价函数依赖集

(I) 不存在冗余函数依赖 (不能去掉一个后和原来等价)

(II) 函数依赖左部不存在冗余属性

(III) 函数依赖右部仅包含一个属性

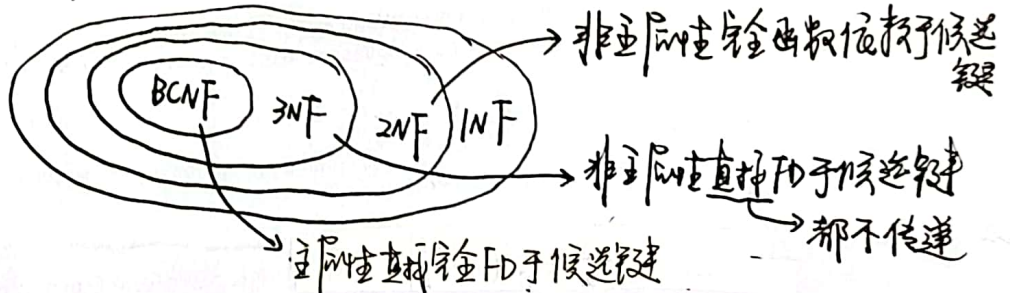
如何计算最小覆盖:

(I) 分解函数依赖, 使右部只有一个属性

(II) 删除函数依赖左部的冗余属性

(III) 删除冗余函数依赖

## (2) 关系模式的范式



① 非主属性: 不参与任何候选键的属性

② 范式之间的关系

1NF

↓ 1 除非主属性时候选键的部分FD

2NF

↓ 1 除非主属性时候选键的传递FD

3NF

↓ 1 除非主属性时候选键的部分FD和传递FD

BCNF

③ 注意: 关系模式的规范化程度并不是一层比一层高, 例如

## (3) 关系模式分解

判断方法: R(U, F) 的分解  $\rho = \{R_1(U_1), R_2(U_2)\}$  是无损分解  $\Leftrightarrow F \models U_1 \cap U_2 \rightarrow U_1 \text{ 或 } F \models U_1 \cap U_2 \rightarrow U_2$

① 两个准则: 无损连接性: 令  $\rho = \{R_1(U_1), \dots, R_n(U_n)\}$  是对关系模式R(U)的一个分解, 若对于R的任意实例r均有  $r = \pi_{U_1}(r) \bowtie \pi_{U_2}(r) \bowtie \dots \bowtie \pi_{U_n}(r)$ , 则是无损的

函数依赖保持性: 设R(U, F)是关系模式, U是属性集, F是函数依赖集. 函数依赖保持性: 对于  $V \subseteq U$ , F在V上的投影是  $F^+$  中所有  $X \rightarrow Y$  且  $X, Y \subseteq V$  的函数依赖  $X \rightarrow Y$  的集合

什么是保持函数依赖的分解?  $\square \Rightarrow \square \square \dots \square$  若  $F_1 \cup F_2 \dots \cup F_n$  覆盖F, 则是

## (三) 物理设计

## 四. 数据库内核

# 1. 存储管理

## 1. 面向磁盘的数据库存储

(1) 为什么不用OS直接管理数据库存储?

当缺页时而由OS的调页, 但OS何时调页是一个全局考量, 会有延迟

(2) 数据库的文件存储

(1) 将一个数据库存储为一个或多个文件。

(2) 每个文件包含多个页。

→ 唯一的编号(页号) → 间接层 → 页的物理地址

(3) 存储管理器:

DBMS的存储管理器负责管理数据库文件

- 记录页中元组的读写
- 记录页中的空闲空间

## 2. 面向行的存储

(1) 字符串的表示

- char(n)型: 若字符数小于n, 则后面用空字符补全
- varchar(n)型: 数组头部先存串长, 再存串的内容

字节串的表示

- BINARY(n): ≤ 0 补全
- VARBINARY(n): ≤

(2) 元组存储: 元组头 + 数据

由元组的所有属性值拼接而成

### 记录元组的元数据

### 存储元组记录

- 指向元组所在表的关系的模式定义的指针
- 元组的长度
- 元组的最后修改时间
- 元组的可见性(并发控制)
- 元组哪些属性值非空。

- 通常按定义关系模式时指定的属性顺序存储属性值
- 每个属性值在元组中存储位置的偏移量是4字节或8字节的倍数。

(当元组是变长元组时的布局: 总长度值和变长属性值分别置于元组两端, 元组头后紧跟指针数组, 指向每个属性值。

(3) 页布局: 页头 + 页数据

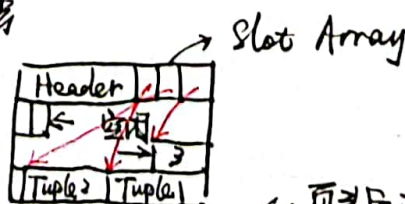
### 记录页的元数据

- 页的大小
- 页的校验和, DBMS版本, 数据库版本号
- 页的可见性(并发控制)

### 组织方法

面向元组的组织方法

日志结构的组织方法



① 面向元组的组织方法 - 分槽页

<1> 每个元组占一个槽

- 槽从后向前布置在页的末尾

- 每个槽的起始位置的偏移量是4B或8B的倍数

<2> 页头后跟槽数组

数组第i个非空元素对应第i个槽的起始位置的偏移量。

<3>槽的元数据存储在页头中

- 槽的数量: 用于确定下一个空闲槽的编号
- 最后一个槽的起始位置的偏移量: 用于确定新元组的插入位置

<4>记录号 (DBMS 为一个表中的每个元组分配的唯一记录号)

- $f_1$ : (页号, 槽号)
- $f_2$ : 唯一的整数作为记录号 (再通过间接层映射成页号, 槽号)

<5> 溢出页、页碎片化 (浪费磁盘空间、增加磁盘 I/O) —— 定期回收

② 日志结构与页号 // 了解

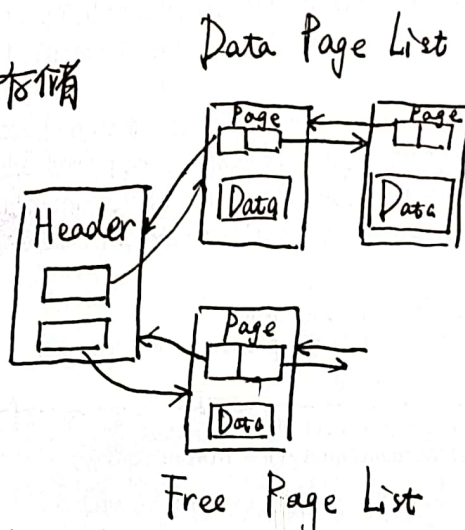
(4) 文件组织

① 堆文件组织: 堆文件中的元组以任意顺序存储

↳ 堆文件中页的组织方法

基于链表的组织方法

- 数据页
- 空闲页
- 头页



基于页目录的堆文件页组织方法: 页目录中记录每个数据页的位置和空闲空间信息

② 顺序/有序文件中的元组按排序键的顺序存储 (除非排序键是主键, 否则用的少)

③ 哈希文件组织

3. 缓冲区管理: 缓冲区管理器负责在 Disk 和 MM 之间复制文件页.

(1) 缓冲区: DBMS 将可用内存区域划分为页数组, 和作为缓冲区  
缓冲区中的页称为页框.

(2) 缓冲区的设计: 页表

↳ 记录 Buffer Pool 中当前有哪些页及它们在内存中的地址  
把页号映射为该页所在页框的地址

① 页框的元数据 -  $pin\_count$ : 页框中的页当前被请求但未释放的次数, 即引用计数 (时时)  
-  $dirty$  (脏位): 该页 Buffer pool 后是否被修改过

② 缓冲区管理器的功能

<1> 请求页 - 若在 Buffer Pool 里:  $pin\_count++$ , 返回页框地址  
不在 - 在池中选一个  $pin\_count=0$  的页框, 加 1, 页区逐旧页 (脏页不写回), 调入新页  
- 没有  $pin\_count=0$  的, 替得直到释放池中页面

<=> 释放页: pin-count --;  
 <=> 修改页: 把 dirty 置为真.

用时间戳记录 Buffer Pool 中每个页框的最后访问时间, 替换时选择 pin-count=0 且最旧的页 (维护一个堆)

③ 页替换策略

- 最近最少使用 (LRU)
- 时钟替换策略 (CLOCK)
- 先进先出 (FIFO)
- 最近最多使用 (MRU)

(3) 缓冲区管理器有预取机制

(4) 缓冲区与 VM 比较

- ① 相同点: 访问内存不够用, 调页, 页替换
- ② 不同点: ① DBMS 更准确地预测页的访问顺序 (更加选择被替换的页, 更加预取)  
 ② DBMS 必须具备强制写回的能力

4. 数据库的存储/柱状存储

(1) 把关系的每个属性单独存储  
 - 当查询仅涉及少量属性时, 不需要读无关属性, 减少 I/O  
 - 表中相邻字节来自于同一个属性, 缓存效率高  
 - 压缩效率高, SIMD 指令向量化计算

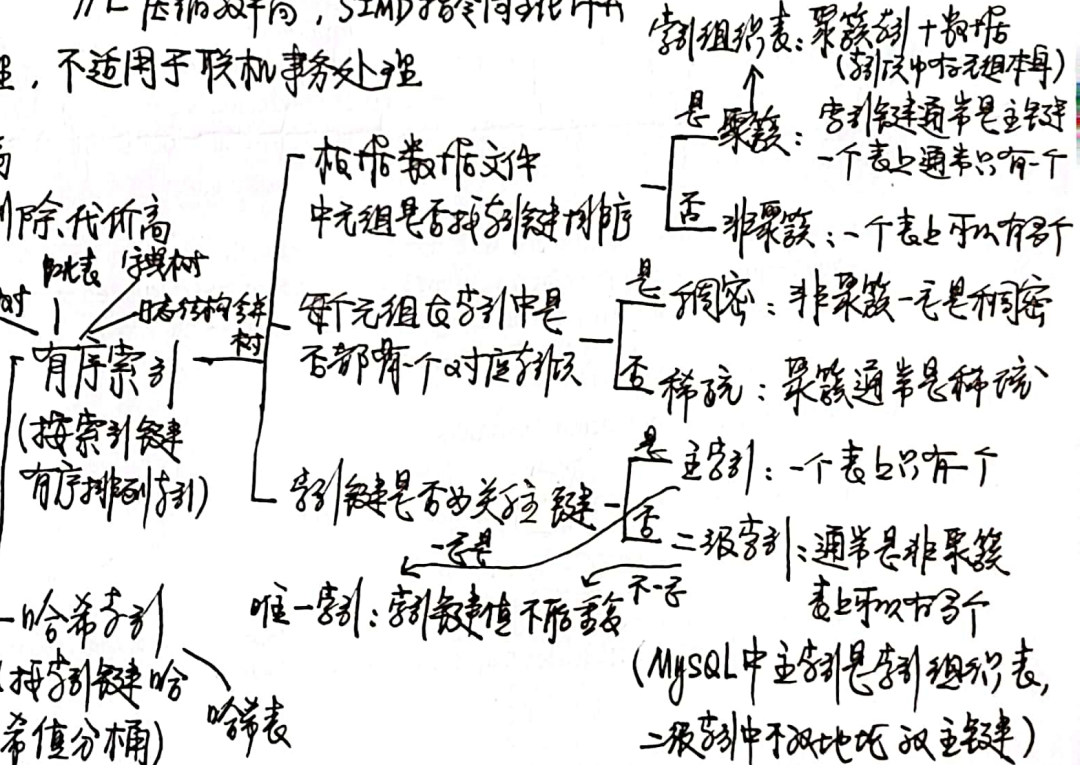
(2) 适合于联机分析处理, 不适用于联机事务处理

- (1) 元组重构代价高
- (2) 元组修改和删除代价高
- (3) 解压代价高

(二) 索引结构

1. 索引的分类

(1) 按照索引的实现方式



(3) 有序索引中的概念 (索引键)

索引键: 索引根据一组属性来定位元组, 索引记录了元组的索引键值与元组地址的对应关系  
 索引项: 索引中的 (键值, 地址) 对, 按索引键值排序

(3) 哈希索引中的概念

h: 哈希函数  
 键为 k 的索引项位于编号为 h(k) 的桶。  
 只支持等值查找

#### (4) 索引的SQL语句

- ① 创建主索引: 声明主键时自动建立
- ② 创建二级索引: create index 索引名 on 表名 (索引键)
- ③ 创建唯一索引: 建表时用 Unique 声明会自动创建唯一索引, 或者在②中 index 前加 unique
- ④ 创建外键索引: 声明外键时自动
- ⑤ 删除二级索引: DROP INDEX 索引名 on 表名
- ⑥ 删除主索引: Drop index 'Primary' on 表名 (删除后重新组织元组)

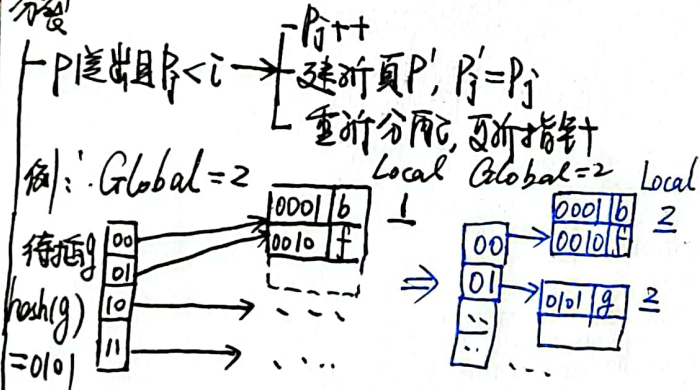
#### 2. 基于哈希的索引结构: 哈希表 (外在哈希表)

静态哈希表: 桶的数是固定不变的  
 动态哈希表: 桶的数可变, 每个桶中索引项所占在大约 1/2 页

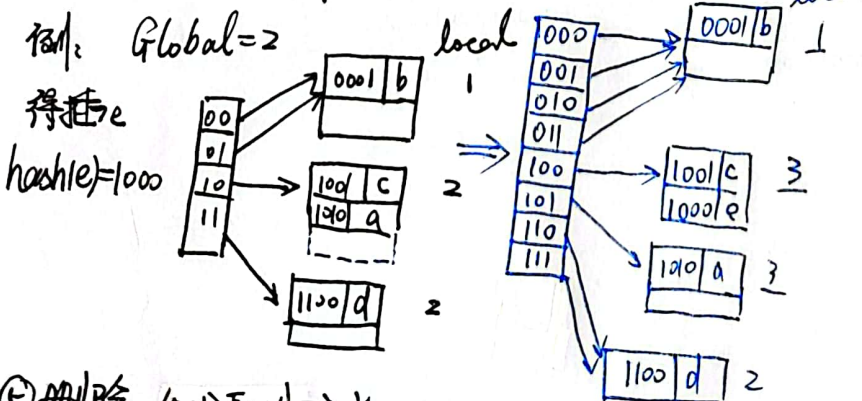
##### 1) 可扩展哈希表

- ① 用二进制全局标识, 一共有  $2^i$  个桶
- ② 每个桶无溢出页, 页存得下可以共用  $\Rightarrow$  该页  $i$
- ③ 每个页记录一个局部标识  $j$ , 该页中前  $2^j$  个  $hash(k)$  相同

插入索引项时: 先找往哪插, 不够插就分裂



Global = 2, Local = 2



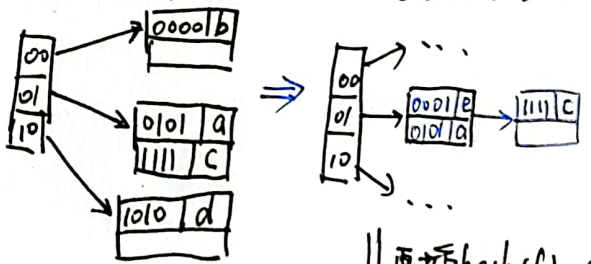
⑤ 删除: 先找再删无合并

#buckets 桶数  
 #entries 记录数 线性哈希表

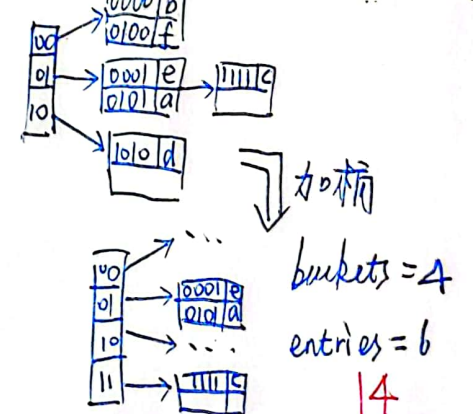
- ① 一共  $n$  个桶, 每个桶中指针指向页链表
- ② 设一页中能存  $b$  个索引项, 一个桶一页则线性哈希表最多存  $nb$  个项 ( $nb < 1$ )
- ③ 先看桶有几位, 然后看哈希值的后这些位决定在哪个桶, 如果没有这个桶就少看一位
- ④ 插入, 先找往哪插, 页数为 0, 是否  $nb < 1$  是则生成, 否则加一个桶, 重新分配

例: 待插入  $hash(e) = 0001$ ,  $\theta = 0.85$

#buckets = 3, #entries = 4  
 buckets = 3, entries = 5  $\leq 0.85 \times 2 \times 3$



再插入  $hash(f) = 0100$   
 bucket = 3, entries = 6 加桶





### (3) 可扩展Hash表

VS 稠密Hash表

桶数  $2^i$   
 是否有溢出 无  
 哈希方案  $hash(k)$  的前  $i$  位  
 页分配中 页发生溢出  
 添加桶方法  $i++$ , 桶翻倍

$n$   
 有  
 $hash(k) \bmod 2^m$  或  $hash(k) \bmod m$   
 $entries > 0.6n$   
 桶数  $+1$

### 3. B+ 树

(1) 性质:

⌈ 路平孩子树

- 除根节点外每个节点至少半满:  $m/2 - 1 \leq keys \leq m - 1$
- 每个节点恰如放入  $i$  个页
- 叶节点包含一个通常按索引键排序的索引项数组和指向右叶的指针

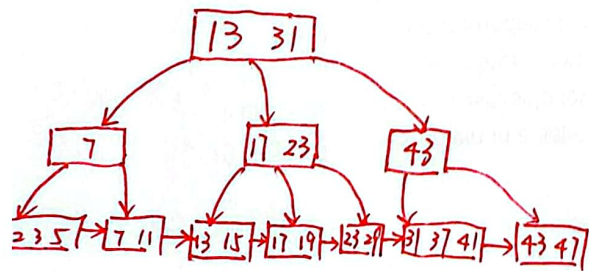
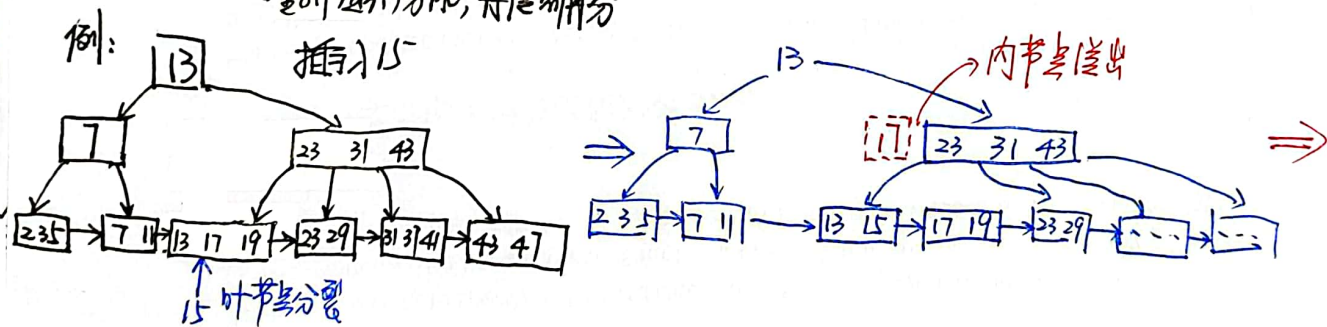
(2) 查找

- 单值查找: 先找在哪个叶上, 再在这个叶上找
- 区间查找: 查找  $[L, U]$ , 先找  $\geq L$  的索引项的索引, 扫描其右侧,  $\leq U$  输出, 否则继续

(3) 插入

- 找到往哪插后插入, 不溢出结束, 溢出则分裂
- 叶节点分裂: 右半溢出页右边再建一个叶, 溢出页的右一半给新页。  
 [ 新页中的最小键值插入父节点, 重新组织分配, 再溢出再分 ]
- 内节点分裂: 建一个新节点, 同溢出节点的右半部分给新节点 (中间键值不传, 传上一级节点)  
 [ 重新组织分配, 再溢出再分 ]

若根节点溢出, 则直接  
 ↑ 增加一个新  
 根节点



### (4) 删除

① 先找到再删除, 若不满足"半满"再处理

② 处理叶节点

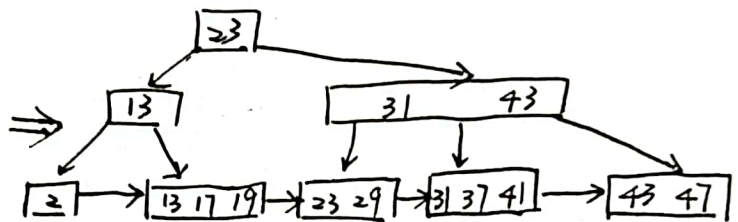
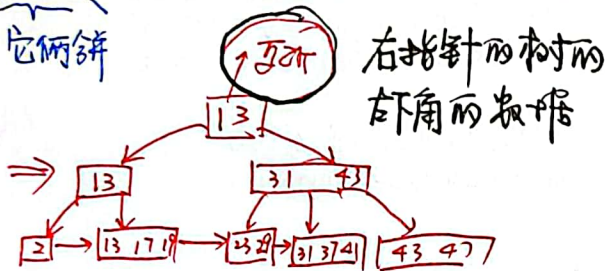
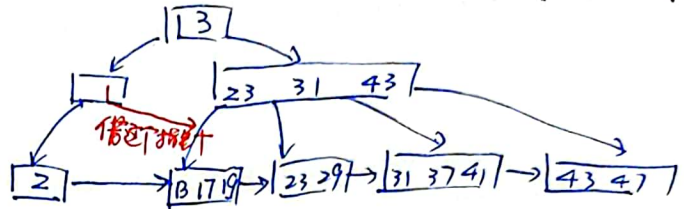
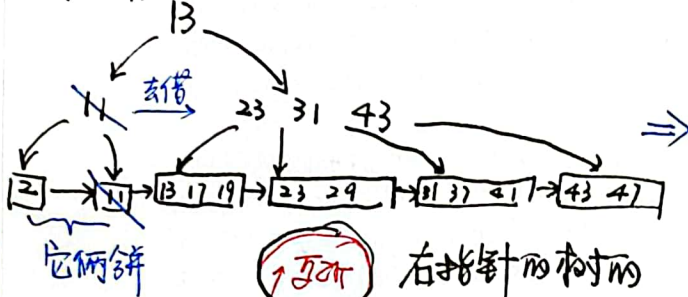
内节点

反对、重新分配

与右兄弟合并, 则从父节点中  
删除指向左兄弟的指针和键  
与右兄弟合并, 则从父节点中  
删除指向右兄弟的指针和键

删除后判断父  
节点是否半满, 否  
则递归处理  
(若根直接删)

例: 删除 11



### (三) 查询执行与优化

SQL 查询

↓ SQL 解析器

关系代数表达式

↓ 查询优化器

— 查询优化

查询执行计划

↓ 查询执行器

— 查询执行

查询结果

### 第一部分: 查询执行

#### 1. 基本记号与约定

$N$ : 关系  $R$  的元组数

$M$ : 缓冲区中可用内存页数

$B$ : 每块最多存  $B$  个元组

$B(R)$ :  $R$  的块数  $\lceil N/B \rceil$

$T(R)$ : 关系  $R$  的元组数

$V(R, A)$ : 关系  $R$  的属性集  $A$  的不同值的个数

算例分析时的约定

输出结果时产生的 I/O 不计入算法的 I/O 代价

输出缓冲区不计入可用内存页数。

## 2. 查询执行时的各个算法

算法大名	算法小名	算法说明	I/O	对M的要求
外存排序	两趟路外归并排序	把R划分成 $\lfloor B(R)/M \rfloor$ 个归并段, 每次调入内存M块(一个归并段)内排序使段内有序后回磁盘, 调入下一个归并段的第一块后开始归并	$3B(R)$	$B(R) \leq M^2$ $(\frac{B(R)}{M} \leq M)$
	三路路外归并排序	$B(R) > M^2$ 时使用	$(2m-1)B(R)$ (m为路数)	
选择操作	基于扫描的选择算法	依次扫描每一块并记录, 符合条件的输出	$B(R)$	$M \geq 1$
不考重量的投选算法	基于哈希的选择算法	前提条件: 选择条件的形式是 $k=V$ 类R采用哈希文件组织形式即 $k$ 是R的哈希键 桶指 $hash(V)$ 桶无组结果在桶内 在这个桶里找关键字值等于V的元组	$\lceil \frac{B(R)}{V(R,K)} \rceil$ (近m)	$M \geq 1$
	基于索引的选择算法	前提条件: 选择条件的形式是 $k=V$ 或 $I \leq k \leq U$ , $k$ 上有索引 直接按索引上搜索	聚簇 $\lceil \frac{B(R)}{V(R,K)} \rceil$ 非聚簇 $\lceil \frac{T(R)}{V(R,K)} \rceil$	$M \geq 1$
考重操作 S(R)	一趟考重算法	把每一块读进来建内存查找结构 读进来的块看它在结构里, 不在就输去并加到查找结构里, 否则跳过	$B(R)$	$M-1 \geq B(S(R))$
	基于排序的考重算法	思想同外排, 归并时相同的只输出一个, 丢弃其他	$3B(R)$	$B(R) \leq M^2$
	基于哈希的考重算法	哈希分桶到 $(M-1)$ 个桶里, 逐桶Disk 对每个桶用一趟考重	$3B(R)$	$B(R) \leq (M-1)^2$ $(\frac{B(R)}{M-1} \leq M-1)$
集合差操作	一趟集合差算法	$R-S$ , $S$ 很小, 在 $(M-1)$ 页能装下 在这 $(M-1)$ 页构造S的内存查找结构, 对R逐块读入后逐条去结构中找, 找不到则输出	$B(R)+B(S)$	$B(S) \leq M-1$
	基于哈希的集合差算法	把R和S分别哈希到 $(M-1)$ 个桶里, 桶号相同的进行一趟集合差, 最后并起来	$3[B(R)+B(S)]$	$B(S) \leq (M-1)^2$
	基于排序的集合差算法	先创建归并段, 归并时去判断元组是否在R且不在S	$3[B(R)+B(S)]$	$B(R)+B(S) \leq M^2$

算法大名	算法小名	算法说明	I/O	对M的要求
连接操作的执行 $R(X,Y) \bowtie S(Y,Z)$	一趟连接算法 (若内存查找结构是哈希表, 则为哈希连接)	假设 $B(S) \leq B(R)$ , 以 $S, Y$ 在 $(M-1)$ 页中建立内存查找结构, 然后按 $R$ 的每一块每一页, 符合条件的就连接起来输出。	$B(R) + B(S)$	$B(S) \leq M-1$
	基于元组的嵌套循环连接	两层循环, 一条一条比较。 for $S$ 的每个元组 $s$ do for $R$ 的每个元组 $r$ do if ... then 输出 外关系 ← 内关系	$\frac{T(S)(T(R)+1)}{T(S)+T(S)T(R)}$	$M \geq 2$
	基于块的循环嵌套连接	for 外关系 $S$ 中的每 $M+1$ 块 do 建查找结构 for 内表 $R$ 中的每块 do * 读入缓冲区 for 每个元组 $r$ do for 查找结构中的与 $r$ 连接的元组 $s$ do 连接, 输出 假设 $B(S) \leq B(R)$	$B(S) + \frac{B(S)}{M-1} B(R)$	$M \geq 2$
	排序归并连接	创建归并段, 归并	$3[B(R) + B(S)]$	$B(R) + B(S) \leq M^2$
	Grace 哈希连接	先把 $R$ 和 $S$ 分别哈希到 $(M-1)$ 个桶里, 逐桶做一趟连接算法, 最后并起来	$3[B(R) + B(S)]$	$B(S) \leq (M-1)^2$
	索引连接	假设 $S$ 上有 $T$ 的索引。把 $R$ 逐块读入缓冲区, 逐条元组去和索引上看是否匹配连接。	果索引: $B(R) + \frac{T(R)B(S)}{V(S,Y)}$ 非索引: $B(R) + \frac{T(R)T(S)}{V(S,Y)}$	$M \geq 2$

### 3. 查询计划的执行方法

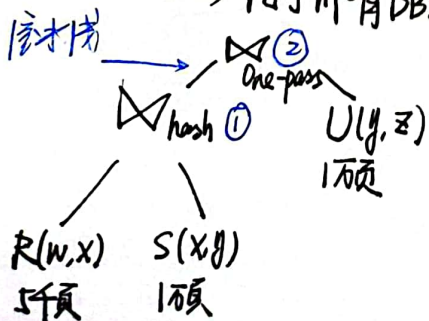
(1) 物化执行: 自底向上执行, 中间操作的结果写入临时关系件。

→ 缺点: 物化临时关系件增加了查询的代价; 用户获得查询结果的时间延迟 ↑

(2) 流水线执行/火山模型: 一个操作的结果直接传给流水线下一个操作。

→ 几乎所有DBMS都采用该方法

执行时 ① 哈希的过物用101页, 原输入缓冲, 100页分桶, 逐桶连接时用50页给R建查找结构, 1页做输入缓冲, 剩下50页做输出缓冲



$M=101$   
 $R, S, U$  上均无索引且无序  
 $B(RMS) \leq 50$

I/O  $B(U)=10000$  ② 一趟连接, 用  $B(RMS)+1$  页内存 18

I/O  $B(R)+B(S) = 45000$

如何实现流水线? 用迭代器模型。

## 第二部分: 查询优化

### 1. 基本概念

- (1) 查询优化的两个阶段
- 逻辑查询优化: 生成更高效的逻辑查询计划 关联表达式
  - 物理查询优化: 生成优化的物理查询计划 带有"如何执行"的关联表达式
- ↓
- 实际上DBMS并不区分

### 2. 逻辑查询优化

(1) 基于代价的查询优化 (所有DBMS采用的找代价最低的计划的方法)

- 计划枚举: 生成各种查询计划
- 代价计算
- 关联表达式分解
- 连接顺序的选择

### (2) 等价关系代数表达式

① 既满足交换律又满足结合律的关系代数操作:  $\times \cap \cup$

(注意: 选择只满足交换律, 不满足结合律)

② 关于选择等价交换与选择下推

可以尽早过滤掉与结果无关的元组

$$\langle 1 \rangle \sigma_{\theta_1 \wedge \theta_2}(R) = \sigma_{\theta_1}(\sigma_{\theta_2}(R)) = \sigma_{\theta_2}(\sigma_{\theta_1}(R))$$

原则:  $\langle 1 \rangle$  选择度高的选择优先做

$$\sigma_{\theta_1 \vee \theta_2}(R) = \sigma_{\theta_1}(R) \cup \sigma_{\theta_2}(R)$$

$\langle 2 \rangle$  把复杂的选择操作分解再下推。

$$\langle 2 \rangle \sigma_{\theta}(R \cup S) = \sigma_{\theta}(R) \cup \sigma_{\theta}(S)$$

$\langle 3 \rangle$  当可以向各个分支下推时更优先:

$$\langle 3 \rangle \sigma_{\theta}(R - S) = \sigma_{\theta}(R) - S = \sigma_{\theta}(R) - \sigma_{\theta}(S)$$

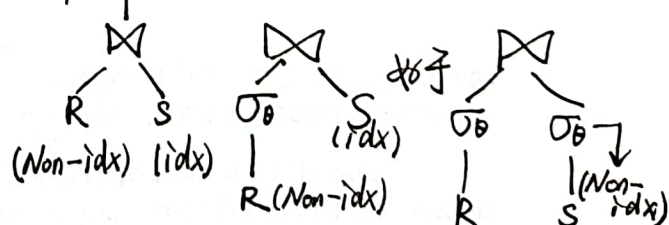
(I) 索引会影响选择下推的效率

$$\langle 4 \rangle \sigma_{\theta}(R \cap S) = \sigma_{\theta}(R) \cap S = R \cap \sigma_{\theta}(S) = \sigma_{\theta}(R) \cap \sigma_{\theta}(S)$$

(II) 选择操作结果上无索引

$$\langle 5 \rangle \sigma_{\theta}(R \times S) = R \times_{\theta} S, \text{ 若 } S \text{ 没有 } \theta \text{ 中的属性, 则 } \sigma_{\theta}(R \times S) = \sigma_{\theta}(R) \times S$$

例:



$$\langle 6 \rangle \sigma_{\theta_1}(R \times_{\theta_2} S) = R \times_{\theta_1 \wedge \theta_2} S, \text{ 若 } S \text{ 没有 } \theta_1 \text{ 中... 则 } \sigma_{\theta_1}(R \times_{\theta_2} S) = \sigma_{\theta_1}(R) \times_{\theta_2} S$$

$\langle 7 \rangle$  如果R和S均包含 $\theta$ 中使用的所有属性

$$\sigma_{\theta}(R \times S) = \sigma_{\theta}(R) \times S = R \times \sigma_{\theta}(S) = \sigma_{\theta}(R) \times \sigma_{\theta}(S)$$

③ 关于投影等价交换与投影下推

降低元组的大小

$$\langle 1 \rangle \pi_{L_1}(\pi_{L_2}(R)) = \pi_{L_1}(R) \quad (L_1 \subseteq L_2)$$

$$\langle 2 \rangle \pi_L(\sigma_{\theta}(R)) = \pi_L(\sigma_{\theta}(\pi_M(R)))$$

有些情况下, 投影下推会降低查询优化的机会: 如投影操作的结果上无索引。

M中既包含L中的属性又包含 $\theta$ 中使用的属性

$$\langle 3 \rangle \pi_L(R \cup S) = \pi_L(R) \cup \pi_L(S)$$

Pay attention:  $\pi_L(R \cap S) \neq \pi_L(R) \cap \pi_L(S)$   
 $\pi_L(R - S) \neq \pi_L(R) - \pi_L(S)$

<4>  $\pi_L(R \times S) = \pi_L(\pi_M(R) \times \pi_N(S))$   
 $M = L \cap R, N = S \cap L$

<5>  $\pi_L(R \bowtie S) = \pi_L(\pi_M(R) \bowtie \pi_N(S))$   
 $M = L \cap R$  以及选择属性(R中)  
 $N = S \cap L$  以及选择属性(S中)

(3) 逻辑查询计划的代价核算

- ① 逻辑查询计划的代价用执行计划过程中产生的中间结果的元组数来度量。
- ② 物理查询计划的代价比逻辑查询计划的代价更可靠。  
 DBMS实际上直接枚举物理查询计划，不区分与优化。

如何度量 cost?

(4) 基数估计

- ① 目的: 估计查询结果的元组数
- ② 要求: 准确、易计算、逻辑一致
- ③ 数据库系统中记录着基数估计所需的DB统计信息
- ④ 基本假设: 关系中每个属性取值均服从均匀分布  
 关系的所有属性相互独立。

单调性: 输入 ↑ 结果估计值 ↑  
 顺序无关性

$T(R)$ : 关系R的元组数  
 $V(R, A)$ : 关系R的属性A的不同值的个数

(5) 各种操作的基数估计:

<1>  $T(R \times S) = T(R)T(S)$

<2>  $T(\pi_L(R)) = \begin{cases} \text{不含去重 } T(R) \\ \text{含去重 } V(R, L) \end{cases}$

<3>  $S = \sigma_{A=c}(R)$

$T(S) = T(R) / V(R, A)$

$S = \sigma_{A \neq c}(R)$

$T(S) = T(R)$  或  $T(S) = T(R) - \frac{T(R)}{V(R, A)}$

$S = \sigma_{A > c}(R)$  或  $S = \sigma_{A < c}(R)$

$T(S) = T(R) / 3$  或  $T(S) = T(R) / 2$   
 若知道A的最佳可用M可整理

(I) <4>  $S = \sigma_{\theta_1 \wedge \theta_2}(R) = \sigma_{\theta_1}(\sigma_{\theta_2}(R)) = \sigma_{\theta_2}(\sigma_{\theta_1}(R))$

$T(S) = T(R) f_1 f_2$ , 其中  $f_i = \begin{cases} 1/V(R, A) & \theta_i: A = C \\ 1 - 1/V(R, A) & \theta_i: A \neq C \\ 1/3 & A < C \text{ or } A > C \end{cases}$

(II)  $S = \sigma_{\theta_1 \vee \theta_2}(R) = \sigma_{\theta_1}(R) \cup \sigma_{\theta_2}(R) = R - \sigma_{\neg \theta_1 \wedge \neg \theta_2}(R)$

$T(S) = T(R) - T(R)(1-f_1)(1-f_2)$

(III)  $S = \sigma_{\neg \theta}(R) = R - \sigma_{\theta}(R)$

$T(S) = T(R) - T(\sigma_{\theta}(R))$

<6> 二路自然连接:  $R \bowtie S$

基本假设: 连接属性值集合包含属性: 设连接属性为k, 若  $V(R, k) \leq V(S, k)$ , 则  $R.k \subseteq S.k$   
 非连接属性值集合保留假设: 对R中非连接属性A, 有  $V(R \bowtie S, A) = V(R, A)$

(I) R和S只有一个连接属性Y:  $T(R \bowtie S) = \frac{T(R)T(S)}{\max(V(R, Y), V(S, Y))}$

(II) R和S有两个连接属性X, Y:  $T(R \bowtie S) = \frac{T(R)T(S)}{\max(V(R, X), V(S, X)) \max(V(R, Y), V(S, Y))}$

(III) 推广到n个:  $\frac{T(R)T(S)}{\max(V(R, X), V(S, X)) \dots \max(V(R, Y), V(S, Y))}$

<7> 集合操作的基数估计

原因:

$T(R \cup S) = \frac{1}{2} (\max(T(R), T(S)) + T(R) + T(S))$

$\max(T(R), T(S)) \leq T(R \cup S) \leq T(R) + T(S)$

$T(R - S) = T(R) - \frac{1}{2} T(S)$

$T(R) - T(S) \leq T(R - S) \leq T(R)$

$T(R \cap S) = \min(T(R), T(S)) / 2$

$0 \leq T(R \cap S) \leq \min(T(R), T(S))$

$T(R \cap S) = T(R \bowtie S)$

$R \cap S = R \bowtie S$

<8> 差操作的基数估计

$T(S(R)) = (T(R) + 1) / 2 \approx T(R) / 2$

$1 \leq T(S(R)) = (T(R) + 1) / 2 \approx T(R) / 2$

$T(S(R)) = \min((T(R) + 1) / 2, V(R, A_1) V(R, A_2) \dots V(R, A_n))$   $1 \leq T(S(R)) \leq V(R, A_1) V(R, A_2) \dots V(R, A_n)$

(5) 用直方图近似数据分布

- 等宽: 属性值各区间的宽度相同, 各区间内属性值的出现次数不同
- 等高: ... 不同 ... 频率相同
- 压缩: 单独记录少数高频出现的属性值, 用等宽或等宽近似其他属性值。

等高更常用 - 包含高频属性值的区间窄, 对低频属性值的频率估计更准确  
 - 区间边界由数据确定, 不用用户指定。

(6) 连接顺序的优化

- ① RMS中 R是右关系, S是左关系
- 一趟连接: 左 → 构建关系; 右 → 探测关系
  - 嵌套循环连接: 左 → 外关系; 右 → 内关系
  - 索引连接: 右 → 索引关系

② 连接树: 左深连接树, 右深连接树, 任意树

- 只考虑这个 - 左深连接树比全部连接树少很多
- 通常更好 (形成完全通路的流水串)
- 对内存的需求低
- 不需要多次构建中间关系

动态规划例题, 课件

### 3. 物理查询优化 (选执行算法+执行模型) (查询执行)

(1) 确定选择操作的物理执行算法: 本质是确定最高效的存取路径.

#### ① 索引扫描+过滤

$\sigma_{\theta_1 \wedge \theta_2}(R)$

<1> 什么时候能用: 索引支持 $\theta_1$ 上的查找

<2> 怎么用? 根据 $\theta_1$ 索引取元组ID, 取元组, 用 $\theta_2$ 过滤. (存取路径)

<3> 什么时候能用: 索引是聚簇索引或非聚簇索引但表中 $\theta_2$ 的元组很少

#### ② 索引合并

$\sigma_{\theta_1 \wedge \theta_2}(R)$

<1> 什么时候能用: 索引分别支持 $\theta_1$ 和 $\theta_2$ 上的查找

<2> 怎么用(存取路径): 两个索引上分别找出 $\theta_1$ 和 $\theta_2$ 中的元组ID, 求交得所有满足的元组ID集合, 取元组

<3> 什么时候能用:  $I_1$ 和 $I_2$ 都是非聚簇, 满足条件 $\theta_1$ 和 $\theta_2$ 的都多, 但同时满足的少

#### ③ 反用索引

$\sigma_L(R)$   $L$ 为该选择操作的反选择操作所涉及到的 $R$ 的 $L$ 属性的集合

<1> ...: 索引支持 $\theta$ 上的查找, 是覆盖 $L$ 的覆盖索引

<2> ...: 在索引上查询是 $L$ 中的索引项后直接取出该元组向 $L$ 的投影

<3> ...: 关于 $R$ 上的聚簇索引不支持选择 $\theta$ , 索引是覆盖索引

#### ④ 顺序扫描 $\sigma_L(R)$

<1> 什么时候用: 没有索引支持 $\theta$ 上的查找

### (2) 确定连接操作的执行算法

适用条件

- 一趟连接: 左关系可以全读入缓存中的可用页面

- 索引连接: 左关系较小, 右关系在连接属性上建有索引

- 排序归并连接: 至少有一个关系已按连接属性排序; 多个关系在相同连接属性上做多路连接也适合. eg  $R(a,b) \bowtie S(a,c) \bowtie T(a,d)$

- 哈希连接: 上面三种都不适用时

- 嵌套循环连接: 缓冲区的面数特别少时 (要求其中最少)

### (四) 并发控制与故障恢复

核心概念: 事务(transaction): 数据库上执行的一个或多个操作构成的序列

事务的ACID性质:  $\sqrt{\text{原子性, 一致性, 隔离性, 持久性}}$   $\rightarrow$  并发控制体系

程序执行

故障恢复



# 第一部分：并发控制

1. 调度：一个或多个事务的重新操作序列

(1) 串行调度：一个调度中不同事务的操作没有交叉

(2) 调度的正确性：单独执行每个事务都会将DB从一种一致状态变为另一种一致状态

- 任意串行调度都能保证DB的一致性
- 非串行调度可能会破坏数据库的一致性

- 脏写：事务T<sub>2</sub>提交前，T<sub>1</sub>写过的记录又被T<sub>2</sub>写
- 脏读：事务T<sub>2</sub>提交前，T<sub>1</sub>写过的记录被T<sub>2</sub>读取
- 不可重复读：一个事务在两次读之间没有修改，但是读取的值不同。
- 幻读

(3) 等价调度：如果两个调度在VDB实例上的效果都相同，则这两个调度等价

(4) 可串行化调度：一个调度等价于一个串行调度，则它是可串行化调度

- ↓ 无脏写、脏读、不可重复读、幻读
- 缺点：并发度低，降低隔离级别可以提高并发度

(5) 有隔离级别：从上至下隔离级别增大

名称	说明	脏写	不可重复读	幻读
读未提交 (Read Uncommitted)	未提交事务的修改对其他事务可见	✓	✓	✓
读提交 (Read Committed)	只有已提交事务的修改对其他事务可见	✗	✓	✓
可重复读 (Repeatable Read)	若事务不写X，则从开始到该事务结束的所有数据都于启动时所见	✗	✗	✓
可串行化 (Serializable)		✗	✗	✗

(6) 支持可串行化隔离级别的DBMS实施的都是冲突可串行化

① 冲突：两个操作属于不同事务，涉及相同对象，至少有一个操作是写

- 写-写冲突 → 脏写
- 写-读冲突 → 脏读
- 读-写冲突 → 不可重复读

② 冲突等价：如果两个调度涉及相同事务的相同操作，且每一个对冲突的操作在两个调度中的次序都相同，则这两个调度等价。

③ 冲突可串行化调度：如果一个调度冲突等价于一个串行调度，则该调度是冲突可串行化调度。

- 判断是否
  - 定义
  - 调用次序及非冲突操作看能不能转换成串行调度
  - 优先图：无环 ⇔ 冲突可串行化

## 2. 并发控制协议 — 基于锁的

(1) 锁：控制访问同一数据的一种手段

① 事务只有获得了对象的锁才有访问权限，如果事务请求锁但没有得到则等待直到得到

③事务完成对某个对象的操作后必须释放其锁

(2) 锁的类型

① 共享锁/S锁 — 读权限

② 互斥锁/X锁 — 读写权限

13) 锁的相容性:  $S \downarrow X$ ,  $S \downarrow S$ ,  $X \downarrow S$ ,  $X \downarrow X$ , 只有S是允许的  
 ↓  
 读冲突    写读冲突    写写冲突

例: 使用共享锁和互斥锁的并发事务执行

T <sub>1</sub>	T <sub>2</sub>	调度器的动作
X-LOCK(A)		A上的互斥锁授予T <sub>1</sub>
r(A)		
	X-LOCK(A)	T <sub>2</sub> 申请A上的互斥锁, 不给
w(A)		
UNLOCK(A)	r(A)	释放A上的互斥锁
	w(A)	A上的互斥锁授予T <sub>2</sub>
S-LOCK(A)		T <sub>1</sub> 申请A上的共享锁, 不给
r(A)		释放A上的互斥锁
UNLOCK(A)		A上的共享锁授予T <sub>1</sub>
		释放A上的共享锁

弊端: 不是冲突可串行化调度, 怎么办?

(4) 两阶段锁协议(2PL)

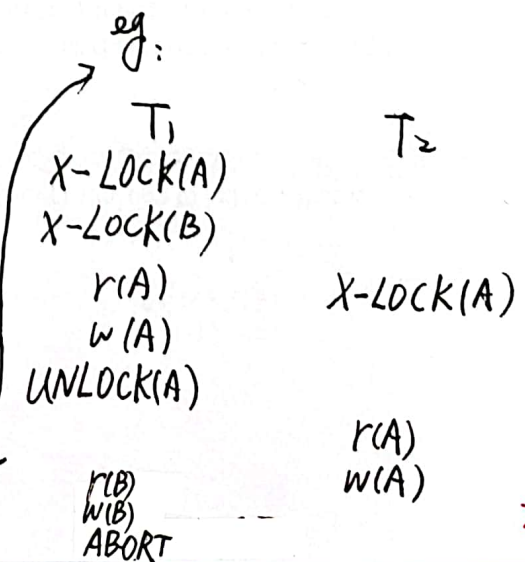
— 增长阶段: 事务向锁管理器请求需要的锁

— 收缩阶段: 事务释放它获得的锁, 但不能再请求加锁

核心: 解锁之后不能再加锁。

例: 基于2PL的并发事务执行

T <sub>1</sub>	T <sub>2</sub>	调度器的动作
X-LOCK(A)		A上的锁(互斥)授予T <sub>1</sub>
r(A)		
	X-LOCK(A)	T <sub>2</sub> 申请A上的互斥锁, 不给
w(A)		
r(A)		
UNLOCK(A)	r(A)	释放A上的锁
	w(A)	A上的互斥锁授予T <sub>2</sub>
	UNLOCK(A)	释放A上的互斥锁



2PL的优点: 能够保证可串行化  
 缺点: [ 2PL面临级联中止的问题  
 [ 2PL可能导致死锁

# (5) 解决级联冲突 —— 两阶段锁协议 (2PL)

① 严格调度: 一个调度中任意事务在结束前由它写入的DB对象的值没有被其他事务读过或修改过, 则该调度称为严格调度。

↳ 不会生级联冲突

## ② SS2PL协议

└ 增长阶段: 与2PL增长阶段相同

└ 收缩阶段: 事务结束时, 释放它获得的全部的锁。

eg:

$T_1$	$T_2$
X-LOCK(A)	
X-LOCK(B)	
r(A)	X-LOCK(A)
w(A)	
r(B)	
w(B)	
UNLOCK(A)	
UNLOCK(B)	
ABORT	

③ SS2PL协议能保证生成严格的冲突可串行化调度

## (6) 死锁问题

① 如果一组事务中每个事务都在等待其他事务释放锁, 则这组事务形成死锁

### ② 死锁的处理

在给定时间内无事务执行, 认为发生

└ 死锁检测: 检测是否发生, 发生则解决 超时检测

└ 死锁预防: 预防死锁的发生 等待图

eg:

$T_1$	$T_2$
X-LOCK(A)	
r(A)	
	S-LOCK(B)
	r(B)
w(A)	S-LOCK(A)
X-LOCK(B)	

### ③ 等待图检测

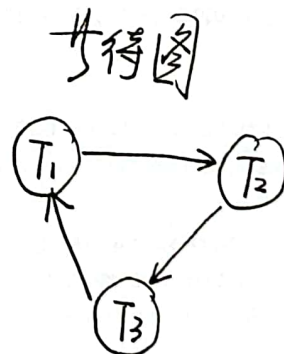
└ 图中的节点代表事务

└ 如果事务 $T_i$ 在等待 $T_j$ 释放锁, 则事务 $T_i$ 死锁 ⇔ 图中有环

└ 则有从 $T_i$ 到 $T_j$ 的有向边。

eg:

$T_1$	$T_2$	$T_3$
S-LOCK(A)	X-LOCK(B)	S-LOCK(C)
r(A)	w(B)	r(C)
S-LOCK(B)	X-LOCK(C)	X-LOCK(A)
(拒绝)	(拒绝)	(拒绝)



在等待图中的死锁解除: 从等待图中选择一个事务作牺牲, 中止该事务

- └ 事务的年龄/启动时间
- └ 事务的进展
- └ 事务持有锁的数量
- └ 需要级联中止的事务数量
- └ 事务作牺牲的“代价”, 防止饿死

#### ④ 死锁预防

- <1> 事务启动时, DBMS为事务分配一个唯一且固化的优先级, 开始越早优先级越高
- <2> Wait-Die 规则 —— "尊老爱幼"
  - 事务 $T_i$ 请求事务 $T_j$ 拥有的锁而无法获得时, 如果 $T_i$ 比 $T_j$ 优先级高, 则 $T_i$ 等待。

**尊老** 否则 $T_i$ 中止( $T_i$ 重启后, 其优先级保持不变) **爱幼**

#### <3> Wound-Wait 规则 —— "论资排辈"

事务 $T_i$ 请求事务 $T_j$ 拥有的锁而无法获得时, 如果 $T_i$ 比 $T_j$ 优先级高, 则在 $T_j$ 中止( $T_j$ 重启后, 其优先级保持不变) 否则 $T_i$ 等待

### 第二部分: 故障恢复

#### 1. 故障时的两类情况:

- 已提交的事务对DB的修改未全部持久化到磁盘 —— 破坏事务的持久性。—— Redo
- 已中止的事务对DB的修改已部分持久化的磁盘 —— 破坏事务的原子性。—— Undo

#### 2. DBMS如何运用undo和redo取决于DBMS如何管理 Buffer Pool

- 只undo
  - 只redo
  - undo+redo
  - 都不用
- 缓冲池策略

- (1) 对未提交事务所做的修改是否刷写到磁盘 ——  刷写 steal
- (2) 对是否强制事务在提交前必须将其所做的修改写回磁盘 ——  强制 force
- 不强制 no-force

#### (3) 四种组合

缓冲池效率高	Steal + Force undo	Steal + No-Force redo+undo
缓冲池效率低	No-Steal + Force I/O效率低	No-steal + No-Force redo I/O效率高

效率最差, 但无需做故障恢复, 实现简单, 缺点是要求缓冲池容量足够大存下来提交事务的修改

#### 3. 预写式日志 (WAL)

(1) DBMS维护一个日志文件, 用于记录事务对数据库的修改。

- (2) WAL协议: 事务启动时 <tid, BEGIN> 事务ID
- 事务提交时 <tid, COMMIT> → 旧值 → 新值
- 事务修改对象A时 <tid, oid, before, after> A的ID

- (3) 事务的分类
- 已提交事务: 既有  $\langle T, \text{BEGIN} \rangle$ , 又有  $\langle T, \text{COMMIT} \rangle$
  - 不完整事务: 只有  $\langle T, \text{BEGIN} \rangle$ , 而没有  $\langle T, \text{COMMIT} \rangle$
  - 已中止事务: 既有  $\langle T, \text{BEGIN} \rangle$ , 又有  $\langle T, \text{ABORT} \rangle$

→ 在恢复的过程中既不需要 Redo, 也不需要 Undo

→ 如果恢复完一个事务, 但没完全取消, 追加一个  $\langle T, \text{ABORT} \rangle$  到日志

#### (4) 三种WAL协议

① Undo Logging: WAL + STEAL + FORCE

事务提交的真正标志

⟨1⟩ 顺序  $\langle T_i, \text{BEGIN} \rangle$  写入磁盘 →  $\langle T_i, A, \text{before}, \text{after} \rangle$  写入磁盘 → A 写入磁盘 →  $\langle T_i, \text{COMMIT} \rangle$  写入磁盘

⟨2⟩ 恢复方法: 从后向前扫描日志

遇到  $\langle T, \text{COMMIT} \rangle$  和  $\langle T, \text{ABORT} \rangle$  都标记为已提交/中止事务, 无需 Undo

遇到  $\langle T, A, \text{before}, \text{after} \rangle$  若 T 不是完整事务, 则把 Disk 上 A 值恢复为 before

遇到  $\langle T, \text{Begin} \rangle$  时, 若 T 不完整, 则追加  $\langle T, \text{ABORT} \rangle$

② Redo Logging: WAL + NO-STEAL + NO-FORCE

⟨1⟩ 顺序 ...  $\langle T_i, \text{COMMIT} \rangle$  写入磁盘 → A 写入磁盘

⟨2⟩ 恢复方法: 从前向后扫描两遍日志

第一遍:  $\langle T, \text{COMMIT} \rangle$  和  $\langle T, \text{ABORT} \rangle$  标记为需要 Redo 和无需 Redo

第二遍:  $\langle T, A, \text{before}, \text{after} \rangle$  若 T 已提交, 则覆盖 A 为 after

不完整的事务追加  $\langle T, \text{ABORT} \rangle$

③ Redo + Undo Logging: WAL + STEAL + NO-FORCE (几乎所有DBMS都采用的策略)

⟨1⟩ 顺序: ...  $\langle T_i, \text{COMMIT} \rangle$   $\xrightarrow{\text{NO-FORCE}}$  A 写入磁盘

或 A 写入磁盘,  $\xrightarrow{\text{STEAL}}$   $\langle T_i, \text{COMMIT} \rangle$

⟨2⟩ 恢复方法: Redo 阶段和 Undo 阶段同步

#### 4. 检查点

(1) WAL 的问题: 日志永远在变大, 故障恢复时要扫描日志, 恢复时间越来越长

(2) DBMS 定期设置检查点

① 将日志刷写到磁盘

② 根据预定的策略, 将脏页写到磁盘

③ 故障恢复时, 只扫描到最新检查点的日志

(3) 模糊检查点

① 检查点开始: 向日志中写入  $\langle \text{BEGIN CHECKPOINT}(T_1, \dots, T_n) \rangle$

其中  $T_i$  是检查点开始时的活跃事务  
→ 尚未提交或中止的事务.

② 检查点中间: 根据缓冲池策略, 把脏页写到 Disk.

┌ 若采用 STEAL, 则将全部脏页写到磁盘  
└ 否则, 只将已提交事务所做的修改写到磁盘

③ 检查点结束: 向日志中写入  $\langle \text{END CHECKPOINT} \rangle$ , 并将日志刷写到磁盘

┌ 若 NO-FORCE, 则写完就可结束检查点

└ 否则, 要等  $T_i$  全部提交后才能结束.

(4) STEAL+NO-FORCE 下的故障恢复

① Redo 阶段

┌ 从前向后扫  
└ 从最近的完整检查点的  $\langle \text{BEGIN CHECKPOINT}(T_1, \dots, T_n) \rangle$  开始扫

而 Redo 的  $T_i \in$  这个集合

② Undo 阶段

┌ 从后向前扫  
└ 向日志中最近的完整检查点为  $\langle \text{BEGIN CHECKPOINT}(T_1, T_2, \dots, T_n) \rangle$   
...  
 $\langle \text{END CHECKPOINT} \rangle$

扫描到  $T_1, \dots, T_n$  中的最早的事务  $T_i$  的日志记录  $\langle T_i, \text{Begin} \rangle$  为止

而 Undo 的  $T_i \in$  这个集合

## 后记

在复习和准备讲义的整个过程中，越深入去理解数据库课程的内容，就越发现它的博大精深和我知识的匮乏。作为我们数据科学与大数据技术的核心课程，它的重要性不言而喻，希望我的这份指南能够帮助大家更好地理解 and 掌握数据库的精髓与内涵。

课程的知识含量毋庸置疑非常大，但路虽难，行则将至；事虽难，做则必成。我相信通过努力复习，大家都能攻克这座难关。

最后，祝愿志在 90+ 的同学都能如愿以偿，祝愿 60 分万岁的同学都能顺利通过。

哈工大计算学部金牌讲师团 杨明达

2023 年 6 月 10 日

## 致谢

感谢张宇非同学和崔志阳同学指出讲义中的错误。

感谢所有 2103501 班同学对我的大力支持。