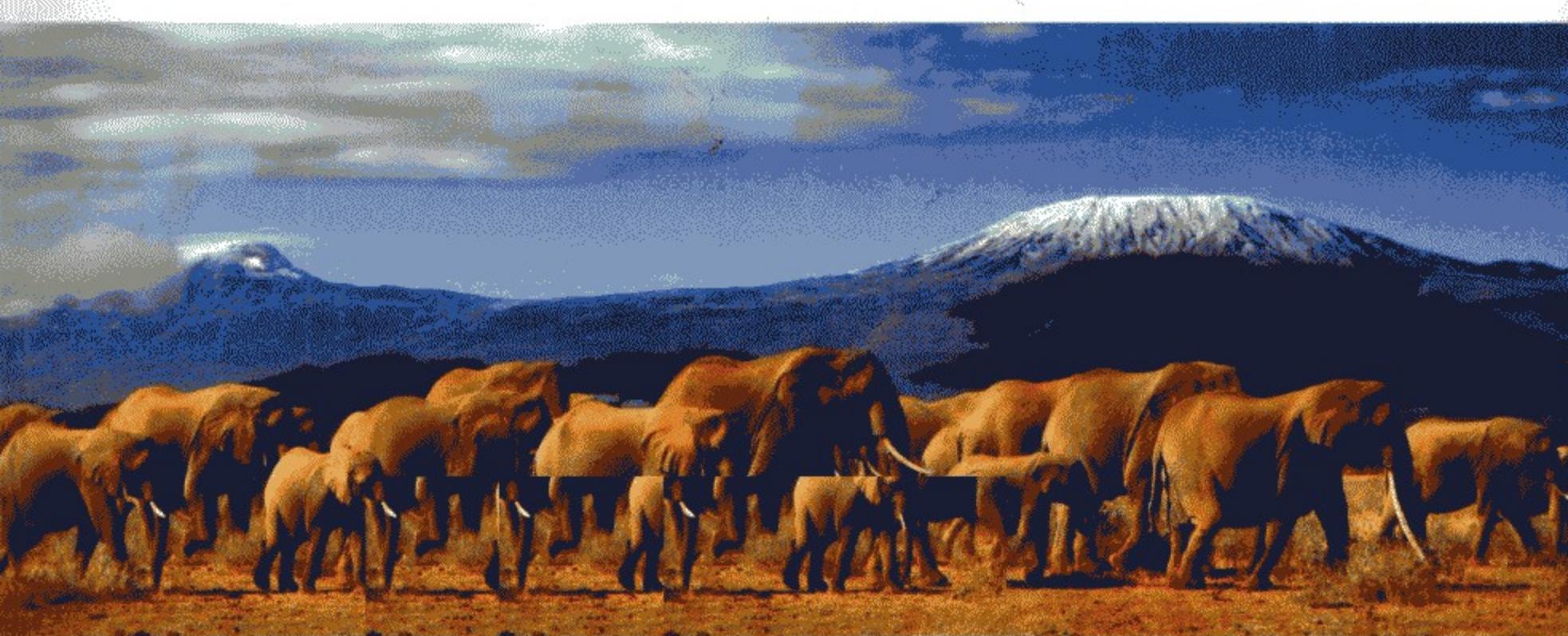


系统、深入地剖析开源数据库PostgreSQL的实现机制和工作原理
作者跟踪、研究PostgreSQL源代码十年的呕心之作

HZ BOOKS
华章科技

PostgreSQL

数据库内核分析



彭智勇 彭煜玮 ◎编著



机械工业出版社
China Machine Press

PostgreSQL 数据库内核分析

随着信息化进程的加快，数据库管理系统的新的应用领域不断出现，现有的数据库管理系统产品已经不能满足日益出现的新需求。很多用户开始采用开源数据库管理系统或者在其上定制、开发满足自己需求的数据库管理系统。开源数据库管理系统已经脱离纯研究层面，吸引大批技术开发人员、研究者以及其他技术爱好者的注意。

本书是全面分析开源数据库管理系统PostgreSQL的开山之作，引导读者轻松、透彻地理解数据库管理系统的内部运行机制，揭示实际数据库管理系统的运行过程。读者更可以通过阅读本书深入地理解、认识数据库或者验证数据库的新技术，进而有助于读者基于 PostgreSQL定制数据库系统、开发数据库内核或数据库管理系统。

本书的特点

- 本书按照PostgreSQL的体系结构，从存储、索引、查询编译、查询执行、并发控制以及安全几个方面切入，全面介绍PostgreSQL各种机制的运行原理。
- 为了清晰阐述复杂的运行机理，书中穿插了大量的原理图、程序流程图进行辅助讲解，使数据库的运行过程一目了然。
- 本书没有单纯地介绍数据库原理或逐行分析源代码，而是从数据库设计者的角度，阐释数据库内部各个模块之间如何配合实现各种功能。

本书对PostgreSQL的分析工作基于PostgreSQL 8.4.1，该版本的源代码可以从`ftp://ftp-archives.postgresql.org/pub/source/v8.4.1`下载。

客服热线:(010) 88378991, 88361066
购书热线:(010) 68326294, 88379649, 68995259
投稿热线:(010) 88379604
读者信箱:hzjsj@hzbook.com

华章网站 <http://www.hzbook.com>

网上购书：www.china-pub.com



上架指导：计算机/数据库

ISBN 978-7-111-35905-0



9 787111 359050

定价：79.00元

数据库技术

PostgreSQL

数据库内核分析

彭智勇 彭煜玮 ◎编著



机械工业出版社
China Machine Press

第 1 章

PostgreSQL 系统概述

1.1 PostgreSQL 简介及发展历程

PostgreSQL 是一种先进的对象 - 关系数据库管理系统（ORDBMS），它不仅支持关系数据库的各种功能，而且还具备类、继承等对象数据库的特征。它是目前功能最强大、特性最丰富和结构最复杂的开源数据库管理系统，其中有些特性甚至连商业数据库都不具备。这个起源于加州大学伯克利分校（UCB）的数据库研究计划，现在已经衍生成一项国际开发项目，并且拥有广泛的用户群。PostgreSQL 主要运行在 Unix 和 Linux 操作系统上（从 8.0 版本开始推出了 Windows 平台上的版本），并且免费开放源代码，用户可以在其官方网站 www.postgresql.org 上下载各种安装程序和文档。

PostgreSQL 的发展历程见证了数据库理论和技术的发展历程，PostgreSQL 由 UCB 计算机科学教授 Michael Stonebraker 于 1986 年创建。在此之前，Stonebraker 教授领导了关系数据库 Ingres 研究项目（Ingres 项目的源代码可以从 www.ingres.com 免费获取），1982 年他离开 UCB 并将 Ingres 商业化使之成为 Relational Technologies（RT）公司的一个产品。而后该公司被 Computer Associations（CA）公司收购。2004 年，CA 在开源许可下发布了 Ingres release 3，并继续开发销售 Ingres。而 Stonebraker 教授在返回 UCB 后开始了一项 post-Ingres 计划，该计划致力于解决基于关系模型的数据库管理系统产品的局限性，这即是 Postgres（但还不是现在的 PostgreSQL）的开端。

从 1986 年开始，Stonebraker 教授发表了一系列论文，引入对象关系理念，探讨了新的数据库的结构设计和扩展设计。1988 年，他提出了 Postgres 的第一个原型设计，1989 年 6 月发布了版本 1，1990 年 6 月发布了带有重写后的规则系统的版本 2。1991 年发布了版本 3，在版本 3 里改进了规则系统，增加了对多种存储系统支持的能力，并且改进了查询引擎。1993 年，Postgres 用户开始剧增，并且特性需求急剧增加。在做了一些代码清理后发布了版本 4，之后 Postgres 项目正式终止。随后，Stonebraker 再次创业，成立 Illustra 公司提供对 Postgres 的商业支持，Illustra 在 1997 年被 Informix 收购，而 Stonebraker 成为 Informix 的 CTO，Informix 由于财务问题在 2001 年被 IBM 收购。

尽管 Postgres 计划终止了，但 BSD 许可证（UCB 在其下发行的 Postgres）却使开放源代码开发者获得副本并进一步开发系统。1994 年，两个 UCB 的研究生，Andrew Yu 和 Jolly Chen，增加了一

个 SQL 语言解释器来替代早先的基于 Ingres 的 QUEL 系统，建立了 Postgres95。代码随后被发行到互联网上。1996 年，该计划被重新命名为 PostgreSQL，以反映数据库的新查询语言 SQL，来自世界各地的数据库开发者和志愿者通过互联网协作起来，发行了 PostgreSQL 的第一个版本 6.0，并且一直维护着这套软件。自此以后，PostgreSQL 开始持续稳定地发布新版本，在新版本中有很多改进。2005 年 1 月 19 日，版本 8.0 发行，从这个版本开始，PostgreSQL 以原生的方式（即不需要模拟中间层的支持）开始支持 Windows 操作系统。

从 Michael Stonebraker 教授开始，全世界无数著名的数据库专家和优秀的黑客为 PostgreSQL 的发展做出了杰出的贡献，使 PostgreSQL 项目充满活力，不断向前发展，并使得 PostgreSQL 成为目前最好的开源数据库管理系统之一。PostgreSQL 的发展历程如图 1-1 所示。

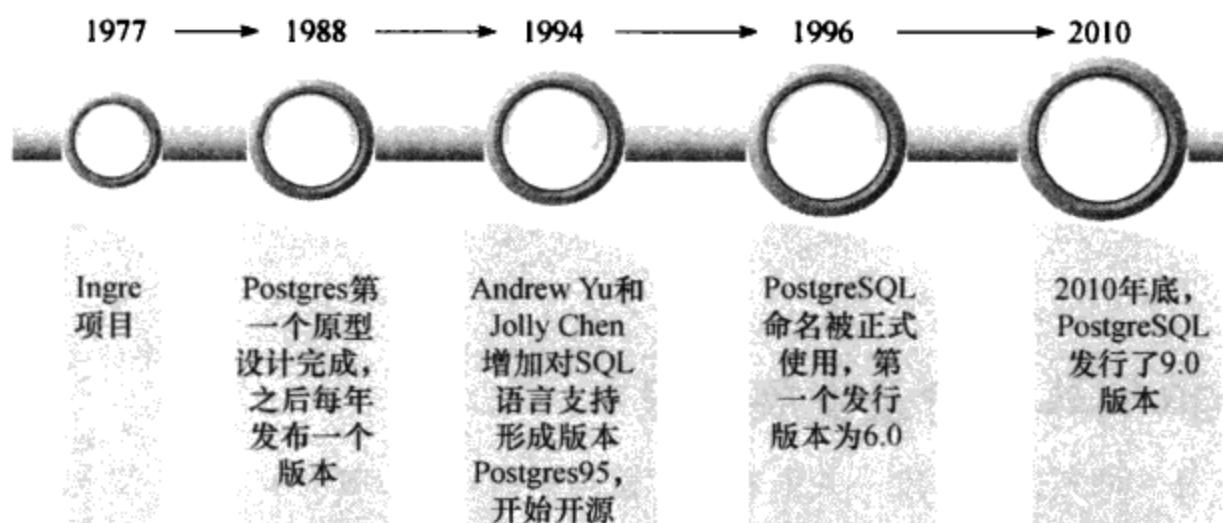


图 1-1 PostgreSQL 发展路线图

1.2 PostgreSQL 的特性

PostgreSQL 是一种几乎可以运行在各种平台上的免费的开放源码的对象关系数据库管理系统，拥有与企业级数据库相媲美的特性，如完善的 SQL 标准支持、多版本并发控制、时间点恢复、表空间机制、异步复制、嵌套事务、在线/热备份、一个复杂的查询优化器、预写日志容错技术。它支持国际字符集、多字节字符编码、Unicode，并且对格式化、排序、大小写敏感提供本地化支持。PostgreSQL 在管理大数据量方面有良好的可扩展性，对并发用户管理具有自适应性。现在已经出现具有管理超过 4 万亿字节数据能力的实用版本产品。

- **开放特性：**PostgreSQL 内置了丰富的数据类型，如任意精度的数值、无限制长度的文本、几何图元、IP 地址、数组等；同时还允许用户定义基于正规 SQL 类型的新类型，让数据库自身理解复杂数据，自定义类型中还可以包含继承关系。用户可以为数据库内几乎所有的对象定义新的类型，如索引、操作符（可重载现有操作符）、聚集函数、数据域、数据类型转换、会话（编码转换）等。
- **可编程性：**PostgreSQL 同样拥有大量的编程接口供用户开发使用，如 ODBC、JDBC（Java）、Libpq（C/C++）等。
- **可定制性：**PostgreSQL 拥有广泛的编程语言支持来实现函数功能，包括内置的 PL/pgSQL 过程语言，PL/Perl、plPHP、PL/Python、PL/Ruby、PL/Tcl 等脚本语言，以及 Java、C/C++ 等高级编程语言。函数的输出是一系列行类型的集合，可以在查询中当做表来使用，函数

也可以被定义成以创建者或者调用者的身份运行。在其他的数据库产品中，函数也会被称为“存储过程”。

- 索引手段：用户可以自定义索引方法或者使用内置 B-Tree 索引、Hash 表索引、GiST 索引、GIN 索引。GiST 索引不是某种特定的索引类型，而是一种通用索引基础结构，可以在这种结构上实现很多不同的索引策略。PostgreSQL 同时还支持如下功能：反向索引检索、表达式索引、部分索引、位图索引扫描。
- 多种身份认证方式：PostgreSQL 中可以使用数据库用户/角色、操作系统、PAM、Kerberos 等方式，根据主机配置文件（`pg_hba.conf`）中的设置执行对应的身份认证。

1.3 PostgreSQL 的应用

PostgreSQL 提供经济有效的、易于部署的复杂数据管理基础设施，在工业界得到广泛应用，在许多国际化大公司的应用中取得了良好的效果。

1. Sony 在线娱乐网站

Sony 在线娱乐网站（SOE）经营着很多知名游戏，因此其数据库需要快速处理大量的用户数据，其原有数据库系统为几十个 Oracle 9i RAC 集群。SOE 分析业务后发现，Oracle 数据库过于昂贵，且许可证限制过多，缺乏灵活性，最终 SOE 将其数据库系统转换为稳定、可扩展、高性能的使用 PostgreSQL 技术的 EnterpriseDB 数据库。

2. Hi5 社交网站

Hi5 从 2003 发展开始已经成为世界上最大的社交网站之一，有超过 200 个国家的八千万注册用户，每个月用户访问量超过 5600 万，是在 Alexa 全球排名前 20 的网站；Hi5 拥有一个运行在几百台服务器节点上的商业 PostgreSQL OLTP 集群，这也是世界上最大的 PostgreSQL 集群，运行着 Hi5 社交网络的所有服务，包括用户数据、配置信息、图片、评论等数据，并提供高效的检索。另外，FortiusOne、Florists' Transworld Delivery、NNT 等公司都选择廉价而高效的 PostgreSQL 作为其数据管理基础设施，取得了良好的应用效果。

PostgreSQL 经过多年的发展，得到了学术界和工业界的充分认可。获得 2008 Developer.com 编辑选择的数据库工具方向的年度产品，2000 年、2003 ~ 2006 年荣获 Linux Journal 杂志编辑评选的“最佳数据库”奖，2004 年获 ArsTechnica 最佳服务器应用奖，2002 年获 Linux New Media 杂志编辑评选的“最佳数据库”奖，1999 年获 Linux World 杂志评选的“最佳数据库”奖等多项荣誉。

1.4 PostgreSQL 代码结构

PostgreSQL 是遵照一个和 BSD 开源协议类似的协议 PostgreSQL License 发布。该许可证可以在 PostgreSQL 源代码的 `COPYRIGHT` 文件中找到。从该协议可以看到，开发人员或者商业组织只要遵循该协议，便可以自由地使用 PostgreSQL，可以完全控制这些第三方代码，在必要的时候可以加以修改或者二次开发。

PostgreSQL 源代码包含 3400 多个文件（截至 8.4.1 版本），主要程序由 C 语言编写，包括十几个大型模块，定义了几百个主要的数据结构和上万个函数。PostgreSQL 源代码结构清晰，每一个子目录都对应一个模块，其中主要目录（模块）及用途如下。

- **Bootstrap**: 用于支持 Bootstrap 运行模式，该模式主要用来创建初始的模板数据库。
- **Main**: 主程序模块，它负责将控制权转到 Postmaster 进程或 Postgres 进程。
- **Postmaster**: 监听用户请求的守护进程，并控制 Postgres 进程的启动和终止。
- **Libpq**: C/C++ 的库函数，处理与客户端间的通信，几乎所有的模块都依赖它。
- **Tcop**: Postgres 服务进程的主要处理部分，它调用 Parser、Optimizer、Executor、和 Commands 中的函数来执行客户端提交的查询。
- **Parser**: 编译器，将 SQL 查询转化为内部查询树。
- **Optimizer**: 优化器，根据查询树创建最优的查询路径和查询计划。
- **Executor**: 执行器，执行来自 Optimizer 的查询计划。
- **Commands**: 执行非计划查询的 SQL 命令，如创建表命令等。
- **Catalog**: 系统表操作，包含用于操作系统表的专用函数。
- **Storage**: 管理各种类型的存储系统（如磁盘、闪存等）。
- **Access**: 提供各种存取方法，支持堆、索引等对数据的存取。
- **Nodes**: 定义系统内部所用到的节点、链表等结构，以及处理这些结构的函数。
- **Utils**: 各种支持函数，如错误报告、各种初始化操作等。
- **Regex**: 正规表达式库及相关函数，用于支持正规表达式处理。
- **Rewrite**: 查询重写，根据规则系统对查询进行重写。
- **Initdb**: 初始化数据库集簇。
- **TSearch**: 全文检索。
- **Psql**: 数据库交互工具。
- **Port**: 平台兼容性处理相关的函数。

PostgreSQL 内部模块之间的调用关系如图 1-2 所示。

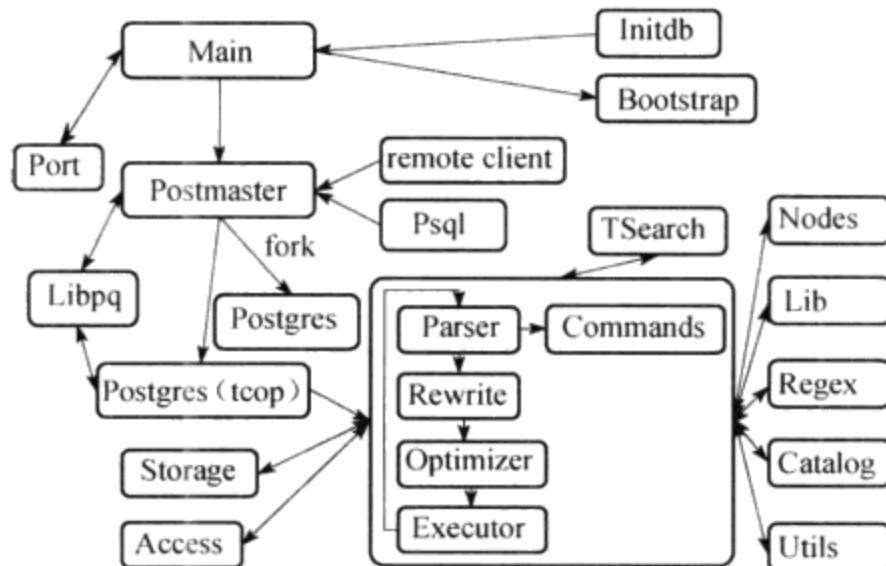


图 1-2 PostgreSQL 模块调用关系

1.5 安装 PostgreSQL

在开始使用 PostgreSQL 之前，必须先安装它。在 Windows 下安装 PostgreSQL 比较简单，按照安装程序的提示一步步进行即可，这里不再赘述。在 Linux 下，则需要从源代码进行安装，过程如下：

- 1) 从 PostgreSQL 官方网站下载最新源代码，通常是 .tar.gz 的压缩格式。
- 2) 通过 tar 或者图形化软件解压源码包，假设解压到 /usr/local/src/pgsrc 目录下。
- 3) 从源代码编译安装 PostgreSQL 需要用到 C 语言开发库（包括 gcc、ld 等），默认配置的情况下还需要 libreadline、zlib 和 xml2 这几个开发库的支持，在 Ubuntu 等 Linux 发行版下可以通过 apt-get 安装上述开发库。
- 4) 在源代码目录下运行 configure 脚本，默认情况下 PostgreSQL 将被安装在 /usr/local/pgsql 目录下。
- 5) 在源代码目录下分别执行 make 和 make install 命令即可从源代码编译 PostgreSQL 并将相关文

件拷贝到目标目录下。

6) 初始化数据库集簇:

```
cd/usr/local/pgsql
mkdir data (创建数据库文件的存放目录)
chown postgres:postgres data
su postgres
cd bin
./initdb --no-locale -D ../data
```

7) 启动数据库服务器:

```
./pg_ctl start -D ../data
```

8) 创建和访问数据库 test:

```
./createdb test
./psql test
```

1.6 PostgreSQL 数据库命令

数据库安装完成后，在 bin 目录下包含了 PostgreSQL 的工具命令。这些命令包括数据库初始化、创建、删除数据库，启动、停止数据库，连接数据库，备份、恢复数据库等丰富的工具。PostgreSQL 中各种命令的功能介绍如表 1-1 所示。

表 1-1 PostgreSQL 常用数据库命令

命令名	常用参数	功能介绍
initdb	-D 数据集簇位置 -E 新建数据库默认编码 -X 事务日志位置	创建一个用于存储数据库的 PostgreSQL 数据目录，并创建预定义的模板数据库 template0 和 template1，生成共享目录表 catalog；此程序通常只在安装 PostgreSQL 时运行一次。如 initdb -D ../data
createuser	-c 指定用户的最大连接数 -d 新用户可以创建数据库 -l 新用户具有登录权限 -U 连接的用户名（不是要创建的用户）	创建一个新的 PostgreSQL 的用户（和 SQL 语句：CREATE USER 相同），如 createuser -d -l -U postgres pguser
dropuser	-U 连接数据库的用户名	删除用户，如 dropuser pguser
createdb	-E 数据库的编码 -l 数据库 locale 设置 -O 数据库所有者	创建一个新的 PostgreSQL 的数据库（和 SQL 语句 CREATE DATABASE 相同），如 createdb test
dropdb	-i 删除之前提示确认操作 -U 连接数据库的用户名	删除数据库，如 dropdb test
pg_dump	-f 备份文件名称 -F 备份文件格式，custom、tar、plain text -Z 压缩等级 -a 只备份数据、不含模式 -b 备份包含大对象 -C 包含创建数据库命令	将 PostgreSQL 数据库导出到一个备份文件，如 pg_dump test >/home/pgsql/backup/1.bak
pg_dumpall	-f 备份文件名称 --lock-wait-timeout 有表锁时若超时则备份失败	将所有的 PostgreSQL 数据库导出到一个备份文件

(续)

命令名	常用参数	功能介绍
pg_restore	-d 连接数据库名称 -f 备份文件名称 -F 备份文件格式 (c, t) -a 只恢复数据，忽略模式 -C 创建目标数据库 -I 恢复有名称的索引	从一个由 pg_dump 或 pg_dumpall 程序导出的备份文件中恢复 PostgreSQL 数据库，如 pg_restore -d test /home/pgsql/backup/1.bak
vacuumdb	-a 清理所有数据库 -d 清理指定名称数据库	清理和分析一个 PostgreSQL 数据库，它是客户端程序 psql 环境下 SQL 语句 VACUUM 的 shell 脚本封装，二者功能完全相同
pg_ctl		启动、停止、重启 PostgreSQL 服务，如 pg_ctl -D .../data start, pg_ctl -D .../data stop, pg_ctl -D .../data restart
postgres		PostgreSQL 单用户模式的数据库服务
postmaster		PostgreSQL 多用户模式的数据库服务
psql	-c 执行单个命令并退出 -d 连接的数据库名字 -f 从文件中执行命令 -l 将文件中的命令当成一个任务执行 -h 数据库所在主机 IP -p 数据库监听端口号 -U 连接数据库用户名 -w 不需要密码 -W 强制使用密码	一个基于命令行的 PostgreSQL 交互式客户端程序。例如：psql -h 192.168.0.12 -p 6666 -U postgres test

使用 PostgreSQL 数据库命令行交互工具 psql 登录数据库后，可以执行 SQL 命令或者 psql 提供的元命令，使用 \? 命令可以查看 psql 的所有元命令和功能说明。登录数据库后命令行提示符为“数据库名#”，如：“test#”。psql 常用元命令如表 1-2 所示。

表 1-2 psql 常用元命令

元命令名	功能说明
\?	查看所有可以使用的元命令和说明信息
\o	将查询结果保存到文件或输出到 shell 命令
\l	列出所有数据库的名称、所有者、编码等信息
\q	退出
\c	连接到某个数据库
\dt	列出所有表
\d	显示数据库对象的模式
\di	列出所有索引
\i	执行文件中的命令
(sql);	执行标准的 SQL 语句，如： test#create table table1 (userID varchar (25), name varchar (25)); 创建表 test#insert into table1 ('1', 'pgtest'); test#select * from table1;

PostgreSQL 中的各种工具命令都是以独立程序的形式存在的，在源代码目录中和数据库核心代码分别放在不同的子目录中。本书的分析对象是 PostgreSQL 数据库的核心代码，各种工具的分析并不是本书的重点，有兴趣的读者可以尝试自行了解相关内容。

第2章

PostgreSQL的体系结构

前面介绍过，PostgreSQL数据库是一种几乎可以运行在各种平台上的免费的开放源码的对象关系数据库，它是一种以关系数据库和SQL为基础，扩展了抽象数据类型，从而具备面向对象特性的数据库。PostgreSQL数据库由连接管理系统（系统控制器）、编译执行系统、存储管理系统、事务系统、系统表五大部分组成，其组成结构和关系如图2-1所示。连接管理系统接受外部操作对系统的请求，对操作请求进行预处理和分发，起系统逻辑控制作用；编译执行系统由查询编译器、查询执行器组成，完成操作请求在数据库中的分析处理和转化工作，最终实现物理存储介质中数据的操作；存储管理系统由索引管理器、内存管理器、外存管理器组成，负责存储和管理物理数据，提供

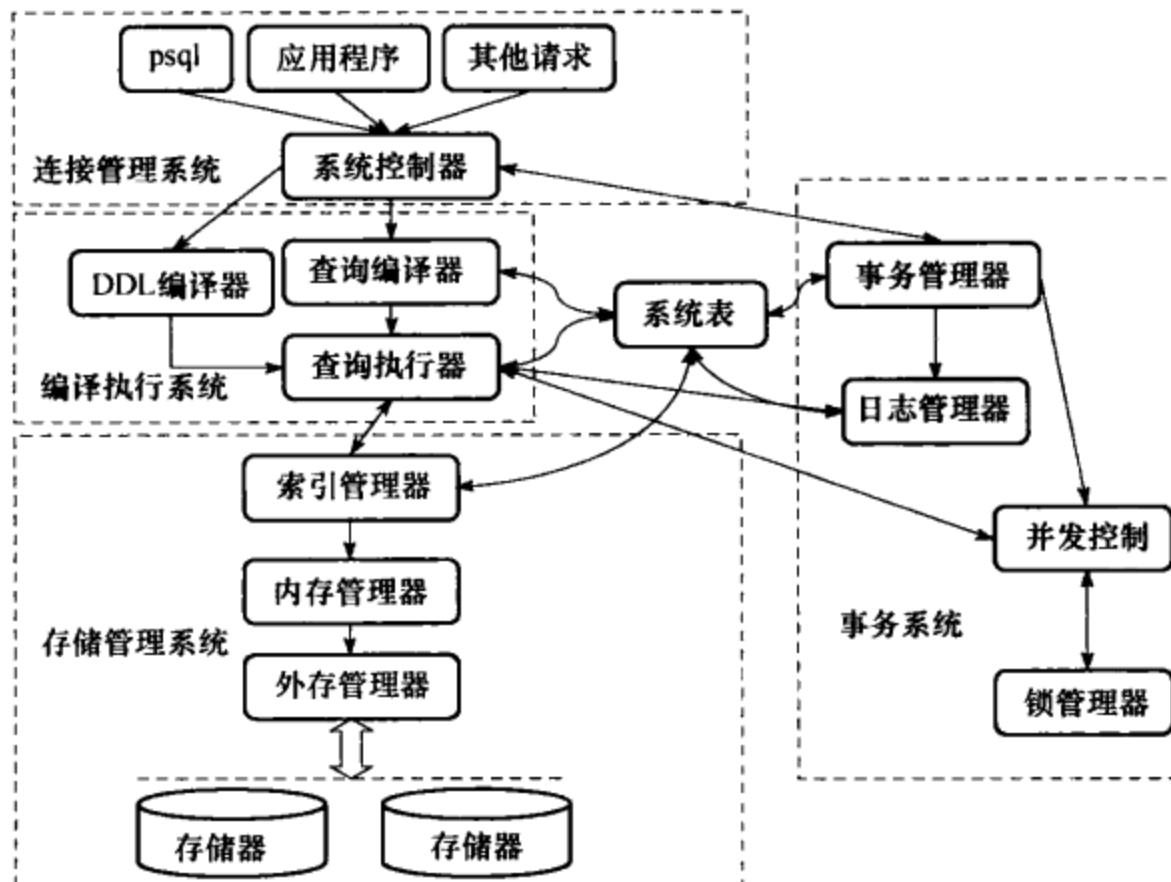


图2-1 PostgreSQL体系结构

对编译查询系统的支持；事务系统由事务管理器、日志管理器、并发控制、锁管理器组成，日志管理器和事务管理器完成对操作请求处理的事务一致性支持，锁管理器和并发控制提供对并发访问数据的一致性支持；系统表是 PostgreSQL 数据库的元信息管理中心，包括数据库对象信息和数据库管理控制信息。系统表管理元数据信息，将 PostgreSQL 数据库的各个模块有机地连接在一起，形成一个高效的数据管理系统。

2.1 系统表

在关系数据库中，为了实现数据库系统的控制，必须提供数据字典的功能。数据字典不仅存储各种对象的描述信息，而且存储系统管理所需的各种对象的细节信息。从内容来看，数据字典包含数据库系统中所有对象及其属性的描述信息、对象之间关系的描述信息、对象属性的自然语言含义以及数据字典变化的历史（即数据库的状态信息）。数据字典是关系数据库系统管理控制信息的核心，在 PostgreSQL 数据库系统中，系统表扮演着数据字典的角色。

系统表是 PostgreSQL 数据库存放结构元数据的地方，它在 PostgreSQL 中表现为存放有系统信息的普通表或者视图。用户可以删除然后重建这些表、增加列、插入和更新数值，然而由用户去修改系统会导致系统信息的不一致性，进而导致系统控制紊乱。正常情况下不应该由用户手工修改系统表，而是由 SQL 命令关联的系统表操作自动维护系统表信息。比如，创建数据库语句（CREATE DATABASE）会向 pg_database 系统表插入一行，并且在磁盘上创建该数据库。

PostgreSQL 的每一个数据库中都有自己的一套系统表，其中大多数系统表都是在数据库创建时从模板数据库中拷贝过来的，因此这些系统表里的数据都是与所属数据库相关的。只有少数系统表是所有数据库共享的（比如 pg_database），这些系统表里的数据是关于所有数据库的。

由于系统表保存了数据库的所有元数据，所以系统运行时对系统表的访问是非常频繁的。为了提高系统性能，在内存中建立了共享的系统表 CACHE，使用 Hash 函数和 Hash 表提高查询效率，这些内容将在第 3 章详细介绍。

系统表功能的实现代码包含系统表定义文件和系统表绑定函数实现文件。分别位于如下位置：

- 在 src/include/catalog 目录下有若干个以“pg_xxx”开头的 .h 文件，它们相应地定义了名为 pg_xxx 的系统表的数据结构，其中的 indexing.h 文件定义了所有的系统表索引，toasting.h 文件定义了所有系统表的 TOAST 表（TOAST 表用于存放普通表中的超长属性值，将在第 3 章介绍）。
- 在 src/backend/catalog 目录下的 pg_xxx.c 文件中定义了对 pg_xxx 进行相关操作的函数，其中的 indexing.c 文件定义了四个操作系统表索引的函数，toasting.c 文件定义了四个操作系统表的 TOAST 表的函数。

2.1.1 主要系统表功能及依赖关系

在 PostgreSQL 8.4.1 中，共有 42 张系统表和 17 张系统视图，系统视图是建立在基本系统表之上的。本节将对主要的系统表进行介绍，其他系统表将在后续章节中分别予以介绍。

1. pg_namespace

系统表 pg_namespace 用于存储命名空间。命名空间是 SQL92 模式下层的结构：每个名字空间有独立的关系、类型等集合，但并不会相互冲突。PostgreSQL 的名字空间层次是：数据库 . 模式 .

表·属性。

当要访问一个对象时，PostgreSQL 会按以下名字空间顺序进行搜索：

- 特殊名字空间（special），仅用于创建模式。
- 临时表的名字空间（TEMP）。
- 系统表的名字空间。

`pg_namespace` 中每一个元组都对应一个名字空间，每一个名字空间都被分配一个 OID（对象标识符，用于在整个数据库系统中唯一地标识一个数据库对象，包括数据库、表、视图、索引等，在 2.2 节中会进一步介绍）作为唯一标识，并且存储在对应元组的隐藏属性（PostgreSQL 中每个元组都有几个用户不可见的属性，用于记录一些系统级的信息）中。每个元组（名字空间）包含的属性如表 2-1 所示。

表 2-1 `pg_namespace` 属性表

属性名	数据类型	注释
<code>nspname</code>	<code>NameData</code> （长度为 64 的字符数组）	用于存放命名空间的名称
<code>nspowner</code>	<code>Oid</code>	用于表示该命名空间的所有者，由于系统中每一个用户也会分配一个 <code>Oid</code> 作为唯一标识，因此这里 <code>nspowner</code> 中只存储所有者的 <code>Oid</code>
<code>nspacl</code>	<code>aclitem</code> 类型的变长数组	其中存放对于该名字空间的访问权限列表，关于 ACL（访问控制列表）将在第 8 章中详细分析

2. `pg_tablespace`

`pg_tablespace` 存储表空间信息，将表放置在不同的表空间有助于实施磁盘文件布局。`pg_tablespace` 在整个数据集簇里只有一份，也就是说同一个数据集簇内的所有数据库共享一个 `pg_tablespace` 表，而不是每个数据库都有自己的 `pg_tablespace` 表。

PostgreSQL 里的表空间允许数据库管理员在文件系统里定义代表数据库对象的文件的存放位置。通过使用表空间，管理员可以控制一个 PostgreSQL 中数据的磁盘布局，即可以通过表空间将 PostgreSQL 系统的数据分布在不同的磁盘位置上。这样做有两个用处。第一，如果初始的集簇所在的分区或者卷用光了空间，而又无法扩展该分区或卷的空间，那么可以通过表空间去利用另一个分区的空间；第二，表空间允许管理员根据数据库对象的使用模式安排数据位置，从而优化性能。比如，一个使用很频繁的索引可以放在非常快并且可靠的磁盘上（如一种非常贵的固态设备）。而一个存储归档的数据、很少使用的或者对性能要求不高的表可以存储在一个相对便宜但速度比较慢的磁盘系统上。

`pg_tablespace` 中每一个元组都对应一个表空间，每一个表空间都被分配一个 OID 作为唯一标识，并且存储在对应元组的隐藏属性中。每个元组（表空间）包含的属性如表 2-2 所示。

表 2-2 `pg_tablespace` 属性表

属性名	数据类型	注释
<code>spcname</code>	<code>NameData</code>	存储表空间名称
<code>spcowner</code>	<code>Oid</code>	表示该表空间的所有者，通常为创建表空间的用户
<code>splocation</code>	<code>Text</code> （实际是一种变长的字符串数据类型）	用于存储表空间的物理位置（操作系统目录路径）
<code>spacl</code>	<code>aclitem</code> 类型的变长数组	存放对于该表空间的访问控制列表

3. pg_database

`pg_database` 中存放了当前数据集簇中数据库的信息，它也是一个在整个集簇范围内共享的系统表。该表中每一个元组就表示集簇中的一个数据库，每一个数据库都被分配一个 OID 作为唯一标识，并且存储在对应元组的隐藏属性中。每个元组（数据库）包含的属性如表 2-3 所示。

表 2-3 `pg_database` 属性表

属性名	数据类型	注释
<code>datname</code>	<code>NameData</code>	表示数据库的名称
<code>datdba</code>	<code>Oid</code>	数据库的拥有者（DBA）
<code>encoding</code>	<code>int4</code>	该整数表示数据库的字符集编码，比如 0 表示 SQL_ASCII 编码、6 表示 UTF8 编码
<code>datcollate</code>	<code>NameData</code>	数据库的 LC_COLLATE 设置
<code>datctype</code>	<code>NameData</code>	数据库的 LC_CTYPE 设置
<code>datistemplate</code>	<code>bool</code>	该数据库是否可以作为创建其他数据库的模板
<code>dataallowconn</code>	<code>bool</code>	数据库是否允许连接，如果取值为假表示该数据库不允许任何人连接，这个属性可用来保护 <code>template0</code>
<code>datconnlimit</code>	<code>int4</code>	数据库允许的最大并发连接数，如果值为 -1 表示对并发连接没有限制
<code>datlastsysoid</code>	<code>Oid</code>	数据库中使用的最后的 OID，在使用 <code>pg_dump</code> 时特别有用
<code>datfrozenxid</code>	<code>TransactionId</code>	数据库中所有在这个事务 ID 之前的事务都被替换为一个永久“冻结”的 ID，该属性用来跟踪该数据库是否需要进行 <code>vacuum</code> （清理）操作
<code>dattablespace</code>	<code>Oid</code>	数据库所在的表空间
<code>datconfig</code>	text 类型动态数组	数据库定义的 GUC 参数
<code>Datacl</code>	<code>aclitem</code> 类型的变长数组	存放对于该数据库的访问权限列表

4. pg_class

`pg_class` 存储表及与表类似结构的数据库对象信息，包含索引、序列、视图、复合数据类型、TOAST 表等。每一个对象都在 `pg_class` 中表示为一个元组，并且每一个对象都会被分配一个 OID 作为唯一标识，该 OID 作为该元组的一个隐藏属性存储。`pg_class` 中每个数据库对象包含的属性见表 2-4。

表 2-4 `pg_class` 属性表

属性名	数据类型	注释
<code>relname</code>	<code>NameData</code>	数据库对象的名称
<code>relnamespace</code>	<code>Oid</code>	数据库对象所属名称空间的 OID
<code>reltype</code>	<code>Oid</code>	对象的行类型对应的数据类型在 <code>pg_type</code> 系统表中的 OID，对于索引该字段值为 0
<code>relowner</code>	<code>Oid</code>	对象的所有者
<code>relam</code>	<code>Oid</code>	对于索引对象，该字段显示索引方式的 OID，如 B-Tree、Hash 等
<code>relfilenode</code>	<code>Oid</code>	对象所在的磁盘文件名称
<code>reltablespace</code>	<code>Oid</code>	对象所在的表空间

(续)

属性名	数据类型	注释
relpages	int4	对象的磁盘页面数，该属性被计划器使用，由 VACUUM、ANALYZE 或者一些 DDL 命令更新
reltuples	float4	表中的元组数，该字段被计划器使用，由 VACUUM、ANALYZE 或者一些 DDL 命令更新
reloastrelid	Oid	如果表关联了 TOAST 表，这里记录 TOAST 表的 OID（TOAST 表将在第 3 章进行介绍）
reloastidxit	Oid	可能存在的 TOAST 表的索引的 OID
relhasindex	bool	表明该表是否有索引
relisshared	bool	表明该表是否是数据集簇所共享的，只有部分系统表的这个属性设置为真
relistemp	bool	说明该表是否为临时表，临时表只能由创建它的会话访问
relkind	char	数据库对象的类型，r 表示表，I 表示索引，S 表示序列，v 表示视图，c 表示复合类型，t 表示 TOAST 表
relnatts	int2	表中用户定义属性的数目，隐藏属性不统计在内
relchecks	int2	CHECK 约束的数目
relhasoids	bool	是否为表的元组生成 OID，当设置为真时，为表的每行生成 OID
relhaspkey	bool	该表是否定义了主键
relhasrules	bool	该表是否定义了规则
relhastriggers	bool	在该表上是否定义了触发器
relhassubclass	bool	该表是否有继承表
relfrozenxid	xid	指定一个事务 ID，在表中所有在这个事务 ID 之前的事务都被替换为一个永久“冻结”的 ID，该属性用来跟踪该表是否需要进行 vacuum（清理）操作
relacl	aclitem 类型的变长数组	存放对于该数据库对象的访问权限列表
reloptions	text 类型的变长数组	数据库对象的访问方法选项，每一个元素是形为“keyword = value”的字符串对

5. pg_type

pg_type 存储数据类型信息。基本数据类型和枚举类型由 CREATE TYPE 创建，域类型由 CREATE DOMAIN 创建，复合数据类型在表创建时自动创建。pg_type 中每一个元组对应一个数据类型，每个元组包含的属性如表 2-5 所示。

表 2-5 pg_type 属性表

属性名	数据类型	注释
typname	NameData	数据类型名
typnamespace	Oid	数据类型所属命名空间的 OID
typlen	int2	数据类型的字节长度
typbyval	bool	是否要求内部设置一个默认值
typtype	char	b 表示基本类型，c 表示复合类型，d 表示域，e 表示枚举类型，p 表示伪类型
typcategory	char	用于数据类型转换的分类
typispreferred	bool	如果 typcategory 字段为其首选类型转换则设置为真
typisdefined	bool	类型为已定义时设置为真，为占位符时设置为假

(续)

属性名	数据类型	注释
typdelim	char	对于数组数据类型，切分输入的分隔符
typrelid	Oid	为复合数据类型时，设置为一个 pg_class 的 OID
typelem	Oid	如果不为 0，则该字段引用 pg_type 中的另外一个数据类型，当前数据类型可以看成是由 typelem 指定的数据类型构成的数组
typarray	Oid	如果不为 0，则该属性引用 pg_type 中的另外一个数据类型，被引用的数据类型是一个“真正的”数组，而本属性则是数组中的一个元素
typinput	regproc	输入转换函数文本模式
typoutput	regproc	输出转换函数文本模式
typreceive	regproc	输入转换函数二进制模式
typsend	regproc	输出转换函数二进制模式
typmodin	regproc	输入调节函数
typmodout	regproc	输出调节函数
typanalyze	regproc	自定义分析函数
typalign	char	不同数据类型存储时的对齐方式，c 表示字符对齐，s 表示短整型对齐（2 字节），i 表示整数对齐（4 字节），d 表示双精度浮点数对齐（8 字节）
typstorage	char	记录变长类型的存储方式，主要用于第 3 章要介绍的 TOAST 技术，这个属性有四种取值：p、e、m、x，分别对应 TOAST 技术中的四种变长类型的存储策略
typnotnull	bool	用于域类型的非空的类型约束
typbasetype	Oid	如果类型为域，则标识该域类型的基本类型
typtypmod	int4	域类型用这个属性记录其基本类型的 typmod
typndims	int4	域类型为数组时，数组的维度
typdefaultbin	text	只用于域类型，记录类型的默认值表达式是通过 nodeToString 函数转换过的形式
typdefault	text	如果为空，表示类型没有关联的默认值。如果 typdefaultbin 非空，typdefault 中必须记录有 typdefaultbin 中默认值表达式的可读版本。但如果 typdefaultbin 为空而 typdefault 非空，则 typdefault 中记录的是该类型的默认值

6. pg_attribute

pg_attribute 存储表的属性信息，对于数据库中表的每个属性都有一个元组。pg_attribute 中每个元组（属性）包含的属性如表 2-6 所示。

表 2-6 pg_attribute 属性表

属性名	数据类型	注释
attrelid	Oid	该属性所属的表 OID
attnname	NameData	属性的名称
atttypid	Oid	属性的数据类型所对应的 OID
attstattarget	int4	控制分析器对该属性累计统计详细信息的登记
attlen	int2	属性类型对应的类型长度
attnum	int2	属性在表中的编号
attndims	int4	属性为数组类型时，数组的维度
attcacheoff	int4	在存储器中总是设置为 -1，被加载到内存中的元组描述符后，该属性用于缓存这个属性在元组描述符中的偏移量，主要用于加快属性的查找

(续)

属性名	数据类型	注释
atttypmod	int4	记录在表创建时指定的与类型相关的数据（例如，一个 VARCHAR 属性的最大长度）。对于绝大多数情况来说，这个属性都被设置为 -1
atthyval	bool	属性类型对应的 pg_type 的 typbyval 字段值
attstorage	char	属性类型对应的 pg_type 的 typstorage 字段值
attalign	char	属性类型对应的 pg_type 的 typalign 字段值
attnonnull	bool	描述一个非 NULL 的约束
atthasdef	bool	属性有一个默认值
attisdropped	bool	属性被删除设置为真
attislocal	bool	属性是当前定义而非继承
attinhcount	int4	属性的直接父节点数
attacl	aclitem 类型的变长数组	存放对于该属性的访问权限列表

7. pg_index

pg_index 存储索引的具体信息，其元组包含的属性见表 2-7。

表 2-7 pg_index 属性表

属性名	数据类型	注释
indexrelid	Oid	索引对应的 pg_class 条目的 OID
indrelid	Oid	创建索引的表的 OID
indnatts	int2	索引的属性的数目
indisunique	bool	唯一索引时设置为真
indisprimary	bool	主键索引时设置为真
indisclustered	bool	表最后在该索引上聚簇
indisvalid	bool	对查询有效设置为真
indcheckxmin	bool	若为真，则在查询时，需要比较该元组的 xmin 值与查询事务的 xmin 值的大小。若前者比后者小，才能使用该索引进行查询；否则不能使用该索引
indisready	bool	可以插入数据时设置为真
indkey	Int2vector	记录表被索引的属性
indclass	oidvector	记录索引中每个属性对应的操作符 OID
indoption	int2vector	给索引中的每个属性的标志位信息，具体含义由索引类型确定，是新增属性。目前对于所有未排序的索引该值都为 0，已排序的有 DESC 等值
indexprs	text	保存创建索引条件的表达式树，对于每一个索引属性，若索引条件不是简单的引用，则存在表达式树。该属性存储的表达式树主要用来在对索引进行插入更新操作时计算键值
indpred	text	若索引为部分索引，则该属性保存部分索引的表达式树，在 indexprs 中则不再保存。当对部分索引进行查找、插入更新时，需要使用该属性判断新增的元组是否需要添加到索引中

关键系统表之间存在的相互依赖关系，如图 2-2 所示。

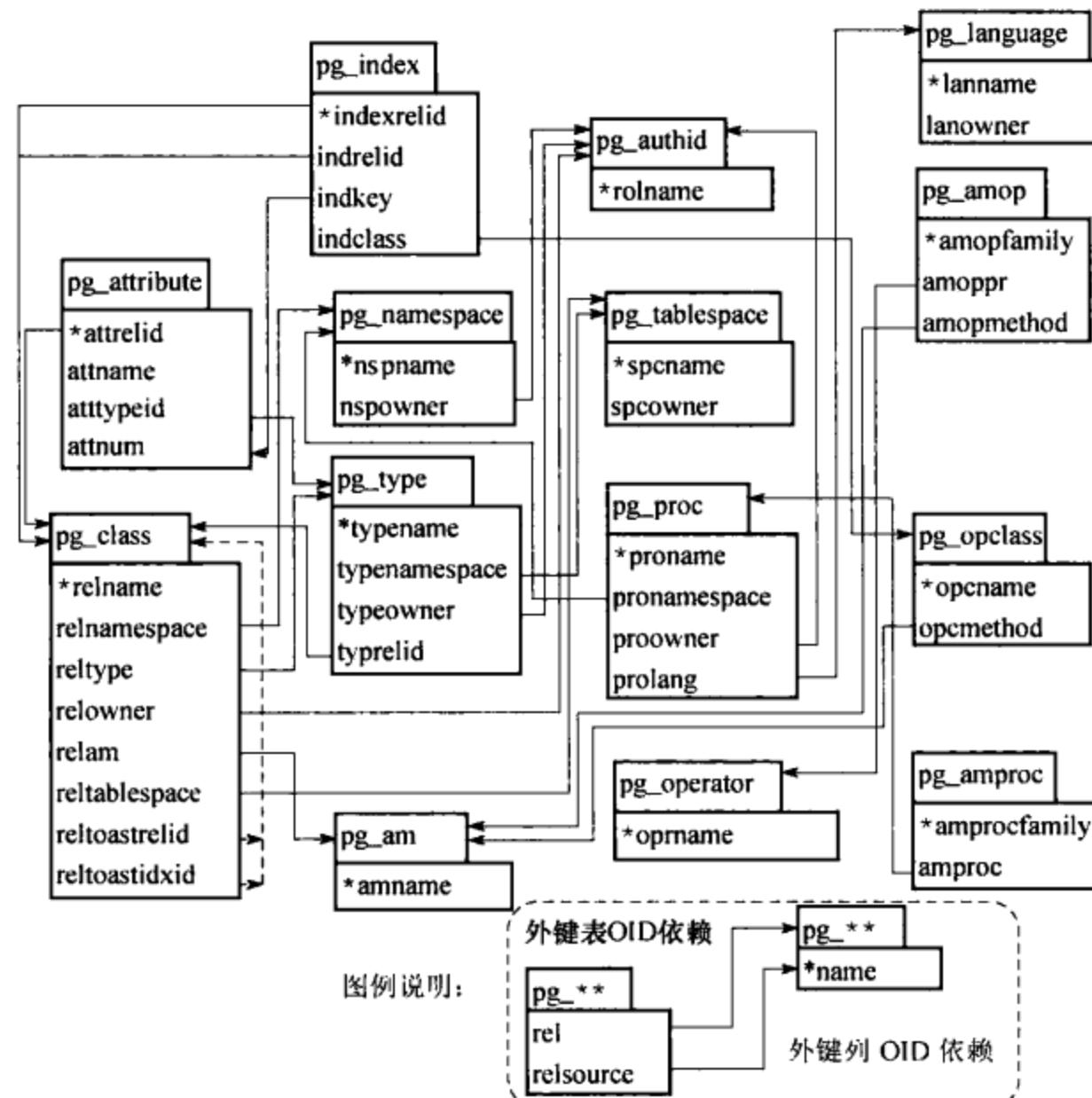


图 2-2 关键系统表之间的相互依赖关系

2.1.2 系统视图

除了系统表之外，PostgreSQL 还提供了一系列内置的视图，这些视图也是在初始化数据库集群的时候读取脚本创建的。系统视图提供了查询系统表和访问数据库内部状态的方法。表 2-8 列出了大部分系统视图及其用途。

表 2-8 系统视图

视图名	用途	视图名	用途	视图名	用途
pg_cursors	打开的游标	pg_roles	数据库角色	pg_timezone_abbrevs	时区缩写
pg_group	数据库用户的组	pg_rules	规则	pg_timezone_names	时区名
pg_indexes	索引	pg_settings	参数设置	pg_user	数据库用户 (无口令属性)
pg_locks	当前持有的锁	pg_shadow	数据库用户	pg_views	视图
pg_prepared_statements	预备语句	pg_stats	规划器统计		
pg_prepared_xacts	预备事务	pg_tables	表		

系统表是 PostgreSQL 数据库系统运行控制信息的来源，是数据库系统的核心组成部分。虽然用户可以操作，但为维护系统表信息的一致性，系统表将由系统统一维护。在 PostgreSQL 数据库安装

完后，需要先进行初始化数据库操作（`initdb`），生成模板数据库和相应的目录、文件信息，系统表即在此阶段生成。而用户数据库及其系统表都是从模板数据库进行复制生成的。初始化数据库的过程将在下节详细介绍。

2.2 数据集簇

PostgreSQL 安装完成后，在做任何其他事情之前，必须先使用 `initdb` 程序初始化磁盘上的数据存储区，即数据集簇，由 PostgreSQL 管理的用户数据库以及系统数据库总称为数据集簇。在 PostgreSQL 的实现中，数据库就是磁盘上一些文件的集合，只不过这些文件有特定的文件名、存储位置等，并且有些文件之间会相互关联。默认情况下，PostgreSQL 的所有数据都存储在其数据目录里，这个数据目录通常会用环境变量 PGDATA 来引用，后文中将会用 PGDATA 来指代数据目录。

在 PostgreSQL 中，对象标识符（OID）用来在整个数据集簇中唯一地标识一个数据库对象，这个对象可以是数据库、表、索引、视图、元组、类型等。PostgreSQL 提供了 `Oid` 数据类型来表示 OID，它实际上是一个无符号整数。

OID——对象标识符

OID 通常是从 1 开始分配，但在初始化数据集簇时，会先将一部分 OID 分配给系统表、系统表元组、系统表上的索引等数据库对象，这一部分 OID 可以在系统表所对应的头文件中找到。同时，为了给后续版本留下扩展的余地，初始化数据集簇时还会预留一部分 OID 资源。这样，在系统运行时可分配的 OID 资源实际是从 16384 开始的。在 PostgreSQL 源代码 `src/include/catalog` 子目录下有一个 shell 脚本 `unused_oids` 用来输出当前版本中预分配和预留的 OID 的使用情况。

数据集簇中的每一个数据库都对应于系统表 `pg_database` 中的一个元组，该元组 OID 属性中记录的就是分配给该数据库的 OID。同样，对于数据库中的对象如表、索引、视图、类型等，也对应于系统表 `pg_class` 中的一个元组（包括 `pg_class` 本身），这个元组的 OID 属性中记录的 OID 也就是该对象所分配的 OID。

对于用户表的元组，是可以在创建用户表时选择是否具有 OID 属性的。如果在 `CREATE TABLE` 语句中使用了 `WITH OIDS` 选项，则该用户表中插入的每一个元组都将被分配一个 OID。否则，默认状态下用户表的元组是没有 OID 属性的。

OID 的分配由系统中的一个全局 OID 计数器来实现，每次需要分配新的 OID 时，就从该计数器中取出当前的 OID，然后该计数器将会加 1。OID 分配时会采用互斥锁加以锁定以避免多个要求分配 OID 的请求获得同一个 OID。

初始化数据集簇包括创建包含数据库系统所有数据的数据目录、创建共享的系统表、创建其他的配置文件和控制文件，并创建三个数据库：模板数据库 `template1` 和 `template0`、默认的用户数据库 `postgres`。以后用户创建一个新数据库时，`template1` 数据库里的所有内容（包括系统表文件）都会拷贝过来，因此，任何在 `template1` 里面安装的内容都自动拷贝到之后创建的数据库中。`template0` 和 `postgres` 都是通过拷贝 `template1` 创建的。

对于某个具体的数据库，在 `PGDATA/base` 里都对应有一个子目录，子目录的名字是该数据库在系统表 `pg_database` 里的 OID。

每个表和索引都存储在其所属数据库目录下的独立文件里，以该表或者该索引的 `filenode` 号命

名，该号码记录在该表或索引在系统表 pg_class 中对应元组的 relfilenode 属性中。

在表或者索引超过 1GB 之后，它就被分裂成多个 1GB 大小的段。第一个段的文件名和 filenode 相同，随后的段命名为 filenode.1, filenode.2……这样的策略避免了在某些有文件大小限制的平台上可能出现的问题。

如果一个表的有些属性要存储相当大的数据，那么就会有个与之相关联的 TOAST 表，用于存储无法在数据行中放置的超大外置数据。表对应的 pg_class 元组的 reltoastrelid 属性记录了它的 TOAST 表 OID。

在 PostgreSQL 中，默认情况下会将数据文件存放在 PGDATA 指定的目录下，但如果 PGDATA 所在磁盘空间不足或者用户出于磁盘性能的考虑，需要为系统增加新的物理存储位置（比如在另外一个磁盘上的目录），则可以使用“表空间”来进行扩展。表空间从物理意义上来说就是一个新的磁盘目录（当然这个目录可以和 PGDATA 同处一个磁盘或者不同磁盘），指定放在表空间中的数据库对象的物理文件都存放在表空间对应的目录中。如果某个数据库中的数据库对象被指定放在一个表空间中，那么在表空间的目录中会以该数据库的 OID 为名称创建一个目录，该数据库中属于该表空间的对象的物理文件都放在这个目录中。每个用户定义的表空间在 PGDATA/pg_tblspc 目录里面都有一个符号链接，它指向表空间的物理目录，该符号链接用表空间的 OID 命名。系统默认创建的表空间 pg_default 和 pg_global 并没有通过符号链接的方式指向其物理目录，而是直接对应 PGDATA/base 和 PGDATA/global。

PGDATA 中还保存有数据集簇的配置文件和其他子目录，各子目录及文件的用途说明可以参见表 2-9。

表 2-9 PGDATA 中的子文件和目录

名称	类型	用途
PG_VERSION	文件	一个包含 PostgreSQL 主版本号的文件
base	目录	包含每个数据库目录，数据库目录以数据库的 OID 编号命名，其中名为 1 的目录对应模板数据库 template1
global	目录	包含整个集簇共享的全局表，比如 pg_database
pg_clog	目录	包含事务提交状态数据的子目录
pg_multixact	目录	包含多重事务状态数据的子目录（用于共享的行锁）
pg_stat_tmp	目录	包含统计子系统所需临时文件的子目录
pg_subtrans	目录	包含子事务状态数据的子目录
pg_tblspc	目录	包含指向表空间的符号链接的子目录
pg_twophase	目录	包含用于预备事务的状态文件的子目录
pg_xlog	目录	包含 WAL（预写日志）文件的子目录
postmaster.opts	文件	记录服务器上一次启动时使用的命令行参数
postmaster.pid	文件	一个锁文件，记录着当前的守护进程 Postmaster 的进程号和共享内存段 ID，在服务器关闭之后此文件将被删除
postgresql.conf	文件	主要配置文件，除基于主机的访问控制和用户名映射之外的其他用户可设置参数都保存在这个文件中
pg_hba.conf	文件	基于主机的访问控制文件，保存对客户端认证方式的设置信息
pg_ident.conf	文件	用户名映射文件，定义了操作系统用户名和 PostgreSQL 用户名之间的对应关系，这些对应关系会被 pg_hba.conf 用到

2.2.1 initdb 的使用

initdb 是在使用 PostgreSQL 之前用于初始化数据集簇的程序，它负责创建数据库目录、系统表、模板数据库。

initdb 命令语法如下：

```
initdb[option...] --pgdata|-D directory
```

其中主要的 option 选项见表 2-10。

表 2-10 常用的 option 选项

选项	长选项	意义
-A METHOD	--auth = METHOD	指定本地连接的默认用户认证方式
-D DATADIR	--pgdata = DATADIR	数据目录的路径，必须是对当前用户可读写的空目录，也可以使用环境变量 PGDATA 来指定而省略本选项
-E ENCODING	--encoding = ENCODING	指定默认数据库编码方式
-U NAME	--username = NAME	指定数据库超级用户名
-W	--pwprompt	提示超级用户设置口令
-d	--debug	以调试模式运行，可以打印出很多调试消息
-L directory	无	指定输入文件（比如 postgres.bki）位置

initdb 把用户指定的选项转换成对应的参数，通过外部程序调用的方式执行 postgres 程序。postgres 程序在这种方式下将进入 bootstrap 模式创建数据集簇，并读取后端接口 postgres.bki 文件来创建模板数据库。下面将先介绍 postgres.bki 文件，然后介绍 bootstrap 模式的运行机制。

2.2.2 postgres.bki

后端接口 postgres.bki 文件是在编译的过程中由 /src/backend/catalog 目录下的脚本 genbki.sh 读取 /src/include/catalog 目录下的以 .h 结尾的系统表定义文件（包括系统表索引和 TOAST 表定义文件）创建，并且通常存放在安装树的 share 子目录下。在 pg_*.h（星号表示对应系统表的名称，每一个这样的头文件对应一个系统表的结构定义）的头文件中包含如下内容的定义：

- 定义 CATALOG 宏，用于以统一的模式去定义系统表的结构以及用以描述系统表的数据结构，如系统表 pg_class 的定义通过 CATALOG (pg_class, 1259) 来表现。
- 通过宏 DATA (x) 和 DESCR (x) 来定义 insert 操作，这样的 insert 操作可能会有多个，用于定义系统表中的初始数据。

模板数据库 template1 是通过运行在 bootstrap 模式的 postgres 程序读取 postgres.bki 文件创建的。BKI 文件是一些用特殊语言写的脚本，这些脚本使 PostgreSQL 后端能够理解，且以特殊的“bootstrap”模式来执行之，这种模式允许在不存在系统表的零初始条件下执行数据库函数，而普通的 SQL 命令要求系统表必须存在。因此 BKI 文件仅用于初始化数据集簇。

BKI 命令的格式如下：

(1) create [bootstrap] [shared_relation] [without_Oids] tablename tableoid (name1 = type1

[, name2 = type2 , ...])

创建一个名为 tablename 并且 OID 为 tableOid 的表，表的字段在圆括弧中给出。Bootstrap 模式支持下列字段类型：bool、bytea、char (1 字节)、name、int2、int4、regproc、regclass、regtype、text、Oid、tid、xid、cid、int2vector、Oidvector、_int4 (数组)、_text (数组)、_Oid (数组)、_char (数组)、_aclitem (数组)。只有在创建完 pg_type 并且填充了合适的数据之后才可以创建包含其他类型字段的表。

- 如果指定了 bootstrap 选项，那么将只创建表而不向 pg_class 等系统表里面插入任何记录。因此这样的表将无法被普通的 SQL 操作访问，直到系统表中的初始记录用硬方法 (BKI 的 insert 命令) 插入。这个选项用于创建诸如 pg_class 等核心的系统表。
- 如果指定了 shared_relation 选项，那么该表将作为整个数据集簇的共享表创建 (比如创建 pg_database)。
- 如果指定了 without_Oids 选项，则该表将不会具有 OID 字段。

(2) open tablename

打开一个名为 tablename 的表，准备插入数据。在 BKI 被执行时，每一个时刻都只有一个表被打开，因此使用 open 命令时，任何当前已经打开的表都会被关闭。

(3) insert [OID = Oid_value] (value1 value2...)

以 value1、value2 等作为字段值插入一个元组，如果该表具有 OID 属性，则把 Oid_value 作为被插入元组的 OID。如果 Oid_value 为零或者省略了 OID 关键字，将为该插入元组分配系统中可用的下一个 OID。

(4) close [tablename]

关闭当前被打开的表。由于每一时刻只有一个表被打开，因此可以省略 tablename。

(5) declare [unique] index indexname indexOid on tablename using amname (opclass1 name1 [, ...])

在名为 tablename 的表上用 amname 访问方法创建一个 OID 为 indexOid 且名为 indexname 的索引。索引属性是 tablename 表的 name1、name2 等属性，在属性上使用的操作符类分别是 opclass1、opclass2 等。该命令仅仅是创建一个索引的结构，索引的内容 (即索引项) 并不由该命令填充。

(6) declare toast toasttableOid toastindexOid on tablename

为名为 tablename 的表创建一个 TOAST 表。该 TOAST 表的 OID 为 toasttableOid，其索引的 OID 为 toastindexOid。与 declare index 命令一样，该命令也不会填充 TOAST 表的索引。

(7) build indices

对使用 (5)、(6) 两种命令创建的索引进行内容填充。

特别需要说明的是，在有一些基本的系统表 (关键表) 被创建并初始化其数据之前，不能使用 open 命令打开非关键系统表。其中关键表包括：pg_class、pg_attribute、pg_proc 和 pg_type。根据前面对这几个系统表用途的介绍，我们可以很容易理解为什么只有在关键表建立之后才能使用 open 命令，因为这几个关键表存储了所有系统表的模式信息，如果它们没有被建立，open 命令不可能在其中找到要打开的系统表的模式信息。由于关键表不能用 open 打开进行填充，BKI 提供了带 bootstrap 选项的 create 命令，该命令可以在创建关键表后自动将其打开。同样，“declare index” 和

“declare toast” 命令也不能在它们所需要的系统表（比如前者需要 pg_index 等）被创建并填充之前使用。

在整个源代码被编译时，genbki.sh 脚本会被调用，它将从每一个 pg_*.h 文件中读取系统表定义、系统表的初始化数据、系统表上的索引等信息，然后分别将其转换为对应的 BKI 命令，最后将所有的 BKI 命令写入到 postgres.bki 文件中，该文件的内容如下：

- 1) 一个“create bootstrap”命令，用于创建其中一个关键表。
- 2) 一个或多个 insert 命令，用于填充步骤 1 创建的关键表中的数据。
- 3) 一个 close 命令，用于关闭步骤 1 创建的关键表。
- 4) 重复步骤 1~3 创建和填充其他关键表。
- 5) 一个不带 bootstrap 选项的 create 命令，用于创建一个非关键表。
- 6) 一个 open 命令，打开非关键表。
- 7) 一个或多个 insert 命令，填充非关键表所需要的数据。
- 8) 一个 close 命令，关闭上面打开的非关键表。
- 9) 重复创建其他非关键表。
- 10) 一个或多个“declare index”命令，用于定义索引。
- 11) 一个“build indices”命令，用于实际建立上一步所定义的索引。

2.2.3 initdb 的执行过程

执行 initdb 程序时，将从 initdb.c 文件中的 main 函数开始执行，main 函数的执行流程如图 2-3 所示。initdb 执行时将按照顺序执行下列工作：

- 1) 根据用户输入的命令行参数获取输入的命令名。
- 2) 设置系统编码为 LC_ALL，查找执行命令的绝对路径并设置该路径。
- 3) 设置环境变量（pg_data 等），获取系统配置文件的源文件路径（postgres.bki、postgresql.conf.sample 等文件），并检查该路径下各文件的可用性。
- 4) 设置中断信号处理函数，对终端命令行 SIGHUP、程序中断 SIGINT、程序退出 SIGQUIT、软件中断 SIGTERM 和管道中断 SIGPIPE 等信号进行屏蔽，保证初始化工作顺利进行。
- 5) 创建数据目录，以及该目录下一些必要的子目录，如 base、global、base/1 等。
- 6) 测试当前服务器系统性能，由测试结果创建配置文件 postgresql.conf、pg_hba.conf、pg_ident.conf，并对其中定义的参数做一些设置。
- 7) 在 bootstrap 模式下创建数据库 template1，存储在数据目录的子目录 base/1 中。
- 8) 创建系统视图、系统表 TOAST 表等，复制 template1 来创建 template0 和 postgres，这些操作都用普通的 SQL 命令来完成。
- 9) 打印操作成功等相关信息，退出。

initdb 是 PostgreSQL 中一个独立的程序，它的主要工作就是对数据集簇进行初始化，创建模板数据库和系统表，并向系统表中插入初始元组。在这以后，用户创建各种数据库、表、视图、索引等数据库对象和进行其他操作时，都是在模板数据库和系统表的基础上进行的。

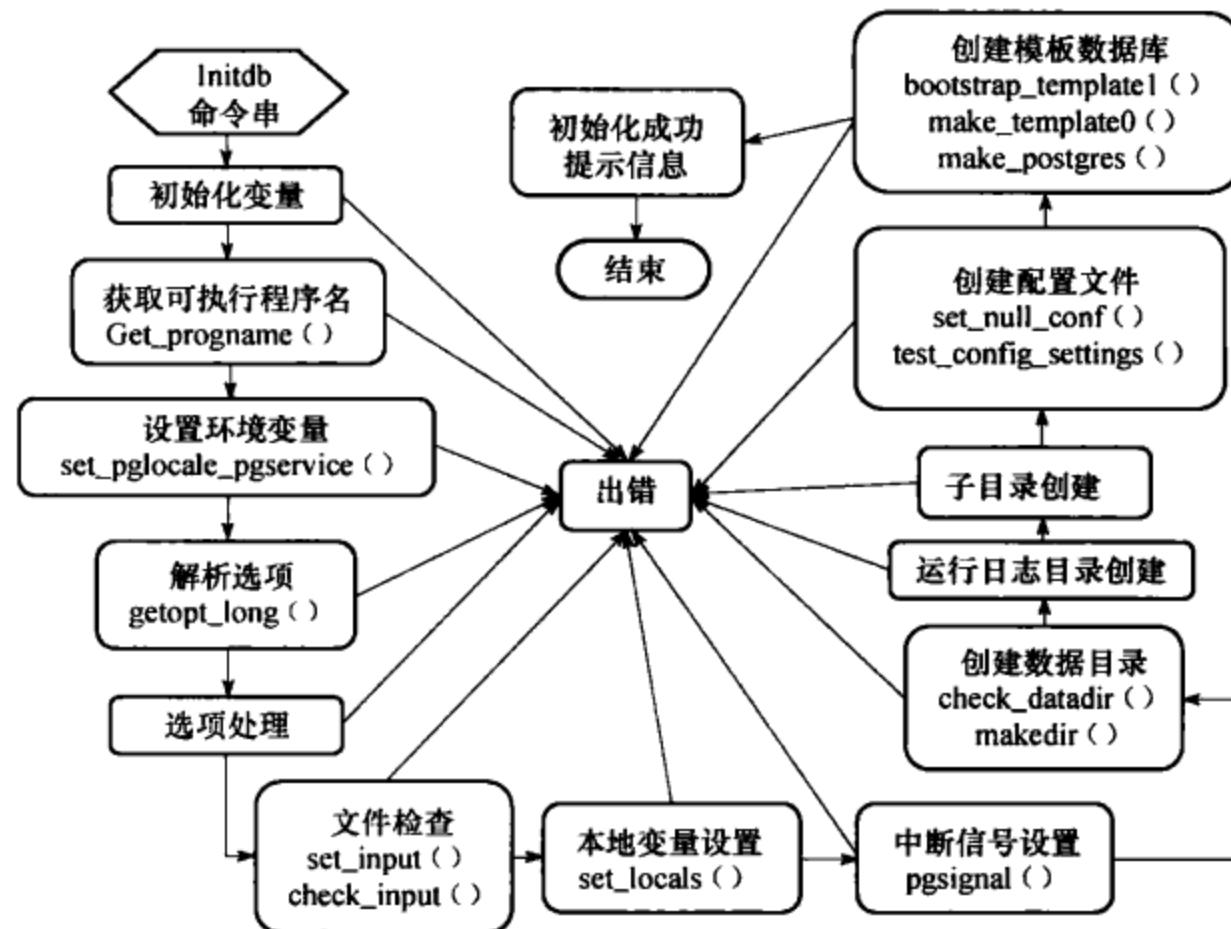


图 2-3 initdb 的运行流程

2.2.4 系统数据库

在创建数据集簇之后，该集簇中默认包含三个系统数据库 template1、template0 和 postgres。其中 template0 和 postgres 都是在初始化过程中从 template1 拷贝而来的。

template1 和 template0 数据库用于创建数据库。PostgreSQL 中采用从模板数据库复制的方式来创建新的数据库，在创建数据库的命令中可以用“-T”选项来指定以哪个数据库为模板来创建新数据库。

template1 数据库是创建数据库命令默认的模板，也就是说通过不带“-T”选项的命令创建的用户数据库是和 template1 一模一样的。template1 是可以修改的，如果对 template1 进行了修改，那么在修改之后创建的用户数据库中也能体现出这些修改的结果。template1 的存在允许用户可以制作一个自定义的模板数据库，在其中用户可以创建一些应用需要的表、数据、索引等，在日后需要多次创建相同内容的数据库时，都可以用 template1 作为模板生成。

由于 template1 的内容有可能被用户修改，因此为了满足用户创建一个“干净”数据库的需求，PostgreSQL 提供了 template0 数据库作为最初始的备份数据，当需要时可以用 template0 作为模板生成“干净”的数据库。

而第三个初始数据库 postgres 用于给初始用户提供一个可连接的数据库，就像 Linux 系统中一个用户的主目录一样。

上述系统数据库都是可以删除的，但是两个模板数据库在删除之前必须将其在 pg_database 中元组的 datistemplate 属性改为 FALSE，否则删除时会提示“不能删除一个模板数据库”。

2.3 PostgreSQL 进程结构

PostgreSQL 系统的主要功能都集中于 Postgres 程序，其入口是 Main 模块中的 main 函数，在初始化数据集簇、启动数据库服务器时，都将从这里开始执行。Main 模块主要的工作是确定当前的操作系统平台，并据此做一些平台相关的环境变量设置和初始化，然后通过对命令行参数的判断，将控制转到相应的模块中去。图 2-4 是 main 函数的调用流程。

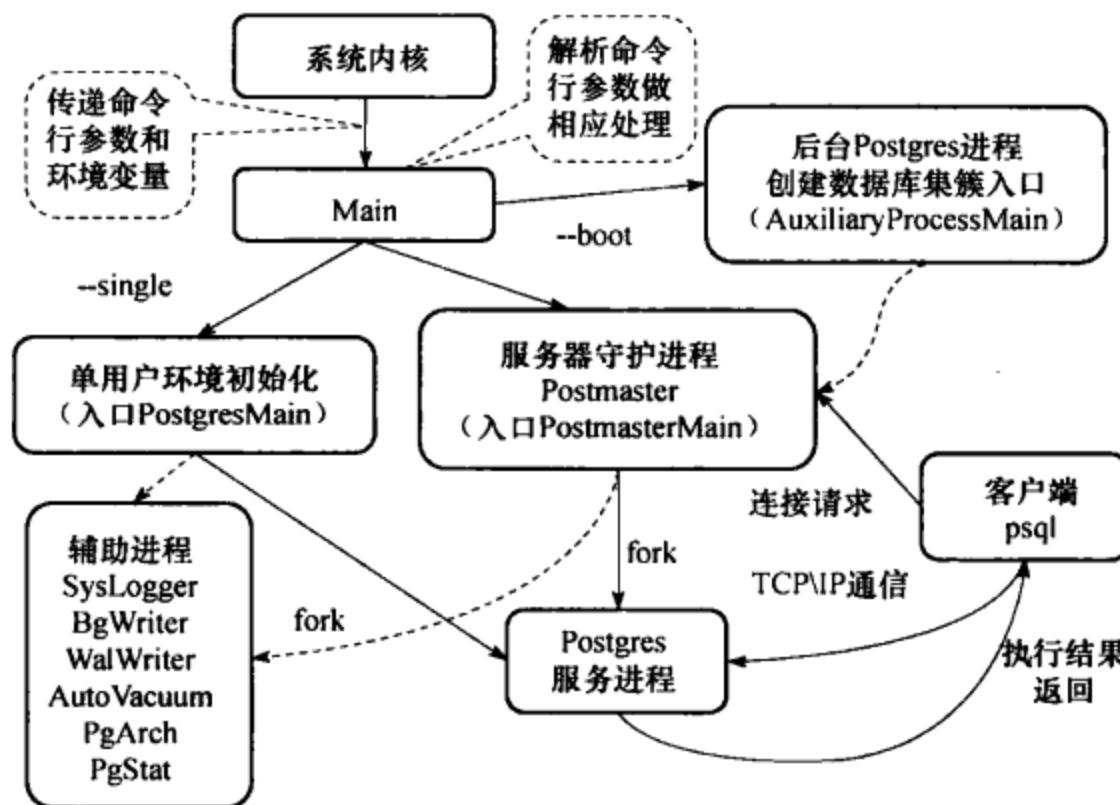


图 2-4 PostgreSQL 系统主函数 main 的流程

PostgreSQL 使用一种专用服务器进程体系结构，其中，最主要的两个进程就是守护进程 Postmaster 和服务进程 Postgres。从本质上来说，Postmaster 和 Postgres 都是通过载入 Postgres 程序而形成的进程，只是在运行时所处的分支不同而已。守护进程 Postmaster 负责整个系统的启动和关闭。它监听并接受客户端的连接请求，为其分配服务进程 Postgres。服务进程 Postgres 接受并执行客户端发送的命令。它在底层模块（如存储、事务管理、索引等）之上调用各个主要的功能模块（如编译器、优化器、执行器等），完成客户端的各种数据库操作，并返回执行结果。

Postmaster 和 Postgres 程序

在 Unix 或 Linux 系列的系统下，Postmaster 仅仅是 Postgres 的一个符号链接；而在 Windows 系统下，Postmaster 是 Postgres 的一个拷贝。所以 PostgreSQL 系统几乎所有的核心功能都是由 Postgres 程序完成的。

PostgreSQL 守护进程 Postmaster（单用户模式下的 Postgres 进程）除为用户连接请求分配后台 Postgres 服务进程外，还将启动相关的后台辅助进程。守护进程 Postmaster 在完成基本运行环境初始化、创建接受用户请求的监听端口后，顺序启动如下系统辅助进程：SysLogger（系统日志进程）、PgStat（统计数据收集进程）、AutoVacuum（系统自动清理进程）。在守护进程 Postmaster 进入到循

环监听中时启动如下进程：BgWriter（后台写进程）、WalWriter（预写式日志写进程）、PgArch（预写式日志归档进程）。这些辅助进程的用途在 2.5 节有详细介绍。

PostgreSQL 采用 C/S 模式，系统为每个客户端分配一个服务进程。前端应用欲访问某一数据库时，就调用接口库（比如 ODBC、libpq）把用户的请求通过网络发给守护进程 Postmaster。Postmaster 将启动一个新的服务进程 Postgres 为用户服务，此后前端进程和服务进程不再通过 Postmaster 而是直接进行通信。也就是说，Postmaster 总是监听用户连接请求并为用户分配服务进程 Postgres，而 Postgres 则负责为客户端执行各种命令。

2.4 守护进程 Postmaster

完成数据集簇初始化后，用户可以启动一个数据库实例来运行数据库管理系统，多用户模式下一个数据库实例由数据库服务器守护进程 Postmaster 来管理。它是一个运行在服务器上的总控进程，负责整个系统的启动和关闭，并且在服务进程出现错误时完成系统的恢复。它管理数据库文件、监听并接受来自客户端的连接请求，并且为客户端连接请求 fork 一个 Postgres 服务进程，来代表客户端在数据库上执行各种命令。同时 Postmaster 还管理与数据库运行相关的辅助进程。用户可以使用 postmaster、postgres 或者 pg_ctl 命令启动 Postmaster。

Postmaster 就像一个处理客户端请求的调度中心。当客户端程序需要对数据库进行操作时，首先会发出一个起始消息给 Postmaster 进行请求。Postmaster 将根据这个起始消息中的信息对客户端的身份进行验证，如果身份验证通过，Postmaster 就为该客户端新建一个服务进程 Postgres。随后 Postmaster 将与客户端的交互工作转交给 Postgres 服务进程，由 Postgres 来完成客户端所需要的数据库操作。Postmaster 的主要作用见图 2-5。

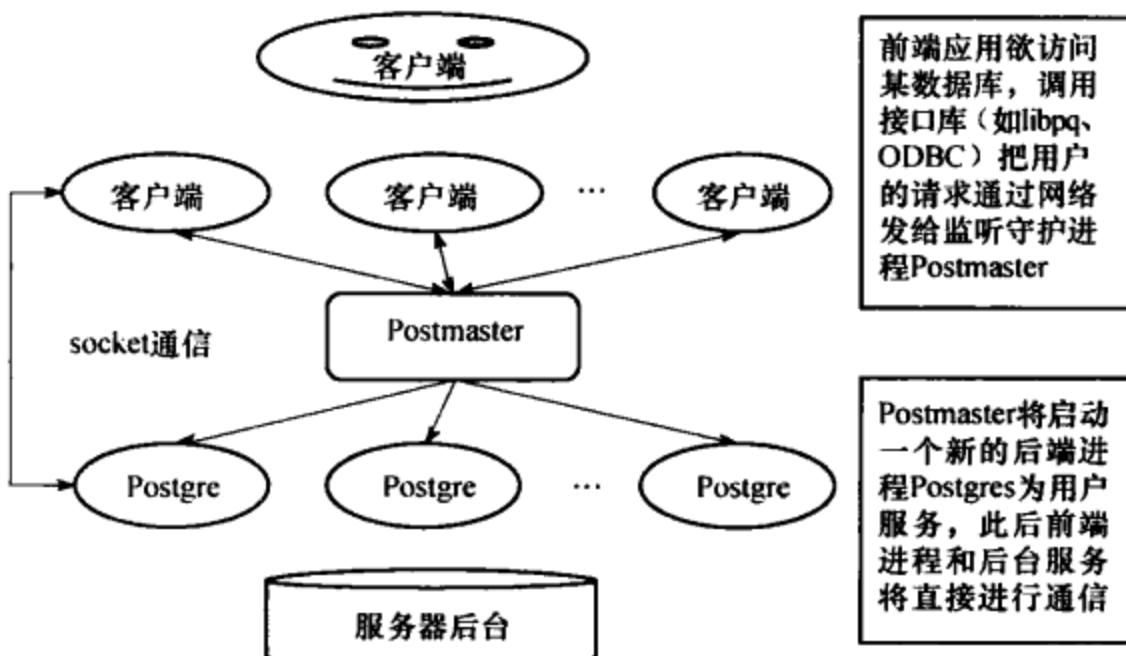


图 2-5 PostgreSQL 请求 – 响应模型

Postmaster 也负责管理整个系统范围的操作，例如中断等操作，Postmaster 本身不进行这些操作，它只是指派一个子进程在适当的时间去处理它们。同时它要在数据库崩溃的时候重启系统。Postmaster 进程在起始时会建立共享内存和信号库，Postmaster 及其子进程的通信就通过共享内存和信号来实现。这种多进程设计使得整个系统的稳定性更好，即使某个后台进程崩溃也不会影响系统

中其他进程的工作，Postmaster 只需要重置共享内存即可从单个后台进程的崩溃中恢复。

Postmaster 文件夹下面的代码主要用来创建一个服务器进程 Postmaster，用来监听用户的连接请求，并 fork 子进程来处理客户端请求，此外还包含了用来进行统计、建立系统日志的代码。相应的代码位置在 src/backend/postmaster 文件夹中。包括下列文件：

- Postmaster 进程源文件 postmaster.c
- 统计数据收集进程的源文件 pqstat.c
- 预写式日志归档进程的源文件 pgarch.c
- 后台写进程的源文件 bgwrite.c
- 系统日志进程的源文件 syslogger.c
- 系统自动清理进程的源文件 autovacuum.c

Postmaster 的入口函数是位于其源文件 postmaster.c 中的函数 PostmasterMain，其流程如图 2-6 所示。下面将对流程的各部分作详细介绍。

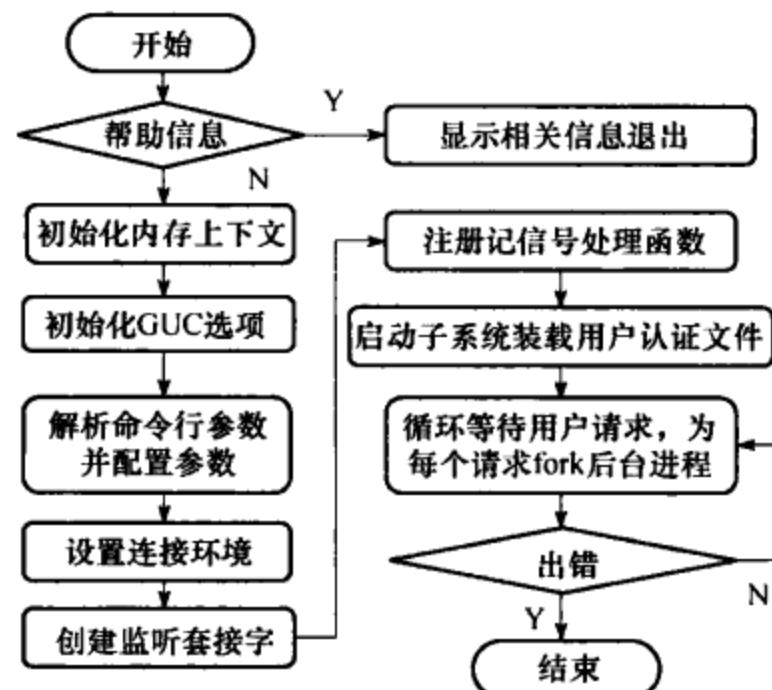


图 2-6 Postmaster 守护进程执行流程

2.4.1 初始化内存上下文

PostgreSQL 7.1 以前的版本在处理大量以指针传值的查询时一直存在着内存泄漏的问题，直到查询结束才能将内存收回。为了解决这个问题，从版本 7.1 开始，系统实现了新的内存管理机制，这样使得运行时大多数内存分配操作在各种语义的内存上下文（MemoryContext）中进行。内存上下文释放时将会释放其中分配的所有内存，这样即使某些内存没有被任何指针指向或忘记了释放，我们都可以通过释放内存上下文来避免这些内存泄漏。这一机制也使内存管理更加方便，开发人员不必再费尽心思地处理内存释放的工作。关于 MemoryContext 的具体分析和介绍请参见 3.3.1 节。

程序首先调用 MemoryContextInit 创建 TopMemoryContext 和 ErrorContext。然后调用 AllocSetContextCreate 以 TopMemoryContext 为根节点创建 PostmasterContext，最后将全局指针 CurrentMemoryContext 指向 PostmasterContext。这些内存上下文的具体含义如下：

- TopMemoryContext：在 TopMemoryContext 中分配的内存直到系统退出时才会释放。例如：TopMemoryContext 中存放了所有打开的文件描述符、内存上下文的控制节点等。它是所有内存上下文的树根。
- ErrorContext：这是错误恢复处理的永久性内存环境，恢复完毕则重设。
- PostmasterContext：这是 Postmaster 正常工作的内存环境，由它通过 fork 函数产生的子进程将会删除这个环境。

2.4.2 配置参数

在初始化内存环境之后，需要配置 Postmaster 运行时所需的各种参数。GUC（Grand Unified Configuration）模块实现了多种数据类型（目前有 boolean、int、float、string 四种）的变量配置。这些参数可能会由不同的进程在不同的时机进行配置，系统会根据既定的优先权来确定什么情况下的

配置可以生效。

参数共有六种类型（通过枚举类型 GucContext 定义），并且只能在合适的环境下进行配置：

- PGC_INTERNAL：参数只能通过内部进程设定，用户不能设定。
- PGC_POSTMASTER：参数只能在 Postmaster 启动时通过读配置文件或处理命令行参数来配置。
- PGC_SIGHUP：参数只能在 Postmaster 启动时配置，或当我们改变了配置文件并发送信号 SIGHUP 通知 Postmaster 或 Postgres 的时候进行配置。
- PGC_BACKEND：参数只能在 Postmaster 启动时读配置文件设置，或由客户端在进行连接请求时设置。已经启动的后台进程会忽略此类参数的改变。
- PGC_USERSET：可以在任何时候配置。
- PGC_SUSET：参数只能在 Postmaster 启动时或由超级用户通过 SQL 语言（SET 命令）进行设置。

下面这个枚举类型数据结构用于描述参数的来源，按照优先级从低到高的顺序排列，一个设置起作用当且仅当先前的设置优先级比当前的设置的优先级低或者相等（参见数据结构 2.1）。

数据结构 2.1 GUCSource

```
typedef enum
{
    PGC_S_DEFAULT,           //参数设为默认值
    PGC_S_ENV_VAR,           //参数通过环境变量得到
    PGC_S_FILE,              //读配置文件得到
    PGC_S_ARGV,               //从 Postmaster 命令行得到
    PGC_S_DATABASE,           //数据库安装时指定
    PGC_S_USER,                //用户指定
    PGC_S_CLIENT,              //通过客户连接请求传送过来的数据包指定
    PGC_S_OVERRIDE,             //在特定情况下用来强制设定为默认值
    PGC_S_INTERACTIVE,          //仅仅是作为交互式来源和非交互式来源之间的分隔
    PGC_S_TEST,                  //仅用于测试
    PGC_S_SESSION,              //通过 SET 命令设置
} GucSource;                 //参数来源，按优先级递增的顺序排列
```

每一种数据类型的 GUC 参数都由两部分组成：共性部分和特性部分。共性部分的数据结构如数据结构 2.2 所示。

数据结构 2.2 config_generic

```
struct config_generic
{
    const char *name;           //参数名
    GucContext context;          //参数类型
    enum config_group group;      //用于根据功能对参数进行分组
```

```

const char *short_desc;           //参数简短描述
const char *long_desc;           //参数详细描述
int         flags;               //参数标志
enum config_type vartype;        //参数值的数据类型
int         status;               //参数状态
GucSource reset_source;          //参数值为 reset_value 时参数的来源
GucSource source;                //当前参数来源
GucStack *stack;                 //当修改发生时，用来保存旧值，以支持回滚
char        *sourcefile;          //配置所在的源文件
Int         sourceline;           //在源文件中的行号
};


```

其中参数的数据类型 config_type 有五种：PGC_BOOL（布尔型）、PGC_INT（整型）、PGC_REAL（实数）、PGC_STRING（字符串）和 PGC_ENUM（枚举），枚举型是在 8.3 版本开始新加入的参数类型。

每一种具体的数据类型的参数都有其特性的数据结构，特性数据结构中的第一个项都是一个指向其共性数据结构的指针。整数型参数的特性部分如数据结构 2.3 所示。

数据结构 2.3 config_int

```

struct config_int
{
    struct config_generic gen;           //参数的共性数据结构
    int      *variable;                  //参数当前被设置的值（所有特性结构都有此项）
    int      boot_val;                  //参数的初始值
    int      reset_val;                 //重新设置时，如果用户未指定，则取该参数值
    int      min;                      //参数值的下界
    int      max;                      //参数值的上界
    GucIntAssignHook assign_hook;       //函数指针，用来设置 reset_val
    GucShowHook show_hook;              //目前没有用到
};


```

Postmaster 配置参数的基本过程包括图 2-7 所示的三个步骤。

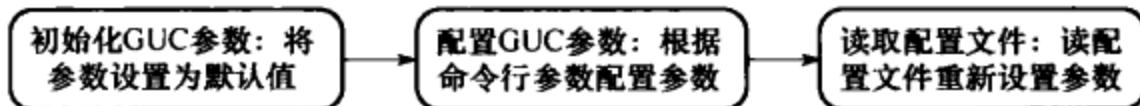


图 2-7 Postmaster 配置参数的过程

1. 初始化 GUC 参数

Postmaster 将首先调用 InitializeGUCOptions 函数将参数设置为默认值：

- 1) 首先调用 build_guc_variables 函数来统计参数个数并分配相应的 config_generic 类型的全局指针数组 guc_variables 以保存每个参数结构体的地址，并且对该数组进行排序。由于参数是通过全局

静态数组 `ConfigureNamesBool`、`ConfigureNamesInt`、`ConfigureNamesReal`、`ConfigureNamesString`、`ConfigureNamesEnum` 存储的，因此在 `build_guc_variables` 函数中只需要遍历相应的数组，统计参数的个数并将参数结构体中 `config_generic` 域的参数 `vartype` 设置为相应的参数类型。当遍历完所有参数后，根据总的参数个数分配 `config_generic` 指针数组 `guc_vars`，然后再次遍历静态参数数组，将每个参数结构的首地址保存到 `guc_vars` 数组中（这里分配的数组个数为当前参数总数的 1.25 倍，主要是为了方便以后参数的扩充）。接着将全局变量 `guc_variables` 也指向 `guc_vars` 数组。最后通过快速排序法把 `guc_variables` 按照参数名进行排序。

2) 接下来将每个参数设置为默认值。对于 `guc_variables` 中的每个参数，`initializeGUCOptions` 函数先将其 `config_generic` 域中的 `status` 设置为 0，将 `reset_source`、`tentative_source`、`source` 设置为 `PGC_S_DEFAULT` 表示默认；`stack`、`sourcefile` 设置为 `NULL`；然后根据参数值 `vartype` 的不同类型分别调用相应的 `assign_hook` 函数（如果该参数设置了该函数），`assign_hook` 函数用来设置 `boot_val`，最后将 `boot_val` 赋值给 `reset_val` 和 `variable` 指向的变量，通过这样一系列的步骤就将参数设置为了默认值。

3) 通过系统调用 `getenv` 来获得环境变量 `PGPORT`、`PGDATESTYLE`、`PGCLIENTENCODING` 的值，不为空则调用 `SetConfigOption` 函数来设置这三个变量对应的参数的值。

4) 最后，检测系统的最大安全栈深度，如果这个深度值大于 100KB 且不超过 2MB，则用它设置 `max_stack_depth` 参数。

2. 配置 GUC 参数

如果用户启动 Postmaster 进程时通过命令行参数指定了一些 GUC 的参数值，那么 Postmaster 需要从命令行参数中将这些 GUC 参数的值解析出来并且设置到相应的 GUC 参数中，这一部分的代码在 `Postmaster.c` 文件的 509 ~ 674 行中。根据命令行设置参数主要是通过 `getopt` 和 `SetConfigOption` 这两个函数来完成的。

用 `getopt` 解析命令行参数

`getopt` (`nargc`, `nargv`, `ostr`) 函数主要用来解析命令行参数，这里主要用到了 3 个全局变量 `optind`、`optopt` 和 `optarg`：

- `optind` 用来保存当前将要解析的参数在参数数组中的下标，其初始值为 1，最后的值应该等于参数个数 `argc`。
- `optopt` 用来指向正被解析到的参数选项。
- `optarg` 则用来保存 `optopt` 选项对应的参数值。

每当解析到一个参数选项 `optopt` 时，将检查 `optopt` 是否出现参数 `ostr` 中，如果没有出现，则说明该参数为非法参数，出错；如果该参数在 `ostr` 中出现，且后面紧跟符号 “：“，则说明该参数将不带参数值，这时候我们将 `optarg` 赋值为 `NULL`，然后将 `optind` 加 1 以进行下一次处理，函数本身返回 `optopt`；如果上述情况都没有发生，那么说明该参数选项将带参数。这时候又有两种情况：一种是参数直接跟在参数选项后面，另一种是参数和参数选项之间用空格隔开。对于第一种情况，只需要将 `optarg` 指向 `optopt` 后面的一个字符即可；对于第二种情况，则将 `optind` 加 1，然后将 `argv` [`optind`] 赋值给 `optarg`。上述工作都完成后增加 `optind` 的值，然后函数返回 `optopt`。

对于 getopt 返回的每一个参数选项及其参数值，通过一个 switch 语句根据参数选项的不同分别调用 SetConfigOption 函数设置相应的参数。

SetConfigOption 函数的第一个参数为参数名；第二个参数为参数值，其值存放在 getopt 函数返回的 optarg 字符串中；第三个参数为参数类型；最后一个参数为参数来源。由于在这里 Postmaster 是在处理命令行参数，所以这里的参数类型和参数来源分别设置为 PGC_POSTMASTER 和 PGC_S_ARGV。

SetConfigOption 函数是通过调用 set_config_option (const char *name, const char *value, GucContext context, GucSource source, bool isLocal, bool changeVal) 函数来实现的，其中最后 2 个参数统一设置为 false 和 true。该函数首先从 guc_variables 指向的参数数组中搜索参数名为 name 的参数，如果没找到则出错；否则将找到的参数的结构体中 GucContext 的值与传过来的参数 context 比较，判断在当前的上下文中参数是否可以设置，如果不能设置的话就报错，否则再将参数结构体中的 GucSource 与传过来的参数 source 进行比较，判断当前操作的优先级是否大于或等于先前的优先级，如果大于或等于先前优先级的话则根据具体参数值的类型将 value 转化为相应数据，然后设置参数结构体中的相应数据项即可。

3. 读取配置文件

当完成了命令行参数的设置之后，接着读配置文件重新配置参数。需要注意的是，在配置文件中设置的参数都不能修改之前通过命令行已经设置的参数，因为其优先级没有通过命令行设置的优先级高。

这个过程主要是调用 SelectConfigFiles (const char *userDoption, const char *progname) 函数来实现的，其中第一个参数是通过命令行设置的用户的数据目录，如果没有设置会通过环境变量 PGDATA 找到；第二个参数为程序名，主要用于错误处理。

该函数首先在数据目录下找到配置文件，然后调用词法分析程序解析文件。对于解析到的每个参数及其参数值，调用 SetConfigOption 来完成参数的修改。

通过上述三个步骤设置完参数后还要检验参数的合法性。比如，数据目录的用户 ID 应该等于当前进程的有效用户 ID、数据目录应该禁止组用户和其他用户的一切访问、缓冲区的数量至少是允许连接的进程数的两倍并且至少为 16，等等。如果一切合法，则将当前目录转入数据目录，然后进行后续的操作。

在创建监听套接字之前，Postmaster 需要保证当前只有一个 Postmaster 在运行且没有任何独立后台进程运行，这是通过 CreateDataDirLockFile 函数来完成的。该函数通过调用 CreateLockFile 函数在数据目录中创建锁文件 postmaster.pid，每次 Postmaster 或独立后台进程 Postgres 启动时都会在数据目录中创建这个独一无二的文件（创建文件的模式中置 O_EXCL 标志位）。因此可以通过尝试去创建该文件来检测当前是否还有别的 Postmaster 在运行。若创建成功则说明当前没有其他 Postmaster 或者 Postgres 正在运行，写入自己的 pid；若创建失败，从文件中取出创建该文件的进程的 pid（若其中的 pid < 0 说明有一个独立后台进程 Postgres 在运行；若 pid > 0 说明有一个 Postmaster 在运行），然后通过 kill 系统函数检测该进程是否依然存在，若还存在则当前的启动过程就会退出。还有一种特殊情况，就是先前的 Postmaster 被强行终止，留下了孤立的后台进程，此时可以通过检查 PostgreSQL 专用的共享内存段是否仍然在使用来判断，若仍然有进程使用则退出。若前面的情况都没有发生，就删除这个文件，然后再次创建该文件并写入当前进程的 pid 和数据目录并返回。

当确定了只有自己在运行时，还将调用 RemovePgTempFiles 函数删除 PGDATA/base/pgsql_tmp 中的临时文件 pgsql_tmp* 以及非默认表空间中的临时文件。

2.4.3 创建监听套接字

PostgreSQL 允许通过网络（TCP、Bonjour 等）或 UNIX 本地机来访问数据库，本节主要介绍一下如何创建 TCP 套接字。在介绍创建套接字之前，先介绍以下主要数据结构。

- 字符串 ListenAddresses：该字符串中存储的是服务器的 IP 地址，如果是多接口主机的话，IP 地址之间用“,” 隔开，“*”表示所有可用 IP 接口，缺省值为“localhost”。
- 整型数组 ListenSocket [MAXLISTEN]：该数组用来保存与服务器上某个 IP 地址（一个服务器可能有多个网卡，因此会设置多个 IP 地址）绑定的监听套接字描述符，初始时全为 -1。MAXLISTEN 的值为 64。
- struct addrinfo 结构体：该结构体用来保存我们调用 getaddrinfo 系统函数返回与协议无关的套接字时需要关心的信息，如要监听的服务器 IP 地址、服务器端口等相关信息。其结构如数据结构 2.4 所示。

数据结构 2.4 addrinfo

```
struct addrinfo
{
    int             ai_flags;           //AI_PASSIVE, AI_CANONNAME
    int             ai_family;          //地址族协议
    int             ai_socktype;        //套接口类型
    int             ai_protocol;        //0 或者 IPv4、IPv6 选项
    size_t          ai_addrlen;         //套接口地址结构长度
    struct sockaddr *ai_addr;          //套接口地址指针
    char            *ai_canonname;       //主机规范名
    struct addrinfo *ai_next;          //指向下一个 struct addrinfo 的指针
};
```

Postmaster 首先调用 SplitIdentifierString 函数解析字符串 ListenAddresses，得到多个服务器 IP 地址并构成一个 List（参见数据结构 2.5）。然后取出每一个服务器地址，调用函数 StreamServerPort 在该地址上创建一个监听套接口，并将返回的套接字描述符保存在整型数组 ListenSocket [MAXLISTEN] 的第一个值为 -1 的单元中。

数据结构 2.5 List

```
typedef struct List
{
    NodeTag      type;               //T_List, T_IntList, or T_OidList
    int          length;              //链表长度
    ListCell     *head;               //链表的头节点
    ListCell     *tail;               //链表的尾节点
} List;
```

List 是 PostgreSQL 中常用的一种链表结构，它的节点类型为 ListCell（参见数据结构 2.6），在 ListCell 中可以包含指向任何数据结构的指针。正是因为 List 和 ListCell 的独特设计，使得 List 能够用于建立任何数据类型的链表，在后面的内容中我们还将见到很多使用 List 的情况。

数据结构 2.6 ListCell

```
struct ListCell
{
    union
    {
        vOid      *ptr_value;           //T_List 时指向该节点所对应的数据结构
        int       int_value;          //T_IntList 时记录整数值
        Oid       Oid_value;          //T_OidList 时记录 OID 值
    } data;
    ListCell  *next;                //指向链表中的下一个节点
};
```

StreamServerPort 函数在创建监听套接字的时候，首先会填写 addrinfo 结构体。这里由于创建的是服务器套接字，所以将 ai_flags 设置为 AI_PASSIVE 表示该套接字为被动打开，ai_socktype 置为 SOCK_STREAM 表示将要创建的是个字节流的套接口，ai_family 将会根据具体情况设置，可以为 AF_INET、AF_INET6 和 AF_UNIX，分别表示 IPv4、IPv6 和 UNIX 套接字。完成这些设置之后就可以调用系统函数 getaddrinfo 来根据服务器地址、服务器端口号对初始 addrinfo 结构体进行补充，填写其中的套接口地址等关键信息，从而得到一个完整的套接口，最后依次调用 socket、bind、listen 系统函数就可以创建一个监听套接字了。

socket 函数先创建一个 socket 对象，bind 函数则将由 getaddrinfo 函数补充完整的 addrinfo 结构体变量和 socket 对象进行绑定，最后通过 listen 函数对绑定的套接口进行监听，返回一个监听套接字。处于监听状态的流套接字将维护一个客户端连接请求队列，listen 函数的第二个参数设定了该队列最多容纳的客户端连接请求数。如果当一个连接请求到来时，队列已满，那么客户端将收到一个 WSAECONNREFUSED 错误。

接下来，将调用 reset_shared 函数创建用于进程间通信的共享内存和信号量，在这之后装载用户认证文件 pg_hba.conf 和 pg_ident.conf，初始化 BackendList 后台进程双向链表，并创建 postmaster.opts 文件。

BackendList 是一个活动的后台进程表，主要用于跟踪和了解目前有多少子进程，并且在必要的时候给它们发送适当的信号。该进程表主要记录的是由 Postmaster 通过 fork 产生的 Postgres 进程，特殊的子进程例如 Startup 和 BgWriter 不在这个表中。

创建监听套接字的代码可参见 Postmaster.c 文件中的第 800 ~ 900 行。

2.4.4 注册信号处理函数

信号是操作系统响应某些错误状况而产生的事件，它可以明确地由一个进程发给另外一个进程，用这种办法传递信息或协调操作行为。进程可以自定义信号处理函数来处理信号，PostgreSQL 系统充分利用了这一点。

首先要声明的是进程有权选择响应或屏蔽信号（SIGKILL 和 SIGSTOP 不能屏蔽）。为了方便起见，Postmaster 定义了三个信号集：BlockSig 是要屏蔽的信号集；UnBlockSig 是不希望屏蔽的信号集；AuthBlockSig 是在进行用户连接认证时需要屏蔽的信号集，它们都是位向量。

在设置响应信号的处理函数之前，要通过 PG_SETMASK 函数把这些信号全部屏蔽，然后通过 pqsignal 函数为感兴趣的信号设置处理函数，例如 pqsignal（SIGHUP, SIGHUP_handler）注册了 SIGHUP 信号的处理函数 SIGHUP_handler。下面分析几个主要的信号处理函数的功能。

1. 信号处理函数 SIGHUP_handler

当配置文件发生改变时产生 SIGHUP 信号。Postmaster 进程在收到 SIGHUP 信号时重读配置文件 postgresql.conf（注意这是在 PGC_SIGHUP 环境下），然后向子进程发同样的信号（Postmaster 维护着一个子进程队列 BackendList，以方便向子进程发消息），并重新装载 pg_hba.conf 和 pg_ident.conf 文件（这两个文件主要用来进行客户端认证）。

2. 信号处理函数 pmdie

pmdie 处理三种信号：SIGTERM、SIGINT 和 SIGQUIT。这三种信号分别对应三种不同的系统关闭方式：

- Smart Shutdown：Postmaster 将等待所有子进程完成当前的任务后再安全关闭系统，对应于 SIGTERM 信号。
- Fast Shutdown：Postmaster 将向所有子进程发出 SIGTERM 信号，等待子进程接收到这个信号回滚当前事务并退出后再安全关闭系统，对应于 SIGINT 信号。
- Immediate Shutdown：Postmaster 将向所有子进程发 SIGQUIT 信号，子进程接收到这个信号马上退出，同时系统非正常关闭，对应于 SIGQUIT 信号。

全局变量 Shutdown 用于表示当前系统所处的关闭状态，可以有三种状态：NoShutdown、SmartShutdown 和 FastShutdown。三种状态值分别对应整数 0、1、2。

(1) SIGTERM 信号处理

若系统已经处于 SmartShutdown 或 FastShutdown（之前已经收到了 SIGTERM 信号），说明 Postmaster 已经在准备关闭系统了，此时不予处理，直接退出信号处理函数。否则将 Shutdown 置为 SmartShutdown 状态，表示正以这种方式关闭系统。如果系统当前正处于正常运行、归档恢复模式或者一致恢复模式，需要分别向自动清理相关的子进程以及 WalWriter 子进程发送 SIGTERM 信号，并且将系统状态置为等待在线备份结束。然后调用函数 PostmasterStateMachine 来根据当前状态机的状态来等待在线备份结束以及子进程退出。

(2) SIGINT 信号处理

若系统已经处于 FastShutdown 说明已经在准备关闭系统，直接退出信号处理函数。否则首先将 Shutdown 设为 FastShutdown。如果系统当前正处于归档恢复模式，则说明当前只有 BgWriter 是活动的，将当前状态设置为等待后台进程退出；如果处于正常运行、等待在线备份结束、等待后台进程结束或者一致恢复模式，则以发送 SIGTERM 信号的方式关闭所有的后台进程。最后仍然是调用函数 PostmasterStateMachine 来根据当前状态机的状态来等待在线备份结束以及子进程退出。

(3) SIGQUIT 信号处理

这是极其暴力的做法，Postmaster 将给 PostgreSQL 系统中所有的活动进程发信号 SIGQUIT，然后自己也退出。当收到 SIGQUIT 信号时，首先通过 SignalChildren 函数向 BackendList 中的所有进程发

送 SIGQUIT 信号，然后依次检查 BgWriter、WalWriter、AutoVacuum、PgArch、PgStat 等辅助进程的存在状态，如果辅助进程存在，则调用 signal_child 函数向它们发送 SIGQUIT 信号。完成向各子进程发送信号的工作后，Postmaster 调用 ExitPostmaster 函数退出。

3. 信号处理函数 reaper

当系统中有子进程退出的时候，子进程就会给 Postmaster 发送一个 SIGCHLD 信号。Postmaster 收到 SIGCHLD 信号后调用 reaper 函数清理退出的子进程。

- 若这个退出的子进程是启动系统的子进程且属于异常退出，那么调用 ExitPostmaster 退出；否则先将 StartupPID 置为 0（标志启动结束）、FatalError 置为 false（标志系统已恢复），然后重新装载权限检测所需要的文件。最后检查其他后台辅助进程是否存在，如果不存在还需要启动它们。
- 如果是 BgWriter 进程退出，首先将 BgWriterPID 置为 0。检查全局变量 pmState 是否为 PM_SHUTDOWN（表明系统状态是等待 BgWriter 做一个关闭检查点），是则向 PgStat 发送 SIGUSER2 信号要求它进行最后一次归档并退出，接着关闭 PgStat 进程；否则认为 BgWriter 是意外退出，调用 HandleChildCrash 来处理子进程崩溃（设置 BgWriterPID 为 0 等工作）。
- 如果是 WalWriter 或者 AutoVacuum 进程退出，这里将会忽略正常退出的情况（因为这两个进程会周期性地启动，执行完工作后又正常退出，且正常退出时会主动完成必要的清理），如果是意外退出则同样调用 HandleChildCrash 来进行崩溃处理。
- 如果是 PgArch 进程异常退出，系统正处于 PM_RUN（系统正常运行）状态且允许进行归档，则调用 pgarch_start 函数重新启动 PgArch 进程。否则将 pmStat 改为 PM_WAIT_DEAD-END 表示正在等待“死亡”的子进程退出。
- 若是 PgStat 进程异常退出，且 pmStat 处于 PM_RUN 的状态，则调用 pgstat_start 函数重新启动 PgStat 进程。
- 如果是 SysLogger 进程意外退出，则调用 SysLogger_Start 函数重启 SysLogger。

2.4.5 辅助进程启动

在 PostgreSQL 中，守护进程 Postmaster 负责整个系统的启动和关闭。在一个数据库管理系统实例中，除了守护进程 Postmaster 和服务进程 Postgres 外，还包括一些其他后台进程，用于专门负责某项具体的工作。在这些辅助进程出现问题时，Postmaster 进程重新产生出现问题的辅助进程。在 Postmaster 的创建过程中会首先启动 SysLogger 日志进程，并完成 PgStat 进程、AutoVacuum 进程的初始化工作，而在 Postmaster 的监听循环中检测辅助进程的状态，并新建或者重新创建这些辅助进程。

1. SysLogger 辅助进程

Postmaster 进程调用 SysLogger_Start 函数启动 SysLogger 子进程。SysLogger 是 8.0 以后新加的特性，它通过一个管道从 Postmaster、所有后台进程以及其他子进程那里收集所有的 stderr 输出，并将这些输出写入到日志文件中。在 postgresql.conf 配置文件中设置了日志文件的大小和存在时间，若当前的日志文件达到这些限制时，就会被关闭，并且一个新的日志文件会被创建。这些日志文件使用一种内部命名模式，文件名中包含了创建时间和当前的 Postmaster 的进程号 pid，日志文件存储于在 postgresql.conf 中设置的子目录下。SysLogger 辅助进程的工作原理将在 2.5 节中详细介绍。

2. 辅助进程初始化

SysLogger 辅助进程启动完后，Postmaster 开始对辅助进程 PgStat 进程、AutoVacuum 进程进行初始化操作，为进程分配必要的资源。

在 PgStat 进程的初始化过程中，主要完成用于发送和接收统计消息的 UDP 端口创建和测试，UDP 端口的创建过程与前面描述的 TCP 端口创建流程相同，但是 Socket 端口类型改为 SOCK_DGRAM。在端口创建完成后发送一个字节长的测试消息，当测试成功时，将创建的 Socket 端口设置为非阻塞 IO 模式，以保证在消息收集器出现故障时发送的消息能被抛弃而不会导致进程因为等待返回而阻塞；当测试失败时，发送错误消息，并关闭之前创建的 UDP 端口释放资源，设置 track_counts 选项为关闭。track_counts 是一个 GUC 选项，它标志了 PgStat 是否会进行统计信息的收集。显然，如果无法创建 PgStat 所需的 UDP 端口也就肯定无法进行统计信息收集。

在 AutoVacuum 进程的初始化过程中，只完成 AutoVacuum 启用选项和 track_counts 启用选项的检查工作。AutoVacuum 启用选项标志变量 autovacuum_start_daemon 会根据 GUC 配置进行设置，如果值为真则表示启用。在启动 AutoVacuum 进程时需要设置 track_counts 选项为启用，如果 AutoVacuum 启动选项标志为“启动”但是 track_counts 选项为关闭，则会报告错误信息。

辅助进程的初始化完成后，Postmaster 启动一个辅助启动过程（函数 StartupProcess）进行数据库的启动操作。其他辅助进程的启动将在循环等待连接 ServerLoop 中检查和启动。辅助进程的处理流程和实现细节将在 2.5 节中详细介绍。

2.4.6 装载客户端认证文件

注册完信号处理函数之后，将逐行读取 data 目录下的 pg_hba.conf 和 pg_ident.conf 两个配置文件的内容到链表变量中，以用于控制客户端认证。

其中，pg_hba.conf 是基于主机认证的配置文件。pg_hba.conf 文件的常用格式为每行一个记录。空白行将被忽略，#开头的注释也被忽略。一个记录是由若干用空格和/或制表符分隔的字段组成。如果字段用引号包围，那么它可以包含空白，记录不能跨行存在。每个记录声明一种连接类型、一个客户端 IP 地址范围（如果和连接类型相关的话）、一个数据库名、一个用户名、对匹配这些参数的连接使用的认证方法。第一条匹配上的记录将用于执行认证。如果选择了一条记录而且认证失败，那么将不再考虑后面的记录。同时，如果没有匹配的记录，那么访问将被拒绝。

pg_ident.conf 是基于身份（ident）认证的配置文件。当使用以身份为基础的认证时，在判断了初始化连接的操作系统用户名后，PostgreSQL 判断该用户名是否可以以所请求的数据库用户的身份连接。这个判断是由跟在 pg_hba.conf 文件里的 ident 关键字后面的身份映射控制的。有一个预定义的身份映射是 sameuser，表示任何操作系统用户都可以以同名数据库用户进行连接（如果后者存在的话），其他映射必须手工创建。非 sameuser 的身份映射定义在身份映射文件 pg_ident.conf 里。身份映射文件包含下面通用的格式：

```
map-name ident-username database-username
```

map-name 是在 pg_hba.conf 中引用这个映射的名称。另外两个域分别说明哪个操作系统用户被允许以哪个数据库用户的身份进行连接。同一个 map-name 可以重复用于在多个映射里声明更多的用户映射。操作系统用户和数据库用户之间的映射是多对多的关系，即一个操作系统用户可以映射

为多个数据库用户，一个数据库用户也可以映射为多个操作系统用户。

2.4.7 循环等待客户连接请求

PostgreSQL 系统采用客户端/服务器模式 (C/S)，系统为每个客户端分配一个后台服务进程 Postgres。一次 PostgreSQL 会话由下列相关的进程（程序）组成：

- 服务进程叫 Postgres，它接受来自客户端应用与数据库的连接，并且代表客户端在数据库上执行操作。
- 客户端进程，是指运行在客户端与 Postmaster 或 Postgres 进程交互的进程，比如在客户端运行的 psql 进程。

和典型的 C/S 应用一样，这些客户端和服务器可以在不同的主机上。这时它们通过网络连接通信。当客户端欲访问某一数据库时，就调用接口库（比如 libpq、ODBC）把用户的请求通过网络连接发给守护进程 Postmaster。Postmaster 将启动一个新的服务进程 Postgres 为用户服务，此后客户端进程和服务进程不再通过 Postmaster 而是直接进行通信。Postmaster 总是监听用户连接请求并为用户分配服务进程 Postgres，而 Postgres 则负责为客户端执行各种命令。从第 3 章开始将对它们的实现机制进行详细介绍。

在完成前述工作之后，Postmaster 将调用 ServerLoop 函数来循环等待客户端的连接请求，该函数的主体是一个死循环，它的主要功能就是在监听到用户的连接请求后建立与该用户的连接，然后通过调用 fork 函数复制出一个 Postgres 进程为该用户服务。

在 ServerLoop 函数中，首先调用 initMasks 函数初始化 Postmaster 所关心的读文件描述符集（通过扫描 ListenSocket 数组将里面的监听套接口字描述符加入到读文件描述符集中）。

接下来进入到一个死循环中，循环中首先调用 select 系统函数等待客户端提出连接请求。当 select 函数正常返回后，扫描 ListenSocket 数组，检查数组中的每一个监听套接字描述符，如果发现该套接字上有用户提出连接请求，那么将调用 ConnCreate 函数来创建一个 Port 结构体。该函数首先从 CurrentMemoryContext 中创建一个 Port 结构体，然后调用 Accept 系统函数接受由该监听套接字所维护的用户连接请求队列中的第一个连接请求，将监听套接字转化成一个已连接的套接字，并将返回的已连接套接字和用户端套接口地址结构都填入到 Port 的相应属性中。对于 TCP 连接，还要为连接设置超时等连接选项。Port 结构的部分定义如数据结构 2.7 所示。

成功创建 Port 结构体之后，调用 BackendStartup 启动后台进程接替 Postmaster 与用户进行连接，然后 Postmaster 关闭与用户的连接，释放 Port 结构体，检查并确保辅助进程如 BgWriter、SysLogger 等正常运行，完成之后就进入下一次循环，继续等待用户连接请求的到来。

BackendStartup 函数将根据 Port 结构体来启动一个 Postgres 服务进程，为用户提供服务。函数首先调用 PostmasterRandom 函数为即将创建的子进程生成一个随机数，并分配一个 Backend 结构。然

数据结构 2.7 Port

```
typedef struct Port
{
    int      sock;           // 已连接的套接口字描述符
    SockAddr laddr;          // 服务器端的地址信息
    SockAddr raddr;          // 客户端的地址信息
    char    *database_name;  // 要连接的数据库名
    char    *user_name;       // 连接的用户名
    .....
} Port;
```

后调用 fork 函数复制出一个 Postgres 子进程。

在子进程中将首先调用 BackendInitialize 函数获取用户端的相关信息（客户端的主机名、端口、用户名、数据库名等），然后将相关数据填入到 Port 结构中，再通过 ClientAuthentication 函数进行用户认证，检查用户端的地址、要连接的数据库、使用的用户名在 pg_hba.conf 文件中是否被允许，并根据 pg_hba.conf 中设置的认证方法与客户端进行相应的交互（见 8.2 节）。接下来该子进程将调用 BackendRun 来设置 Postgres 运行时的命令行参数，并且进行一些初始化工作，最后将调用 PostgresMain 函数进入 Postgres 的执行入口。

而在父进程也就是 Postmaster 中，会将随机数和子进程 pid 填入到 Backend 结构中，然后将该结构体加入到 Postmaster 维护的 BackendList 链表中。这样，Postmaster 才能通过 BackendList 链表来监控其子进程。

在 ServerLoop 循环中，除等待用户连接请求建立 Postgres 服务进程外，将检查辅助进程的状态，并根据需要启动辅助进程。

Postmaster 在数据库中扮演着管理员的角色，为用户的请求准备一个干净的环境。用户发送一个开始消息，Postmaster 即可根据该消息建立一个后端 Postgres 服务进程处理用户的操作，之后用户直接与 Postgres 进程通信而无需通过 Postmaster 中转。Postmaster 则继续监听新用户的请求。在整个数据库的生存周期中，Postmaster 都只执行管理工作，通过派生出各种辅助进程来执行相关功能，这个分离的结构极大地增强了系统的容错恢复特性。

2.5 辅助进程

PostgreSQL 的各个辅助进程完成特定的功能，支撑 PostgreSQL 系统运行和管理工作。在 Postmaster 进程中，为每个辅助进程设置了一个全局变量来标识该进程的进程号，分别为 SysLoggerPID、BgWriterPID、WalWriterPID、AutoVacPID、PgArchPID、PgStatPID。当这些变量值为 0 时，表示相应的进程尚未启动。只有系统日志辅助进程 Syslogger 在 Postmaster 进入监听循环之前完成启动操作，在监听循环中检测其状态进行重启操作；其他辅助进程的进程号在进入监听循环之前都被设置为 0，因此它们的初次启动和重新启动操作都会在 Postmaster 的监听循环中执行。在监听循环 Server-Loop 中，每一次循环时都会检查各个辅助进程号所对应的全局变量值是否为 0，并根据系统中进程运行状态机标志变量 pmState 的状态值来启动相应的辅助进程。

BgWriter 和 WalWriter 辅助进程采用相同的模式启动，设置辅助进程的类型 type，调用 StartChildProcess (type) 进入复制子进程 (fork) 操作流程，复制子进程成功后在父进程中直接返回，在子进程中封装传递参数 av、ac，关闭父进程中继承的端口，并切换内存上下文，进入统一的子进程处理函数 AuxiliaryProcessMain (ac, av)。在其中根据 type 类型启动各辅助进程。其他的辅助进程使用其对应的独立启动入口。下面将详细介绍每个辅助进程的处理流程和实现原理。

2.5.1 SysLogger 系统日志进程

日志信息是数据库管理员获取数据库系统运行状态的有效手段。在数据库出现故障时，日志信息是非常有用的。把数据库日志信息集中输出到一个位置将极大方便管理员维护数据库系统。然而，日志输出将产生大量数据（特别是在比较高的调试级别上），单文件保存时不利于日志文件的操作。因此，在 SysLogger 的配置选项中可以设置日志文件的大小，SysLogger 会在日志文件达到指

定的大小时关闭当前日志文件，产生新的日志文件。在 `postgresql.conf` 里可以配置日志操作的相关参数：

- `log_destination`: 配置日志输出目标，根据不同的运行平台会设置不同的值，Linux 下默认为 `stderr`。
- `logging_collector`: 是否开启日志收集器，当设置为 `on` 时启动日志功能；否则，系统将不产生系统日志辅助进程。
- `log_directory`: 配置日志输出文件夹。
- `log_filename`: 配置日志文件名称命名规则。
- `log_rotation_size`: 配置日志文件大小，当前日志文件达到这个大小时会被关闭，然后创建一个新的文件来作为当前日志文件。

此外，`postgresql.conf` 中还提供了其他配置参数，可以根据需要进行设置。

系统日志辅助进程的入口位置为 `SysLogger_Start` 函数。在 `SysLogger_Start` 函数第一次运行时，首先将调用系统函数 `pipe` 创建管道 `syslogPipe` 用于接收 `stderr` 的输出。接下来创建日志目录和文件，然后调用 `fork_process` 函数 `fork` 一个子进程来运行 `syslogger`。成功创建子进程后，在父进程即 `Postmaster` 中用 `dup2` 系统函数将 `stdout` 和 `stderr` 重定向到日志管道的写入端，然后关闭日志管道的写入端和日志文件（`Postmaster` 不需要再来写日志文件，`stdout` 和 `stderr` 的输出将直接通过管道写入日志文件），最后返回子进程的 PID。而在刚创建的子进程中，则关闭已创建的监听套接字，丢弃与 `Postmaster` 的共享内存的联系，最后调用 `SysLoggerMain` 函数真正进入日志进程的运行。

在 `SysLoggerMain` 函数体中，首先将对系统日志进程相关的全局变量进行赋值初始化，如记录自己的进程号、获取当前时间、进程状态标志等，若是重新启动 `SysLogger` 进程则要对 `stderr` 和管道进行重新设置。在进入系统日志进程处理流程主循环之前，需要注册信号处理函数，包括 `SIGHUP` 信号设置读取配置文件、`SIGUSR1` 信号设置切换日志请求；记录当前日志文件所在目录、日志文件名称、本次日志文件的生成时间，并调用 `set_next_rotation_time` 函数设置下一次日志周期（即关闭当前日志文件并创建新的日志文件）的时间点。完成准备工作后，程序将进入一个无限循环，它是 `SysLogger` 的工作主体。

每一次循环将首先判断是否收到了 `SIGHUP` 信号，如收到则调用 `ProcessConfigFile` 函数重新处理配置文件，并设置相应的参数的值。之后判断配置参数 `Log_directory`、`Log_filename` 和当前日志文件信息是否相同，即日志文件目录和名称是否发生修改，若修改则设置全局变量 `rotation_requested` 为真，表示需要进入新的日志周期。接下来判断日志文件是否达到下一次周期时间、是否达到规定的日志文件大小限制，如果满足日志切换条件，则设置 `rotation_requested` 为真。完成上述检查后，判断 `rotation_requested` 是否为真，是则调用 `logfile_rotate` 函数创建一个新日志文件，并关闭旧的日志文件。

接下来将调用 `select` 系统函数监听日志管道的读取端，当监听到发生了数据读取事件时，调用 `piperead` 函数从日志管道的读取端中读取一定长度的数据到临时缓冲区 `logbuffer` 中。如果读取数据的长度大于 0，则调用 `process_pipe_input` 函数对读取的数据进行处理，完成后跳出本次循环进入下一次循环；如果长度等于 0 则表明已经到达了日志管道末端，当前已没有任何进程维持着管道写入端，那么退出 `SysLogger` 进程。

process_pipe_input 函数

process_pipe_input 函数解析以 chunk (块) 为单位传送日志信息的日志管道协议，它采用 save_buffer 类型的全局数组 saved_chunks [20] 保存最多 20 个进程的日志信息 chunk (一个 save_buffer 保存一个进程的 chunk) 直到接收完最后一个 chunk，从而构成一个完整的消息后调用 write_syslogger_file 函数把该消息写入到日志文件中。

2.5.2 BgWriter 后台写进程

BgWriter 是 PostgreSQL 中在后台将脏页写出到磁盘的辅助进程，引入该进程主要为达到如下两个目的：首先，数据库在进行查询处理时若发现要读取的数据不在缓冲区中时要先从磁盘中读入要读取的数据所在的页面，此时如果缓冲区已满，则需要先选择部分缓冲区中的页面替换出去。如果被替换的页面没有被修改过，那么可以直接丢弃；但如果要被替换的页已被修改，则必需先将这页写出到磁盘中后才能替换，这样数据库的查询处理就会被阻塞。通过使用 BgWriter 定期写出缓冲区中的部分脏页到磁盘中，为缓冲区腾出空间，就可以降低查询处理被阻塞的可能性。其次，PostgreSQL 在定期作检查点时需要把所有脏页写出到磁盘，通过 BgWriter 预先写出一些脏页，可以减少设置检查点 (CheckPoint，数据库恢复技术的一种) 时要进行的 IO 操作，使系统的 IO 负载趋向平稳。通过 BgWriter 对共享缓冲区写操作的统一管理，避免了其他服务进程在需要读入新的页面到共享缓冲区时，不得不将之前修改过的页面写出到磁盘的操作。不过，当 BgWriter 无法维护足够的干净共享缓冲区时，其他服务进程仍然可以自行完成将脏页写回磁盘的操作。BgWriter 同时也负责处理所有的检查点，它也会定期地发出一个检查点请求，当然也可以由其他进程通过信号要求 BgWriter 执行一个检查点。

BgWriter 是 PostgreSQL 8.0 以后新加的特性，但在 8.2 以前版本中，使用 BgWriter 需要管理员进行很复杂的配置。在 PostgreSQL 8.4 中，数据库配置文件 postgresql.conf 中与 BgWriter 相关的配置选项有 3 个：bgwriter_delay、bgwriter_lru_maxpages、bgwriter_lru_multiplier。系统每隔 bgwriter_delay 指定的时间启动 BgWriter。BgWriter 从后向前扫描缓冲区的 LRU 链表，写出至多 bgwriter_lru_multiplier * N 个脏页，并且不超过 bgwriter_lru_maxpages 值的限制。其中 N 是最近一段时间在两次 BgWriter 运行期间系统新申请的缓冲页数。在 BgWriter 参数的配置中，如果 BgWriter 过于频繁地将脏页写出，则经常被更新的数据页很可能被一次又一次地写出到磁盘上，反而增加了数据库的 IO 次数，进而导致系统性能下降。另一方面，若 BgWriter 写周期过长，又不能起到优化数据库写 IO 操作的作用。因此确定 BgWriter 以什么速率将脏页写出才能达到最佳效果需要综合考虑系统的实际运行状态，默认将 bgwriter_delay 设置为 200 毫秒，bgwriter_lru_maxpages 设置为 100，bgwriter_lru_multiplier 设置为 2.0。

BgWriter 辅助进程在 Postmaster 中启动，入口为 StartBackgroundWriter 函数，BgWriter 和 WalWriter 两个辅助进程采用相同的模式启动。BgWriter 实际的工作函数是 BackgroundWriterMain。其处理流程如图 2-8

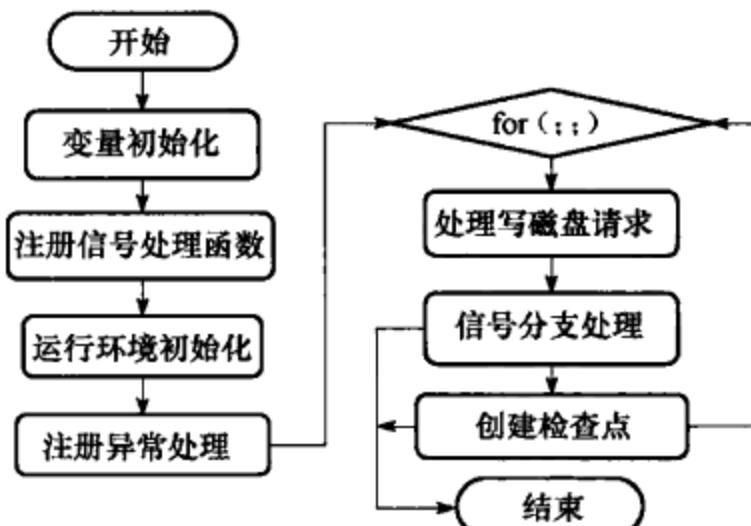


图 2-8 BgWriter 处理流程

所示。

- 1) 变量初始化：定义局部变量，完成部分全局变量的赋值操作，记录进程 PID。
- 2) 注册信号处理函数：SIGHUP 信号，在信号响应函数中设置 `got_SIGHUP = true`；SIGINT 信号，在信号响应函数中设置 `checkpoint_requested = true`；信号 SIGUSR2，在信号响应函数中设置 `shutdown_requested = true`；最后注册 SIGQUIT 用于快速退出处理。在后续的其他辅助进程中，信号的处理方式与此相同，都是在响应的信号处理函数中设置对应的标志变量为 `true`。在循环处理中根据标志变量的取值来执行相应的处理操作。
- 3) 运行环境初始化：通过 `ResourceOwnerCreate` 函数创建一个名称为“Background Writer”的资源跟踪器（见 3.6 节）。然后为 `BgWriter` 创建运行内存上下文，并将运行环境切换到新创建的内存上下文中。
- 4) 注册异常处理：结合系统调用 `setjmp` 和 `longjmp` 完成异常处理，实现进程的错误处理过程。在 `setjmp` 设置的异常处理结构中，会向系统日志中写错误信息、清理正在运行的缓冲区 I/O、释放初始化中创建的资源跟踪器，同时清理事务相关资源、更新 ckpt 状态标识 `checkpoint` 执行状态。最后清理 `ErrorContext` 并重置 `BgWriter` 内存上下文、关闭所有存在的 `SmgrRelation` 对象。使用 `setjmp` 和 `longjump` 是 C 语言编程中常用的一种错误恢复方法，有兴趣的读者可以参考相关的文献。
- 5) 处理写磁盘请求：调用 `AbsorbFsyncRequests` 函数，处理 `fsync` 请求（内存中的数据刷新到磁盘文件中）队列，将请求发送到本地的 SMGR（见第 3 章），该函数必须在执行检查点创建工作时执行。
- 6) 处理信号分支：在每次循环中检测信号对应的标志变量进行相应的处理。`got_SIGHUP` 代表重新读取配置文件；`shutdown_requested` 表示关闭数据库并退出 `BgWriter`；`checkpoint_requested` 表明有创建检查点请求，设置 `do_checkpoint` 标志，将检查点请求计数器加 1。
- 7) 创建检查点：如果没有创建检查点请求或者创建请求点时间间隔未到，则调用函数 `BgBufferSync` 按照 `BgWriter` 相关的配置参数把“脏”缓冲块刷回磁盘中。否则执行创建检查点操作。检查点的创建分为检查点类型设置、执行检查点创建操作、检查检查点是否创建成功三个阶段。与检查点创建相关的状态标志集合为 `ckpt_started`、`ckpt_failed`、`ckpt_done`、`ckpt_active`、`ckpt_performed`，与检查点类型相关的标志集合为 `do_checkpoint`（创建检查点）、`do_restartpoint`（创建恢复点）、`flags`。检查点的创建步骤如下：
 - ①当有创建检查点信号或者是到达创建检查点时间间隔时，设置 `do_checkpoint` 标志为需要创建检查点，否则调用函数 `BgBufferSync` 按照 `BgWriter` 相关的配置参数把“脏”缓冲块刷回磁盘中。
 - ②根据系统运行状态，以及在发送创建检查点请求时设置的检查点标志类型，来确定是创建检查点还是重启点。当系统处于恢复状态且不是 WAL 恢复的结束时刻时，将设置创建重启点 `do_restartpoint` 标志，完成检查点类型设置阶段。
 - ③在执行创建检查点阶段，将根据 `do_restartpoint` 的状态创建检查点或者重启点。在检查点创建后将直接设置 `ckpt_performed` 为真；而在重启点中将根据创建检查点成功与否标志来设置 `ckpt_performed`。
 - ④当 `ckpt_performed` 为真时，即检查点或者重启点创建成功时，将最近检查点创建时间设置为当前时间。否则设置最近检查点创建时间为 15 秒后再试。
 - ⑤完成上述工作后，如果需要将进行 XLog 日志切换操作，最后将休眠配置中设定的时间或中途被信号唤醒继续执行上述过程。

检查点创建流程的一次处理过程如图 2-9 所示。

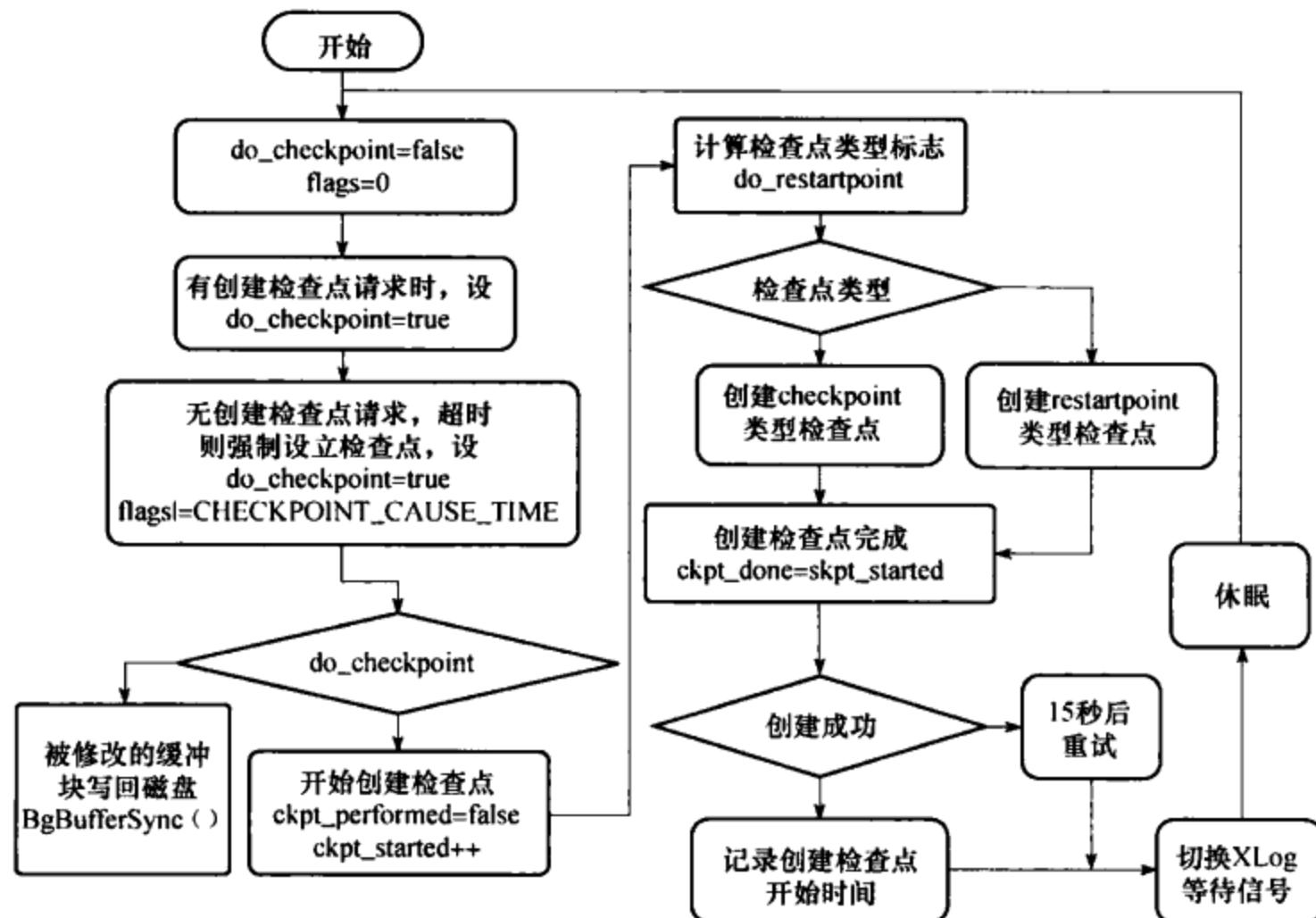


图 2-9 检查点创建流程

BgWriter 后台写进程最理想的情况是后台写进程负责刷回所有的缓冲区。但是，如果后台写进程不能保证有足够的干净的缓冲区情况下，常规后台进程仍然有权刷回缓冲区。

2.5.3 WalWriter 预写式日志写进程

预写式日志 WAL (Write Ahead Log, 也称为 Xlog) 的中心思想是对数据文件的修改必须是只能发生在这些修改已经记录到日志之后，也就是先写日志后写数据。如果遵循这个过程，那么就不需要在每次事务提交的时候都把数据块刷回到磁盘，因为在出现崩溃的情况下可以用日志来恢复数据库。使用 WAL 主要的好处就是显著地减少了写磁盘的次数，因为在日志提交的时候只需要把日志文件刷新到磁盘，而不是事务修改的所有数据文件。在多用户环境里，许多事务的提交可以用日志文件的一次 fsync 来完成。而且日志文件是顺序写的，因此同步日志的开销远比同步数据块的开销要小。WalWriter 是 Postgres 8.3 以后才新加入的新特性，它避免了其他服务进程在事务提交时需要同步地写入预写式日志到磁盘，也使得事务提交记录不是在提交时同步地写入磁盘，而是在一个已知的预先设置的时间异步地写入。同 BgWriter 一样，其他服务进程在 WalWriter 出错时也允许直接进行预写日志写操作。

WAL 日志文件存放在数据集簇中的 pg_xlog 目录里。它是作为一个段文件的集合存储的，每个段 16MB，并分割成若干页，每页 8KB。日志记录头格式在 xlog.h 里描述。日志内容取决于它记录的事件的类型。一个段文件的名字由 24 个十六进制字符组成，分为三个部分，每个部分由 8 个十

六进制字符组成。第一部分表示时间线，第二部分表示日志文件标号，第三部分表示日志文件的段标号。时间线由 1 开始，日志文件标号和日志文件的段标号由 0 开始，所以系统中的第一个事务日志文件是 000000010000000000000000，第二个事务日志文件是 000000010000000000000001。目前这些数字不能循环使用（不过要把所有可用的数字都用光也需要非常长的时间）。

WAL 的缓冲区和控制结构在共享内存里，它们是用轻量的锁保护的，对共享内存的需求由缓冲区数量决定，默认的 WAL 缓冲区大小是 8 个 8KB 的缓冲区（即 64KB）。出于安全考虑，可以将日志文件和数据文件分别存储在不同的磁盘上，可以通过把 pg_xlog 目录移动到另外一个位置，然后在数据集簇里原来的位置创建一个指向新位置的符号链接来实现。

在 PostgreSQL 数据库的系统配置文件 postgresql.conf 中有如下参数可以配置 WAL 的属性：

- fsync：该参数直接控制日志是否先写入磁盘。默认值是 ON（先写入），表示更新数据写入磁盘时系统必须等待 WAL 的写入完成。可以配置该参数为 OFF，表示更新数据写入磁盘完全不用等待 WAL 的写入完成。
- synchronous_commit：参数配置是否等待 WAL 完成后才返回给用户事务的状态信息。默认值是 ON，表明必须等待 WAL 完成后才返回事务状态信息；配置成 OFF 能够更快地反馈回事务状态。
- wal_sync_method：WAL 写入磁盘的控制方式，默认值是 fsync，可选用值包括 open_datasync、fdatasync、fsync_writethrough、fsync、open_sync。open_datasync 和 open_sync 分别表示在打开 WAL 文件时使用 O_DSYNC 和 O_SYNC 标志；fdatasync 和 fsync 分别表示在每次提交时调用 fdatasync 和 fsync 函数进行数据写入，两个函数都是把操作系统的磁盘缓存写回磁盘，但前者只写入文件的数据部分，而后者还会同步更新文件的属性；fsync_writethrough 表示在每次提交并写回磁盘会保证操作系统磁盘缓存和内存中的内容一致。
- full_page_writes：表明是否将整个 page 写入 WAL。
- wal_buffers：用于存放 WAL 数据的内存空间大小，系统默认值是 64K，该参数还受 wal_writer_delay、commit_delay 两个参数的影响。
- wal_writer_delay：WalWriter 进程的写间隔时间，默认值是 200 毫秒，如果时间过长可能造成 WAL 缓冲区的内存不足；时间过短将会引起 WAL 的不断写入，增加磁盘 I/O 负担。
- commit_delay：表示一个已经提交的数据在 WAL 缓冲区中存放的时间，默认值是 0 毫秒，表示不用延迟；设置为非 0 值时事务执行 commit 后不会立即写入 WAL 中，而仍存放在 WAL 缓冲区中，等待 WalWriter 进程周期性地写入磁盘。
- commit_siblings：表示当一个事务发出提交请求时，如果数据库中正在执行的事务数量大于 commit_siblings 值，则该事务将等待一段时间（commit_delay 的值）；否则该事务则直接写入 WAL。系统默认值是 5，该参数还决定了 commit_delay 的有效性。

WalWriter 辅助进程在 Postmaster 中启动，启动函数为 StartWalWriter，WalWriter 的实际工作函数为 WalWriterMain，该函数的处理流程如图 2-10 所示。

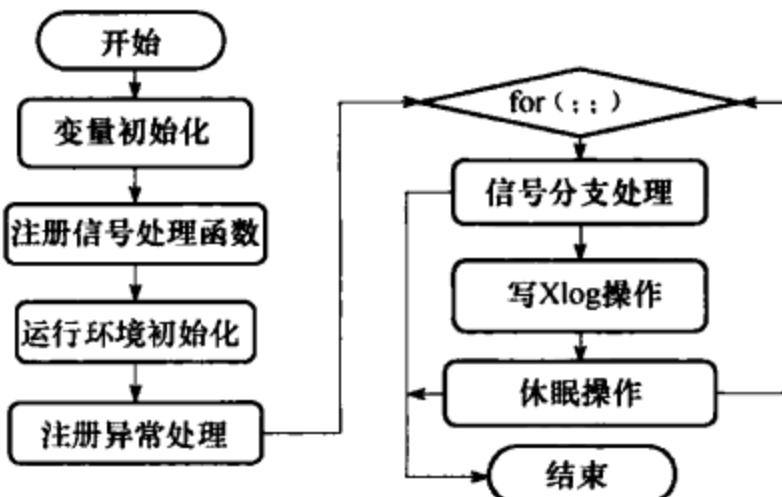


图 2-10 WalWriter 预写日志进程处理流程

从图 2-10 可以看到，WalWriter 的处理流程和 BgWriter 非常相似，只是 WalWriter 定期写磁盘的是存放预写式日志的 WAL 缓冲区，而 BgWriter 定期写入磁盘的是存放普通数据的共享缓冲区。WalWriter 的具体执行流程这里不再赘述，感兴趣的读者可以参考 BgWriter 流程的分析来阅读代码。

结合 WAL 日志和数据文件可以实现 PostgreSQL 数据库的在线备份和恢复。使用这种备份恢复方法时，我们可能要经常性地把数据文件、WAL 日志文件保存到另外一个存储设备上。数据库文件拷贝和日志归档文件可以用于灾难恢复，每次做归档以后，过时的日志文件就可以删除。PostgreSQL 提供了一种 dump 方式备份数据库文件，完成此工作的 pg_dump 工具存放在安装目录的子目录 bin 下。

2.5.4 PgArch 预写式日志归档进程

PostgreSQL 从 8.x 版本开始提出了 PITR (Point-In-Time-Recovery) 技术，支持将数据库恢复到其运行历史中任意一个有记录的时间点。除 2.5.3 节中所述的 WalWriter 外，PITR 的另一个重要的基础就是对 WAL 文件的归档功能。PgArch 辅助进程的目标就是对 WAL 日志在磁盘上的存储形式 (Xlog 文件) 进行归档备份。

PostgreSQL 在数据集簇的 pg_xlog 子目录中始终只使用一个 WAL 日志文件，这个日志文件记录数据库中数据文件的每个改变。从逻辑上来看，PostgreSQL 数据库会产生一个无限长的顺序的 WAL 记录序列。PostgreSQL 在物理上把这个 WAL 记录序列分割成多个 WAL 文件（每个 WAL 文件为一个 WAL 段），通常每个段的大小为 16MB（在编译 PostgreSQL 时可以通过编译选项改变这个大小）。每个段文件的名字是一个数字，用来反映它们在 WAL 序列中的位置。即使不进行 WAL 归档，PostgreSQL 也会创建一些 WAL 段文件。但是只会使用其中一个来记录 WAL 日志，如果当前使用的 WAL 段文件超过了大小限制，则会关闭当前段文件，然后把另外一个可重复使用的段文件作为当前段文件来使用。能够被重复使用的段文件必须保证其中的内容都在最后一次检查点之前产生并保证其中的内容都写入磁盘中。在这种情况下，被保存下来的将只有部分 WAL 日志，因此也不能实现任意历史时间点的恢复。为实现 PITR，需要在 WAL 段文件被重用时进行归档备份操作，把将被重用的 WAL 段中的日志记录保存到其他位置。这样归档日志加上当前日志就可以形成连续的 WAL 日志记录。为了给数据库管理员提供最大的灵活性，PostgreSQL 不对如何归档做任何假设，而是让管理员提供一个 shell 命令来拷贝一个完整的 WAL 段文件到备份存储位置。该命令可以就是一个 cp 命令，或者是一个复杂的 shell 脚本，所有的操作都由管理员决定。

在 postgresql.conf 中与预写式日志归档相关的属性有：

- archive_mode：表示是否进行归档操作，默认值为 off（关闭）。
- archive_command：由管理员设置的用于归档 WAL 日志的命令。
- archive_time：表示归档周期，在超过该参数设定的时间时强制切换 WAL 段，默认值为 0（表示禁用该功能）。

为允许归档，需要把 postgresql.conf 配置文件中的 wal_level 参数设置为“archive”或“hot_standby”，archive_mode 参数设置为“on”，并为 archive_command 命令指定一个 shell 命令。在用于归档的命令中，预定义变量“%p”用来指代需要归档的 WAL 全路径文件名，“%f”表示不带路径的文件名（这里的路径都是相对于当前工作目录的路径）。每个 WAL 段文件归档时将调用 ar-

chive_command 所指定的命令。当归档命令返回 0 时，PostgreSQL 就会认为文件被成功归档，然后就会删除或循环使用该 WAL 段文件。否则，如果返回一个非零值，PostgreSQL 会认为文件没有被成功归档，便会周期性地重试直到成功。

为了标识各个段文件的状态，PostgreSQL 在数据集簇的 pg_xlog/archiver_status 目录下记录了每一个 WAL 段文件的状态文件，状态文件的前缀与段文件同名，以表示时间顺序的整数形式命名。状态文件后缀为 .ready 或者 .done，代表段文件的归档状态，分别代表“就绪”和“已完成”两种模式。PgArch 进程会找到所有状态为“就绪”的段文件，找到状态文件后，若用户设置了归档命令，PgArch 进程将归档命令解析后交由系统的 shell 函数 system(3) 执行。文件归档成功后会把 pg_xlog/archiver_status 目录下相应的状态文件后缀修改为 .done。

PgArch 进程的启动由函数 pgarch_start 负责。在该函数中执行 fork 操作创建 Postmaster 的子进程 PgArch，在新创建的子进程中关闭从 Postmaster 进程中复制的网络连接端口，同时关闭与父进程的共享内存之间的联系，最后进入 PgArch 的工作函数 PgArchiverMain，其处理流程如图 2-11 所示。

可见，PgArch 进程的执行流程和 BgWriter、WalWriter 大同小异。PgArch 进程在工作循环中检测到需要进行归档时会按照前面介绍的方式找到“就绪”的 WAL 段文件，然后对其进行归档。完成一次归档之后，PgArch 将休眠直至下一次需要归档时。

对 WAL 日志的归档给管理员提供了一种新的数据库备份策略，这种备份策略组合了文件系统备份与 WAL 文件的备份。在恢复时，首先恢复数据文件，然后重放 WAL 日志到指定的时间点。应用这种备份恢复策略，在开始的时候并不需要一个非常完美的一致性备份。任何备份内部的不一致都会被日志的重放动作修改正确。因此，我们不需要文件系统快照的功能，只需要 tar 或者类似的归档工具。另外，WAL 文件的分段归档也把连续的备份简化为了分段的备份。这个功能对大数据库特别有用，因为大数据库的完全备份可能并不方便。如果持续把 WAL 文件重放给其他装载了同样的基础备份文件的机器，就有了一套“热备份”系统：在任何点我们都可以启动第二台机器，而它拥有近乎当前的数据库拷贝。和简单的文件系统备份技术一样，这个方法只能支持整个数据集簇的恢复。

2.5.5 AutoVacuum 系统自动清理进程

在 PostgreSQL 数据库中，对表元组的 UPDATE 或 DELETE 操作并未立即删除旧版本的数据，表中的旧元组只是被标识为删除状态，并未立即释放空间。这种处理对于获取多版本并发控制是必要的，如果一个元组的版本仍有可能被其他事务看到，那么就不能删除元组的该版本。当事务提交后，过期元组版本将对事务不再有效，因而其占据的空间必须回收以供其他新元组使用，以避免对磁盘空间增长的无休止的需求，此时对数据库的清理工作通过运行 VACUUM 来实现。从 PostgreSQL 8.1 开始，PostgreSQL 数据库引入一个额外的可选辅助进程 AutoVacuum（系统自动清理进程），自动执行 VACUUM 和 ANALYZE 命令，回收被标识为删除状态记录的空间，更

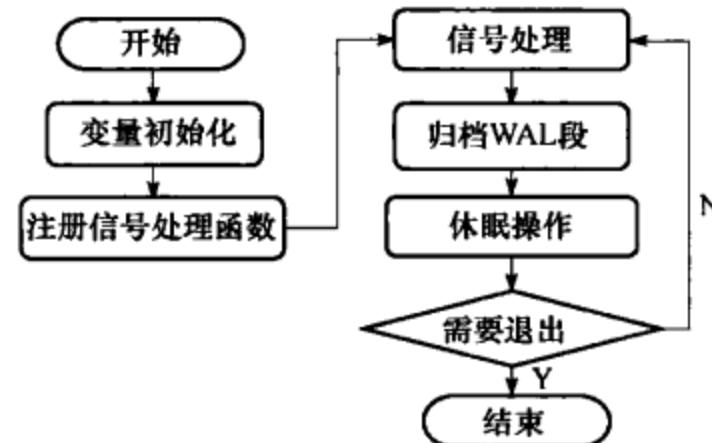


图 2-11 PgArch 预写日志归档进程处理流程

新表的统计信息。

在 PostgreSQL 数据库系统配置文件中，与系统自动清理相关的主要相关参数如下：

- autovacuum：是否启动系统自动清理功能，默认值为 on。
- autovacuum_max_workers：设置系统自动清理工作进程的最大数量。
- autovacuum_naptime：设置两次系统自动清理操作之间的间隔时间。
- autovacuum_vacuum_threshold 和 autovacuum_analyze_threshold：设置当表上被更新的元组数的阈值超过这些阈值时分别需要执行 vacuum 和 analyze。
- autovacuum_vacuum_scale_factor 和 autovacuum_analyze_scale_factor：设置表大小的缩放系数。
- autovacuum_freeze_max_age：设置需要强制对数据库进行清理的 XID 上限值。

AutoVacuum 系统自动清理进程中包含两种不同的处理进程：AutoVacuum Launcher 和 AutoVacuum Worker。AutoVacuum Launcher 进程为监控进程，用于收集数据库运行信息，根据数据库选择规则选中一个数据库，并调度一个 AutoVacuum Worker 进程执行清理操作。在 AutoVacuum Launcher 进程中，选择数据库的规则如下：首先由于数据库事务 XID 是 32 位整数且递增分配，当超过最大值时会从头开始计数使用，而事务 XID 的大小表示事务开始的时间，事务 XID 重新计数使用会使数据库中部分事务数据丢失，因此当 XID 超过配置的 autovacuum_freeze_max_age 时，强制对该数据库进行清理并更新事务 XID；其次，若无强制清理操作，则选择数据库列表中最早未执行过自动清理操作的数据库。Launcher 进程会定时（或者被信号驱动）选择数据库并调度 Worker 进程去执行清理工作。

AutoVacuum 中的 AutoVacuum Worker 进程执行实际的清理任务，Launcher 进程中维护有 Worker 进程列表。Worker 进程列表由三种不同状态的进程列表构成，即空闲的 Worker 进程列表、正在启动的 Worker 进程、运行中的 Worker 进程列表。Launcher 进程在不同状态之间的切换实现了 Worker 进程的调度工作。首先，在初始化阶段创建的运行内存上下文中，创建长度为 autovacuum_max_workers 的空闲 Worker 进程描述信息列表，而正在启动 Worker 进程和运行中的 Worker 进程的列表被置为空。如果 Launcher 进程需要一个 Worker 进程，空闲 Worker 进程列表不为空且当前没有正在启动中的 Worker 进程，则开始一个启动 Worker 进程的操作，即向 Postmaster 进程发送启动消息，从空闲 Worker 进程列表中取出一个进程描述信息，设置为启动中状态。Launcher 进程中只允许存在一个启动中状态的 Worker 进程，启动中的 Worker 进程如果超时（超时时间由 autovacuum_naptime 设置）将被取消并重新开始启动 Worker 进程的循环。如果 Worker 启动成功，将启动成功的 Worker 进程信息添加到运行中的 Worker 进程列表中。运行中的 Worker 进程即连接上根据规则选中的数据库。

在启动成功的 Worker 进程连接数据库成功后，将遍历该数据库中的表，根据对表的清理规则选择要执行的表和在该表上执行的操作。对表的操作分为 VACUUM 和 ANALYZE 两种，对选中的表如果上次 VACUUM 之后的过期元组的数量超过了“清理阈值”（vacuum threshold），那么就清理该表，清理阈值是定义为：

$$\text{清理阈值} = \text{清理基本阈值} + \text{清理缩放系数} * \text{元组数}$$

这里的清理基本阈值是 autovacuum_vacuum_threshold，清理的缩放系数是 autovacuum_vacuum_scale_factor，元组的数目可以从统计收集器里面获取。这是一个部分精确的计数，由每次 UPDATE 和 DELETE 操作更新。

如果表上次被执行 ANALYZE 操作之后，其中过期元组的数量超过了“分析阈值”(analyze threshold)，那么就分析该表更新表统计信息，分析阈值的定义与清理阈值相似，定义如下：

$$\text{分析阈值} = \text{分析基本阈值} + \text{分析缩放系数} * \text{元组数目}$$

缺省的阈值和伸缩系数都是从 postgresql.conf 里面取得的。不过，我们可以以每个表独立设置的方式覆盖它，方法就是在系统表 pg_autovacuum 里输入信息。pg_autovacuum 表中一个元组可以用来记录一个需要自动清理的表及其清理设置，AutoVacuum 进程将使用其中的清理设置来清理该表。如果没有特别设置该表的清理设置，AutoVacuum 将使用全局设置。

1. AutoVacuum Launcher 进程

通常，在入口函数 StartAutoVacLauncher 中执行 fork 操作创建 Postmaster 的子进程 AutoVacuum Launcher，在新创建的子进程中执行体中关闭从 Postmaster 进程中复制出的网络连接端口，同时进入进程 AutoVacuum Launcher 的执行体函数 AutoVacLauncherMain，其处理流程如图 2-12 所示。

下面对其中几个重要的步骤进行说明：

1) 构建数据库列表：调用函数 rebuild_database_list 完成，其步骤如下：

①建立一个 Hash 表，其中每一个元素代表一个数据库，记录了该数据库的 OID (adl_datid)、启动 worker 的时间戳 (adl_next_worker) 以及一个评分值 (adl_score)。初始时该 Hash 表中没有元素。

②将 pg_database 平面文件（在 PGDATA/global 目录下）中的数据库构成一个链表，链表中的每

一个节点代表一个数据库，其中包括数据库的 OID、名称、该数据库的统计信息等。

③调用 pgstat_fetch_stat_dbentry 来填充每个节点的统计信息。

④对每一个统计信息不为空的数据库，在 Hash 表中搜索该数据库，如果没有找到则将该数据库加入到 Hash 表中，并且将该数据库的 adl_score 设置为该数据库被加入到 Hash 表中时的顺序号。

⑤将 Hash 表中的数据库按照 adl_score 值升序的顺序依次加入到全局变量 DatabaseList 所指向的链表中，并设置每一个数据库的 adl_next_worker 值。其中第一个数据库的 adl_next_worker 值设为当前时间，之后的每一个数据库的 adl_next_worker 的值都比前一个增加 millis_increment。millis_increment 的值由 autovacuum_naptime 参数值除以 Hash 表中数据库的个数来设定。

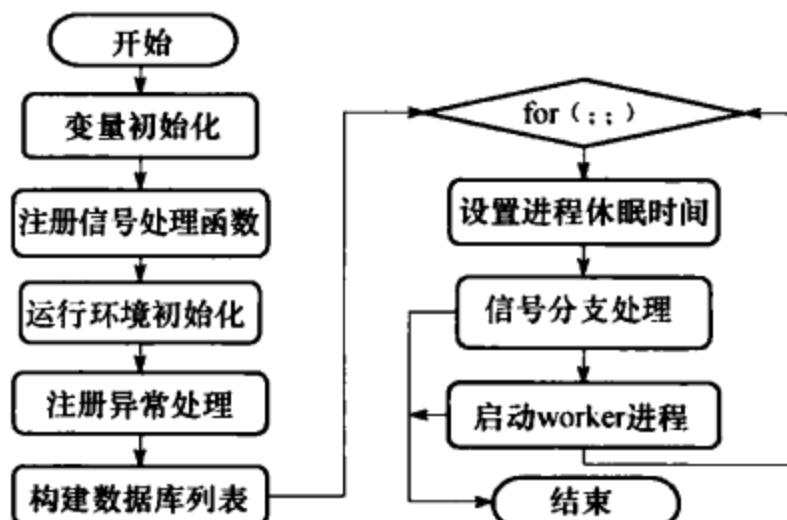


图 2-12 AutoVacuum Launcher 处理流程

为什么使用平面文件？

在 PostgreSQL 的数据集簇中，提供了两个平面文件：pg_database 和 pg_auth。这两个文件分别记录了 pg_database 系统表和 pg_authid 系统表中的部分信息。如果有些还没有启动完毕的后台进程需要访问这两个系统表的内容，它们将会使用两个平面文件来进行信息的获取，这是由于进程还未完全启动时是无法连接到数据库并读取相关系统表内容的。

在“构建数据库列表”这一步骤由于并未连接到数据库，因此只能用平面文件来替代系统表 pg_database。

2) 设置进程休眠时间：根据空闲 Worker 和数据库列表来计算休眠的时间，当所有 Worker 进程都在运行时要设置一个较长的休眠时间。而当 Worker 进程退出时可以唤醒休眠，同时休眠也可以被其他信号中断。

3) 信号处理分支中 got_SIGUSR1 信号通知有 Worker 进程退出或者 Postmaster 通知有 Worker 启动失败。若是 Postmaster 通知 Worker 启动失败，则给 Postmaster 重新发送启动 Worker 进程的消息。

4) 启动 Worker 进程：如果当前有一个 Worker 正在启动中，则再休眠一会儿等待该 Worker 启动完成。如果可以开始启动一个新的 Worker，则进行以下判断：

①如果数据库列表不为空，则检查 DatabaseList 尾部数据库的 adl_next_worker 参数，如果早于当前时间（表示该数据库早就应该被处理）则启动 Worker 进程。

②数据库列表为空时，立即启动 Worker 进程。

2. AutoVacuum Worker 进程

AutoVacuum Worker 进程的入口为 launch_worker 函数，在该入口处调用 do_start_worker 创建 worker 进程，并且返回连接数据库的 OID。如果返回的 OID 为有效的数据库 OID，则遍历数据库列表找到该 OID 对应的数据库在数据库列表中对应的节点，更新该节点的 adl_next_worker 域值，并将该节点移动到数据库列表的头部。如果遍历数据库列表没有对应于该 OID 的节点，则调用 rebuild_database_list 重建数据库列表。创建 worker 进程的函数体 do_start_worker 处理流程如图 2-13 所示。

AutoVacuum Worker 进程的处理流程和 AutoVacuum Launcher 进程的处理流程基本类似，选择要进行清理的数据库的规则如前所述：选中数据库后遍历数据库中的表，根据表的统计信息计算清理阈值和分析阈值，来确定是否要对表执行相应的操作。

在系统进行自动清理的同时，用户可以使用安装目录 bin 文件夹下的 vacuumdb 或者 vacuumlo 工具对数据进行手动清理工作。vacuumdb 工具清理数据库并对数据库执行分析操作，vacuumlo 工具清理数据库中无效的大对象。

2.5.6 PgStat 统计数据收集进程

PgStat 辅助进程是 PostgreSQL 数据库系统的统计信息收集器，它专门负责收集数据库系统运行中的统计信息，如在一个表和索引上进行了多少次插入与更新操作、磁盘块的数量和元组的数量、每个表上最近一次执行清理和分析操作的时间，以及统计每个用户自定义函数调用执行的时间等。由于统计数据收集给查询处理增加了一些负荷，所以可以把系统配置为收集信息，也可以配置为不收集信息。系统表 pg_statistic 中存储了 PgStat 收集的各类统计信息，另外在数据库簇的目录下有与统计信息收集器相关的文件：global 子文件夹下的 pgstat.stat 文件用于保存当前全局的统计信息；pg_stat_tmp 文件则是 PgStat 进程和各个后台进程进行交互的临时文件所在地。

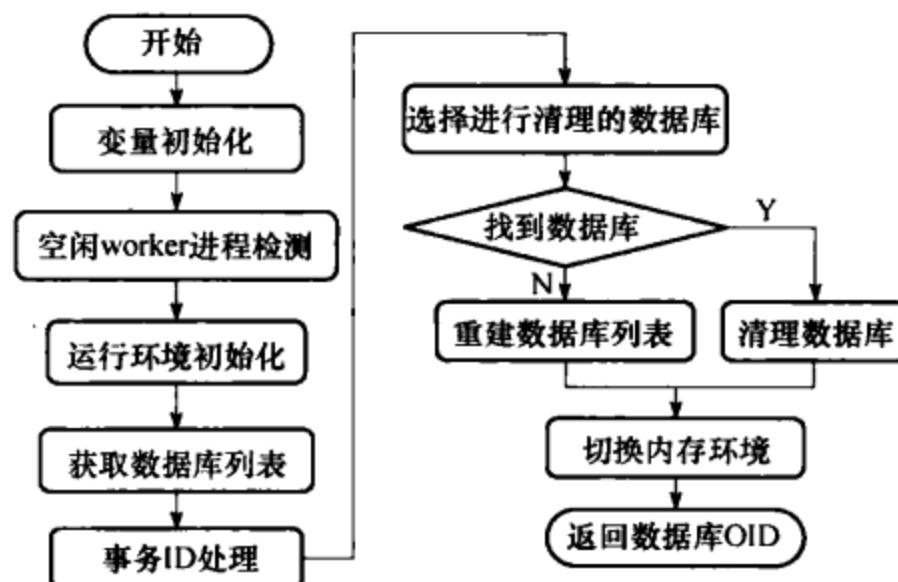


图 2-13 AutoVacuum Worker 进程处理流程

PgStat 辅助进程收集的统计信息主要用于查询优化时的代价估算。在 PostgreSQL 的查询优化过程中，查询请求的不同执行方案是通过建立不同的路径（Path）来表达的。在生成了许多符合条件的路径之后，从中选择出代价最小的路径转化为一个计划，这个计划将被传递给执行器执行。因此优化器的核心工作就是建立许许多多的路径，然后从中找出最优的路径。造成同一个查询请求有不同路径的主要原因是：表不同的访问方式（如顺序访问（Sequential Access）、索引访问（Index Access），PostgreSQL 中还有可能使用 TID 直接访问元组）；表间不同的连接方式（嵌套循环连接（Nest-loop join）、归并连接（Merge Join）、Hash 连接（Hash Join））；表间不同的连接顺序（左连接（Left-join）、右连接（Right-join）、布希连接（Bushy-join））。而评价路径优劣的依据是用系统表 pg_statistic 中的系统统计信息估计出的不同路径的代价。

在 PostgreSQL 数据库系统配置文件 postgresql.conf 中与 PgStat 相关的配置选项有：

- track_activities：表示是否对会话中当前执行的命令开启统计信息收集功能，该参数只对超级用户和会话所有者可见，默认值为 on（开启）。
- track_counts：表示是否对数据库活动开启统计信息收集功能，由于在 AutoVacuum 自动清理进程中选择清理的数据库时，需要数据库的统计信息，因此该参数默认值为 on。
- track_function：表示是否开启函数的调用次数和调用耗时统计。
- track_activity_query_size：设置用于跟踪每一个活动会话的当前执行命令的字节数，默认值为 1024，只能在数据库启动后设置。

统计信息功能设置为启动状态时，在 Postmaster 进程启动数据库过程完成后，将通过 pgstat_init 函数对 PgStat 辅助进程进行初始化操作，创建用于和后台进程之间通信的 UDP 端口。在 PgStat 辅助进程中调用 select 函数（或者 poll）来监听 UDP 端口上的数据变化，当后台进程通过 UDP 端口向统计收集进程发送统计消息时将触发监听过程，PgStat 进程根据消息类型进入相应的处理接口。

在 PgStat 辅助进程中定义了可以处理的消息类型，不同的消息类型在 PgStat 进程中定义了对应的描述结构（数据结构 2.8）。

数据结构 2.8 StatMsgType

```
typedef enum StatMsgType
{
    PGSTAT_MTYPE_DUMMY,           //空信息
    PGSTAT_MTYPE_INQUIRY,         //通知收集器写统计文件消息
    PGSTAT_MTYPE_TABSTAT,         //发送表和缓冲访问统计消息
    PGSTAT_MTYPE_TABPURGE,        //发送无效表的消息
    PGSTAT_MTYPE_DROPDB,          //发送删除的数据库信息消息
    PGSTAT_MTYPE_RESETCOUNTER,    //通知收集器重置技术器消息
    PGSTAT_MTYPE_AUTOVAC_START,   //发送一个数据库即将被清理的消息
    PGSTAT_MTYPE_VACUUM,          //发送 VACUUM 或者 VACUUM ANALYZE 执行完消息
    PGSTAT_MTYPE_ANALYZE,         //发送 ANALYZE 执行完消息
    PGSTAT_MTYPE_BGWRITER,        //bgwriter 通知更新统计信息消息
    PGSTAT_MTYPE_FUNCSTAT,         //发送函数使用统计信息消息
    PGSTAT_MTYPE_FUNCPURGE       //发送无效函数的消息
} StatMsgType;
```

对于不同类型的消息，有对应的数据结构表示。这里以删除数据库消息（数据结构 2.9）为例说明。

数据结构 2.9 PgStat_MsgDropdb

```
typedef struct PgStat_MsgDropdb
{
    PgStat_MsgHdr m_hdr;           //消息头部信息，记录消息的类型和长度
    Oid            m_databaseid;   //要删除的数据库的 OID
} PgStat_MsgDropdb;
```

其中，`m_databaseid` 是要删除的数据库的 OID，而 `m_hdr` 则是一个消息头（数据结构 2.10）。`m_hdr` 中记录了当前消息的类型以及整个消息的长度。

`PgStat` 进程被初始化之后，将进入 `pgstat_start` 函数。在 `pgstat_start` 中，首先会检查 `PgStat` 进程的初始化函数 `pgstat_init` 是否成功创建 UDP Socket 端口，如果创建失败，则退出 `PgStat` 进程；否则创建 `PgStat` 进程。最后进入 `PgStat` 进程的 `PgstatCollectorMain` 函数。其处理流程如图 2-14 所示。

数据结构 2.10 PgStat_MsgHdr

```
typedef struct PgStat_MsgHdr
{
    StatMsgType m_type;          //消息类型
    int         m_size;          //消息长度
} PgStat_MsgHdr;
```

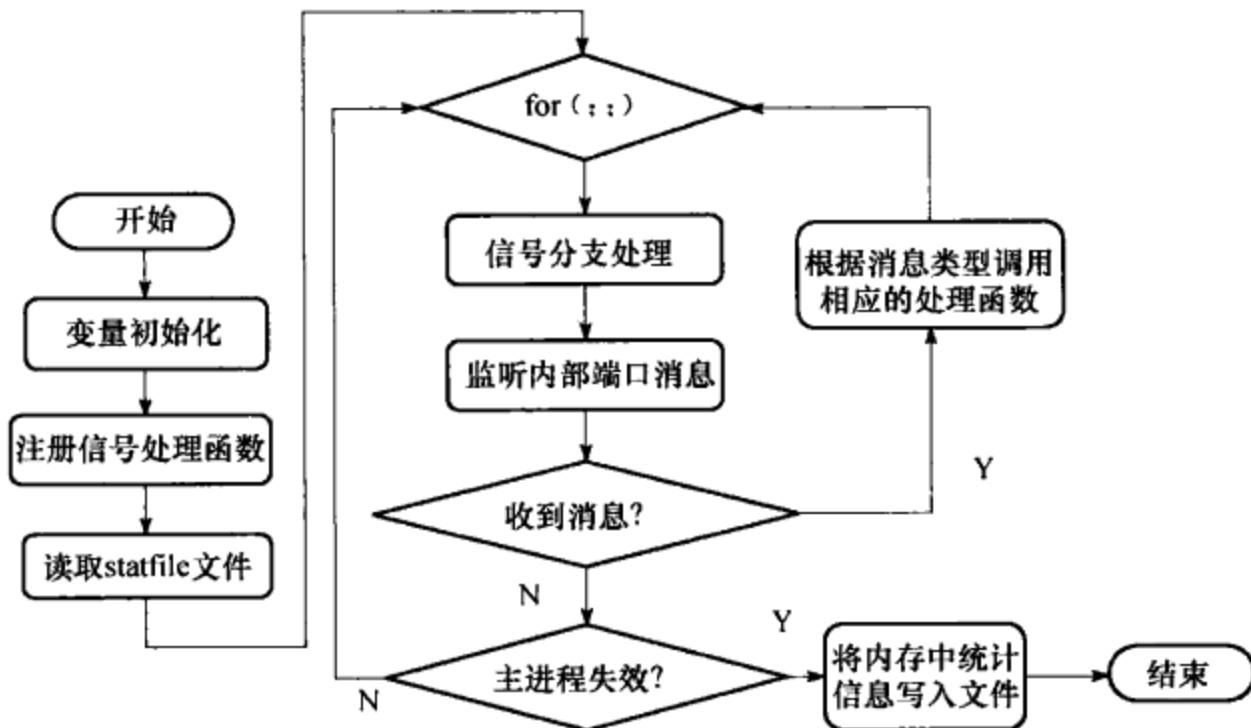


图 2-14 PgstatCollectorMain 处理流程

下面对几个重要步骤加以说明：

- 1) 在初始化阶段，会读取默认位置（数据集簇中的 global 子目录下）的统计信息文件 `pgstat.stat`。从统计文件中读取数据库、表、函数的统计信息，构建数据库统计信息的 Hash 表。若读取失败或是第一次启动，则各种统计信息计数器被设置为 0。
- 2) 如果最近更新统计文件时间小于最近申请读统计文件的时间，需要将统计信息临时文

件中最新的统计信息写入 pgstat.stat 文件中，并设置最近更新统计文件时间和最近申请读统计文件时间为当前系统时间；最近申请读统计文件的时间由 PGSTAT_MTYPE_INQUIRY 消息设置。

3) 在 PgStat 辅助进程的工作循环中，将在创建的 UDP Socket 端口上使用 select（或者 poll）进行监听。当端口上有可读数据时（即有其他后台进程发送了统计消息），则从端口中读取消息，对正确的消息（即消息中包含的消息类型是正确的）根据消息中的消息头包含的消息类型进入相应的消息处理分支，而错误的消息直接忽略掉。

4) 如果 UDP Socket 端口监听超时，将检查 Postmaster 状态，若 Postmaster 失效或者 PgStat 接收到退出消息，则退出 PgStat 进程。在退出前要将数据库、表、函数的统计信息写入 pgstat.stat 文件，并删除 pg_stat_tmp 中的临时文件。

对于 PgStat 收集的统计信息，系统提供了部分标准视图以供管理员查看；同时，还提供了统计信息相关的操作函数，供管理员来定义视图。需要注意的是，使用的统计信息并非实时的，每个活动的独立数据库服务仅仅在处理空闲时才发送信息到统计信息收集器，因此查询或者事务仍然在处理中并未更新当前的统计信息。在设置了 track_activities 参数时，当前查询信息在统计信息收集器中是实时的。

2.6 服务进程 Postgres

Postgres 进程是实际的接受查询请求并调用相应模块处理查询的 PostgreSQL 服务进程。它直接接受用户的命令进行编译执行，并将结果返回给用户。如此循环，直到用户断开连接。用户的命令分为两种：一种是查询命令，即插入、删除、更新和选择四种命令。另一种是非查询命令，如创建/删除表、视图、索引等命令。服务进程 Postgres 根据不同的命令类型选择不同的策略进行处理。

从前面 Main 主程序的执行流程图（见图 2-3）中可以看出，Postgres 的启动方式有两种：一种是在 Postmaster 监控下，动态地被 Postmaster 创建为用户服务，这是一种多用户的服务方式，也是比较通常的启动方式；另外一种是不经过 Postmaster 以单用户模式直接启动，为单一用户提供服务，这种方式由 --single 选项启动。在这种模式下，Postgres 服务器进程必须自己完成初始化内存环境、配置参数等操作，而这些操作在多用户模式下是由 Postmaster 服务器进程完成的。

Postgres 进程的主要源代码文件位于 src/backend/tcop 文件夹下，主要文件包括：

- 服务进程的源代码文件 postgres.c，它是 Postgres 的入口文件，负责管理查询的整体流程。
- 对于查询命令进行处理的源代码文件 pquery.c，它执行一个分析好的查询命令。
- 对于非查询命令进行处理的源代码文件 utility.c，它执行各种非查询命令。
- dest.c 中的代码主要处理 Postgres 和远端客户的一些消息通信操作，并负责返回命令的执行结果。

Postgres 进程的入口是位于 postgres.c 文件中的 PostgresMain 函数，其工作流程如图 2-15 所示。下面将对 PostgresMain 函数进行详细介绍。

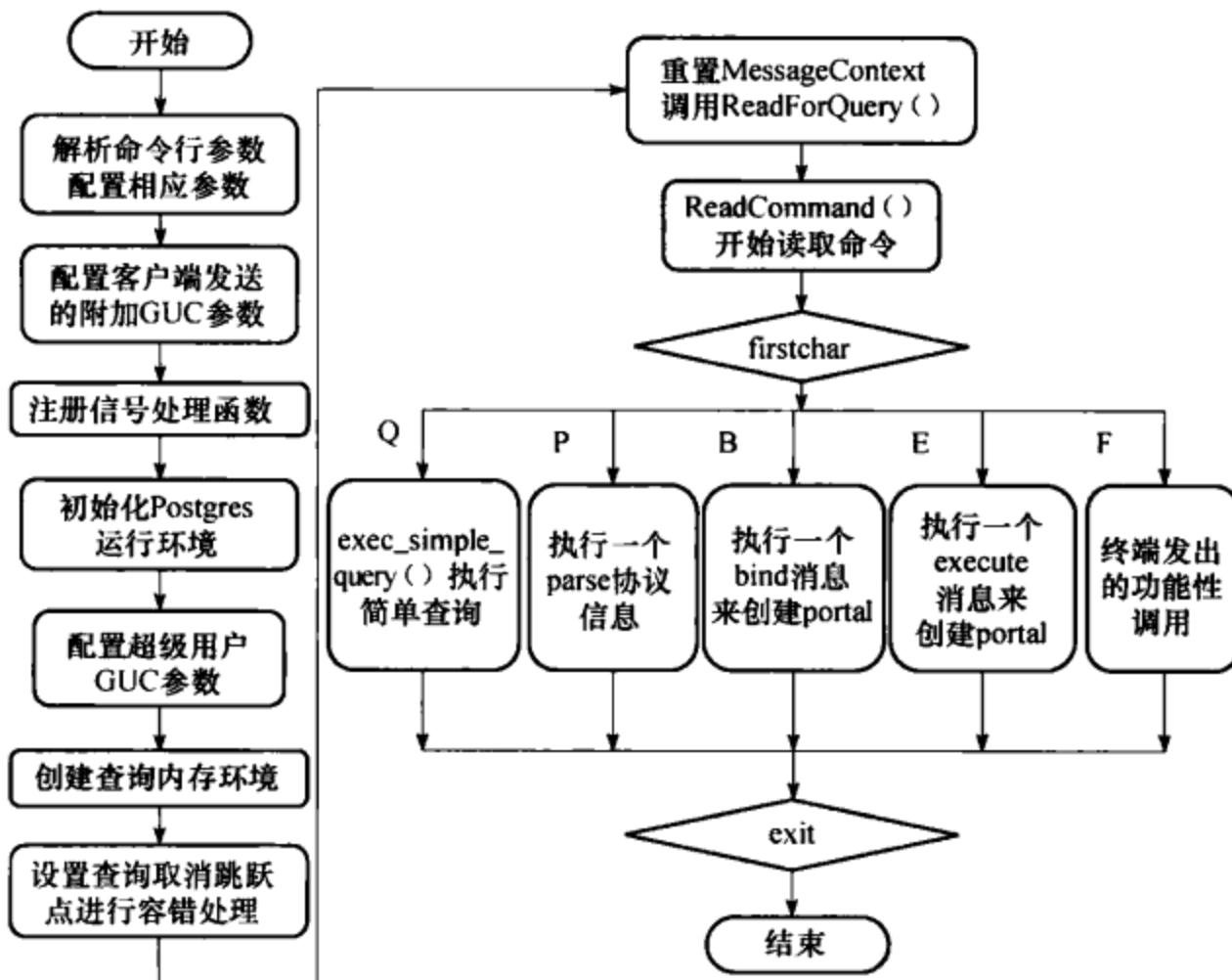


图 2-15 PostgresMain 函数流程图

2.6.1 初始化内存环境

前面讲到，Postgres 进程有两种运行方式。如果是在多用户模式下被 Postmaster 动态创建，则这一部分的工作由 Postmaster 进程来完成，在 Postgres 进程中此部分跳过。事实上，由于 Postmaster 进程和 Postgres 进程的特殊关系，Postgres 进程中会多次出现这样的情况。全局布尔变量 IsUnderPostmaster 表明进程的运行状态，该变量值为 true 表明运行在多用户模式下。Postgres 中内存环境的初始化也是调用 MemoryContextInit 来完成，这里不再赘述。

2.6.2 配置运行参数和处理客户端传递的 GUC 参数

和 Postmaster 一样，这里的配置参数也包括三个步骤，即将参数设置为默认值、根据命令行参数配置参数、读配置文件重新设置参数。其中，第一步和第三步在多用户模式下（即 IsUnderPostmaster 为 true 时）已由 Postmaster 完成，而在单用户模式下这三步都和 Postmaster 完全相同，故不再重复。命令行方式的参数设置优先级是高于缺省设置的，将会覆盖缺省设置。命令行参数有很多，和 Postmaster 一样，在 Postgres 中仍使用 while + switch 控制结构读入参数并进行配置。

之后 Postgres 还将从 Port 结构得到客户端传递的 GUC 选项，然后根据 GUC 选项的具体情况调用 SetConfigOption 进行设置。

2.6.3 设置信号处理和信号屏蔽

Postgres 中信号的注册方式和 Postmaster 类似，因此这里只对 Postgres 中处理的信号及其意义进

行介绍。Postgres 中信号和对应的信号处理函数如表 2-11 所示。

其中几个重要的信号处理函数的意义如下：

- **SIGHupHandler**: 当配置文件发生改变时产生此信号。服务进程在收到 SIGHUP 信号时重读配置文件 postgresql.conf，并重新装载 pg_hba.conf 和 pg_ident.conf 文件。
- **StatementCancelHandler**: 在收到 SIGINT 信号后调用此函数，终止正在进行的查询操作。若此时进程正在退出则什么也不做（proc_exit_in-progress 为真），否则置标志位 QueryCancelPending、InterruptPending 为真，表明准备处理查询取消中断。如果现在可以中断的话，先阻碍其他中断，调用 LockWaitCancel 函数判定是否在等待锁（若是则关闭计时器，将自己从等待队列删除，将进程信号量置 0），调用 ProcessInterrupts 由 elog 函数退回到在处理查询之前设置的跳跃点位置处理下一个查询。
- **die**: 该函数处理 SIGTERM 信号，用来终止当前的事务。若此时进程正在退出则什么也不做；否则置标志位 ProcDiePending 为真，表明准备处理进程退出中断，如果现在可以中断的话，先阻碍其他中断，调用 LockWaitCancel 函数判定是否在等待锁（若是则关闭计时器，将自己从等待队列删除，将进程信号量置 0），调用 ProcessInterrupts 由 elog 函数退出（进程退出的优先级大于取消查询）。
- **quickdie**: 处理 SIGQUIT 信号。首先屏蔽其他信号，然后结束正在进行的工作并退出。
- **handle_sig_alarm**: 该函数处理 SIGALRM 信号，这是由进程等待锁的时间超时引发的。如果存在死锁，就将自己从锁等待队列中退出，唤醒自己并在 PROC 结构中将错误类型置为 STATUS_ERROR。

2.6.4 初始化 Postgres 的运行环境

基本参数和信号量处理函数初始化完成之后，将进行 Postgres 进程运行环境的初始化工作。Postgres 进程首先会检查 DataDir 变量，确保给定的 DataDir 字符串是一个格式正确的数据目录路径。在多用户模式下，此工作由 Postmaster 进程完成。在有效数据目录路径下，检查 PG_VERSION 文件中的版本信息是否与当前版本的程序兼容。在版本兼容的前提下，将当前工作目录转到 DataDir 字符串表示的目录，以方便 Postmaster 进程及其他后台进程使用相对路径访问数据目录。

路径设置完成后，Postgres 的初始化工作分为两个阶段：首先调用 BaseInit 函数来完成基本初始化，之后调用 InitPostgres 来完成 Postgres 的初始化。之所以将 BaseInit 和 InitPostgres 划分为两个阶段，是因为与 XLog 相关的初始化工作必须在 InitPostgres 之前。

在 BaseInit 中，第一步调用 InitCommunication 创建共享内存和信号量并进行初始化，第二步调用 DebugFileOpen 初始化 input/output/debugging 文件描述符，第三步调用 Init FileAccess 初始化文件访问，第四步调用 Smgrinit 初始化或者关闭存储管理器，最后调用 InitBufferPoolAccess 初始化共享

表 2-11 Postgres 信号及其处理函数

信号	信号处理函数
SIGHUP	SigHupHandler
SIGINT	StatementCancelHandler
SIGTERM	die
SIGQUIT	quickdie
SIGALRM	handle_sig_alarm
SIGPIPE	SIG_IGN
SIGUSR1	CatchupInterruptHandler
SIGUSR2	NotifyInterruptHandler
SIGFPE	FloatExceptionHandler
SIGCHLD	默认

缓冲区存储器。在整个 BaseInit 阶段，将完成 Postgres 进程的内存、信号量以及文件句柄的创建和初始化工作。

在 InitPostgres 中，以数据库名或者数据库 OID、用户名为参数，如果参数给出的是数据库的 OID，该函数还将把对应的数据库名返回给调用者，在 bootstrap 模式不需要参数。在 InitPostgres 中获取数据库的路径并设置该路径的环境变量，完成 PGPROC 结构的填充并将其添加到 ProcArray，即可使得其对其他后台可见。填充 PGPROC 结构后，开始一系列的初始化操作，包含初始化后端缓冲池、Xlog 访问、关系 Cache、系统表 Cache、查询计划 Cache、Portal 管理器、状态收集器以及退出码的设置。完成初始化工作后，即开始一个新的事务并创建相应的锁对象，再次检查数据库对象保证数据操作安全。

2.6.5 创建内存上下文并设置查询取消跳跃点

Postgres 首先创建一个名为 MessageContext 的内存上下文，该内存上下文用于存储从前端发送过来的消息中的查询命令，以及在查询过程中产生的中间数据（例如在简单查询模式时产生的分析树和计划树），每当 PostgresMain 进行下一次循环（即进行新的查询）时该内存上下文将被重设，即所有经过 MessageContext 分配的内存块将被释放。

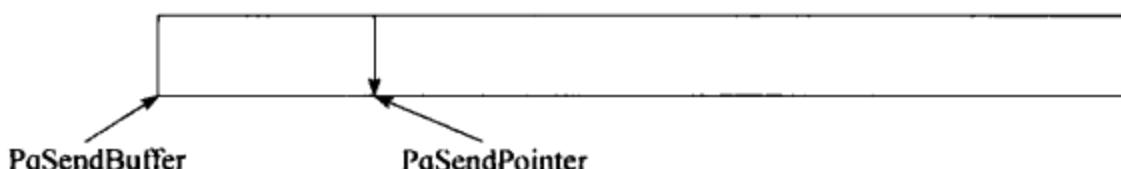
创建完成后将调用 sigsetjmp（系统调用函数）设置跳跃点，当客户端取消一次查询请求或发生错误时，将通过调用 siglongjmp 利用全局变量 PG_exception_stack（该变量是指向跳跃点的指针）从这个点退出当前事务然后重新开始查询，在错误恢复期间不允许被中断，也不接受任何客户端取消查询的请求（正在取消）。

PostgreSQL 使用的这种错误恢复的方式是 C 语言程序中进行错误处理的一种常见方式，具体细节可以参考 C 语言的相关资料。

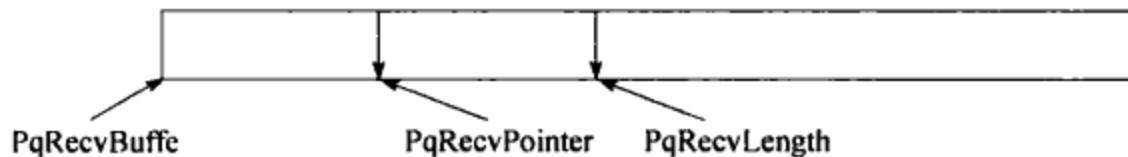
2.6.6 循环等待处理查询

完成上述工作后，Postgres 可以真正接受客户端查询并处理，系统开始运行。在介绍具体实现之前，先介绍一下在服务进程与用户进程通信的过程中所涉及的数据结构：

- 共享缓冲区 PqSendBuffer：这是一个大小为 8192 字节的输出缓冲区，用来接收服务器发送给客户端的消息。全局变量 PqSendPointer 是一个字符指针，表示在 PqSendBuffer 缓冲区中在此地址之前的消息正等待发送给客户端，以后新加入的消息放在 PqSendPointer 所指向的位置。



- 共享缓冲区 PqRecvBuffer：它也是一个大小为 8192 字节的输入缓冲区，用来接受客户端发给服务器的消息。全局变量 PqRecvLength 是一个字符指针，表示下一个发送过来的消息在缓冲区中应存放的位置，PqRecvPointer 表示在此偏移地址之前的消息已处理，PqRecvPointer 和 PqRecvLength 之间的消息还未被处理。初始化的时候 PqRecvLength 和 PqRecvPointer 大小均为 0。



Postgres 使用一种基于消息的协议用于服务器和客户端之间通信。所有通信都是通过一个消息流进行的。消息的第一个字节标识消息类型，接下来四个字节声明消息剩下部分的长度（这个长度包括长度域自身，但不包括消息类型字节）。剩下的消息内容由消息类型决定。主要有以下几种消息：

- 启动消息：要开始一个会话，客户端需要打开一个与服务器的连接并且发送一个启动消息。这个消息包括用户名以及用户希望与之连接的数据库；它还标识要使用的协议版本。然后服务器就使用这些信息以及配置文件的内容（比如 `pg_hba.conf`）来判断这个连接是否可以接受，以及需要什么样的额外的认证。
- 简单查询：客户端发送一条简单查询消息给服务器请求开始一个查询。该消息包含一个用文本字符串表达的 SQL 命令（或者一组命令）。服务器根据查询命令字符串的内容发送一条或者更多条响应消息给客户端，并且最后是一条 `ReadyForQuery` 响应消息。该消息通知客户端它可以安全地发送新命令了（也表示服务器已经处理完之前发送的命令）。
- 扩展查询：扩展查询协议把前面的简单查询分割成若干个步骤，准备的步骤可以多次复用以提高效率。另外，还可以获得额外的特性。在扩展的协议里，前端发送一个 `Parse` 消息，它包含一个文本查询字符串，还有一些有关参数占位符的数据类型的信息，以及一个最终准备好的语句对象的名字。
- 函数调用：函数调用允许客户端请求对数据库 `pg_proc` 系统表中的一个任意函数的直接调用，但是用户必须具有在该函数上的执行权限。
- 取消正在处理的请求：客户端可以向服务器发送 `CancelRequest` 消息来取消正在处理的请求。服务器收到这类消息后将处理这个请求然后关闭连接。出于安全原因，对取消请求消息不做直接的响应。
- 终止：通常终止过程的方法是客户端发送一条 `Terminate`（终止）消息并且立刻关闭连接。一旦收到消息，服务器将马上关闭连接并且退出。

在第一次进入循环时，服务进程首先会释放上次查询（循环）时的内存，并为新的查询（循环）分配内存，准备好查询执行环境。接着调用 `ReadyForQuery` 函数给客户端发送一条消息告诉客户端它已经准备好接收查询了。接下来服务进程就可以调用 `ReadCommand` 函数开始从客户端接收消息了。

客户端可以有两种方式递交请求：一种是通过网络连接，这时候我们调用 `SocketBackend` 函数接收客户端请求；另一种是客户端和服务器在同一台机器上，这时候我们调用 `InteractiveBackend` 函数从交互终端读取客户端请求。

在网络连接的方式下，`SocketBackend` 首先检查 `PqRecvBuffer` 中是否存在未处理的消息，如果没有则先要通过网络连接从客户端读取消息到 `PqRecvBuffer` 中，然后取出第一条未处理的消息的消息类型。接着 `SocketBackend` 根据消息类型设置相应的全局变量，例如当消息类型为 Q 的时候表明该消息为简单查询，这时候会将全局变量 `doing_extended_query_message` 设置为 `false`，然后函数从 `PqRecvBuffer` 中读取一条完整的消息到 `inBuf` 中。最后函数返回消息的消息类型。

在终端方式下，InteractiveBackend 负责从终端读取用户输入的一条完整的消息到 inBuf 中，此时消息的消息类型设置为 Q，表示一个简单消息。

根据前面函数返回的消息类型，服务进程就可以调用相应的函数来具体处理这些消息了。例如，对类型为“Q”的查询（简单查询），Postgres 会调用 exec_simple_query 函数来处理，事实上，几乎所有我们用到的 SELECT、UPDATE、DELETE 和 INSERT 语句都是从这里开始执行的。

2.6.7 简单查询的执行流程

一般的数据操作语言（DML）命令都是作为简单查询来处理的，系统调用 exec_simple_query 函数来执行简单查询，该函数的执行流程如图 2-16 所示。

exec_simple_query 函数是 Postgres 后台实际执行 DML 语句的函数。它解读用户输入的字符串形式的命令并检查合法性，最终执行并返回结果。《数据库系统实现》^①中所提到的各个主要的数据库模块都是在这个函数中调用的，包括编译器、分析器、优化器和执行器，它们相互协作完成各种实际的处理工作，构成了服务进程的主体。这些模块的内部都十分复杂，下面大致说明它们的主要功能和代码位置。

- 编译器：编译器是主流程中的第一个模块。它的作用是扫描用户输入

的字符串形式的命令，检查其合法性，并将其转换为 Postgres 定义的内部数据结构。Postgres 为每一条 SQL 命令都定义了相应的 C 语言结构体，用来存放命令中的各种参数。编译器是利用著名的 lex 和 yacc 工具编写的，其入口为 pg_parse_query 函数。它的代码位于 src/backend/paser 目录下的 scan.l 和 gram.y 文件中。

- 分析器：分析器接收编译器传递过来的各种命令数据结构（语法树），对它们进行相应的处理，最终转换为统一的数据结构 Query。如果是查询命令，在生成 Query 后进行规则重写（rewrite）。重写部分的入口是 QueryRewrite 函数，代码位于 src/backend/rewrite 目录下。分析器的入口是 parse_analyze 函数，其代码位于 src/backend/paser 目录下。
- 优化器：优化器接收分析器输出的 Query 结构体，进行优化处理后，输出执行器可以执行的计划（Plan）。一个特定的 SQL 查询可以以多种不同的方式执行，优化器将检查每个可能的查询计划，最终选择运行最快的查询计划。优化器的入口是 pg_plan_query 函数，代码位于 src/backend/optimizer 目录下。
- 执行器：执行非查询命令的入口函数是 ProcessUtility，代码位于 src/backend/tcop/utility.c 中，该函数的主体结构是一个 switch 语句，根据输入的命令类型调用相应的执行函数。执行查询命令的入口函数是 ProcessQuery，代码主要位于 src/backend/executor 目录下。

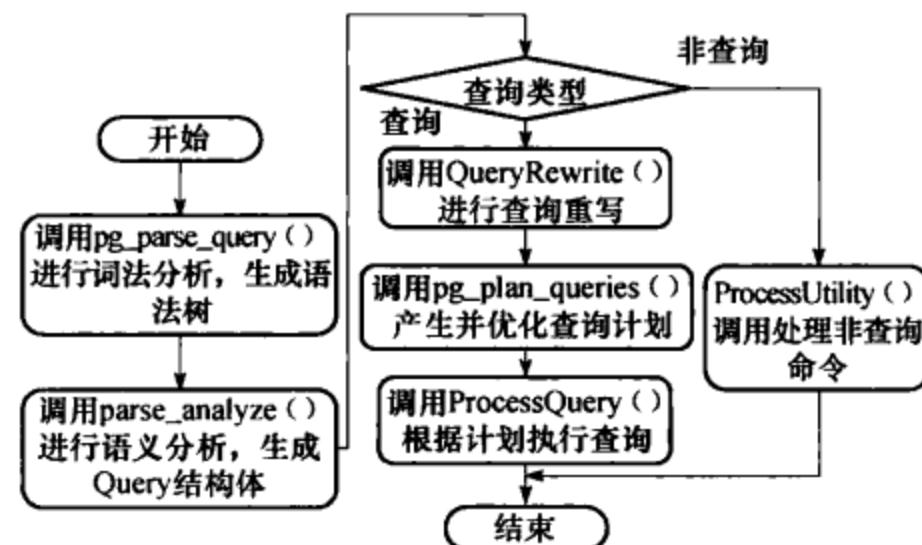


图 2-16 简单查询流程图

^① 该书由机械工业出版社于 2010 年引进出版。——编辑注

Postgres 进程在系统中扮演着一个工作执行者的角色，在单用户或者多用户模式下，客户端请求通过认证后将直接与服务进程 Postgres 进程通信，而无须守护进程干预，只在客户端对应的后台进程出现问题时由守护进程执行容错恢复工作。Postgres 和用户进行交互，执行客户端提交的查询请求和命令，并将执行结果通过网络返回给用户。Postgres 后台进程的运行即实现了 PostgreSQL 的多任务并发执行。

2.7 小结

本章从宏观上对 PostgreSQL 的控制和处理流程进行了简要介绍，说明了各个模块之间是如何协同工作，以使得整个数据库系统能够稳定、正确地处理用户的各种操作和请求的。至于每个模块是如何各司其职，其内部具体是如何运作的，将会在后续的章节进行专门的介绍。

第3章

存储管理

数据库管理系统的任务本质上是向存储设备上写入数据或者从存储设备上读出数据，因此对于一个DBMS来说，存储的管理是一项非常基础和重要的技术。在PostgreSQL中，有专门的模块负责管理存储设备（包括内存和外存），我们称之为存储管理器。存储管理器提供了一组统一的管理外存与内存资源的功能模块，所有对外存与内存的操作都将交由存储管理器处理，可以认为存储管理器是数据库管理系统与物理存取设备的接口。与PostgreSQL其他模块相比，存储管理器处在系统结构的底层，它包含了操作物理存取设备的接口。本章主要介绍存储管理器的功能和实现。

3.1 存储管理器的体系结构

PostgreSQL的存储管理器主要包括两个功能：内存管理和外存管理。除了管理内存和外存的交互外，存储管理器的另一个主要任务是对内存进行统筹安排和规划。

存储管理器的体系结构如图3-1所示。

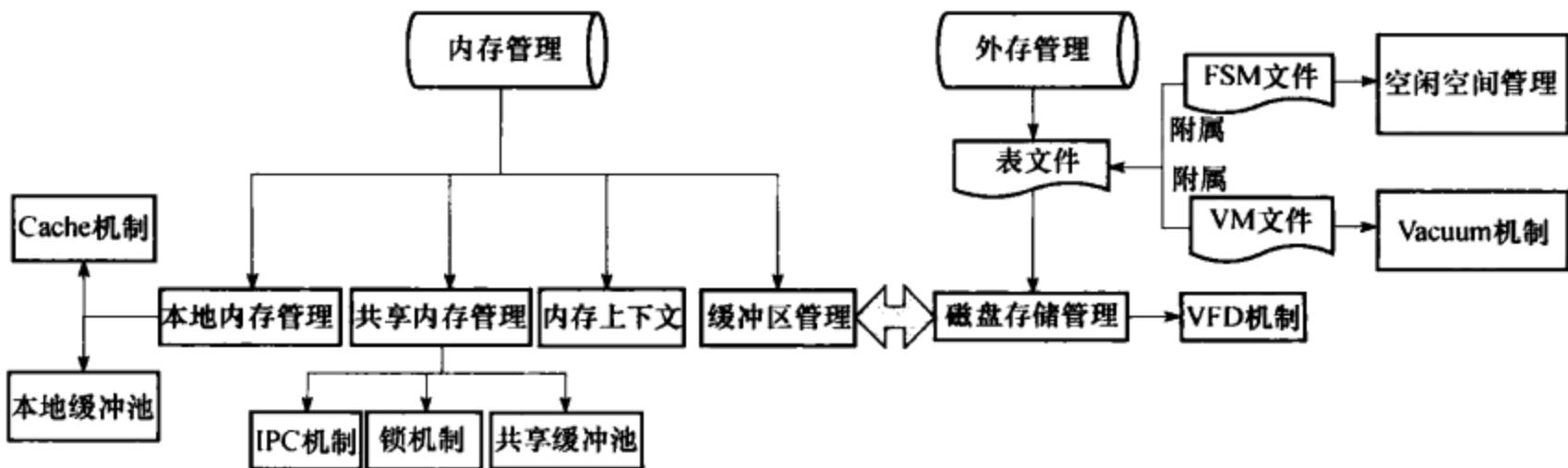


图3-1 存储管理器的体系结构

内存管理包括共享内存的管理以及进程本地内存的管理。在共享内存中存储着所有进程的公共数据，例如锁变量、进程通信状态、缓冲区等。而本地内存为每个后台进程所专有，是它们的工作

区域，存储着属于该进程的 Cache（高速缓存）、事务管理信息、进程信息等。为了防止多个进程并发访问共享内存中数据时产生冲突，PostgreSQL 提供了轻量级锁，用于支持对共享内存中同一数据的互斥访问。PostgreSQL 使用共享内存实现了 IPC（进程间通信）以及无效消息共享，用以支持进程间的相互通信。此外，存储管理器还提供内存上下文（MemoryContext）用于统一管理内存的分配和回收，从而更加有效安全地对内存空间进行管理。

外存的管理包括表文件管理、空闲空间管理、虚拟文件描述符管理及大数据存储管理等。在 PostgreSQL 中，每个表都用一个文件（表文件）存储，表文件以表的 OID 命名。对于超出操作系统文件大小限制（比如 FAT32 限制为 4G）的表文件，PostgreSQL 会自动将其切分成多个文件来存储，并在原表文件名的尾部加上切分文件的顺序号来标识它们。从 PostgreSQL 8.4 版本开始，每个表除了表文件外还拥有两个附属文件：可见性映射表文件（VM）和空闲空间映射表文件（FSM）。前者用于加快清理操作（VACUUM）的执行速度，后者则用于表文件空闲空间的管理。为了避免超过操作系统对每个进程打开文件数的限制，存储管理器使用了虚拟文件描述符机制，使得后台进程可以打开“无限多个”文件。此外，存储管理器还提供了大对象机制以及 TOAST 机制，用以支持大数据存储。PostgreSQL 还在存储管理选择器中封装了对具体存储器的操作接口，以便扩展支持多种存储介质（比如对闪存、光盘的支持）。

PostgreSQL 的存储管理器采用与操作系统类似的分页存储管理方式，即数据在内存中是以页面块的形式存在。每个表文件由多个 BLCKSZ（一个可配置的常量）字节大小的文件块组成，每个文件块又可以包含多个元组。表文件以文件块为单位读入内存中，每一个文件块在内存中形成一个页面块。由于页面块是文件块在内存中的存在形式，因此在后文中如不进行特殊说明也会使用页面来指代文件块。同样，文件的写入也是以页面块为单位。PostgreSQL 采用传统的行式存储，即以元组为单位进行数据的存储。一个文件块中可以存放多个元组，但是 PostgreSQL 不支持元组的跨块存储，每个元组最大为 MaxHeapTupleSize。这样保证了每个文件块中存储的是多个完整的元组。

与操作系统一样，PostgreSQL 在内存中开辟了缓冲区域用于存储这些文件块，我们将其在内存中开辟的缓冲区域称为缓冲池，缓冲池被划分成若干个固定大小（和文件块的尺寸相同，也是 BLCKSZ）的缓冲区，磁盘上的文件块读入内存后被存放在缓冲区中，称之为页面块或者缓冲块。BLCKSZ 的默认值是 8192，因此一个标准缓冲块的大小默认为 8KB。

总体来说，存储管理器的主要任务包括：

- 1) 缓冲池管理：缓冲池在 PostgreSQL 中起缓存的作用。数据库中的事务常常需要频繁地存取数据，为了减少对磁盘的读写，在事务执行时，数据首先将会放入缓冲池中，PostgreSQL 设立了进程间共享的缓冲池（共享缓冲池）以及进程私有的缓冲池（本地缓冲池）。
- 2) Cache 机制：将进程最近使用的一些系统数据缓存在其私有内存中，其级别高于缓冲池。
- 3) 虚拟文件描述符管理：PostgreSQL 通过虚拟文件描述符（Virtual File Descriptor, VFD）来对物理文件进行管理，避免因为操作系统对进程打开文件数的限制出现错误。
- 4) 空闲空间管理：用于快速定位到表文件中的空闲空间以便于插入新数据，从而提高空间利用率。
- 5) 进程间通信机制（IPC）：PostgreSQL 是一个多进程的系统，IPC 用来在多个后台进程之间进行通信和消息的传递，比如使用消息队列来同步进程产生的无效消息，同时 IPC 还提供了对共享内存的管理。

6) 大数据存储管理：提供大对象和 TOAST 机制。大对象机制是一种由用户控制的大数据存储方法，它允许用户调用函数，通过 SQL 语句直接向表中插入一个大尺寸文件（图片、视频、文档等）。而 TOAST 机制则是在用户插入的变长数据超过一定限度时自动触发，用户无法对 TOAST 加以控制。

当一个 PostgreSQL 进程（Postmaster、Postgres 等）从数据库中读写一个元组的时候，对于以上各个功能模块的使用顺序可以使用图 3-2 来描述：

1) 读取一个元组数据时，首先需要获取表的基本信息（如表的 OID、索引信息、统计信息等）及元组的模式信息，这些信息被分散记录在多个系统表中。通常一个表的模式在设定好后的变化频率很低，因此在对同一个表的多个元组操作时，每次都去读取系统表元组来构建模式信息显然是没必要的，也会降低元组操作的效率。为了减少对系统表的访问，在每个进程的本地内存区域设置了两种 Cache，一种用来存储系统表元组，一种用来存储表的基本信息，从而让进程可以以很高的效率快速构建出表的基本信息和元组的模式结构。Cache 的具体实现分析将在 3.3.2 节介绍。

2) 获取表基本信息和元组的模式信息后，接下来就需要读取元组所在的文件块来获得元组的数据。PostgreSQL 进程将首先从共享缓冲池中查找所需文件对应的缓冲块（缓冲块是文件块被载入到缓冲区后的存在形式）。如果能够找到对应的缓冲块，则直接从中取得元组的数据并配合上一步得到的模式信息解析出元组的各个属性值。如果找不到

相应的缓冲块则要读入文件块，然后再从被调入的缓冲块中获取元组数据。PostgreSQL 除了设置共享缓冲池外，还为每个进程设置了本地缓冲池，用于缓存临时表数据及其他进程不可见数据。缓冲池的管理将在 3.3.3 节中介绍。

3) 如果缓冲池中没有包含所需元组的缓冲块，则需要从存储介质中读取，并写入到缓冲区中。缓冲区与存储介质的交互是通过存储介质管理器（SMGR）来进行的。不同存储介质的底层实现各有差异，SMGR 负责统管各种存储介质，并对上层的请求提供统一的接口，对下层则根据不同介质调用相应的介质管理器，因此 SMGR 是 PostgreSQL 外存管理部分的核心。在版本 8.4.1 中仅支持对磁盘的管理，磁盘管理器负责管理所有存储在磁盘上的文件的操作，包括创建、删除、读取等。磁盘管理器支持对大表（超过 2G）的管理。当表文件超过 RELSEG_SIZE 大小限制时，将会扩展一个新的表文件，假定表的 OID 为 12000，则该表的多个表文件被命名为 12000.1、12000.2……依此类推，每一个表文件的尺寸都不超过 RELSEG_SIZE。磁盘管理器将在 3.2.2 节中介绍。

4) 在磁盘管理器与物理文件之间还存在一层虚拟文件描述符（VFD）机制，这是为了防止进程打开的文件数超过操作系统限制而引起不可预知的错误。磁盘管理器操作的是 VFD，对于 PostgreSQL 的进程来说，VFD 的个数是无限的，即 VFD 机制通过合理使用有限个实际文件描述符来满足无限的 VFD 访问需求。而在 VFD 管理模块，通过维持一个 LRU 池来管理实际文件描述符

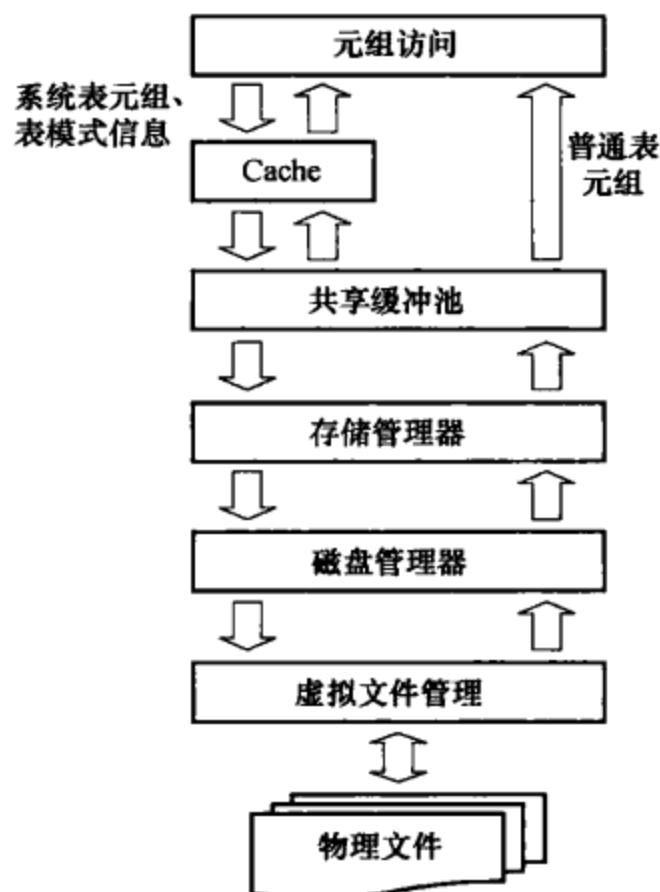


图 3-2 读写元组的过程

(FD)。由于系统允许的进程 FD 个数是有限的，因此当 LRU 池满时需要执行替换机制。VFD 机制将在 3.2.3 节介绍。

5) 在写入元组的时候，需要使用一定的策略找到一个具有合适空闲空间的缓冲块（这个过程和读取元组类似），然后将该元组装填到缓冲块中。系统会在适当的时候将这些被写入数据的缓冲块“刷”回到存储介质中。对于单个进程而言，应尽量将连续插入的数据放在一个缓冲块中，因此在表的基本信息（RelationData）中记录了最近执行插入操作的缓冲块块号。如果其中有足够的空闲空间，就使用该缓冲块，否则需要去寻找一个合适的文件块，将其读入内存并进行插入。显然，如果遍历所有的文件块查询其空闲空间，效率是很低的。为了加快查找的速度，在 PostgreSQL 8.4.1 中，为每个表增加了一个附属文件，即空闲空间映射表（FSM），用于记录每个表文件块的空闲空间大小，通过一定的查找机制和数据组织实现了文件块的快速选择。在 FSM 查找过程中，还提供了一些额外的特性，例如尽量使不同的进程选择不同的缓冲块以提高效率。FSM 的内容将在 3.2.4 节中介绍。

6) PostgreSQL 使用标记删除的方式来处理元组的删除，即对元组作上删除标记，而非从物理上删除元组。元组的物理清除工作将由 VACUUM 机制来完成，VACUUM 操作将会把这类元组所占用的空间置为可用，并且根据不同的 VACUUM 策略会选择是否将这些腾出来的空间还给操作系统。显然，在做 VACUUM 操作的时候去遍历所有文件块查找被删除的元组也是效率很低的，在版本 8.4.1 中同样也为表设计了一个附属文件来加快查找的速度，该文件被称为可见性映射表（VM）。VM 的相关内容在 3.2.5 节中介绍。

7) 除了一般的数据类型外，PostgreSQL 还支持大数据存储，大数据的存储使用 TOAST 机制和大对象机制来实现，前者主要用于变长字符串，后者则主要用于大尺寸的文件。在写入元组时，如果该元组的大小超过一定限制，将会自动触发 TOAST 机制对该元组进行相应的处理，使其大小满足系统限制。而大对象机制则由用户来调用，当用户需要在数据库中存放一个大的文件时，可以调用 PostgreSQL 提供的接口函数创建一个大对象并获得一个大对象 OID，用户通过该大对象 OID 来存储并使用大对象。大数据存储将在 3.2.6 节介绍。

3.2 外存管理

外存管理负责处理数据库与外存介质（在 PostgreSQL 中只实现了磁盘的管理操作）的交互过程。在 PostgreSQL 中，外存管理由 SMGR（主要代码在 smgr.c 中）提供对外存操作的统一接口，其结构如图 3-3 所示。SMGR 负责统管各种介质管理器，会根据上层的请求选择一个具体的介质管理器进行操作。

每个表文件在磁盘中都以一定的结构进行存储，针对磁盘，外存管理模块提供了磁盘管理器和 VFD 机制。在 PostgreSQL 8.4.1 版本中，还为每个表文件创建了两个附属文件，即空闲空间映射表文件（FSM）和可见性映射表文件（VM）。另外，对于大数据存储，PostgreSQL 也提供了两种处理

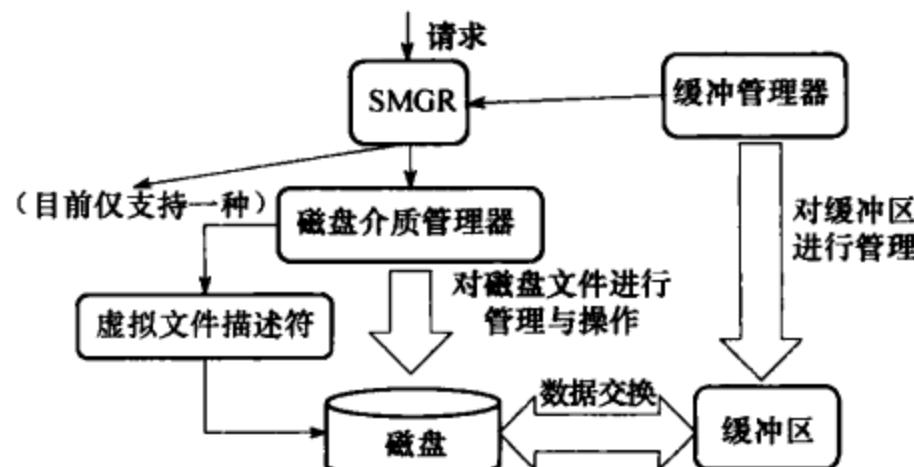


图 3-3 外存管理体系结构

机制。下文将就上述几个模块进行详细介绍。

3.2.1 表和元组的组织方式

在 PostgreSQL 中，同一个表中的元组按照创建顺序依次插入到表文件中（在 3.5 节我们将看到，在进行清理操作清除被删除的元组后，元组也可以以无序的方式插入到具有空闲空间的文件块中）。元组之间不进行关联，这样的表文件称为堆文件。PostgreSQL 系统中包含了四种堆文件：普通堆（ordinary cataloged heap）、临时堆（temporary heap）、序列（SEQUENCE relation，一种特殊的单行表）和 TOAST 表（TOAST table）。前面所述的堆文件就是指普通堆。临时堆的结构与普通堆相同，但临时堆仅在会话过程中临时创建，会话结束会自动删除。序列则是一种元组值自动增长的特殊堆。TOAST 表其实也是一种普通堆，但是它被专门用于存储变长数据。尽管这几种堆文件功能各异，但在底层的文件结构却是相似的：每个堆文件都是由多个文件块组成，在物理磁盘中的存储形式如图 3-4 所示。其中：

1) PageHeaderData 是长度为 20 字节的页头数据，包含该文件块的一般信息，如：

- 空闲空间的起始和结束位置。
- Special space 的开始位置。
- 项指针的开始位置。
- 标志信息，如是否存在空闲项指针、是否所有的元组都可见。

2) Linp 是 ItemIdData 类型的数组，ItemIdData 类型由 lp_ off、lp_ flags 和 lp_ len 三个属性组成。每一个 ItemIdData 结构用来指向文件块中的一个元组，其中 lp_ off 是元组在文件块中的偏移量，而 lp_ len 则说明了该元组的长度，lp_flags 表示元组的状态（分为未使用、正常使用、HOT 重定向和死亡四种状态）。每个 Linp 数组元素的长度为 4 字节，在图 3-4 的例子中，Linp1 指向 tuple1，Linp2 指向 tuple2。

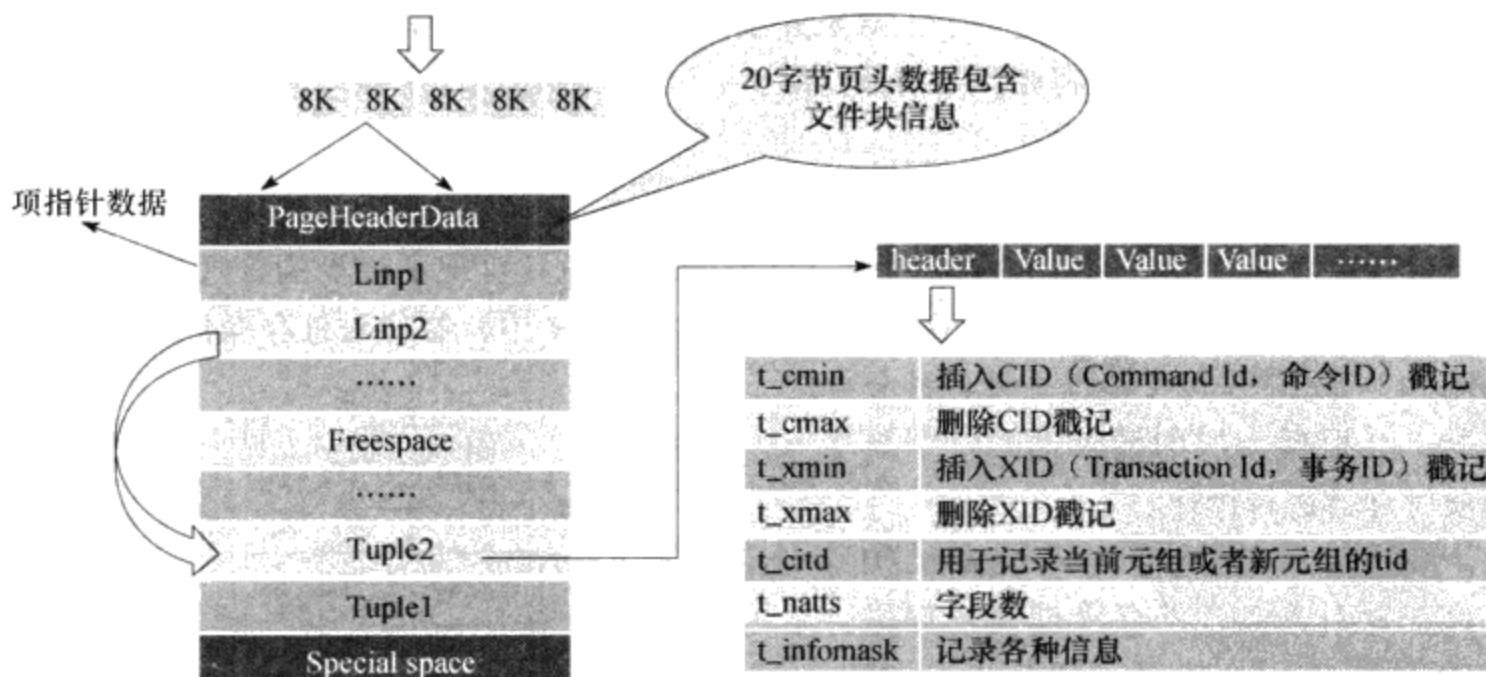


图 3-4 堆文件的物理结构

3) Freespace 是指未分配的空间（空闲空间），新插入页面中的元组及其对应的 Linp 元素都将从这部分空间中来分配，其中 Linp 元素从 Freespace 的开头开始分配，而新元组数据则从尾部开始分配。

4) Special space 是特殊空间，用于存放与索引方法相关的特定数据，不同的索引方法在 Special

space 中存放不同的数据（详见第 4 章）。由于索引文件的文件块结构和普通表文件的相同，因此 Special space 在普通表文件块中并没有使用，其内容被置为空。

元组信息中除了存放元组的实际数据，还存放元组头部信息，该信息通过 `HeapTupleHeaderData` 结构描述（数据结构 3.1），其中记录了操作此元组的事务 ID 和命令 ID 等信息，每个元组都有一个这样的头部信息。出于编程的考虑，PostgreSQL 的源代码中常用指向 `HeapTupleHeaderData` 的结构指针 `HeapTupleHeader` 来访问元组的头部信息。

`HeapTupleHeaderData` 的几个主要字段包括：

1) `t_choice` 是具有两个成员的联合类型：

- `t_heap`: 用于记录对元组执行插入/删除操作的事务 ID 和命令 ID，这些信息主要用于并发控制时检查元组对事务的可见性，详见第 7 章。
- `t_datum`: 当一个新元组在内存中形成的时候，我们并不关心其事务可见性，因此在 `t_choice` 中只需用 `DatumTupleFields` 结构来记录元组的长度等信息。但在把该元组插入到表文件时，需要在元组头信息中记录插入该元组的事务和命令 ID，故此时会把 `t_choice` 所占用的内存转换为 `HeapTupleFields` 结构并填充相应数据后再进行元组的插入。

2) `t_ctid` 用于记录当前元组或者新元组的物理位置（块内偏移量和元组长度），若元组被更新（PostgreSQL 对元组的更新采用的是标记删除旧版本元组并插入新版本元组的方式），则记录的是新版本元组的物理位置。

3) `t_infomask2` 使用其低 11 位表示当前元组的属性个数，其他位则用于包括用于 HOT 技术及元组可见性的标志位。

4) `t_infomask` 用于标识元组当前的状态，比如元组是否具有 OID、是否有空属性等，`t_infomask` 的每一位对应不同的状态，共 16 种状态。

5) `t_hoff` 表示该元组头的大小。

6) `_bits []` 数组用于标识该元组哪些字段为空。

这里需要介绍一下 PostgreSQL 中使用的“HOT 技术”。PostgreSQL 中对于元组采用多版本技术存储，对元组的每个更新操作都会产生一个新版本，版本之间从老到新形成一条版本链（将旧版本的 `t_ctid` 字段指向下一个版本的位置即可）。此外，更新操作不但会在表文件中产生元组的新版本，在表的每个索引中也会产生新版本的索引记录，即对一条元组的每个版本都有对应版本的索引记录。即使更新操作没有修改索引属性，也会在每个索引中都产生一个新版本。这一技术的问题是浪费存储空间，旧版本占用的空间只有在进行 VACUUM 时才能被回收，增加了数据库的负担。

为了解决这个问题，从版本 8.3 开始，使用一种 HOT 机制，当更新的元组同时满足如下条件时（通过 `HeapSatisfiesHOTUpdate` 函数判断）称为 HOT 元组：

数据结构 3.1 `HeapTupleHeaderData`

```
typedef struct HeapTupleHeaderData
{
    union
    {
        HeapTupleFields t_heap;
        DatumTupleFields t_datum;
    } t_choice;
    ItemPointerData t_ctid;
    int16 t_natts;
    uint16 t_infomask;
    uint16 t_infomask2;
    uint8 t_hoff;
    bits8 t_bits[1];
    //元组的具体数据将跟在这个结构后面
} HeapTupleHeaderData;
```

1) 所有索引属性都没被修改（索引键是否修改是在执行时逐行判断的，因此若一条 UPDATE 语句修改了某属性，但前后值相同则认为没有修改）。

2) 更新的元组新版本与旧版本在同一文件块内（限制在同一文件块的目的是为了通过版本链向后找时不产生额外的 I/O 操作从而影响到性能）。

HOT 元组会被打上 HEAP_ONLY_TUPLE 标志，而 HOT 元组的上一个版本则被打上 HEAP_HOT_UPDATED 标志。更新一条 HOT 元组将不会在索引中引入新版本，当通过索引获取元组时首先会找到同一块中最老的版本，然后顺着版本链向后找，直到遇到 HOT 元组为止。因此 HOT 技术消除了拥有完全相同键值的索引记录，减小了索引大小。

在堆中要删除一个元组，理论上有两种方法：

1) 直接物理删除：找到该元组所在文件块，并将其读取至缓冲区。然后在缓冲区中删除这个元组，最后再将缓冲块写回磁盘。

2) 标记删除：为每个元组使用额外的数据位作为删除标记。当删除元组时，只需设置相应的删除标记，即可实现快速删除。这种方法并不立即回收删除元组占用的空间。

PostgreSQL 采用的是第二种方法，每个元组的头部信息 HeapTupleHeader 就包含了这个删除标记位，其中记录了删除这个元组的事务 ID 和命令 ID。如果上述两个 ID 有效，则表明该元组被删除；若无效，则表明该元组是有效的或者说没有被删除的。在第 7 章中将会看到这种方法对多版本并发控制也是有好处的。

3.2.2 磁盘管理器

磁盘管理器是 SMGR 的一种具体实现，它对外提供了管理磁盘介质的接口，其主要实现在文件 md.c 中。磁盘管理器并非对磁盘上的文件直接进行操作，而是通过 VFD 机制来进行文件操作。

在 PostgreSQL 中，凡是需要对存储在磁盘中的表进行磁盘操作（如打开/关闭、读/写等），都是必须与磁盘管理器打交道，由它来统一处理。

对文件的操作是通过磁盘文件描述符（数据结构 3.2）进行的。MdfdVec 的各个字段意义如下：

- mdfd_vfd：该文件所对应的 VFD（见 3.2.3 节 VFD 机制）。
- mdfd_segno：由于较大的表文件会被切分成多个文件（可以称为段），因此用这个字段表示当前这个文件是表文件的第几段。
- mdfd_chain：指向同一表文件下一段的 MdfdVec，通过这个字段可以把表文件的各段链接起来，形成链表。

通过 mdfd_chain 链表使得大于操作系统文件大小（通常为 2GB）限制的表文件更易被支持，因为我们可以将表文件切割成段，使每段小于 2GB。

数据结构 3.2 MdfdVec

```
typedef struct _MdfdVec
{
    File      mdfd_vfd;           // 该磁盘文件对应的 VFD 的编号
    BlockNumber mdfd_segno;       // 该磁盘文件对应的表文件段号
    struct _MdfdVec *mdfd_chain; // 指向同一个表文件的下一段
} MdfdVec;
```

需要注意的是，`mdfd_chain` 为空并不一定表示表仅有一段，可能其他段没有打开。对于一个大表的操作，往往是逐段打开的。在 PostgreSQL 中，使用结构体 `SmgrRelation` 来表示一个被打开的表文件，该结构中就保存了该表文件对应的 `MdfdVec` 链表的头部。所有的 `SmgrRelation` 被组织成一个 Hash 表，该 Hash 表仅对后台进程可见，这样可以通过 Hash 表快速找到表对应的 `MdfdVec` 链表。在 `md.c` 中仅提供了读写文件块的接口函数，其具体的实现则在文件 `fd.c` 中。`fd.c` 实现了 VFD 管理以及通过 VFD 执行对实际物理文件的操作。

3.2.3 VFD 机制

在操作系统中，当一个进程创建或是打开一个文件时，操作系统会为该文件分配一个唯一文件描述符（或者叫文件句柄），并通过该文件描述符来唯一标识和操作该文件。由于每个操作系统都对一个进程能打开的文件数加以限制，因此进程能获得的文件描述符是有限的。对于经常需要打开许多文件的数据库进程来说，很容易就会超过操作系统对于文件描述符数量的限制。为了解决这个问题，在 PostgreSQL 中使用了虚拟文件描述符（VFD）机制，并实现了相应的管理机制。

1. VFD 实现原理

VFD 机制的原理类似连接池，当进程申请打开一个文件时，总是能返回一个虚拟文件描述符，对外封装了打开物理文件的实际操作。所谓的虚拟文件描述符，是指一个叫做 VFD 的数据结构，其中记录了操作系统为文件分配的真实文件描述符。实际上，一个进程能同时打开的文件的最大数仍然为操作系统规定的最大值，但是由于进程使用了 VFD，它自身会觉得可以打开任意多的文件，不再受操作系统的限制。如果多个进程对同一个文件进行操作，那么每个进程将获得一个真实文件描述符，而每个真实文件描述符又对应一个 VFD。因此同一个文件对于不同进程而言，其 VFD 可能是不同的。VFD、真实文件描述符（FD）、文件之间的关系如图 3-5 所示。

在 PostgreSQL 中，一个进程打开的 VFD 都是存储在 `VfdCache` 数组当中，该数组的每一个元素都表示该进程拥有的一个虚拟文件描述符，其数据类型为 `Vfd`（数据结构 3.3）。

`Vfd` 的每个字段意义如下：

- `fd` 记录该 VFD 所对应的真实文件描述符。如果当前 VFD 没有打开文件描述符则其值为 `VFD_CLOSED` (`VFD_CLOSED = -1`)。此外，在下文提到的 `VfdCache` 数组中，其中第一个元素的 `fd` 值恒为 `VFD_CLOSED`，以此来标志 LRU 池的开始。
- `fdstate` 表示该虚拟文件描述符的标记位：
①如果它的第 0 位置 1，即为 `FD_DIRTY`，

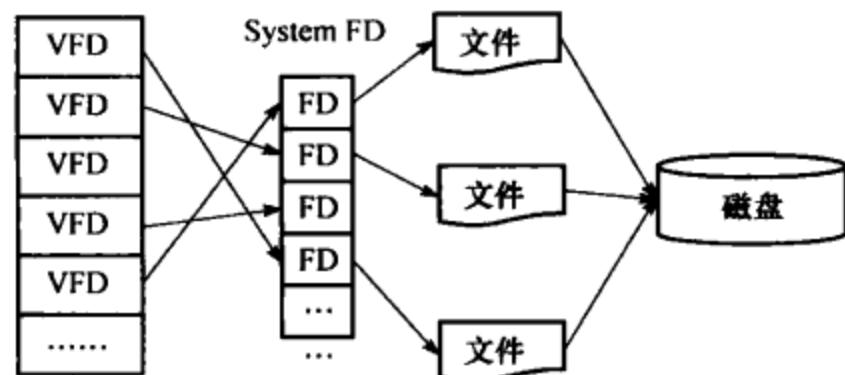


图 3-5 VFD、真实文件描述符、文件之间的关系

数据结构 3.3 `Vfd`

```
typedef struct vfd
{
    int fd;
    unsigned short fdstate;
    SubTransactionId create_subid;
    File nextFree;
    File lruMoreRecently;
    File lruLessRecently;
    off_t seekPos;
    *fileName;
    int fileFlags;
    int fileMode;
} Vfd;
```

表明该文件的内容已被修改过，但还没有被写回到磁盘，则在关闭此文件时要将该文件同步到磁盘里；②如果它的第 1 位置 1，即为 FD_TEMPORARY，表明该文件在关闭时要被删除。

- nextFree 指向下一个空闲的 VFD，其数据类型 File 就是一个整数，表示 VFD 在 VFD 数组中的下标。
- lruMoreRecently 指向比该 VFD 最近更常用的虚拟文件描述符。
- lruLessRecently 指向比该 VFD 最近更不常用的虚拟文件描述符。
- seekPos 记录该 VFD 的当前读写指针的位置。
- fileName 表示该 VFD 对应文件的文件名，如果是空闲的 VFD 则 fileName 为空值。
- fileFlags 表示该文件打开时标志，包括只读、只写、读写等。
- fileMode 表示文件创建时所指定模式。

2. LRU 池

在每一个 PostgreSQL 后台进程中都使用一个 LRU (Last Recently Used，最近最少使用) 池来管理所有已经打开的 VFD，池中每个 VFD 都对应一个物理上已经打开的文件。每一个进程都拥有其私有的 LRU 池和一系列的 VFD，进程需要打开文件时都是从自己私有的 LRU 池中申请 VFD。

当 LRU 池未满时，即进程打开的文件个数未超过系统限制时，进程可以照常申请一个 VFD 用来打开一个物理文件；而当 LRU 池已满的时候，进程需要首先关闭一个 VFD，这样打开新的文件时就不会因为超出操作系统限制而造成不可预料的错误。在 LRU 池中，使用替换最长时间未使用 VFD 的策略。

事实上，进程在 VfdCache 上保持了两个链表，一个是 LRU 池（双向链表），另一个是 FreeList（空闲链表，记录了所有可被分配的 VFD）。前者通过 Vfd 数据结构的 lruMoreRecently 属性和 lruLessRecently 属性来链接，后者则通过 Vfd 数据结构的 nextFree 属性来链接。而 VfdCache [0] 不参与 VFD 分配，它仅用来标识 FreeList 和 LRU 池的链表头部。

当进程需要打开一个文件时，将会为其分配一个 VFD；而关闭文件时则会回收 VFD。VFD 的分配和回收流程如下：

- 1) 进程在打开第一个文件时，将初始化 VfdCache 数组，置其大小为 32，为其中每一个 Vfd 结构分配内存空间，将 Vfd 结构中的 fd 字段置为 VFD_CLOSED，并将所有的数组元素放在 FreeList 上。
- 2) 分配一个 VFD，即从 FreeList 头取一个 VFD，并打开该文件，将该文件的相关信息（包括真实文件描述符、文件名、各种标志等）记录在分配的 VFD 中。
- 3) 若 FreeList 上没有空闲 VFD，则将 VfdCache 数组扩大一倍，新增加的 VFD 放入 FreeList 链表中。
- 4) 关闭文件时，将该文件所对应的 VFD 插入到 FreeList 的头部。

进程获得 VFD 之后，需要检查 LRU 池是否已满，也就是检查当前进程所打开的物理文件个数是否已经达到了操作系统的限制。如果没有超过限制，进程可以使用该 VFD 打开物理文件并将其插入到 LRU 池中；否则需要用到接下来要介绍的 LRU 池替换算法，先关闭一个 VFD 及其对应的物理文件，然后再使用获得的 VFD 来打开物理文件。在介绍 LRU 池替换算法之前先介绍一下 LRU 池的组织方式。

PostgreSQL 将一个进程当前正打开（所谓当前正打开是指操作系统当前确实为该文件分配了真实文件描述符）的所有文件的 VFD 都链成一个环，即 LRU 池，如图 3-6 所示。

图 3-6 中的每个方框代表一个 VFD，即一个 VfdCache 数组元素，方框里的符号表示该 VFD 在

数组中的下标。在 LRU 池中，每一个 VFD 都通过指针链接着两个 VFD，通过指针 lruMoreRecently 链接的是最近更常用的 VFD，通过指针 lruLessRecently 链接的是最近更不常用的 VFD。例如图 3-6 中 a_2 链接着 a_1 和 a_3 。LRU 池中有一个特殊的 VFD 和两个特殊的链接，VfdCache [0] 充当了 LRU 链头部的作用，它永远不会被实际分配给任何文件。一个特殊链接是 LRU 链末尾的 VfdCache [a_n] 通过 lruLessRecently 链接到 VfdCache [0]，另一个特殊链接是 VfdCache [0] 通过 lruMoreRecently 链接到 VfdCache [a_n]，这样系统可以通过 VfdCache [0] 的 lruLessRecently 值找到最近最少使用的文件的 VFD。显然，LRU 池的大小与操作系统对于进程打开文件数的限制是一致的，在 PostgreSQL 的实现中用全局变量 max_safe_fds 来记录该限制数，在 Postmaster 进程的启动过程中将会调用 set_max_safe_fds 函数来检测操作系统限制，并设置 max_safe_fds 的值。

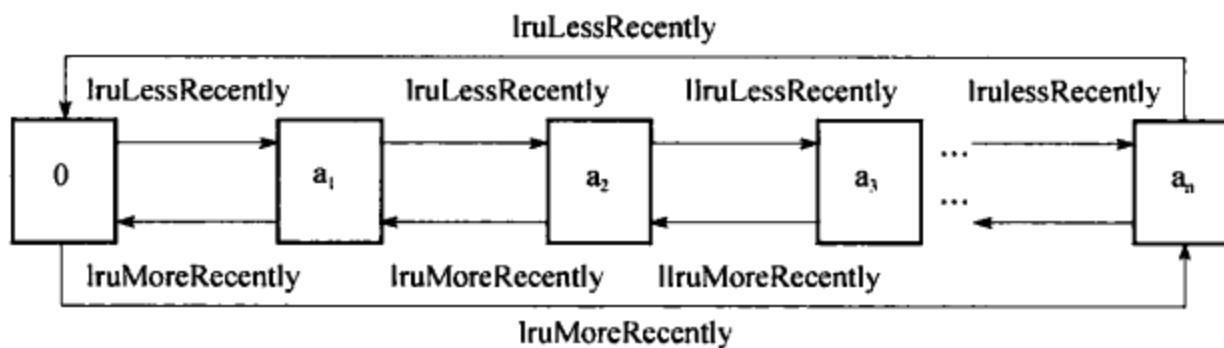


图 3-6 LRU 池

对 LRU 池里的 VFD 的操作主要包括以下三种：

(1) 从 LRU 池删除 VFD

该操作发生在进程使用完一个文件并关闭它时，通过 LruDelete 函数实现。该操作将指定的 VFD 从 LRU 池中删除，并将该 VFD 对应的文件关闭掉。例如，我们要对 Vfd [a_2] 执行 LruDelete 操作，则首先将 Vfd [a_1] 的 lruLessRecently 指向 Vfd [a_3]，将 Vfd [a_3] 的 lruMoreRecently 指向 Vfd [a_1]。这样就将 Vfd [a_2] 从 LRU 池中删除了，如果 Vfd [a_2] 的 fdstate 被置为 FD_DIRTY，则要先将 Vfd [a_2] 对应的文件同步回磁盘，再清掉 FD_DIRTY 位。接着将 Vfd [a_2] 的 seekpos 置为该文件当前读写的指针位置，最后将 Vfd [a_2] 对应的文件关闭掉并将 Vfd [a_2] 中的 fd 置为 VFD_CLOSED。这样 Vfd [a_2] 就变为空闲，还需要将其加入到空闲链表中。

(2) 将 VFD 插入 LRU 池

该操作发生在打开一个新的 VFD 时，通过 LruInsert 函数实现。该操作将指定 VFD 对应的物理文件打开，并将该 VFD 插入到 VFD [0] 之后的位置。例如，我们要插入的 VFD 为 b，首先要根据 b 中的 fd 字段判断对应的物理文件是否已经打开了，如果没有，则根据 b 中的 fileName 打开此文件，并插入到 LRU 池中。插入时要将 b 的 lruMoreRecently 指向 Vfd [0]，Vfd [0] 的 lruLessRecently 指向 b，然后 b 的 lruLessRecently 指向 Vfd [a_1]，Vfd [a_1] 的 lruMoreRecently 指向 b。

(3) 删除 LRU 池尾的 VFD

该操作通过 ReleaseLruFile 函数实现，它将 LRU 池中末尾的那个 VFD 删除。当 LRU 池已满而此时又要打开新的文件时，就需要执行 ReleaseLruFile 操作，将池中末尾的 VFD（最少使用的 VFD）删掉，这样新打开的 VFD 就可以插入到 LRU 中。注意，这里被删除的 VFD 仅仅只是从 LRU 池中脱链并关闭其对应的物理文件，VFD 结构本身并不做其他修改和删除。因为进程后面的操作还可能会用到该 VFD 所对应的物理文件。当再次需要访问一个 LRU 池之外的 VFD 时，需要先根据 VFD 中记

录的文件打开标志打开其对应的物理文件，然后根据 VFD 中记录的读写指针位置将物理文件描述符的读写指针移动到正确的位置，最后还要把该 VFD 重新插入到 LRU 池中。

3.2.4 空闲空间映射表

随着表中不断插入和删除元组，文件块中必然会产生空闲空间（free space）。在插入元组时优先选择将其存放在空闲空间内是利用存储的好方法。但是这需要在该表众多拥有空闲空间的文件块之间进行选择，而遍历所有文件块进行选择的开销将是非常庞大的。在 PostgreSQL 8.4 之前使用了一种全局 FSM 文件来记录所有表文件的空闲空间状况，其缺点是在全局 FSM 文件中只能记录每个表文件一定数量的文件块空闲空间状况，这使得对于空闲空间的管理比较复杂和低效。为了解决这一问题，PostgreSQL 8.4 中采用了一种新的策略，即对于每个表文件（包括系统表在内），同时创建一个名为“关系表 OID_fsm”的文件，用于记录该表的空闲空间大小，称之为`空闲空间映射表文件(FSM)`。例如，一个 OID 为 12000 的表，其空闲空间映射表文件将被命名为 12000_fsm。

为了能够更快地查找合适的空闲空间，FSM 文件不应该太大。因此在 FSM 中存储的并不是实际的文件块空闲空间的大小，而是仅仅用一个字节来记录。该字节的值用于描述对应文件块中空闲空间的范围，取值如表 3-1 所示。

为避免混淆，我们将表中的文件块称为表块，将 FSM 文件中的文件块称为 FSM 块。对于每个表块，不论其中是否有空闲空间，在 FSM 文件中都能找到对应的字节。若表块存在空闲空间，则按照表 3.1 将空闲空间字节映射为一个整数，在 FSM 中用 1 字节记录；若表块没有空闲空间，在 FSM 中使用 0 值来进行记录。

对于任何一个表块，只要找到它在 FSM 文件中对应的字节，根据该字节的值就可以知道这个表块中空闲空间的范围是多少。同样，对于任何一个表块，可以根据其空闲空间的大小计算出它对应的 FSM 字节的取值。比如，对于一个有 N 字节空闲空间的表块，其在 FSM 中记录的值为 $\lfloor (31 + N)/32 \rfloor$ 。事实上，大小在 0 ~ 31 字节的空闲空间将不会被使用，而当申请空闲空间大小在 8164 ~ 8192 字节之间时，系统将为该请求分配 MaxHeapTupleSize 的空闲空间（也就是一个空闲的表块）。

为了实现快速查找，FSM 文件中并不是简单地使用数组顺序存储每个表块的空闲空间，而是使用了树结构。在 FSM 块之间使用了一个三层树的结构，第 0 层和第 1 层为辅助层，第 2 层 FSM 块中用于实际存放各表块的空闲空间值。第 0 层 FSM 块只有一个，作为树根；第 1 层 FSM 块可以有多个，作为第 0 层 FSM 块的子节点；第 2 层 FSM 块也可以有多个，作为第 1 层 FSM 块的子节点。每一个 FSM 块内又构成一个局部的最大堆二叉树，但每一层 FSM 块内的最大堆二叉树又略有不同。

1) 第 2 层 FSM 块中最大堆二叉树的每一个叶子节点表示一个表块的空闲空间值。按照从左至右的顺序，所有第 2 层 FSM 块中的叶子节点排列起来就一一对应了表文件中的每一个表块。也就是说最左边的 FSM 块中左边第一个叶子节点对应表文件的第一块，左边第二个叶子节点对应表文件的第二块，依此类推。

- 2) 第 1 层 FSM 块中的叶子节点从左至右顺序对应第 2 层 FSM 块的根节点。
- 3) 第 0 层 FSM 块中的叶子节点从左至右顺序对应第 1 层 FSM 块的根节点。

表 3-1 FSM 中文件块空闲空间的表示

字节取值	表示的空闲空间范围（单位为字节）
0	0 ~ 31
1	32 ~ 63
.....
255	8164 ~ 8192

从上面的介绍可以看到，第0层和第1层的FSM仅用于快速定位到存在满足需求表块所属的第2层FSM文件块，而只有第2层中的FSM文件块，其块内最大堆二叉树的叶子节点中存储的才是和表块对应的空闲空间值。所有FSM块中的节点从逻辑上构成了一个全局的最大堆二叉树，其中的叶子节点从左至右顺序对应表文件中的文件块。

在单个FSM块中，使用最大堆的二叉树能够保证所有叶子节点的最大值被提升到根节点。这样我们只需要判断FSM块中的二叉树根节点是否满足需求，就可以知道这个块是否存在具有满足要求的叶子节点。

每个FSM块大小默认为8KB，除去必要的文件块头部信息，FSM块中剩下的空间都用来存储块内的二叉树结构，每个叶子节点用一个字节记录，因此FSM块内大约可以保存4000个叶子节点，其他空间均用作构成最大堆二叉树的辅助节点。这里4000并不是一个精确值，但为了表述方便，我们默认每个FSM块中二叉树最多有4000个叶子节点。这样，一个FSM文件中总共可以记录 4000^3 个叶子节点（表块），远远大于 2^{32} ——单个表的最大块数（PostgreSQL中的块号长度为32位，因此单个表最多只能有 2^{32} 个块），这样三层的结构足以记录表文件所有文件块的空闲空间值。其编号规律是从0开始按照深度优先遍历的方式对树中的FSM块进行编号：

1) 第0、1、 $4001+1$ 、 $4001*2+1$ ……号FSM块为FSM文件的辅助块，用于快速查找FSM文件。其中第0号就是根节点（第0层FSM块），第1、 $4001+1$ 、 $4001*2+1$ ……号FSM块是第1层FSM块。

2) 其余的FSM块属于第2层，用于存储表块的空闲空间值。即：

- 第2个FSM块存储对应表第0~4000-1表块的空闲空间值。
- 第3个FSM块存储对应表第4000~4000*2-1表块的空闲空间值。
-

图3-7给出了一个FSM文件磁盘存储与其逻辑结构的映射关系实例图。其中，0号FSM块中的

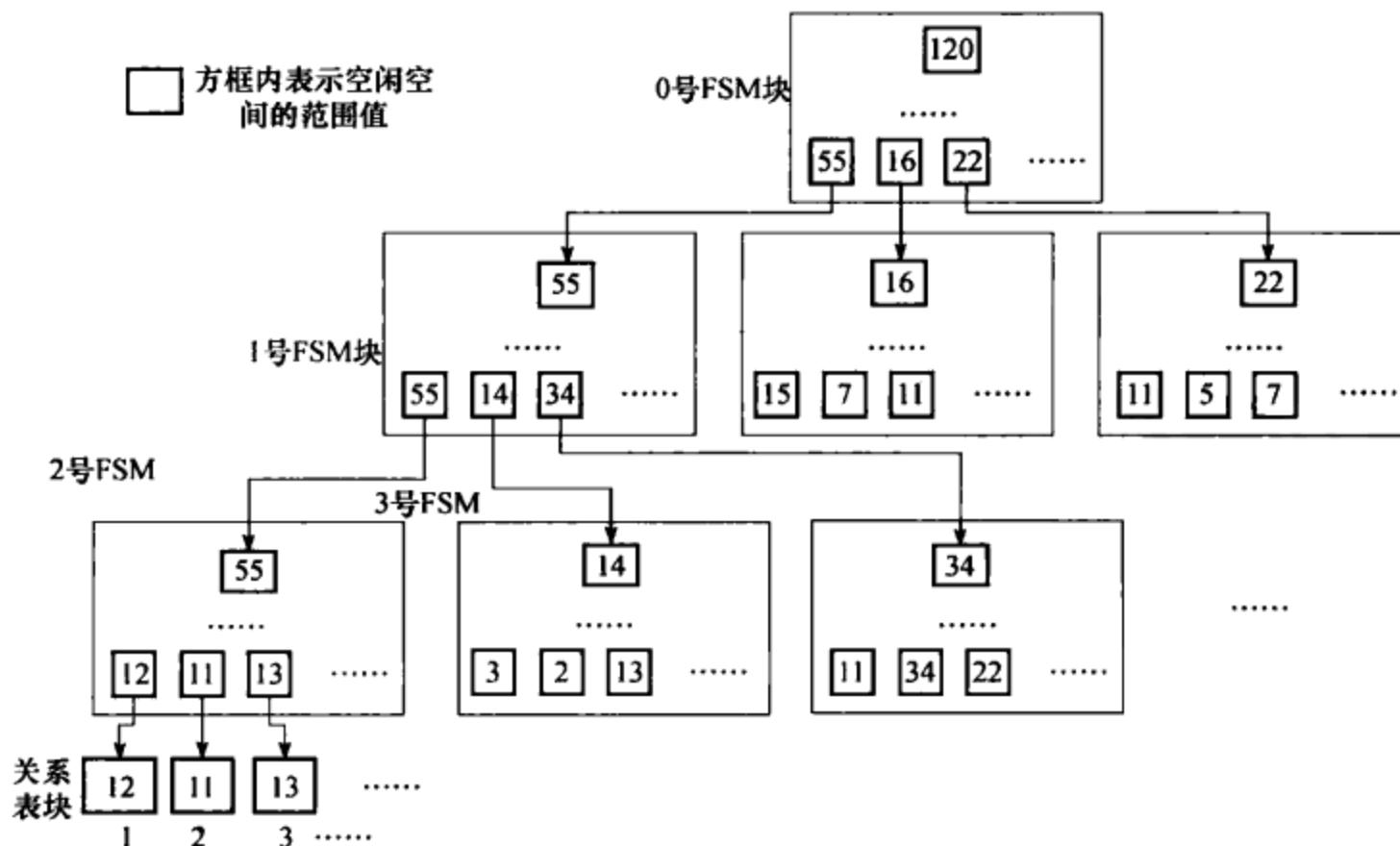


图3-7 FSM物理块与逻辑地址对照

叶子节点记录了第 1 层的所有 FSM 块根节点，例如其第一个叶子节点值为 55，对应 1 号 FSM 块的根节点；第二个叶子节点对应第 4002 号 FSM 块（第 1 层的第 2 个 FSM 块）。通过这样的对应关系，形成了三层树结构，当我们需要找到一个具有指定大小空闲空间的表块时，实际上是要找到一条从第 0 层到某个第 2 层 FSM 块的路径，最后定位到该第 2 层 FSM 块中二叉树的某个叶子节点。由于该叶子节点在所有叶子节点中的顺序就是其对应表块在表文件中的顺序，因此可以通过简单的计算得到对应的表块号，有兴趣的读者可以尝试建立该计算公式。

通过这样的组织结构，可以保证 FSM 文件的第一个文件块中二叉树根节点存储的是所有表块空闲空间最大值。这样，当需要从 FSM 文件中选择一个表块时就可以快速判定是否存在满足要求的表块。下面将就 FSM 文件的创建、查找和调整来进行具体阐述。

1. 创建 FSM

FSM 文件并不是在创建表文件的时候就被创建，而是推迟到需要使用 FSM 文件的时候，即执行 VACUUM 操作时或为了插入元组第一次查找 FSM 文件时才创建。

在创建 FSM 文件时，会预先创建三个 FSM 块：第 0 号为根节点即第 0 层的 FSM 块，第 1 号为第 1 层的第一个节点，第 2 号为第 2 层的第一个节点。在第 2 号 FSM 块内叶子节点中依次存储从第 0 号开始的各表块的空闲空间值，若没有空闲空间或空闲空间太小，则记录为 0。当第 2 号 FSM 块满了之后，将会扩展新的 FSM 块。该工作主要调用函数 fsm_extend 进行实现。每个 FSM 块的物理结构如图 3-8 所示。

在内存中，FSM 块对应的结构为 FSMPagedata（数据结构 3.4）。

可以看到，FSM 块的内存结构和物理存储结构是相对应的。其中 fp_nodes 数组存储了一个最大堆二叉树，从根节点开始每层顺序排列直到叶子节点，每个数组元素是一个整数值。如图 3-9 中的最大堆二叉树将被存储为 {7, 7, 6, 5, 7, 6, 5, 4, 5, 5, 7, 2, 6, 5, 2}。fp_next_slot 值是一个整数，用于提示下一次开始查询二叉树的叶子节点位置。当从 FSM 块中找到一个合适的叶子节点位置（也就是一个满足请求的二叉树叶子节点序号，从 0 开始记，用 slot 表示）时，若该 FSM 块为叶子 FSM 块，则将 fp_next_slot 设置为 slot + 1，否则设置为 slot。这样做有两点好处：第一，当多个进程同时向一个表中插入数据时，结合 fp_next_slot 使用一定的寻径算法（3.2.4 节），可以向请求的进程返回不同的表块，从而避免了对同一个表块的竞争。此外，由于操作系统本身具有预读取和批量写技术，按照表块顺序进行插入将获益于操作系统这一特性。下一节将介绍如何从 fp_next_slot 标记的位置开始查询二叉树。

2. 查找 FSM

函数 fsm_search 利用 FSM 文件查找一个具有指定空闲空间的表块。该函数有两个参数：表的基本信息 rel 和最小空闲空间值 min_cat，该函数将在 rel 中找到一个至少具有 $min_cat * 32$ 字节空闲空间的表块。当然，这个查找过程是利用 FSM 文件来完成的，该函数返回的是一个具有满足条件空闲空间的表块的块号。

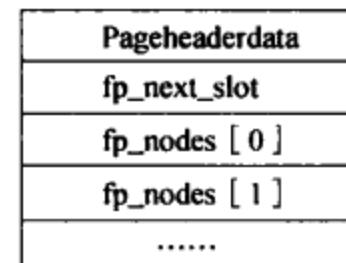


图 3-8 FSM 块的物理结构

数据结构 3.4 FSMPagedata

```
typedef struct
{
    int      fp_next_slot;
    uint8   fp_nodes [ 1 ];
} FSMPagedata;
```

在 `fsm_search` 中为了方便查找 FSM 文件，使用了数据结构 3.5 来表示 FSM 块在树中的位置。

其中，`level` 表示该 FSM 块所处的层，但是这里 `level` 的取值和前面介绍 FSM 块的层号并不一致：第 0 层 FSM 块的 `level` 值为 2，第 1 层 FSM 块的 `level` 值为 1，第 2 层 FSM 块的 `level` 值为 0。`FSMAddress` 的 `logpageno` 字段表示该块是其所在层次中的第 `logpageno` 个 FSM 块。可以认为 `FSMAddress` 是 FSM 块的逻辑地址，可以调用函数 `fsm_logical_to_physical` 从逻辑地址计算出物理地址，也就是该 FSM 块的物理块号。

`fsm_search` 查找的基本思想非常简单，将其看做一个查找最大堆二叉树的过程。其实际执行流程如下：

- 1) 首先初始化块指针（逻辑地址）指向第 0 号 FSM 块。然后循环地执行步骤 2 ~ 4。
- 2) 由块指针表示的逻辑地址计算得到物理地址，并将这个 FSM 块装入缓冲区。
- 3) 如果对应的 FSM 块不存在，则直接返回一个无效块号。
- 4) 调用函数 `fsm_search_avail` 在步骤 2 中装入的 FSM 块中查找一个符合条件的叶子节点。如果能找到这样的叶子节点，`fsm_search_avail` 将返回该叶子节点在块内 `fp_nodes` 数组中的编号，否则将返回 -1。根据这个返回值进行如下判断：

①如果找到的叶子节点编号不等于 -1 且当前的块指针已经到达第 2 层，说明当前的 FSM 块中包含一个具有合适空闲空间的表块的信息，则调用函数 `fsm_get_heap_blk` 计算出表块的块号作为返回值返回。

②如果找到的叶子节点编号不等于 -1 且当前的块指针还没有到达第 2 层，说明当前 FSM 块是第 0 层或第 1 层的辅助块，还需要继续向下搜索，则调用函数 `fsm_get_child` 根据返回的叶子节点编号取得下一个要搜索的 FSM 块，然后继续回到步骤 2 继续循环。

③如果找到的叶子节点编号等于 -1 且当前块指针指向第 0 层，说明整个 FSM 文件中都没有能满足要求的空闲空间信息，因此直接返回无效块号。

④如果找到的叶子节点编号等于 -1 且当前块指针没有指向第 0 层，这种情况有可能是因为更低层的 FSM 块的最大空闲空间值发生变化后没有反映到高层的 FSM 块中（在后面的“调整 FSM”部分中可以找到原因），这时需要更新相关 FSM 块的最大空闲空间值。更新后将块指针重新指向第 0 号 FSM 块，再次从 FSM 文件的根块开始搜索。由于有可能会有其他后台进程同时在更新 FSM 文件，因此有可能重新开始的搜索也会得到同样的结果，在这里会在重试之前设置一个计数器，如果重试超过 10000 次还是得到同样的结果，则直接返回无效块号。

从 `fsm_search` 的流程可以看出，我们通过 FSM 文件有可能找不到满足要求的表块，原因有可能是确实不存在合适的表块，也有可能是因为 FSM 文件中的最大空闲空间值信息没有及时更新。在这种情况下，我们可以重新申请一个新的表块来插入数据，至于 FSM 文件中的信息更新会由 VACUUM 进行。

`fsm_search` 只是反映了在 FSM 块层面上的查找过程，而每一个 FSM 块内最大堆二叉树的搜索过程则由函数 `fsm_search_avail` 实现，该函数从指定的 FSM 块中获取一个大于或等于 `min_cat` 的叶子节点，其流程如下：

- 1) 首先检测 FSM 块中的根节点 (`fp_nodes[0]`)，若根节点无法满足 `min_cat`，则表明该块中

数据结构 3.5 `FSMAddress`

```
typedef struct
{
    int           level;
    int           logpageno;
} FSMAddress;
```

没有满足该请求的节点，返回 -1。

2) 否则，根据 fp_next_slot 的值（fp_next_slot 在初始时为 0，每一次找到合适的叶子节点后会设置 fp_next_slot 的值，见步骤 5），设置当前节点指针 nodeno = 每个 FSM 块中非叶子节点个数 + fp_next_slot。进入以下循环：

①若当前节点满足需求，则跳出循环。

②否则，设置当前节点指针为当前节点右节点（这个右节点是可以循环的，即同一层的最后一个节点的右节点是该层的第一个节点）的父亲节点，并返回步骤 1。

3) 从步骤 2 找到的节点开始向下寻找一条通往叶子节点的路径，该路径上的每个节点的值都满足请求。注意在寻径的过程中，可能会碰到左右子节点都满足要求的情况，这时候优先选择左节点。从该路径达到的叶子节点即为最终的结果，将其序号记录为 slot。

4) 若当前块为叶子块，则设置 fp_next_slot = slot + 1，否则设置 fp_next_slot = slot。

5) 返回 slot，即找到的叶子节点序号。

为了便于理解，下面将通过一个例子来对上述算法进行描述。图 3-9 是一个简单的 FSM 块二叉树。每个方框内的数字代表一个空闲空间值。

如图 3-9 所示，假设 fp_next_slot 指的是最下一层中左起第 7 个节点，记该节点为 node (4, 7)，表示第 4 层第 7 个节点，它的空闲空间值记为 value (4, 7)。如果需要找到一个空闲空间值为 6 的表块，其过程如下：

1) value (1, 1) = 7，大于 6，表明该 FSM 块中有可以满足要求的表块的信息。

2) value (4, 7) = 5，小于 6，则向右并找到其父节点即 node (3, 4)。

3) value (3, 4) = 5，不满足需求，向右扫描到 node (3, 1) 并找到其父节点 node (2, 1)。

4) value (2, 1) = 7，大于 6，满足需求，则从 node (2, 1) 开始，找到路径 node (2, 1) → node (3, 2) → node (4, 4)，即满足条件的是第 4 个叶子节点。

5) 由于 slot 是从 0 开始记数，所以返回的 slot 值为 $4 - 1 = 3$ 。

可见，第 4 个叶子节点是满足我们需求的。若考虑到多层次结构，即当前 FSM 块不在第 2 层，那么需要根据得到的 slot 找到它的第 slot 个 FSM 子块，然后按照以上的步骤在子块中查找，直到找到一个第 2 层 FSM 块并在其中找到一个符合要求的叶子节点。

3. 调整 FSM

当从 FSM 文件中找到了一个具有合适空闲空间的表块并使用它进行了数据插入之后，需要对 FSM 文件中相关信息进行修改，需要调整该表块的空闲空间值，同时对其所属的 FSM 块内的最大堆二叉树进行相应的调整。这个调整的过程由函数 RecordPageWithFreeSpace 完成，其参数包括表块号、表的信息以及表块当前的空闲空间值，主要流程如下所示：

1) 根据表块号，找到其对应的 FSM 块，将其读入缓冲区中。

2) 调用函数 fsm_set_avail 对 FSM 块进行调整：

①若新的空闲空间值与旧值相等，则不做调整。

②如果空闲空间值发生了变化，则将新值赋给相应的叶子节点。

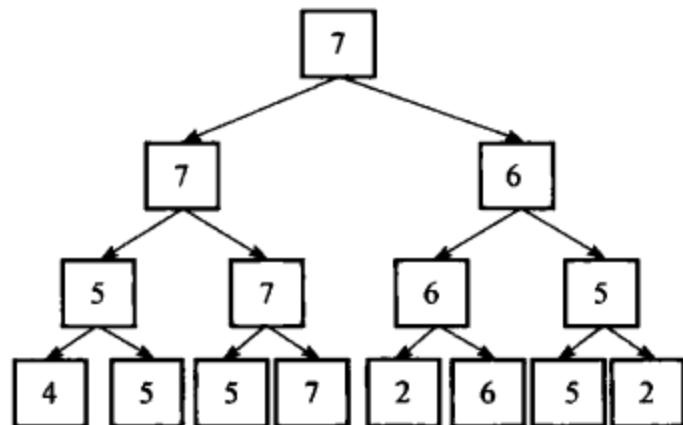


图 3-9 最大堆二叉树

- ③调整除二叉树根节点外的其他节点，使父亲节点的值总是为两个子节点的最大值。
- ④若新值大于树的根节点值时，调用 fsm_rebuild_page 重建块内结构，使其符合最大堆结构。
- 3) 将该 FSM 块标记为脏。

可见，当一个 FSM 块发生变化时，不管修改后的叶子节点值是大于还是小于全局二叉树根节点值，都不会去调整 FSM 块之间的三层树结构，对全局树的调整是在 VACUUM 时完成的（函数 FreeSpaceMapVacuum）。这样做是为了防止频繁调用 FSM 块引起额外的开销。

3.2.5 可见性映射表

PostgreSQL 中为了实现多版本并发控制，当事务删除或更新元组时，并非从物理上删除，而是通过将其标记为无效的方式进行标记删除，最终对这些无效元组的清理操作需要调用 VACUUM 来完成。

为了能够加快 VACUUM 查找包含无效元组的文件块的过程，在 PostgreSQL 8.4.1 中为每个表文件定义了一个新的附属文件——可见性映射表（VM）。VM 中为表的每个文件块设置了一位，用来标记该文件块是否存在无效元组。对包含无效元组的文件块，VACUUM 有两种方式处理，即快速清理（Lazy VACUUM）和完全清理（Full VACUUM）。注意，VM 文件仅在 Lazy VACUUM 操作中被使用到，而 Full VACUUM 操作由于要执行跨块清理等复杂操作，需要对整个表文件进行扫描，这时候 VM 文件的作用并不大。

对于每个表文件，其对应的 VM 文件命名为：“关系表 OID_vm”。对该文件的操作在 visibility-map.c 文件中进行了定义。

与其他文件一样，VM 文件也被划分为若干个文件块（简称 VM 块）。VM 块中除了必要的标记信息外，其他的每一位都对应一个表块，当表块中所有的元组对当前的事务都是可见的时候，表块对应的位才设置为 1。其文件块结构如图 3-10 所示。

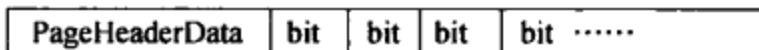


图 3-10 VM 文件块

每个 VM 文件块中能够记录 $\text{size} = (\text{blkbsz} - \text{SizeOfPageHeaderData}) * 8$ 个表块的信息，第一个 VM 块记录第 1 至 size 号表块的信息，第二个 VM 块记录第 size + 1 至 2 * size + 1 号表块的信息，依此类推。

当对某个表块中的元组进行更新或者删除后，那么该表块在 VM 文件中对应位置的标志位将被置为 0。在设置标志位的时候，需要对其对应的 VM 页面加锁。这是为了避免在 VACUUM 判断该页面是否对所有事务可见的同时，其他进程修改该页面，从而导致 VACUUM 清理过程中忽略了此页面。

当标志位为 1 时，VACUUM 操作会忽略扫描对应的表块，所以能大大提高 VACUUM 的效率。由于 VM 文件不跟踪索引，所以对索引的清理操作还是需要进行完全扫描。

当前，VM 文件仅仅是作为一个提示（hint）来加快 VACUUM 的速度，所以即使 VM 文件损坏也仅仅会导致 VACUUM 忽略那些需要清理的页面，而不会对数据产生任何负面影响。

3.2.6 大数据存储

在 PostgreSQL 中提供了两种大数据存储方式：第一种是 TOAST 机制，使用数据压缩和线外存储来实现；第二种是大对象机制，使用一个专门的系统表来存储大对象数据。下面将就这两种机制

进行详细的介绍，并对这两种机制进行比较。

1. TOAST

TOAST (The Oversized-Attribute Storage Technique, 超尺寸字段存储技巧) 是 PostgreSQL 提供的一种存储大数据的机制，我们先介绍 TOAST 的实现机制。

(1) TOAST 实现原理

在 PostgreSQL 中只有一部分数据类型支持 TOAST，这些数据类型必须有变长（代码内部常称为 varlena 类型）的表现形式，比如 TEXT 类型。只有在准备向支持 TOAST 的属性中存储超过 BLCKSZ/4 字节（通常是 2KB）的数据时，TOAST 机制才会被触发。TOAST 机制会试图对将要存储的数据进行压缩和（或）线外存储数据（即把数据存储在其他的表中），直到数据比 BLCKSZ/4 字节短，或者无法得到更好的结果的时候才停止。不管是否触发 TOAST 机制，表的目标属性中都会存有数值，该数值的前 2 位表示数据的存储方式：

- 如果前两位都是 0，那么数值是该数据类型一个普通的未 TOAST 的数值。
- 如果前两位为 01，那么表示该数据被压缩过，使用前必须先解压缩。剩下的 30 位表示压缩后的数据的大小，在这 4 字节后还会附加 32 位来表示压缩前的数据大小。
- 如果第一位为 1，这时候数据头部仅使用 1 字节。当存储的是短字符串时（小于 128 字节），剩下的 7 位表示字符串的长度；当该字节被设置为 10000000 时，表明这是线外存储的数据，紧随其后使用 1 字节记录指针的大小 (TOAST_POINTER_SIZE)，数据区域则记录 TOAST 指针（数据结构 3.6）。

注意，线外存储时也可能进行压缩，这种情况通过 TOAST 指针中的 va_rawsize 和 va_extsize 来进行比较。

线外存储的数据会保存在称为 TOAST 表的普通表中。如果一个表中有任何一个属性是可以 TOAST 的，那么该表将有一个关联的 TOAST 表，其 OID 存储在表的基本信息（也就是表在 pg_class 中的元组）的 reltoastrelid 属性里。如果没有关联的 TOAST 表，则 reltoastrelid 属性值为 0。

TOAST 机制识别四种不同的存储数据的策略：

- PLAIN：避免压缩或者线外存储。这只是对那些不能 TOAST 的数据类型才有可能。
- EXTENDED：允许压缩和线外存储。这是大多数可以 TOAST 的数据类型的缺省策略。TOAST 机制首先将企图对数据进行压缩，如果元组仍然太大，则进行线外存储。
- EXTERNAL：允许线外存储，但是不允许压缩。使用 EXTERNAL 将令那些在 text 和 bytea 字段上的子串操作更快（但代价是增加了存储空间），因此这些操作是经过优化的，即如果线外数据没有压缩，那么它们能够只抓取需要的部分。
- MAIN：允许压缩，但不允许线外存储。实际上，在这样的数据上仍然可能会进行线外存储，但只是作为没有办法把数据变得更小的情况下的最后手段。

每一种可以 TOAST 的数据类型都有一个缺省策略，但是某个属性的存储策略可以用 ALTER TABLE SET STORAGE 命令修改。

TOAST 机制的主要优点在于可以有效地节省查询时所占的内存空间。通常的查询是用相对比较短键值的主属性来完成，TOAST 数据只是在向用户显示结果时才被取出来，在查询的过程中并不需要检查 TOAST 数据的具体取值。因此，查询时所用到的数据要小得多，并且它的大部分元组都存储在共享缓冲区里，可以不需要任何线外存储。另一方面，在排序时数据量也比较小，因此排序可

以更多地在内存里完成。

(2) TOAST 表结构

如果采用线外存储方式处理 TOAST 数据，则线外数据会被分割成（如果压缩过，在压缩之后）不超过 TOAST_MAX_CHUNK_SIZE 字节的片段（chunk），每个片段都作为独立的元组存储在相关联的 TOAST 表中。每个 TOAST 表的模式都是相同的，其属性如表 3-2 所示。

每一个被线外存储的 TOAST 数据都会被分配一个 OID，通过这个 OID 可以在 TOAST 表中找到属于该 TOAST 数据的所有片段，进而可以重组该数据。

下面马上会看到，TOAST 数据的 OID 是构成 TOAST 指针的一部分。

根据前文所述，一个线外存储的数据可以通过一个存储在 TOAST 属性中的 TOAST 指针找到。该指针的结构如数据结构 3.6 所示。加上头部的长度字，一个 TOAST 指针的总长度是 20 字节，不管它代表的数值的实际长度是多大。

数据结构 3.6 varatt_external

```
struct varatt_external
{
    int32 va_rawsize;           // 原始数据（未 TOAST 过）的大小（包含数据头部信息）
    int32 va_extsize;          // 线外存储（TOAST 过）的数据大小（不包括头部）。
    Oid     va_valueid;        // 线外存储数据的 OID，也是其在 TOAST 表中的唯一标识
    Oid     va_toastrelid;      // TOAST 表的 OID
};
```

从这里可以看到，访问一个线外 TOAST 数据的基本流程是：首先从表的 TOAST 属性中获取 TOAST 指针，然后通过 TOAST 指针找到 TOAST 表，再通过 TOAST 数据的 OID 在 TOAST 表中找到所有的片段并且按序号拼装起来得到 TOAST 数据。具体的过程将在下一小节中介绍。

(3) TOAST 操作

TOAST 机制的实现代码在 src/backend/access/heap/tuptoaster.c 文件中，与 TOAST 相关的主要操作及其实现如图 3-11 所示。

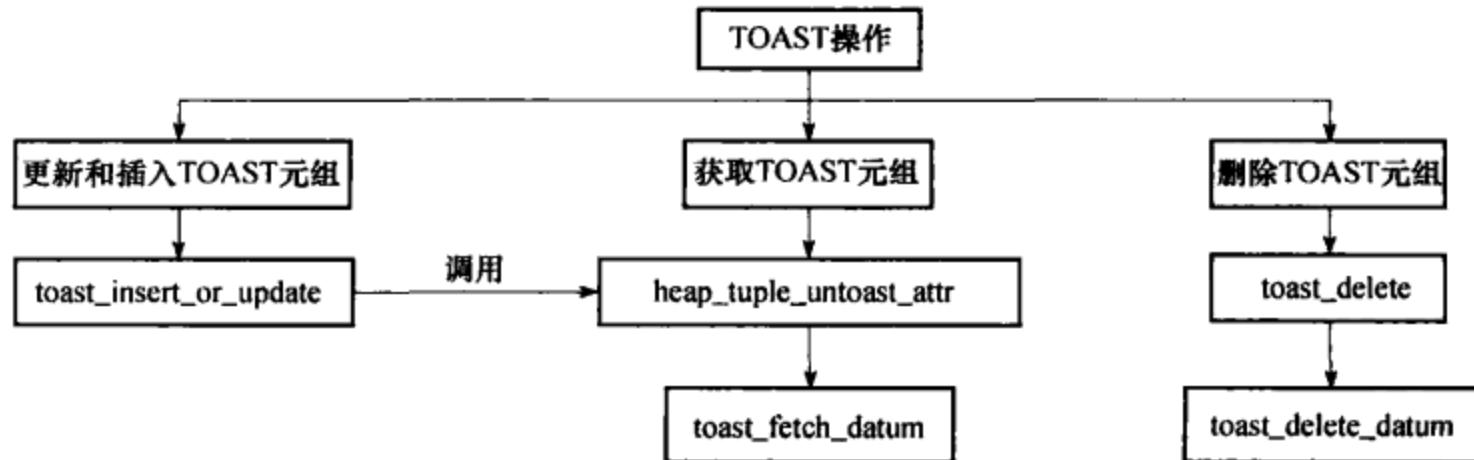


图 3-11 TOAST 操作实现

表 3-2 TOAST 表属性

属性名称	类型	描述
chunk_id	oid	线外存储时为整个 TOAST 数据分配的 OID
chunk_seq	int4	一个序列号，存储该片段在整个 TOAST 数据中的位置
chunk_data	bytea	该片段实际的数据

1) 元组插入/更新时的 TOAST 操作。

在进行元组的插入或更新操作时有可能需要对元组进行 TOAST 处理。这一工作由函数 `toast_insert_or_update` 来完成，该函数的参数中包括新版本元组和旧版本元组。在进行插入或更新操作时会调用该函数对将要插入的新版本元组进行 TOAST 处理，并将该函数返回的结果元组（可能与作为参数传入的新版本元组不同）替代新版本元组插入表中。该函数将通过判断旧版本元组是否为空来确定当前进行的是插入操作还是更新操作。

`toast_insert_or_update` 的主要流程如下：

①首先从传递的新版本元组中提取出属性数组（包括属性的描述和值）等信息。

②检查每一个属性：

- 对于元组更新操作，若该属性已经进行线外存储且更新后的值与旧值不相同，则将该属性标记为需要删除旧版本中该属性的值和对应的 TOAST 表元组。否则将重用原有的线外存储数据于新版本元组中。
- 对于元组插入操作，若表被设置为不允许线外存储或压缩，则设置该属性的 TOAST 动作为“p”，表示采用 plain 策略。

③根据不同的策略，分为四个 while 循环来处理数据：

- 循环一：找到需要并且可以压缩的属性对其进行压缩。循环一会反复执行下面的操作直到新版本元组的大小小于元组进行 TOAST 操作的临界值：

- i) 检查每一个属性，找到其中尺寸最大并且没有进行过 TOAST 处理的属性，如果不存在这样的属性则直接退出循环。
- ii) 如果该属性的 TOAST 策略为“x”（即允许同时采取压缩和线外存储），将首先尝试进行压缩，若成功压缩则获取压缩后的newValue；若压缩失败，则设置其 TOAST 动作为忽略压缩，表明在后续的压缩操作中不考虑这个属性。
- iii) 如果该属性的 TOAST 策略为“e”（即允许线外存储但不做压缩），同样将设置其 TOAST 动作为忽略压缩。
- iv) 再次检查该属性值的尺寸是否超过元组进行 TOAST 操作的临界值，如果超过则表明该属性值过大，将立刻调用函数 `toast_save_datum` 将其存储到线外（前提是 TOAST 表存在），并设置 TOAST 动作为忽略压缩和线外存储。

- 循环二：找到可以进行线外存储但还没有实行的属性，并对其进行线外存储。在 TOAST 表存在的前提下循环二将会反复执行下面的操作，直到新版本元组的大小小于元组进行 TOAST 操作的临界值：

- i) 检查每一个可以线外存储的属性（即其 TOAST 策略为“e”或者“x”），找到其中尺寸最大并且没有进行过线外存储的属性，如果不存在这样的属性则直接退出循环。

- ii) 调用函数 `toast_save_datum` 将其存储到线外，并设置 TOAST 动作为忽略压缩和线外存储。

- 循环三：找到只能线内压缩并且还未实行的属性，并对其进行线内压缩。注意，循环一只是压缩了那些同时允许线外存储和压缩的属性，并没有对只能线内压缩的属性进行处理。循环二将会反复执行下面的操作直到新版本元组的大小小于元组进行 TOAST 操作的临界值：

- i) 检查每一个只能线内压缩的属性（即其 TOAST 策略为“m”），找到其中尺寸最大并且没有进行过压缩或者线外存储的属性，如果不存在这样的属性则直接退出循环。

- ii) 对该属性尝试线内压缩, 如果压缩不成功则设置其 TOAST 动作为忽略压缩。
 - 循环四: 如果经过前面三个循环的处理, 元组仍然超过临界值, 我们只能尝试将那些 TOAST 策略为 “m” 且已经被处理过的属性存储到线外。循环四会逐个将这样的属性存储到线外。如果真的进行到这一步, 实际上所有需要进行 TOAST 处理的属性都已经被存储到了线外, 这个循环就可以确保最终得到的元组不会超过临界值。
- ④使用参数中新元组的信息和 TOAST 处理过的属性值构造一个结果元组。
- ⑤从旧元组中删除没有在新元组中重用的线外存储数据。
- ⑥返回结果元组。

前面介绍过 PostgreSQL 中是不允许元组跨文件块存储的, 因此一个元组的尺寸是有上限的。经过四个循环的处理, 可以保证处理过的元组不会因为过大而无法存入一个文件块中。

另外, 还有必要提一下 TOAST 机制所采用的压缩方法。TOAST 机制实际上采用了一种 LZ 压缩算法, 这是一种无损压缩算法, 该算法在函数 `toast_compress_datum` 中进行了具体实现。简单来说, LZ 压缩算法被认为是基于字符串匹配的算法。例如, 在一段文本中某字符串经常出现, 并且可以通过前面文本中出现过的字符串的指针来表示。当然这个做法的前提是指针长度应该比字符串本身要短。LZ 压缩算法的详情, 读者可以参阅相关文献, 这里就不再赘述。

2) 读取 TOAST 数据。

在查询涉及被 TOAST 过的数据时, 我们需要取出 TOAST 数据并且将其恢复原貌 (TOAST 数据可能被压缩或者线外存储) 并返回给用户。TOAST 数据的读取同样要根据数据所采用的 TOAST 策略来获取, 函数 `heap_tuple_untost_attrs` 将根据 TOAST 数据采用的策略读取出准确的数据。

`heap_tuple_untost_attr` 函数将用于从 TOAST 数据中读出原始数据, 其参数是一个从 TOAST 属性中取出的值。由于 TOAST 属性中存储的值的前 2 位表示该属性对应的 TOAST 数据的存储策略, 因此 `heap_tuple_untost_attr` 可以根据不同的策略来读取 TOAST 数据。该函数的主要流程如下:

①如果数据是线外存储的, 则先调用函数 `toast_fetch_datum` 从 TOAST 表中获取该数据的片段来重组数据。获取线外数据的原理非常简单: 根据 TOAST 属性中存储的 TOAST 指针, 从 TOAST 表中得到所有属于该 TOAST 数据的片段元组, 然后根据元组中记录的片段顺序号将片段拼成 TOAST 数据。接下来对从线外取回的数据进行判断, 如果是经过压缩的则用解压算法进行解压缩返回数据; 否则直接返回数据。

②如果数据是没有线外存储但是经过压缩的, 则解压缩后返回数据。

③如果数据是短数据 (没有被压缩和线外存储), 则直接返回数据。

另外, 还有函数 `heap_tuple_fetch_attr` 用于专门从线外获取一个 TOAST 数据, 但获取的数据仍然可能是被 TOAST 过的 (被压缩过)。实际上 `heap_tuple_fetch_attr` 就是通过调用 `toast_fetch_datum` 函数来获取线外数据的。

按四种 TOAST 策略处理的 TOAST 数据都可以通过以上几个函数或函数组合来读取。

3) 删除 TOAST 数据。

当包含 TOAST 数据的元组被删除时, 其对应的 TOAST 数据也需要被删除。对于采用 PLAIN 和 MAIN 策略的数据来说, TOAST 数据会包含在元组内一起被删除, 不需要做额外的处理。但是对于采用 EXTENDED 和 EXTERNAL 两种策略的数据来说, 除了删除元组, 还需要将线外存储的数据也级联删除, 这种情况下将调用函数 `toast_delete` 来实现删除。`toast_delete` 函数以被删除的元组及其所在

的表作为参数，该函数会对元组中的 TOAST 属性进行检查，如果该属性的数据是线外存储的，则调用 toast_delete_datum 函数删除线外数据。toast_delete_datum 的实现和 toast_fetch_datum 一样简单：根据属性值中存储的 TOAST 指针，在 TOAST 表中删除所有属于该 TOAST 数据对应的片段元组即可。

TOAST 技术保障了 PostgreSQL 中的元组大小足以存放在一个文件块中，正因为如此，TOAST 技术必须被设置成为一种系统自动处理的机制，用户不能加以控制。而且 TOAST 机制主要集中于变长的数据类型。如果用户需要存储文件型的大数据或者要对自己使用的大数据进行主动控制，就需要用到下面要介绍的大对象存储技术。

2. 大对象

随着信息技术的发展，数据库需要存储的数据对象变得越来越大，为了解决大对象的存储，PostgreSQL 提供了大对象存储机制。PostgreSQL 的大对象存储机制可以支持三种数据类型的存储：

- 1) 二进制大对象 (BLOB)：主要用来保存非传统数据，如图片、视频、混合媒体等。
- 2) 字符大对象 (CLOB)：存储大的单字节字符集数据，如文档等。
- 3) 双字节字符大对象 (DBCLOB)：用于存储大的双字节字符集数据，如变长双字节字符图形字符串。

这些类型可以保存诸如音视频、图片以及文本等大数据量的对象，最大可支持 2G 的尺寸。

(1) 大对象存储

在 PostgreSQL 中，所有的大对象都存储在一个叫 pg_largeobject 的系统表中，其在数据库中的 OID 固定为 2613。pg_largeobject 的属性如表 3-3 所示。

一个大对象是使用其创建时分配的 OID 标识的，即 OID 是 PostgreSQL 数据库识别大对象的唯一标记。一个大对象将会被分成若干元组存放在系统表 pg_largeobject 中，每一个元组也称为一个页面。一个元组的大小为 2KB（默认值）。通常一个 20M 以上的文件，将会被存储为 1 万多个元组。这里将一个元组大小设置为 2KB 主要是出于两个目的：第一是更新某个元组时不会浪费太大的空间，第二是可以触发 TOAST 的压缩机制（2KB 是 TOAST 压缩机制的临界点）。

PostgreSQL 对于大对象存储允许松散的存储，可以容忍页面的丢失。大对象内丢失的部分在读取的时候将读做零。例如一个大对象有 6 个页面（编号从 0~5），现在因为某种原因导致 1、2 号页面内容丢失，但 pg_largeobject 表不会将该大对象的这两个页面删除，只是在以后读到这两个页面时会将它们以全 0 方式读出来。

为了加快大对象的查找效率，PostgreSQL 为 pg_largeobject 创建了一个叫 pg_largeobject_loid_pn_index 的索引表，可以利用该索引根据大对象的 OID 及它的 pageno 对大对象的某一页进行快速查找。

(2) 大对象管理

如数据结构 3.7 所示，当我们要对一个大对象进行实际操作时，我们都要先为该大对象创建一个大对象描述符用来记录它的相关信息。

表 3-3 pg_largeobject 系统表属性

属性名称	类型	描述
loid	oid	包含本页的大对象的标识符
pageno	int4	本页在其大对象数据中的页码从零开始计算
data	bytea	存储在大对象中的实际数据

数据结构 3.7 LargeObjectDesc

typedef struct LargeObjectDesc	
{	
Oid id;	
Snapshot snapshot;	
SubTransactionId subid;	
uint32 offset;	
int flags;	
}	LargeObjectDesc;

其中各个字段意义如下：

- `id`: 指的是分配给该大对象的 OID。
- `snapshot`: 用于大对象读写时检查数据可见性的快照；快照中存储当前事务和命令的相关信息以及用于匹配元组的函数指针。
- `subid`: 打开该大对象的子事务 ID 或者当前拥有该大对象的子事务 ID。
- `offset`: 指明该大对象内部当前的读写指针位置。
- `flags`: 指明当要对大对象进行操作时，是以读锁还是写锁的方式打开 `pg_largeobject`。

大对象的操作主要有创建、打开、读取、写入、关闭和删除，其实现都比较简单，下面将逐一介绍。

1) 创建大对象。

创建一个大对象由 `inv_create` 函数实现，该函数具有一个类型为 OID 的参数 `lobjId`，该参数可以用来指定要创建的大对象的 OID。该函数的流程如下：

①如果在参数 `lobjId` 中没有指定 OID（即 `lobjId` 的值为 `InvalidOid`），则先为要创建的大对象分配一个 OID，然后在 `pg_largeobject` 系统表中检查该 OID 是否已经存在。如果 OID 已存在，则报错返回。

②在 `pg_largeobject` 中插入一个元组，该元组的 `loid` 为该大对象的 OID，`pageno` 为 0，`data` 的 size 为 0。然后更新 `pg_largeobject` 上的索引并返回大对象的 OID。

2) 打开大对象。

该操作实际上是由已创建的大对象生成一个大对象操作符，由函数 `inv_open` 实现。该函数有三个参数：大对象的 OID、打开标志以及打开大对象描述符所在的内存上下文。其中，打开标志用于指定打开大对象时对其加锁的方式，取值可以是 `INV_WRITE` 或 `INV_READ`。`inv_open` 函数的流程如下：

①首先在内存上下文中为该大对象创建一个大对象描述符，其中的 `id` 字段置为大对象的 OID，`offset` 置为 0，`subid` 调用 `GetCurrentSubTransactionId` 函数置为当前子事务的 id（顶层事务的子事务 id 是 1，见 7.4.2 节）。

②打开标志若是 `INV_WRITE`，则描述符里的 `flags` 置为 `IFS_WRLOCK | IFS_RDLOCK`；否则描述符里的 `flags` 置为 `IFS_RDLOCK`。

③最后检查在 `pg_largeobject` 中是否有属于该大对象的元组，如果没有属于该大对象的元组则报告错误返回；否则正常返回该大对象的描述符。

3) 读取大对象。

该操作从大对象中读取指定大小的数据放入到给定的内存中，由函数 `inv_read` 实现，该函数有三个参数：大对象描述符、缓冲区以及要读取的字节数。`inv_read` 的返回值为实际读出的字节数。

该函数的流程由一个循环组成，根据大对象描述符中保存的当前读写位置 `offset` 以及大对象页面的大小可以计算出这一次读出的内容所涉及的页面范围，然后在循环中利用大对象 OID 从 `pg_largeobject` 中逐个取出大对象的页面，并将其数据放入缓冲区中。如果遇到某些页面丢失，`inv_read` 将会把该页面以全 0 的方式读出。最后将返回实际读出的字节数。

4) 写入大对象。

该操作向大对象中写入指定大小的数据，由函数 `inv_write` 实现。该函数有三个参数：大对象描

述符、缓冲区以及要写入的字节数。其中，缓冲区中存放的是需要写入的数据，`inv_write` 的返回值为实际写入的字节数。

`inv_write` 的实现和 `inv_read` 的思想类似，只是数据的流向相反。`inv_write` 先从大对象描述符找到该大对象的 OID 和读写指针的位置，然后从读写指针开始将缓冲区中的数据逐页面地写入，写满一页换另一页。在写入的过程中，如果当前读写指针的位置还在大对象已有的页面中，则覆盖原来页的内容；如果读写指针的位置将要超过大对象的最后一页，则新增加页面并继续写数据，直到所有数据都被写入。新增的页面写完后将会以元组形式插入到 `pg_largeobject`。最后将实际写入的字节数返回。

5) 关闭大对象。

关闭大对象实际是指关闭一个被 `inv_open` 打开的大对象描述符，由函数 `inv_close` 实现。该函数会将大对象描述符相关的信息从内存中彻底清除。

6) 删除大对象。

删除大对象是指将该大对象在数据库中的存储部分删除，而不是仅仅清除其描述符，由函数 `inv_drop` 实现。该函数先根据大对象的 OID 找到大对象在 `pg_largeobject` 中的所有元组，并依次将其删除，最终达到删除该大对象的目的。

(3) TOAST 与大对象的比较

TOAST 和大对象虽然都是对大数据进行存储的技术，但两者有一定的区别。

首先，TOAST 用于存储变长的数据，如 `VARCHAR` 变量等。只要表中有变长的数据类型，都会自动地创建 TOAST 表，但只有存入的数据长度超过 2KB，才会触发 TOAST 存储机制。而大对象操作是客户端通过代码或者 SQL 命令调用的。因此，两者的第一个区别在于 TOAST 是可变长数据类型的一种大数据存储机制，属于自动触发机制；而大对象属于用户手动调用机制。

其次，TOAST 中的数据不能丢失，一旦丢失则会报错；而大对象中允许数据丢失，如果丢失，则用 0 替代。

第三，文件不适合使用 TOAST 技术进行存储，需要将文件以二进制方式读出，然后将二进制当作字符串存储到变长数据类型的属性中。而读取时则需要将数据读取出来，然后再转变成文件。而大对象操作不一样，它直接将文件作为一个对象存储到大对象表中，读取的时候也直接读取成一个文件，大对象对文件的存储非常容易通过编程来实现。

最后，TOAST 和大对象操作保存大数据时都是采用了将数据切成片段存储到表中的方式。但 TOAST 提供了线外和压缩两种存储机制，而大对象只是对数据不做处理直接存储。

3.3 内存管理

不管是什么样的数据库管理系统，其存储管理中涉及的问题本质上是一样的：如何减少 I/O 次数。在磁盘上读或写一个块大约要花 10 ~30 毫秒，在这段时间内一台普通的机器或许能执行数万条指令。在通常情况下，读写磁盘所用的时间决定了数据库操作所花费的总时间。因此，要尽可能地让最近使用的文件块停留在内存中，这样就能有效地减小磁盘 I/O 的代价。合理有效的内存管理对于整个 DBMS 的性能起着非常重要的作用。

PostgreSQL 中的内存管理包括对共享内存和本地内存的管理（参见图 3-12）。接下来将对图 3-12 中的模块进行详细介绍。

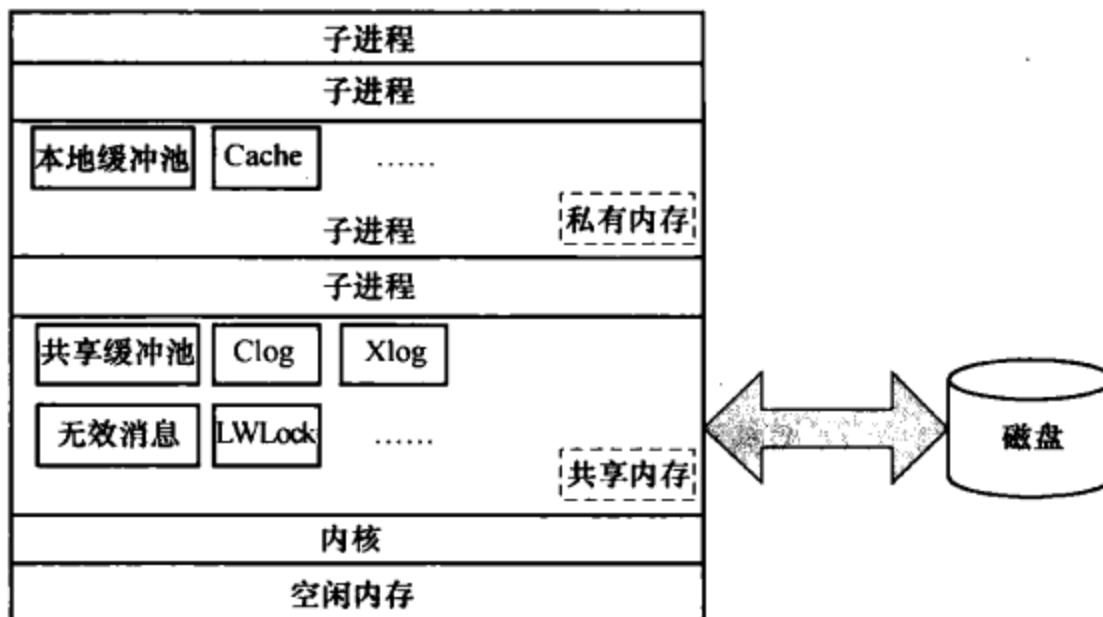


图 3-12 PostgreSQL 内存管理

3.3.1 内存上下文概述

在 PostgreSQL 的旧版本中，常常需要处理大量以指针传值的查询，因而存在着内存泄漏的问题，直到查询结束才能将内存收回。尤其是在处理引用 TOAST 机制的查询时，需要使用大量的内存，使得这个问题尤为明显。为此，从版本 7.1 开始，PostgreSQL 中实现了新的内存管理机制——内存上下文（MemoryContext）。系统中的内存分配操作在各种语义的内存上下文中进行，所有在内存上下文中分配的内存空间都通过内存上下文进行记录。因此可以很轻松地通过释放内存上下文来释放其中的所有内容，而不用费心地去释放其中的每一块内存，使得内存分配和释放更加快捷和可靠。

内存上下文机制实际上借鉴了操作系统的一些概念。我们知道，操作系统为每个进程分配了进程执行环境，进程环境之间不互相影响，由操作系统来对环境进行切换，进程可以在其进程环境中调用一些库函数（如 malloc、free、realloc 等）来执行内存操作。类似的，一个内存上下文实际上就相当于一个进程环境，PostgreSQL 以类似的方式提供了在内存上下文进行内存操作的函数：malloc、pfree、repalloc 等。

1. MemoryContext

PostgreSQL 的每一个子进程都拥有多个私有的内存上下文，每个子进程的内存上下文组成一个树形结构（参见图 3-13），其根节点为 TopMemoryContext。在根节点之下有多个子节点，每个子节点都用于不同的功能模块，例如 CacheMemoryContext 用于管理 Cache、ErrorMemoryContext 用于错误处理，每个子节点又可以有自己的子节点。

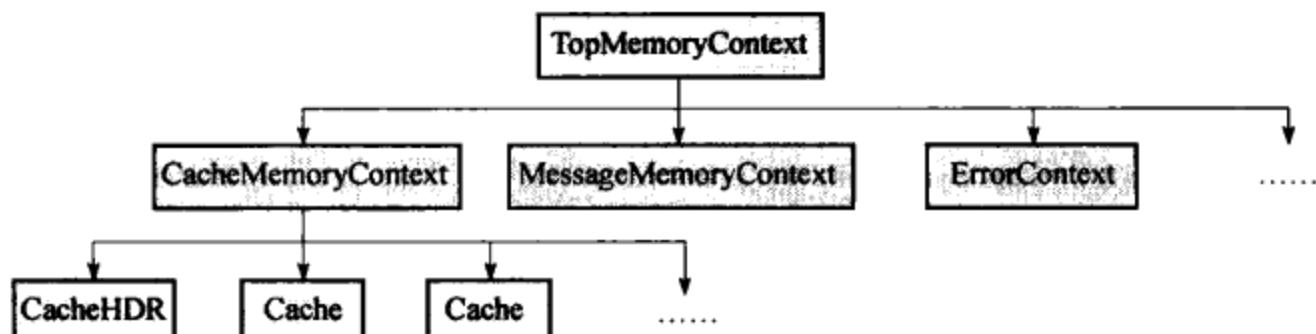


图 3-13 内存上下文树

通过树形结构可以跟踪进程中内存上下文的创建和使用情况，当创建一个新的内存上下文时，将其添到某个已存在的内存上下文下面作为其子节点。在清除内存时从根节点开始遍历内存上下文树可以将其所有节点占用的内存完全释放。

每个内存上下文中都定义了这个内存上下文所占用内存块的具体位置、大小等相关信息以及其他内存上下文之间的关联信息，只要能获得这个内存上下文，就可以获得其子节点的内存使用状态。

在数据结构 3.8 中显示了每个内存上下文的节点信息。从内存上下文的结构可以看到，每一个节点只能有一个父节点，但是可以有多个子节点。

数据结构 3.8 MemoryContext

```
typedef struct MemoryContextData
{
    NodeTag           type;          // 内存节点类型
    MemoryContextMethods *methods;   // 内存处理函数指针
    MemoryContext     parent;        // 父节点指针
    MemoryContext     firstchild;    // 第一个孩子节点指针
    MemoryContext     nextchild;    // 第一个兄弟节点指针
    Char              *name;         // 节点的名称
}MemoryContextData;
typedef struct MemoryContextData * MemoryContext;
```

MemoryContext 中的 methods 字段是一个 MemoryContextMethods 类型（数据结构 3.9），它是由一系列函数指针组成的集合，其中包含了对内存上下文进行操作的函数。对于不同的 MemoryContext 实现，可以设置不同的方法集合。但目前 MemoryContext 中只有 AllocSetContext 一种实现，因此 PostgreSQL 中只有针对 AllocSetContext 的一种操作函数集合，由全局变量 AllocSetMethods 表示。每当创建新的 MemoryContext 时，会将其 methods 字段置为 AllocSetMethods。

数据结构 3.9 MemoryContextMethods

```
typedef struct MemoryContextMethods
{
    void * (*alloc) (MemoryContext context, Size size);           // 分配内存函数
    void (*free_p) (MemoryContext context, void *pointer);        // 释放内存函数
    void * (*realloc) (MemoryContext context, void *pointer, Size size); // 重分配内存
    void (*init) (MemoryContext context);                          // 初始化内存上下文
    void (*reset) (MemoryContext context);                        // 重置内存上下文
    void (*delete) (MemoryContext context);                       // 删除内存上下文
    Size (*get_chunk_space) (MemoryContext context, void *pointer); // 检查内存片段的大小
    bool (*is_empty) (MemoryContext context);                     // 检查内存上下文是否为空
    void (*stats) (MemoryContext context, int level);           // 打印内存上下文状态
    void (*check) (MemoryContext context);                        // 检查所有内存片段
} MemoryContextMethods;
```

全局变量 AllocSetMethods 中指定了 AllocSetContext 实现的操作函数，它们一一对应 MemoryContextMethods 中的操作函数：AllocSetAlloc、AllocSetFree、AllocSetRealloc、AllocSetInit、AllocSetReset、AllocSetDelete、AllocSetGetChunkSpace、AllocSetIsEmpty、AllocSetStats 和 AllocSetCheck。后面我们将看到，对于内存上下文的操作都是通过这些函数来实现的。

在任何时候，都有一个“当前”的 MemoryContext，记录在全局变量 CurrentMemoryContext 里，进程就在这个内存上下文中调用 palloc 函数来分配内存。在需要变换内存上下文时，可以使用 MemoryContextSwitchTo 函数将 CurrentMemoryContext 指向其他的内存上下文。

MemoryContext 是一个抽象类，可以有多个实现，但目前只有 AllocSetContext 一个实现。事实上，MemoryContext 并不管理实际上的内存分配，仅仅是用作对 MemoryContext 树的控制。管理一个内存上下文中的内存块是通过 AllocSet 结构来完成的，而 MemoryContext 仅作为 AllocSet 的头部信息存在，AllocSet 是一个指向 AllocSetContext 结构的类型指针，AllocSetContext 的结构如数据结构 3.10 所示。

数据结构 3.10 AllocSetContext

```
typedef struct AllocSetContext
{
    MemoryContextData header;           // 对应于该内存上下文的头部信息
    AllocBlock blocks;                 // 该内存上下文中所有内存块的链表
    AllocChunk freelist[ALLOCSET_NUM_FREELISTS]; // 该内存上下文中空闲内存片的数组
    bool isReset;                     // 如果为真表示从上次重置以来没有分配过内存
    Size initBlockSize;               // 初始内存块的大小
    Size maxBlockSize;                // 最大内存块大小
    Size nextBlockSize;               // 下一个要分配的内存块的大小
    Size allocChunkLimit;             // 分配内存片的尺寸阈值，该值在分配内存片时会用到
    AllocBlock keeper;                // 保存在 keeper 中的内存块在内存上下文重置时
                                    // 会被保留不做释放
} AllocSetContext;
```

在 AllocSetContext 中，有几个字段需要进一步说明：

- isReset：PostgreSQL 中提供了对内存上下文的重置操作（将在本节的后续部分介绍），所谓重置就是释放内存上下文中所有分配的内存，并将这些内存交还给操作系统。在一个内存上下文被创建时，其 isReset 字段置为 True，表示从上一次重置到当前没有内存被分配。只要在该内存上下文中进行了分配，则将其 isReset 字段置为 False。这样在进行重置时，可以检查内存上下文的 isReset 字段，如果为 True 则表示该内存上下文中没有进行过内存分配，所以不需要进行实际的重置工作，从而提高操作效率。
- initBlockSize、maxBlockSize、nextBlockSize：initBlockSize 和 maxBlockSize 字段在内存上下文创建时指定，且在创建时 nextBlockSize 会置为与 initBlockSize 相同的值。nextBlockSize 表示下一次分配的内存块的大小，在进行内存分配时，如果需要分配一个新的内存块，则这个新内存块的大小将采用 nextBlockSize 的值。在后面的介绍中将会看到，有些情况下需要将下一次要分配的内存块的大小置为上一次的 2 倍，所以 nextBlockSize 是会变大的。当对内存上下文进行重置时，需要将 nextBlockSize 恢复到初始值，也就是 initBlockSize，所

以 initBlockSize 充当了内存块初始大小的备份值。虽然 nextBlockSize 可以增大，但也并不能无限制地增加，maxBlockSize 字段指定了内存块可以到达的最大尺寸。

- allocChunkLimit：内存块内会分成多个称为内存片的内存单元，在分配内存片时，如果一个内存片的尺寸超过了宏 ALLOC_CHUNK_LIMIT 时，将会为该内存片单独分配一个独立的内存块，这样做是为了避免日后进行内存回收时造成过多的碎片。由于宏 ALLOC_CHUNK_LIMIT 是不能在运行时更改的，因此 PostgreSQL 提供了 allocChunkLimit 用于自定义一个阈值。如果定义了这个字段的值，则在进行超限检查时会用该字段来替换宏定义进行判断。
- keeper：在内存上下文进行重置时不会对 keeper 中记录的内存块进行释放，而是对其内容进行清空。这样可以保证内存上下文重置结束后就已经包含一定的可用内存空间，而不需要通过 malloc 另行申请。另外也可以避免在某个内存上下文被反复重置时，反复进行 malloc 带来的风险。例如，在执行查询时，每一个元组被获得时，都需要有一个内存上下文用于存放与之相关的信息，而当获取下一个元组时，该内存上下文将被重置后重复使用。通过 keeper 字段保留一个内存块，可以避免每次重置后都进行 malloc 操作。

内存上下文的组织结构如图 3-14 所

示。每个 AllocSet 结构都对应一个内存上下文，AllocSet 所管理的内存区域被分成若干个内存块（block），内存块用 AllocBlockData 结构（数据结构 3.11）表示。每个内存块内又被分成多个称为内存片（chunk）的单元。AllocSet 结构主要包括以下三部分信息：

(1) 头部信息 header

头部信息就是一个 MemoryContextData 结构，header 是进入一个内存上下文的唯一外部接口，事实上管理内存上下文的接口函数都是通过对 header 的管理来实现的。

(2) 内存块链表 blocks

该元素为一个指向 AllocBlockData 结构体（数据结构 3.11）的指针，AllocBlockData 用于表示一个内存块。AllocBlockData 之间通过其 next 字段链接成一个单向链表，而 AllocSet 的 blocks 字段则指向这个链表的头部。

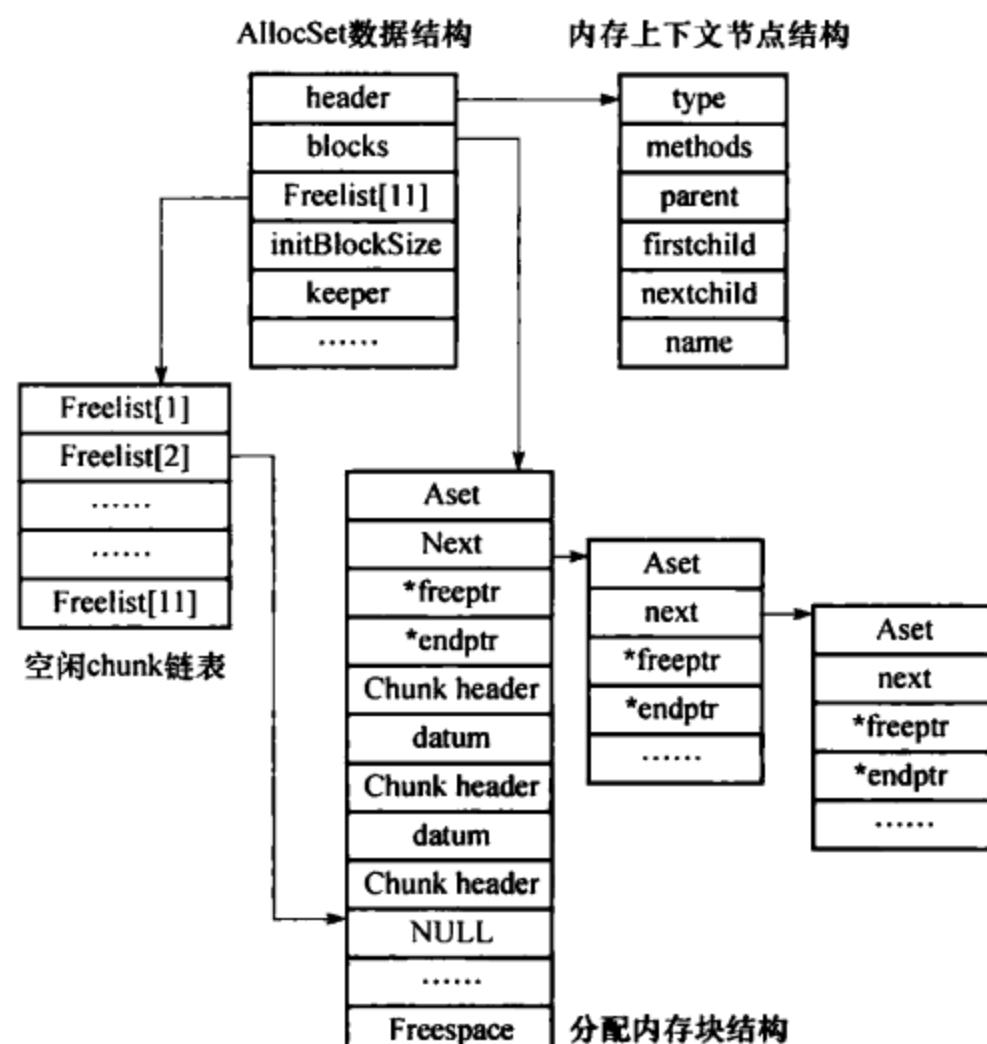


图 3-14 AllocSet 的结构

数据结构 3.11 AllocBlockData

```

typedef struct AllocBlockData
{
    AllocSet      aset;           // 该内存块所在的 AllocSet

```

```

    AllocBlock      next;          //指向下一个内存块的指针
    char           *freeptr;      //指向该块空闲区域的首地址
    char           *endptr;       //该内存块的末地址
}AllocBlockData;

```

AllocBlockData 记录在一块内存区域的起始地址处，这块内存区域通过标准库函数 malloc 进行分配，称为一个内存块。在每个内存块中进行内存分配时产生的内存片段称之为内存片，每个内存片包括一个头部信息和数据区域，其中头部信息包括该内存片所属的内存上下文以及该内存区的其他相关信息，数据区则存储实际数据。内存片的头部信息由数据结构 AllocChunkData 描述（数据结构 3.12），内存片的数据区域则紧跟在其头部信息之后分配。通过 PostgreSQL 中定义的 palloc 函数和 pfree 函数，我们可以自由地在内存上下文中申请和释放内存片，被释放的内存片将被加入到 FreeList 中以备重复使用。

数据结构 3.12 AllocChunkData

```

typedef struct AllocChunkData
{
    void  *aset;          //该内存片所在的 AllocSet,如果内存片为空闲,则用于链接其空闲链表
    Size   size;          //内存片的实际大小,由于内存片都是以 2 的幂为大小进行对齐,因此申请的大小
                          //可能比实际大小要小
    Size   requested_size; //内存片中被使用的空间大小,如果是空闲内存片则置为 0
}AllocChunkData;

```

(3) FreeList 数组

该数组用于维护在内存块中被回收的空闲内存片，这些空闲内存片将被用于再分配。FreeList 数组元素类型为 AllocChunk，数组长度默认为 11（由宏 ALLOCSET_NUM_FREELISTS 定义）。

FreeList 数组中的每一个元素指向一个由特定“大小”空闲内存片组成的链表，这个“大小”与该元素在数组中的顺序有关。比如，FreeList 数组中第 K 个元素所指向链表的每个空闲数据块的大小为 $2^{(k+2)}$ 字节，空闲内存片最小为 8 字节，最大不超过 8K 字节。因此，FreeList 数组中实际上维护了 11 个空闲链表。管理着 11 种“大小”的空闲内存片。从 AllocChunkData 的结构我们可以发现，并没有明显的字段用来将内存片链接到空闲链表。其实 aset 字段具有两种用途，如果一个内存片正在使用，则它的 aset 字段指向其所属的 AllocSet。如果内存片是空闲的，也就是说它处于某个空闲链表中，那么它的 aset 字段指向空闲链表中在它之后的内存片。这样从 FreeList 数组元素所指向的链表头部开始，顺着 aset 字段指向的下一个内存片就可以找到该空闲链表中所有的空闲内存片。

需要注意的是，所有 FreeList 中的内存片的大小都为 2 的指数。当需要申请一块内存时，我们可以迅速地定位到相应的空闲空间链表中。对于一个大小为 size 的内存分配请求，将会在第 K 个空闲链表中为其分配内存片，K 的计算规则是：

- 当 $size < (1 \ll \text{ALLOC_MINBITS})$ 时，K 为 0，其中 ALLOC_MINBITS 系统预定义为 3，即表示 FreeList 中保存的最小字节数为 $2^3 = 8$ 字节。
- 当 $(1 \ll \text{ALLOC_MINBITS}) < size < \text{ALLOC_CHUNK_LIMIT}$ 时，若 $2^{N-1} < size < 2^N$ ，则 $K = N - 3$ ，其中 ALLOC_CHUNK_LIMIT 为 FreeList 数组中所能维持空闲内存片的最大值。

一个内存片的大小不会超过 ALLOC_CHUNK_LIMIT，否则将会被当作一个独立的块来看待，具体的内存分配过程在本节的后续部分中会详细介绍。

这样的分配方式实际上会造成一定的内存浪费，因为大多数时候分配的内存总是略大于请求的内存，但这种策略将为内存片重分配提供便利。

在创建一个 AllocSet 时，我们可以为其指定 allocChunkLimit 字段的值，如果申请分配的内存大小超过这个值，那么将为这次请求分配一个独立的内存块，这个内存块中只存放一个内存片，当该内存片释放的时候，整个内存块将被释放，而不是将内存片加到 FreeList 中。

2. 内存上下文初始化与创建

任何一个 PostgreSQL 进程在使用内存上下文之前，都需要先进行初始化。内存上下文的初始化工作由函数 MemoryContextInit 来完成。在初始化时首先创建所有内存上下文的根节点 TopMemoryContext，然后在该节点下创建子节点 ErrorContext 用于错误恢复处理：

- TopMemoryContext：该节点在分配后将一直存在，直到系统退出时候才释放。在该节点下面分配其他内存上下文节点本身所占用的空间。
- ErrorContext：该节点是 TopMemoryContext 的第一个子节点，是错误恢复处理的永久性内存上下文，恢复完毕就会进行重置。

当初始化完毕，建立根节点和错误恢复子节点后，PostgreSQL 进程就可以开始创建其他的内存上下文。内存上下文的创建由 AllocSetContextCreate 函数来实现，主要有两个工作：创建内存上下文节点以及分配内存块。

内存上下文节点的创建由 MemoryContextCreate 函数来完成，主要流程如下：

- 1) 从 TopMemoryContext 节点中分配一块内存用于存放内存上下文节点（该内存上下文管理的内存块不包括在内），该内存块略大于一个 AllocSet 结构体的大小。

- 2) 初始化 MemoryContext 节点，设置其父节点、节点类型、节点名等相关信息，还要设置该节点的 methods 属性为 AllocSetMethods。

AllocSetContextCreate 函数在调用 MemoryContextCreate 函数完成了内存节点初始化后，将根据用户请求为其预分配内存块。AllocSetContextCreate 包括以下 5 个参数：

- parent：类型为 MemoryContext，表示当前要创建的内存上下文的父节点。
- name：字符串类型，是当前要创建的内存上下文的名称。
- minContextSize、initBlockSize、maxBlockSize：都是 size_t 类型，分别表示内存上下文的最小尺寸、初始内存块尺寸和最大内存块尺寸。

AllocSetContextCreate 函数的处理流程如下：

- 1) 调用 MemoryContextCreate 创建 MemoryContext 结构。
- 2) 将上一步创建的 MemoryContext 结构强制转换为 AllocSet 结构体，填充 AllocSet 其他结构体元素的信息，包括最小块大小、初始化块大小以及最大块大小，并根据最大块大小设置其 allocChunkLimit 的值。
- 3) 如果 minContextSize 超过一定限制（内存块头部信息尺寸和内存片头部信息尺寸之和）时，调用标准库函数，以 minContextSize 为大小分配一个内存块，并初始化块结构体，加入 AllocSet 的内存块链表中。需要注意的是，预分配的内存块将被记入到内存上下文的 keeper 字段中，作为内存上下文的保留块，以便重置内存上下文的时候，该内存块不会被释放。

当完成以上步骤后，我们就可以在该内存上下文中进行内存的分配操作，主要通过 MemoryContextAlloc 和 MemoryContextAllocZero 两个接口函数来完成，其中后者会将分配的内存块内容清零。

3. 内存上下文中内存的分配

在 PostgreSQL 中，内存的分配、重分配和释放都是在内存上下文中进行，因此不再使用 C 语言的标准库函数 malloc、realloc 和 free 来操作，PostgreSQL 实现了 palloc、repalloc 和 pfree 来分别实现内存上下文中对于内存的分配、重分配和释放。

- **palloc:** palloc 是一个宏定义，它会被转换为在“当前”内存上下文中对 MemoryContextAlloc 函数的调用，而 MemoryContextAlloc 函数实际上是调用了“当前”内存上下文的 methods 字段中所指定的 alloc 函数指针。对于目前 PostgreSQL 的实现来说，调用 palloc 实际就是调用了 alloc 指向的 AllocSetAlloc 函数。使用 palloc 分配的内存空间中的内容是随机的，与之相对应的还定义了一个宏 palloc0，后者会将分配的内存中的内容全部置为 0。
- **realloc:** realloc 是一个函数，其参数是一个内存片的指针和新的内存片大小。realloc 将调用内存片所属的内存上下文的 realloc 函数指针，把该内存片调整为新的大小，并返回新内存片的指针。目前，realloc 函数指针对应于 AllocSetRealloc 函数。
- **pfree:** pfree 是一个函数，其参数是一个内存片的指针，pfree 将调用内存片所属的内存上下文的 methods 字段中的 free_p 函数指针来释放内存片的空间。目前，PostgreSQL 中 free_p 指针实际指向 AllocSetFree 函数。

函数 AllocSetAlloc 负责处理具体的内存分配工作，该函数的参数为一个内存上下文节点以及需要申请的内存大小。具体的内存分配流程如下所示：

- 1) 判断需要申请的内存大小是否超过了当前内存上下文中允许分配内存片的最大值（即内存上下文节点的 allocChunkLimit 字段）。若超过，则为其分配一个新的独立的内存块，然后在该内存块中分配指定大小的内存片。接下来将该内存块加入到内存块链表中，最后设置内存上下文的 isReset 字段为 False 并返回内存片的指针。如申请的大小没有超过限制则执行步骤 2。

- 2) 计算申请的内存大小在 FreeList 数组中对应的位置，如果存在合适的空闲内存片，则将空闲链表的指针（Freelist 数组的某个元素）指向该内存片的 aset 字段所指向的地址（在空闲内存片中，aset 字段指向它在空闲链表中的下一个内存片）。然后将该内存片的 aset 字段指向其所属的内存上下文节点，最后返回该内存片的指针。如果空闲链表中没有满足要求的内存片则执行步骤 3。

- 3) 对内存上下文的内存块链表（blocks 字段）的第一个内存块进行检查，如果该内存块中的未分配空间足以满足申请的内存，则直接在该内存块中分配内存片并返回内存片的指针。这里可以看到，在内存上下文中进行内存分配时，总是在内存块链表中的第一个内存块中进行，当该内存块中空间用完之后会分配新的内存块并作为新的内存块链表首部，因此内存块链表中的第一块也称作活动内存块。如果内存块链表中的第一个内存块没有足够的未分配空间则执行步骤 4。

- 4) 由于现有的内存块都不能满足这一次内存分配的要求，因此需要申请新的内存块，但是当前的活动内存块中还有未分配空间，如果申请新的内存块并将之作为新的活动内存块，则当前活动内存块中的未分配空间就会被浪费。为了避免浪费，这里会先将当前活动内存块中的未分配空间分解成个数尽可能少的内存片（即每个内存片尽可能大），并将它们加入到 FreeList 数组中，然后创建一个新的内存块（其大小为前一次分配的内存块的两倍，但不超过 maxBlockSize）并将之作为新

的活动内存块（即加入到内存块链表的首部）。最后在活动内存块中分配一个满足申请内存大小的内存片，并返回其指针。

需要说明的是，AllocSetAlloc 为内存片分配空间时并不是严格按照申请的大小来分配的，而是将申请的大小向上对齐为 2 的幂，然后按照对齐后的大小来分配空间。例如，我们要申请一块大小为 30 字节的空间，则 AllocSetAlloc 实际会为我们分配一块大小为 $2^5 = 32$ 字节的内存片，该内存片对应的 AllocChunkData 中的 size 字段设置为 32，而 requested_size 设置为 30。

4. 内存上下文中的内存重分配

内存重分配由 AllocSetRealloc 函数实现。AllocSetRealloc 将在指定的内存上下文中对参数 pointer 指向的内存空间进行重新分配，新分配的内存大小由参数 size 指定。pointer 所指向的内存中的内容将被复制到新的内存中，并释放 pointer 指向的内存空间。AllocSetRealloc 函数的返回值就是指向新内存空间的指针。整个重分配内存流程如下：

- 1) 由于内存片在分配之初就被对齐为 2 的幂，因此有可能参数 pointer 指向的旧的内存空间的大小本来就大于参数 size 指定的新大小。如果是这种情况，则修改 pointer 所指向的 AllocChunkData 的 requested_size 为新的内存大小并返回 pointer；否则执行步骤 2。

- 2) 若 pointer 所指向的内存片占据一个内存块时，则找到这个内存块并增大这个内存块的空间，即将该内存块的 freeptr 和 endptr 指针都向后移动到 size 所指定的位置。如果 pointer 指向的内存片不是独占一个内存块则转而执行步骤 3。

- 3) 调用 AllocSetAlloc 分配一个新的内存片，并将 pointer 所指向内存片的数据复制到其中。然后调用 AllocSetFree 函数释放旧的内存片，如果是占用一个内存块的内存片则直接释放（这样的内存片所占的内存空间较大，直接释放不会造成过多碎片）；否则将其加入到 FreeList 中以便下次分配（经常释放较小的内存片会造成内存空间碎片化）。

5. 释放内存上下文

释放内存上下文中的内存，主要有以下三种方式：

(1) 释放一个内存上下文中指定的内存片

当释放一个内存上下文中指定的内存片时，调用函数 AllocSetFree。该函数的执行方式如下：

- 1) 如果指定要释放的内存片是内存块中唯一的一个内存片，则将该内存块直接释放。
- 2) 否则，将指定的内存片加入到 FreeList 链表中以便下次分配。

(2) 重置内存上下文

重置内存上下文的工作由函数 AllocSetReset 完成。在进行重置时，内存上下文中除了在 keeper 字段中指定要保留的内存块外，其他内存块全部释放，包括空闲链表中的内存。keeper 中指定保留的内存块将被清空内容，它使得内存上下文重置之后就立刻有一块内存可供使用。

(3) 释放当前内存上下文中的全部内存块

这个工作由 AllocSetDelete 函数完成，该函数释放当前内存上下文中的所有内存块，包括 keeper 指定的内存块在内。但内存上下文节点并不释放，因为内存上下文节点是在 TopMemoryContext 中申请的内存，将在进程运行结束时统一释放。

3.3.2 高速缓存

当数据库访问表时，需要表的模式信息，比如表的列属性、OID、统计信息等。PostgreSQL 将

表的模式信息存放在系统表中，因此要访问表，就需要首先在系统表中取得表的模式信息。对于一个 PostgreSQL 系统来说，对于系统表和普通表模式的访问是非常频繁的。为了提高这些访问的效率，PostgreSQL 设立了高速缓存（Cache）来提高访问效率。Cache 中包括一个系统表元组 Cache（SysCache）和一个表模式信息 Cache（RelCache）。SysCache 中存放的是最近使用过的系统表的元组，而 RelCache 中包含所有最近访问过的表的模式信息（包含系统表的信息）。RelCache 中存放的不是元组，而是 RelationData 数据结构（数据结构 3.13），每一个 RelationData 结构表示一个表的模式信息，这些信息都由系统表元组中的信息构造而来。值得注意的是，两种 Cache 都不是所有进程共享的，每一个 PostgreSQL 的进程都维护着自己的 SysCache 和 RelCache。

数据结构 3.13 RelationData

```

typedef struct RelationData
{
    RelFileNode          rd_node;           // 表的物理标识, 包括表空间、数据库、表的 OID
    struct SMgrRelationData *rd_smgr;       // 表的文件句柄
    .....
    Form_pg_class        rd_rel;            // 表在 pg_class 系统表中对应的元组里的信息
    TupleDesc             rd_att;            // 表的元组描述符, 描述了表的各个属性
    Oid                  rd_id;              // 表的 OID
    List                 *rd_indexlist;      // 表上所有索引的 OID 链表
    Bitmapset            *rd_indexattr;     // rd_indexlist 中各索引用到的属性
    Oid                  rd_oidindex;       // 若 OID 属性上有唯一索引, 这里记录该索引的 OID
    .....
    Form_pg_index         rd_index;          // 如果该表是一个索引表, 这里记录它在 pg_index 中的信息
    .....
} RelationData;

```

1. SysCache

SysCache 主要用于缓存系统表元组。从实现上看 SysCache 就是一个数组，数组的长度为预定义的系统表的个数。在 PostgreSQL 8.4.1 中实现了 54 个系统表，因此 SysCache 数组具有 54 个元素，每个元素的数据结构为 CatCache（数据结构 3.14），该结构体内使用 Hash 来存储被缓存的系统表元组，每一个系统表唯一地对应一个 SysCache 数组中的 CatCache 结构。每个 CatCache 都有若干个（不超过 4 个）查找关键字，这些关键字及其组合可以用来在 CatCache 中查找系统表元组，在初始化数据集簇时会在这些关键字上为系统表创建索引。

数据结构 3.14 CatCache

```

typedef struct catcache
{
    in                  id;                // CatCache 的 ID
    struct catcache   *cc_next;          // 连接到 SysCache 中的下一个 CatCache
    const char         *cc_relname;        // CatCache 对应的系统表名称

```

```

Oid          cc_reloid;           //CatCache 对应的系统表的 OID
Oid          cc_indexoid;        //CatCache 查找关键字上的索引 OID
bool         cc_relissshared;    //对应的系统表是否是数据库间共享的
TupleDesc    cc_tupdesc;         //对应的系统表元组描述符
int          cc_reloidattr;     //系统表中 OID 属性的属性号, 0 表示没有 OID 属性
int          cc_ntup;            //缓存在该 CatCache 中的元组数量
int          cc_nbuckets;       //该 CatCache 中的 Hash 桶数量
int          cc_nkeys;           //查找关键字个数
int          cc_key[4];          //每个关键字在系统表中的属性号
PGFunction  cc_hashfunc[4];     //用于每一个关键字的 Hash 函数
ScanKeyData cc_skey[4];         //为关键字预先填充好的 ScanKey 结构
bool         cc_isname[4];       //关键字的数据类型是否为 name 类型, name 类型用于系统
                                //表中表示名称的属性
Dllist      cc_lists;           //由 CatClist 结构构成的 Dllist 链表, 用于缓存部分匹配
                                //的元组, 每一个 CatClist 结构保存一次部分匹配的结果
                                //元组

#ifndef CATCACHE_STATS
long         cc_searches;        //总共查询该 CatCache 的次数, 不管查询是否成功都会加 1
long         cc_hits;             //命中该 CatCache 的次数
long         cc_neg_hits;        //负元组的命中次数, 如果一次查找在 CatCache 和物理表中
                                //都无法找到匹配的元组, 则会将其作为一个负元组加入到
                                //CatCache 中
long         cc_newloads;         //将 CatCache 中没有的元组载入其中的次数
long         cc_invals;           //该 CatCache 中无效元组的个数
long         cc_lsearches;        //在该 CatCache 中进行部分匹配查询的次数
long         cc_lhits;             //在 cc_lists 中命中的次数
#endif
Dllist      cc_bucket[1];         //由 Hash 桶构成的 Dllist 链表, Hash 桶中实际缓存着系统
                                //表元组
}CatCache;

```

(1) SysCache 初始化

在 Postgres 进程初始化时（在 `InitProgres` 中），将会对 SysCache 进行初始化。SysCache 的初始化实际上是填充 SysCache 数组中每个元素的 CatCache 结构的过程，主要任务是将查找系统表元组的关键字信息写入 SysCache 数组元素中。这样通过指定的关键字可以快速定位到系统表元组的存储位置。

在 `SysCache.c` 文件中已经将所有系统表的 CatCache 信息存储在一个名为 `cacheinfo` 的静态数组中，每个系统表的 CatCache 信息用一个数组元素来描述，其数据类型为 `cachedesc`（数据结构 3.15）。

数据结构 3.15 cachedesc

```

struct cachedesc
{
    Oid      reloid;        //CatCache 对应的系统表 OID
    Oid      indoid;        //CatCache 用到的索引 OID
    int     reloidattr;    //系统表中 OID 属性的属性号
    int     nkeys;          //查询关键字的个数
    int     key[4];         //查询关键字的属性号
    int     nbuckets;       //该 CatCache 需要的 Hash 桶数
};

```

为了便于查找，SysCache 中的 CatCache 通过其 cc_next 字段构成了一个单向链表，其头部用全局变量 CacheHdr 记录，其数据结构如数据结构 3.16 所示。

数据结构 3.16 catcacheheader

```

typedef struct catcacheheader
{
    CatCache *ch_caches; //指向 CatCache 链表的头部
    int      ch_ntup;    //所有 CatCache 中缓存元组的总数
}CatCacheHeader;

```

在 Postgres 进程初始化时，会调用 InitCatalogCache 函数对 SysCache 数组进行初始化，并建立由 CacheHdr 记录的 CatCache 链表。

InitCatalogCache 函数中对 SysCache 的初始化主要分为以下几个步骤：

- 1) 根据 cacheinfo 为 SysCache 数组分配空间，这里将 SysCache 的长度设置为和 cacheinfo 数组相同。
- 2) 循环调用 InitCatcache 函数根据 cacheinfo 中的每一个元素生成 CatCache 结构并放入 SysCache 数组的对应位置中。InitCatcache 每调用一次将处理一个 cachedesc 结构。InitCatcache 将首先确保 CacheMemoryContext 存在（如不存在会创建之），然后切换到 CacheMemoryContext 中。接下来将检查 CacheHdr 是否存在，不存在则先建立 CacheHdr。然后该函数根据 cachedesc 中要求的 Hash 桶的数量为即将建立的 CatCache 结构分配内存，并根据 cachedesc 结构中的信息填充 CatCache 的各个字段。最后将生成的 CatCache 链接在 CacheHdr 所指向的链表的头部。

在 InitCatalogCache 函数中实际只完成了 SysCache 初始化的第一个阶段，在稍后被调用的函数 RelationCacheInitializePhase2（负责 RelCache 的初始化）还将调用 InitCatcachePhase2 进行第二阶段也是最后的 SysCache 初始化工作。InitCatcachePhase2 将依次完善 SysCache 数组中的 CatCache 结构，主要是根据对应的系统表填充 CatCache 结构中的元组描述符（cc_tupdesc）、系统表名（cc_relname）、查找关键字的相关字段（cc_hashfunc、cc_isname、cc_skey）等。

SysCache 数组初始化完成之后，CatCache 内是没有任何元组的，但是随着系统运行时对于系统表元组的访问，CatCache 中的系统表元组会逐渐增多。

(2) CatCache 中缓存元组的组织

CatCache 中对缓存元组的组织如图 3-15 所示。每个 CatCache 中的 cc_bucket 数组中的每一个元

素都表示一个 Hash 桶，元组的键值通过 Hash 函数可以映射到 cc_bucket 数组的下标。每一个 Hash 桶都被组织成一个双向链表（Dlist，见数据结构 3.17），其中的节点为 Dlelem 类型（数据结构 3.18），Dlelem 是一个包装过的缓存元组，其 dle_val 字段指向一个 CatCTup（数据结构 3.19）形式的缓存元组。

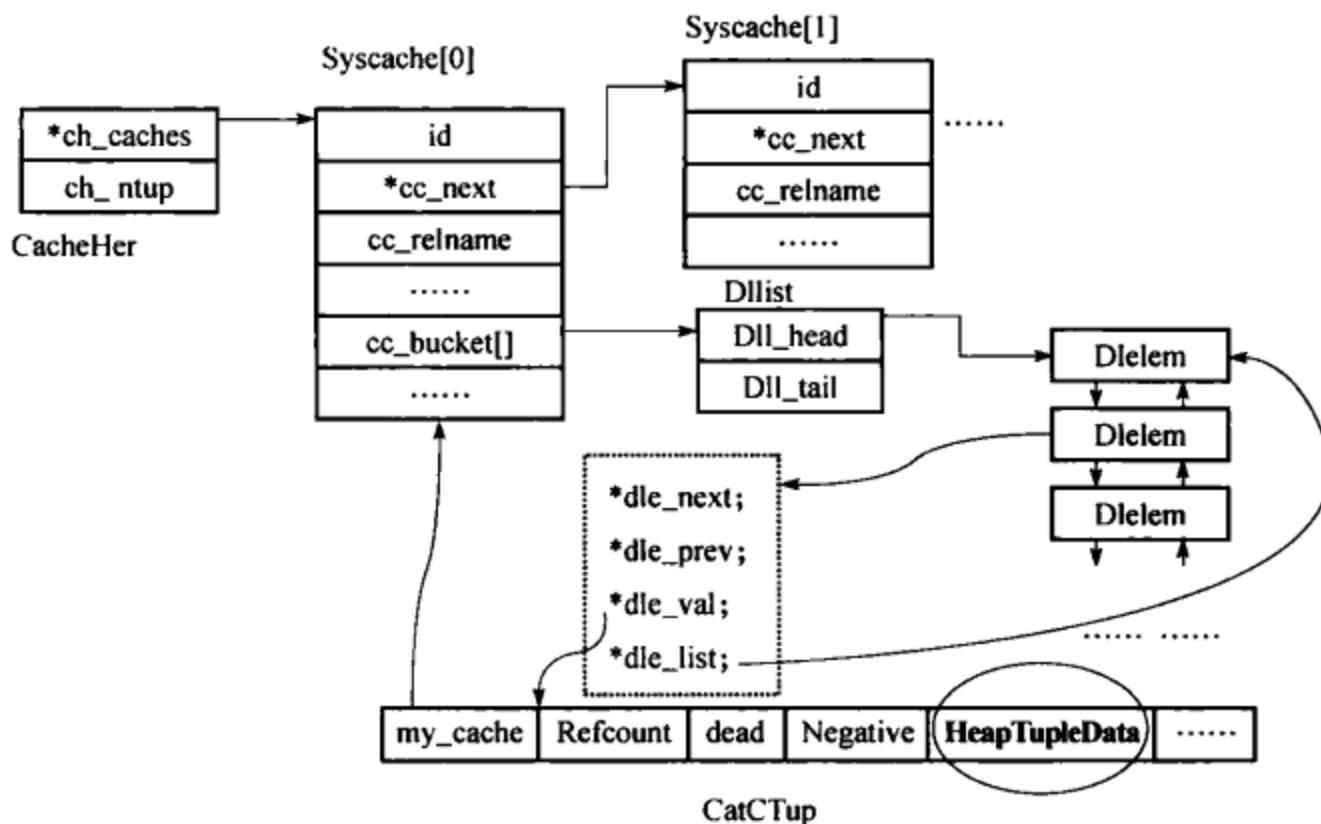


图 3-15 在 SysCache 中查找元组

从图 3-15 可以看到，具有同一 Hash 键值的元组被缓存在同一个 Hash 桶中，每一个 Hash 桶中的缓存元组都被先包装成 Dlelem 结构并链接成一个链表。因此在查找某一个元组时，需要先计算其 Hash 键值并通过键值找到其所在的 Hash 桶，之后要遍历 Hash 桶的链表逐一比对缓存元组。为了尽量减少遍历 Hash 桶的代价，在组织 Hash 桶中链表时，会将这一次命中的缓存元组移动到链表的头部，这样下一次查找同一个元组时可以在尽可能少的时间内命中。Dlelem 结构中的 dle_list 字段用来记录缓存元组所在链表的头部，以方便将该缓存元组移动到链表的头部。

CatCache 中的缓存元组将先被包装成 CatCTup 形式，然后加入到其所在 Hash 桶的链表中。在 CatCTup 中通过

数据结构 3.17 Dlist

```
typedef struct Dlist
{
    Dlelem     *dll_head;           // 指向链表的头部
    Dlelem     *dll_tail;          // 指向链表的尾部
}Dlist;
```

数据结构 3.18 Dlelem

```
typedef struct Dlelem
{
    struct Dlelem *dle_next;      // 指向下一个节点
    struct Dlelem *dle_prev;      // 指向上一个节点
    void         *dle_val;        // 节点的实际内容
    struct Dlist *dle_list;       // 指向链表头部
}Dlelem;
```

`my_cache` 和 `cache_elem` 分别指向该缓存元组所在的 `CatCache` 及 Hash 桶链表中的节点。一个被标记为“死亡”的 `CatCTup` (`dead` 字段为真) 并不会实际从 `CatCache` 中删除，但是在后续的查找中它不会被返回。“死亡”的缓存元组将一直被保留在 `CatCache` 中，直到没有人访问它，即其 `refcount` 变为 0。但如果“死亡”元组同时也属于一个 `CatCList`，则必须等到 `CatCList` 和 `CatCTup` 的 `refcount` 都变为 0 时才能将其从 `CatCache` 中清除。`CatCTup` 的 `negative` 字段表明该缓存元组是否为一个“负元组”，所谓负元组就是实际并不存在于系统表中，但是其键值曾经用于在 `CatCache` 中进行查找的元组。负元组只有键值，其他属性均为空。负元组的存在是为了避免反复到物理表中去查找不存在的元组所带来的 I/O 开销，具体将在下一节中介绍。

数据结构 3.19 CatCTup

```
typedef struct catctup
{
    int             ct_magic;          // 用于标识 CatCTup 数据结构
#define CT_MAGIC      0x57261502        // ct_magic 字段被设置为该值
    CatCache       *my_cache;         // 指向该缓存元组所在的 CatCache
    Dlelem        *cache_elem;        // 指向该缓存元组所在的 Dlelem
    struct catclist *c_list;         // 指向该缓存元组所在的 catclist, 如果该元组不属于任何 cat-
                                    //     clist 则置空值
    int             refcount;         // 对该缓存元组的引用计数
    bool            dead;              // 标记该元组是否已“死亡”
    bool            negative;         // 标记该元组是否为“负元组”
    uint32          hash_value;       // 记录该元组的 Hash 键值
    HeapTupleData tuple;            // 缓存元组的头部信息, 实际的元组数据在接下来的内存中
}CatCTup;
```

(3) 在 `CatCache` 中查找元组

在 `CatCache` 中查找元组有两种方式：精确匹配和部分匹配。前者用于给定 `CatCache` 所有的所有键值，并返回 `CatCache` 中能完全匹配这个键值的元组；而后者只需要给出部分键值，并将部分匹配的元组以一个 `CatCList` 的方式返回。

精确匹配查找由函数 `SearchCatCache` 函数实现，其函数原型如下：

```
SearchCatcache(CatCache* Cache, Datum v1, Datum v2, Datum v3, Datum v4)
```

其中，`v1`、`v2`、`v3` 和 `v4` 都用于查找元组的键值，分别对应该 `Cache` 中记录的元组搜索键。可以看到，`SearchCatcache` 最多可以使用 4 个属性的键值进行查询，4 个参数分别对应该 `CatCache` 数据结构中 `cc_key` 字段定义的查找键。

`SearchCatCache` 需要在一个给定的 `CatCache` 中查找元组，为了确定要在哪个 `CatCache` 中进行查找，还需要先通过 `CacheHdr` 遍历 `SysCache` 中所有的 `CatCache` 结构体，并根据查询的系统表名或系统表 OID 找到对应的 `CatCache`。

`SearchCatCache` 在给定的 `CatCache` 中查找元组的过程如下：

- 1) 对所查找元组的键值进行 Hash，按照 Hash 值得到该 `CatCache` 在 `cc_bucket` 数组中对应的

Hash 桶的下标。

2) 遍历 Hash 桶链找到满足查询需求的 Dlelem，并将其结构体中 dle_val 属性强制转换为 CatCTup 类型，CatCTup 中的 HeapTupleData 就是要查找的元组头部。另外，还要将该 Dlelem 移到链表头部并将 CatCache 的 cc_hits（命中计数器）加 1。

3) 如果在 Hash 桶链中无法找到满足条件的元组，则需要进一步对物理系统表进行扫描，以确认要查找的元组是确实不存在还是没有缓存在 CatCache 中。如果扫描物理系统表能够找到满足条件的元组，则需要将该元组包装成 Dlelem 之后加入到其对应的 Hash 桶内链表头部。如果在物理系统表中找不到要查找的元组，则说明该元组确实不存在，此时构建一个只有键值但没有实际元组的“负元组”，并将它包装好加入到 Hash 桶内链表头部。

从 SearchCatCache 的查找过程可以看到，由于 CatCache 只是一个缓存，因此即使在其中找不到某个元组也不能确定该元组是否存在于系统表中，还需要进一步扫描物理系统表来查找该元组。但是，如果在 CatCache 中为这个不存在的元组放置一个“负元组”则可避免这些额外的开销，因为每次查找同一个不存在的元组时将会得到这个“负元组”，此时即可判定要查找的元组并不存在于系统表中，因而不用进一步去扫描物理系统表确认，从而造成浪费。

SearchCatCache 函数的流程如图 3-16 所示。SearchCatCache 的调用者不能修改返回的元组，并且使用完之后要调用 ReleaseCatCache 将其释放。

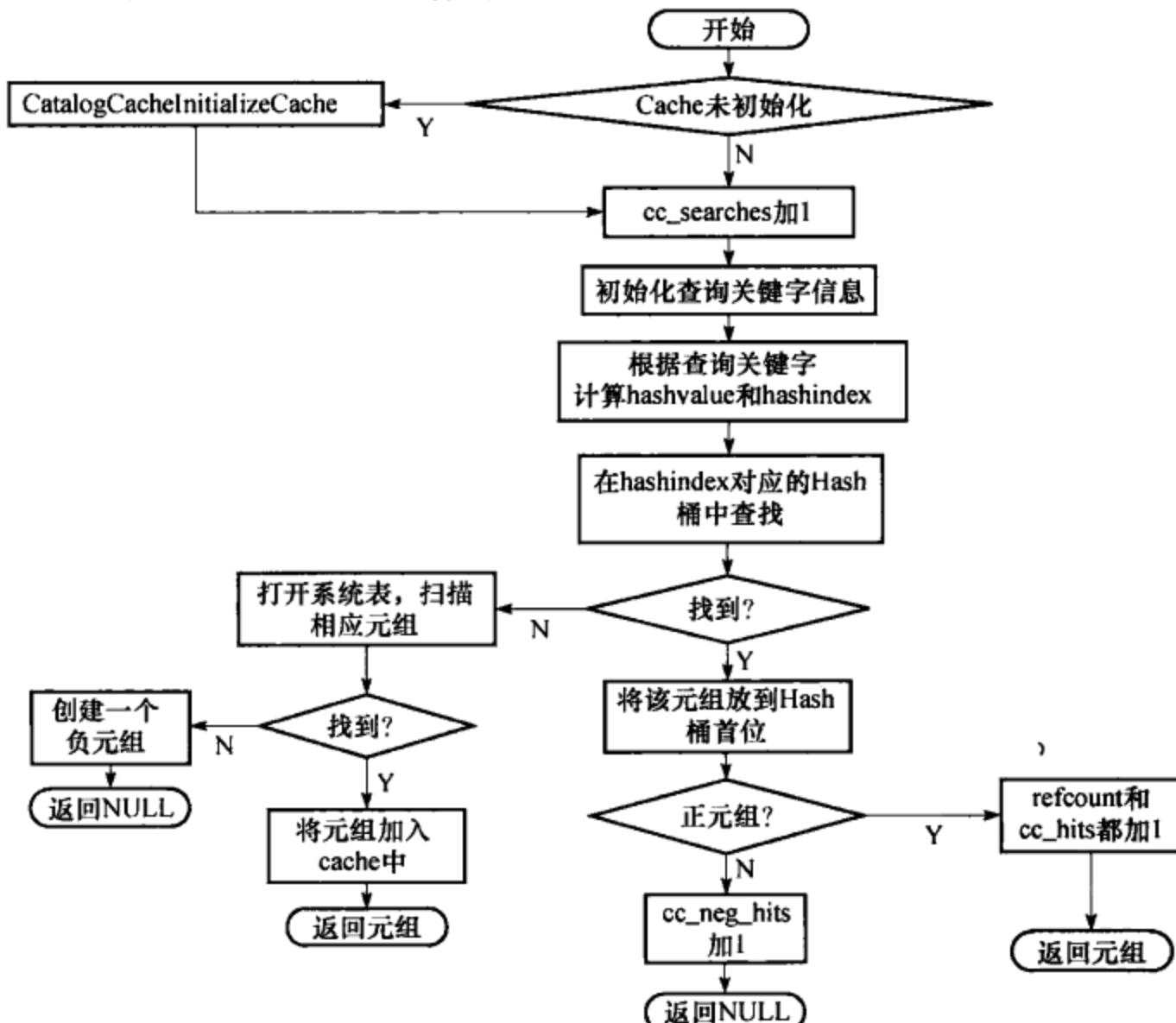


图 3-16 函数 SearchCatcache 流程

在 CatCache 中，部分匹配使用另外一个函数 SearchCatcacheList，该函数产生一个 CatCList 结构（数据结构 3.20），其中以链表的方式存放了在 Cache 中找到的元组。CatCList 中的 tuple 字段记录的是一个“负元组”，它仅仅用来存放该 CatCList 所用到的键值，没有其他用途。CatCList 中所包含的元组实际通过 members 字段表示的变长数据来记录，该数组的实际长度由 n_members 字段记录。

数据结构 3.20 catlist

```

typedef struct catlist
{
    int             cl_magic;      // 用于标识 CatCList 数据结构
#define CL_MAGIC      0x52765103 // cl_magic 字段被设置为该值
    CatCache        *my_Cache;     // 指向所在的 CatCache
    Dlelem          cache_elem;   // 指向 CatCache 中 cc_lists 链表中对应于该 CatCList 的
                                  // Dlelem 结构
    int             refcount;     // CatCList 的引用次数
    bool            dead;         // 当其中任何一个元组死亡时为真
    bool            ordered;      // 所包含的元组是否按索引排序
    short           nkeys;        // 查找所用到的关键字数目
    uint32          hash_value;   // 查找关键字的 Hash 值
    HeapTupleData tuple;        // 只包含键值的负元组
    int             n_members;    // 该 CatCList 中元组的个数
    CatCTup         *members[1];  // 该 CatCList 中元组的数组
}CatCList;

```

SearchCatcacheList 函数也会先计算查找键的 Hash 值，不过该函数首先会在 CatCache 的 cc_lists 字段中记录的 CatCList 链表中查找以前是否缓存了该查找键的结果，该查找过程将使用 CatCList 中 tuple 字段指向的元组与查找键进行 Hash 值比较。如果能够找到匹配该 Hash 值的 CatCList，则 cc_lhits 加 1 并将该 CatCList 移到 cc_lists 所指向链表的头部，然后返回找到的 CatCList。如果在 CatCache 中找不到 CatCList，则扫描物理系统表并构建相应的 CatCList 并将它加入到 cc_lists 所指向链表的头部。

SearchCatcacheList 的流程如图 3-17 所示。同样，SearchCatcacheList 的调用者不能修改返回的 CatCList 对象或者里面的元组，并且使用完之后要调用 ReleaseCatCacheList 将其释放。

2. RelCache

对 RelCache 的管理比 SysCache 要简单许多，原因在于大多数时候 RelCache 中存储的 RelationData 的结构是不变的，因此 PostgreSQL 仅用一个 Hash 表来维持这样一个结构。对 RelCache 的查找、插入、删除、修改等操作也非常简单。当需要打开一个表时，首先在 RelCache 中寻找该表的 RelationData 结构，如果没有找到，则创建该结构并加入到 RelCache 中。

和 SysCache 的初始化类似，RelCache 的初始化同样也在 InitPostgres 函数中进行，同样分为两个阶段：RelationCacheInitialize 和 RelationCacheInitializePhase2。

InitPostgres 会调用函数 RelationCacheInitialize 对 RelCache 进行第一阶段初始化，该函数将为该

进程创建一个 Hash 表，其 Hash 键为表的 OID，并设置 Hash 函数为 oid_hash。Hash 表的创建在函数 hash_create 中实现，该函数将创建一个标准 Hash 表结构体 HTAB。

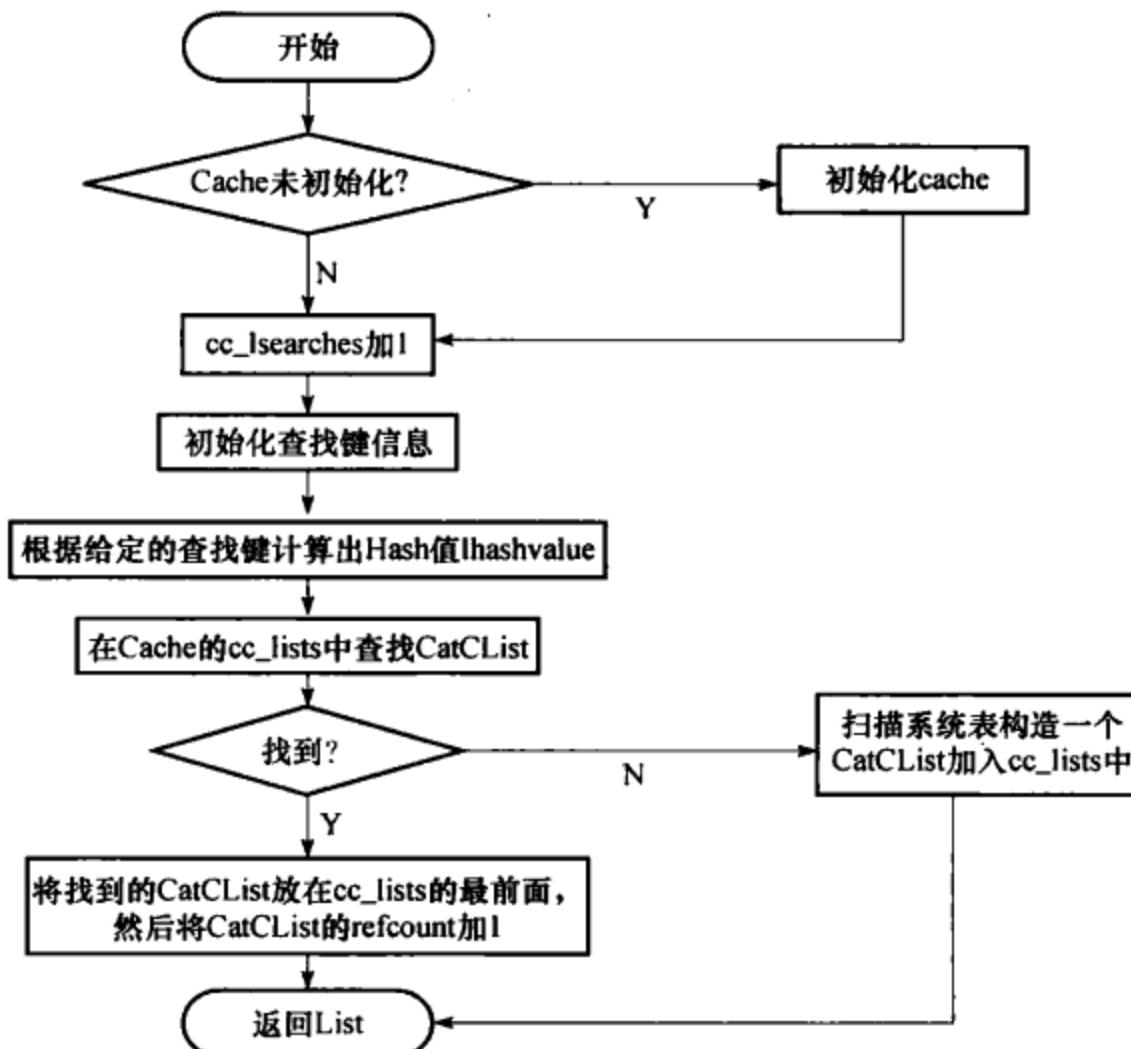


图 3-17 函数 SearchCatcacheList 流程

在完成了 Hash 表的创建后，InitPostgres 将调用 RelationCacheInitializePhase2 进入第二阶段的初始化。该函数将必要的系统表和系统表索引的模式信息加入到 RelCache 中，这个过程通过函数 RelationCacheInitializePhase2 来实现。这个阶段会确保 pg_class、pg_attribute、pg_proc、pg_type 四个系统表及相关索引的模式信息被加入到 RelCache。在 PostgreSQL 中，使用一个文件 pg_internal.init 来记录系统表 RelationData 结构体，若该文件存在且未损坏，则将其内容直接读入 RelCache 中。否则，分别建立 pg_class、pg_attribute、pg_proc、pg_type 及其索引的 RelationData 结构，加入到 RelCache 上的 Hash 表中，并重写 pg_internal.init 文件。

当 RelCache 初始化完成后，我们就可以使用它来查找表的模式信息。RelCache 的主要操作包括：

(1) 插入新打开的表

当打开新的表时，要把它 RelationData 加入到 RelCache 中。该操作通过宏 RelationCacheInsert 来实现：首先，根据关系表 OID 在 Hash 表中找到对应的位置，调用函数 hash_search，指定查询模式为 HASH_ENTER，该模式下若发现 OID 对应的 Hash 桶已存在，则返回其指针；否则创建一个新的空 Hash 桶，返回其指针。然后将返回的指针强制转换为 RelIdCacheEnt（数据结构 3.21），然后把打开表的 RelationData 赋值给 reldesc 字段。

数据结构 3.21 RelIdCacheEnt

```
typedef struct relidCacheent
{
    Oid      reloid;          //表的 OID
    Relation reldesc;         //表的 RelationData 结构指针
}RelIdCacheEnt;
```

(2) 查找 Hash 表

查找 Hash 表通过定义宏 RelationIdCacheLookup (ID, RELATION) 来实现，调用函数 hash_search，指定查询模式为 HASH_FIND，若找到 ID 对应的 RelIdCacheEnt，则将其 reldesc 字段的值赋值给 RELATION；否则，设置 RELATION 为 NULL。

(3) 从 Hash 表中删除元素

从 Hash 表中删除元素通过定义宏 RelationCacheDelete (RELATION) 来实现，调用函数 hash_search，指定查询模式为 HASH_REVOKE，在该模式下，若找到对应的 Hash 桶，则将其删除；否则返回 NULL。

Hash 表实际上扮演了 RelCache 索引的角色，所有对于 RelCache 的查找都是通过 Hash 表辅助进行的。在 PostgreSQL 8.4.1 以前的版本中，还允许使用表名来为系统表的 RelationData 结构建立 Hash 表，但目前这个特性已被删除，而是使用 OID 作为关键字在 Hash 表中查询。

3. Cache 同步

在 PostgreSQL 中，每一个进程都有属于自己的 Cache。换句话说，同一个系统表在不同的进程中都有对应的 Cache 来缓存它的元组（对于 RelCache 来说缓存的是一个 RelationData 结构）。同一个系统表的元组可能同时被多个进程的 Cache 所缓存，当其中某个 Cache 中的一个元组被删除或更新时，需要通知其他进程对其 Cache 进行同步。在 PostgreSQL 的实现中，会记录下已被删除的无效元组，并通过 SI Message 方式（即共享消息队列方式）在进程之间传递这一消息。收到无效消息的进程将同步地把无效元组（或 RelationData 结构）从自己的 Cache 中删除。

当前系统支持三种无效消息传递方式：第一种是使 SysCache 中元组无效，第二种是使 RelCache 中 RelationData 结构无效，第三种是使 SMGR 无效（表物理位置发生变化时，需要通知 SMGR 关闭表文件）。PostgreSQL 使用数据结构 3.22 所示的结构体来存储无效消息。

其中，id 为 0 或正数表示 CatCache，-1 表示 RelCache，-2 表示 SMGR。结构中的字段 cc、rc、sm 分别在 id 取值为 0 或正数、-1、-2 时有有效值，其他情况下为空。当 id 为 0 或正数时，它同时也表示产生该无效消息的 CatCache 的编号。

进程通过调用函数 CacheInvalidateHeapTuple 对无效消息进行注册，主要包括以下几步：

- 1) 注册 SysCache 无效消息。

数据结构 3.22 SharedInvalidationMessage

```
typedef union
{
    int16 id;
    SharedInvalCatcacheMsg cc;
    SharedInvalRelCacheMsg rc;
    SharedInvalSmgrMsg sm;
}SharedInvalidationMessage;
```

2) 如果是对 pg_class 系统表元组进行的更新/删除操作，其 relfilenode 或 reltablespace 可能发生变化，即该表物理位置发生变化，需要通知其他进程关闭相应的 SMGR。这时首先设置 relationid 和 databaseid，然后注册 SMGR 无效消息；否则转而执行步骤 3。

3) 如果是对 pg_attribute 或者 pg_index 系统表元组进行的更新/删除操作，则设置 relationid 和 databaseid，否则返回。

4) 注册 RelCache 无效消息。

当一个元组被删除或者更新时，在同一个 SQL 命令的后续执行步骤中我们依然认为该元组是有效的，直到下一个命令开始或者事务提交时改动才生效。在命令的边界，旧元组变为失效，同时新元组置为有效。因此当执行 heap_delete 或者 heap_update 这样的操作时，不能简单地刷新 Cache。而且，即使刷新了，也可能由于同一个命令中的请求把该元组再次加载到 Cache 中。因此正确的方法是保持一个无效链表用于记录元组的 delete/update 操作。

事务完成后，根据前述的无效链表中的信息广播该事务过程中产生的无效消息，其他进程通过 SI Message 队列读取无效消息对各自的 Cache 进行刷新。当子事务提交时，只需要将该事务产生的无效消息提交到父事务，最后由最上层的事务广播无效消息。

需要注意的是，若涉及对系统表结构的改变，还需要重新加载 pg_internal.init 文件，该文件记录了所有系统表的结构。

3.3.3 缓冲池管理

如果需要访问的系统表元组在 Cache 中无法找到或者需要访问普通表的元组，就需要对缓冲池进行访问。任何对于表、元组、索引表等的操作都在缓冲池中进行，缓冲池的数据调度都以磁盘块为单位，需要访问的数据以磁盘块为单位调用函数 smgrread 写入缓冲池，而 smgrwrite 将缓冲池数据写回到磁盘。调入缓冲池中的磁盘块称为缓冲区、缓冲块或者页面，多个缓冲区组成了缓冲池。

缓冲池管理模块的主要结构如图 3-18 所示。

PostgreSQL 有两种缓冲池：共享缓冲池和本地缓冲池。共享缓冲池主要用作普通可共享表的操作场所；本地缓冲池则用作仅本地可见的临时表的操作场所。

对于缓冲池中缓冲区的管理通过两种机制完成：pin 和 lock。pin 实际就是缓冲区的访问计数器。当进程要访问缓冲区前，对缓冲区加 pin，pin 的数目保存在缓冲区的 refcount 属性中。当 refcount 不为 0 时表明有进程正在访问缓冲区，此时该缓冲区不能被替换。而 lock 机制为缓冲区的并发访问提供了保障，当有进程对缓冲区进行写操作时加 EXCLUSIVE 锁，读操作时加 SHARE 锁，其意义和数据库的锁机制是类似的。

为了便于实现，PostgreSQL 对共享缓冲区的管理基本上采取了静态方式：它在系统配置时规定好了共享缓冲区的总数（1000 个，由全局变量 NBuffers 定义），以后在每次系统启动时，由 Postmaster 或某一独立的 Postgres 从共享内存中分配一片空间用作共享缓冲区，全部（共享）缓冲区构

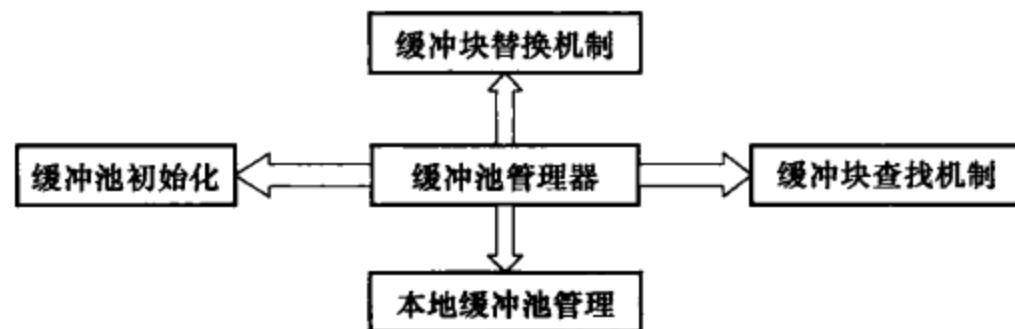


图 3-18 缓冲池模块图

成共享缓冲池。

1. 初始化共享缓冲池

共享缓冲池处于共享内存区域，在系统启动的时候需要对其进行初始化操作，负责这一工作的函数为 InitBufferPool。在共享缓冲池管理中，使用了一个全局数组 BufferDescriptors 来管理池中的缓冲区，其数组元素类型为 BufferDesc（参见数据结构 3.23）。数组元素个数为缓冲区的总数，用于管理所有缓冲区的描述符，可以认为 BufferDescriptors 就是共享缓冲池。另外还使用了一个全局指针变量 BufferBlocks 来存储缓冲池的起始地址。每个缓冲区对应一个描述符结构 BufferDesc（见数据结构 3.23），用于描述相应缓冲区的状态信息。

数据结构 3.23 BufferDesc

```
typedef struct sbufdesc
{
    BufferTag      tag;           // 缓冲区页面 ID, 指明了该缓冲区对应表块的物理信息
    BufFlags       flags;         // 标志位, 表示缓冲区是否为脏等
    uint16          usage_count;   // 记录最近缓冲区使用次数, 用于缓冲区替换
    unsigned        refcount;      // 引用计数
    int             wait_backend_pid; // 用于记录等待所有 Pin 结束的进程 PID
    slock_t         buf_hdr_lock;  // 用于修改上述数据域时加锁保护
    int             buf_id;        // 缓冲区的索引号
    int             freeNext;      // 用于链接 FreeList
    LWLockId        io_in_progress_lock; // 缓冲区和磁盘进行 I/O 时加锁
    LWLockId        content_lock;  // 进程对缓冲区内容访问时加锁
} BufferDesc;
```

BufferDesc 主要用于对缓冲区进行描述，还保存了该缓冲区内缓冲块的信息。下面将就几个主要的字段进行介绍：

- **tag**: 用于标识该缓冲块的物理信息（见数据结构 3.24）。
- **refcount**: 表示当前正在引用该缓冲块的进程数，通过 pin 操作来修改该字段。默认情况下，当一个进程第一次做 pin 操作的时候，**usage_count** 的值加 1。如果一个缓冲区的 **refcount** 大于 0 就说明有进程正在访问它，或者正在进行 IO。
- **buf_id**: 缓冲区的索引号，由于共享缓冲区和本地缓冲区使用同一种描述符，所以它们的索引号的编号规则不同：共享缓冲区的索引号为大于 0 的整数，从 0 开始编号，后续依次加 1；而本地缓冲区的索引号是小于等于 -2 的负整数，从 -2 开始编号，后续的依次减 1。宏定义 BufferDescriptorGetBuffer 负责将 **buf_id** 转换为缓冲区在描述符数组中的下标号。
- **freeNext**: 如果当前缓冲区在空闲链表中，该值表示在当前缓冲区之后的空闲缓冲区描述符在数组中的索引号，通过 **freeNext** 字段可以把所有空闲缓冲区链接起来。
- **wait_backend_pid**: 用于记录一个请求修改缓冲区的进程号。例如，若进程 X 需要删除的元组所在缓冲块有其他进程访问，即 **refcount** 不为 0 时（有其他进程访问该缓冲区），则进程 X 不能物理上删除元组。故系统将该进程的 id 记录在 **wait_backend_pid** 域中，然后对缓冲

块加 pin，并阻塞自己等待其他进程都不再访问该缓冲块（即 refcount 为 0）。同时，当 refcount 为 1 时，最后一个使用该缓冲块的进程释放缓冲区时，会向该缓冲块的 wait_backend_id 进程发送消息。

- io_in_progress_lock：该锁用于缓冲区和磁盘进行 I/O 时，请求这样的缓冲块时需要等到 I/O 结束。
- content_lock：当进程访问缓冲块时，会在 content_lock 上加锁，读访问加 LW_SHARE 锁，写访问加 LW_EXCLUSIVE 锁，此锁可以防止因多个进程对缓冲区访问的冲突而造成数据不一致。

对共享缓冲池初始化的主要流程如图 3-19 所示。

对共享缓冲区的初始化实际上是对两个全局变量进行初始化：第一个是缓冲区描述符数组 BufferDescriptors，该数组存储了所有共享缓冲区的描述符；第二个是指针 BufferBlocks，全局变量 BufferBlocks 指向共享缓冲池头部。

InitBufferPool 函数在共享内存初始化时调用一次（在 Postmaster 或一个单独的 Postgres 进程中），该函数将初始化缓冲区描述符，建立便于查找缓冲区的 Hash 表（见 3.3.3 节）和 FreeList 结构。

在该函数中，最终会通过调用函数 ShmemInitHash 来初始化缓冲区 Hash 表，该函数是 PostgreSQL 在共享内存中创建 Hash 表的通用函数，它会初始化 Hash 表的大小和 Hash 函数，并对 HTAB 进行初始化。

FreeList 的构建与前文 VFD 的 FreeList 构建类似，通过 BufferDescriptors 数组元素的 freeNext 字段来进行连接以形成一个链表。

2. 共享缓冲区查询

为了在共享缓冲池中快速查找缓冲区，在初始化缓冲池阶段系统为其在共享内存中创建了 Hash 表。缓冲区描述符中的 BufferTag 结构体包含一个缓冲块的物理信息，因此具有唯一性。共享缓冲区查询以 BufferTag 作为索引键，查找 Hash 表并返回目标缓冲区在缓冲区描述符数组中的位置。

数据结构 3.24 BufferTag

```
typedef struct buftag
{
    RelFileNode    rnode;        //由表所在的表空间 OID、数据库 OID 及表本身的 OID 构成
    ForkNumber     forkNum;      //枚举类型，标记缓冲区中是什么类型的文件块
    BlockNumber    blockNum;     //块号
}BufferTag;
```

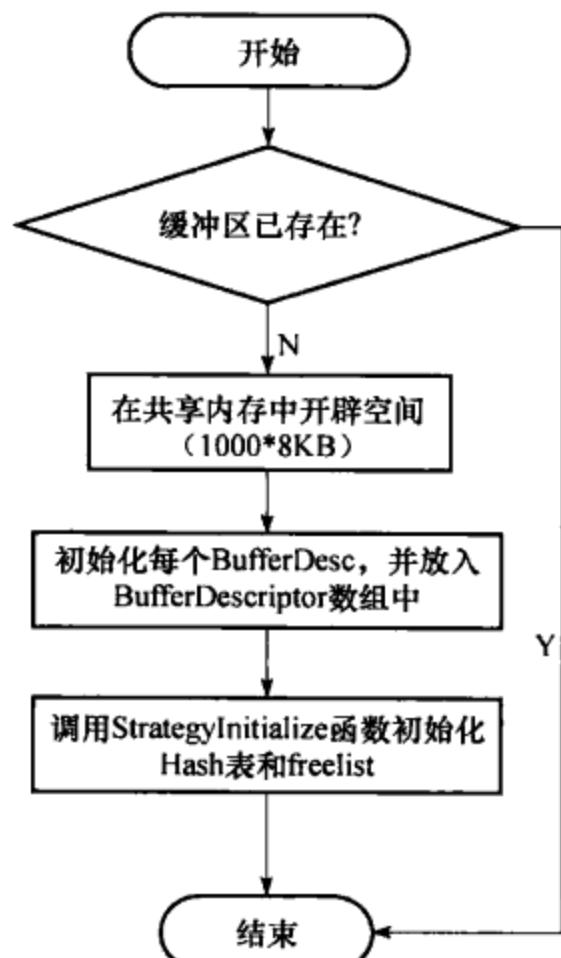


图 3-19 缓冲区初始化

给定一个文件的 RelFileNode 信息、文件类型（用于区别表文件、FSM 文件和 VM 文件）、文件块号，可以唯一标识一个文件块，把这些信息包装成一个 BufferTag 后可以用来查询对应于该文件

块的缓冲区，查询流程如下：

- 1) 首先从 Hash 表中，根据 BufferTag 进行查询，若 Hash 表中存在记录说明请求的文件块已经在缓冲池中，直接返回缓冲区。
- 2) 当 Hash 表中不存在时，需要从缓冲池中找到一个空闲的缓冲区来装入文件块，若存在空闲的缓冲区，则返回该缓冲区；若不存在，则使用替换机制进行替换缓冲区，最后获得一个空闲缓冲区并返回。
- 3) 从磁盘中将文件块读入返回的缓冲区，并将该缓冲区记录到 Hash 表中。

函数 ReadBuffer_common 是所有缓冲区读取的通用函数，该函数定义了本地缓冲区和共享缓冲区的通用读取方法，其基本流程如图 3-20 所示。

ReadBuffer_common 的参数中包括前面提到的组成 BufferTag 的三个信息：表文件的 RelTreeNode、文件类型和块号，其中 RelTreeNode 并没有作为一个单独的参数传入，而是作为表的 SMgrRelation 描述符的一部分传入。该函数使用了一个布尔类型的参数以区分是在本地缓冲池还是在共享缓冲池中查找缓冲区。该函数返回包含所要求文件块的缓冲区。如果所要求的 block-Num 是新的，则扩展该表文件并且分配一个新的块。

ReadBuffer 函数是一个一般性的缓冲区读取函数，该函数是读缓冲区的一个上层接口，它的参数只有两个：要访问的表的 Relation (RelationData 的指针) 结构和要读取的块号。ReadBuffer 会打开相应的 SMgrRelation 并调用到函数 ReadBuffer_common，ReadBuffer 会将文件类型 (BufferTag 中的 forkNum 的值) 默认为 MAIN_FORKNUM (即表文件)。此外，在 PostgreSQL 中还定义了一个扩展函数 ReleaseAndReadBuffer。该函数结合了 ReleaseBuffer 和 ReadBuffer 两个函数的功能，ReleaseBuffer 用于释放加在一个缓冲区上的 pin 锁。两者结合后，将首先判断参数中指定的缓冲区与另外两个参数指定的文件块是否相匹配，若相匹配则直接返回这个缓冲区；否则先解除该缓冲区的 pin 锁，并调用 ReadBuffer 读取指定的文件块。

在 ReadBuffer_common 中，将会调用函数 BufferAlloc 函数来获得指定的共享缓冲区。该函数是执行缓冲区替换策略的核心函数，其流程如图 3-21 所示。

BufferAlloc 函数首先查找 Hash 表，如果缓冲区已经在 Hash 表中，则直接从表中读取。否则执

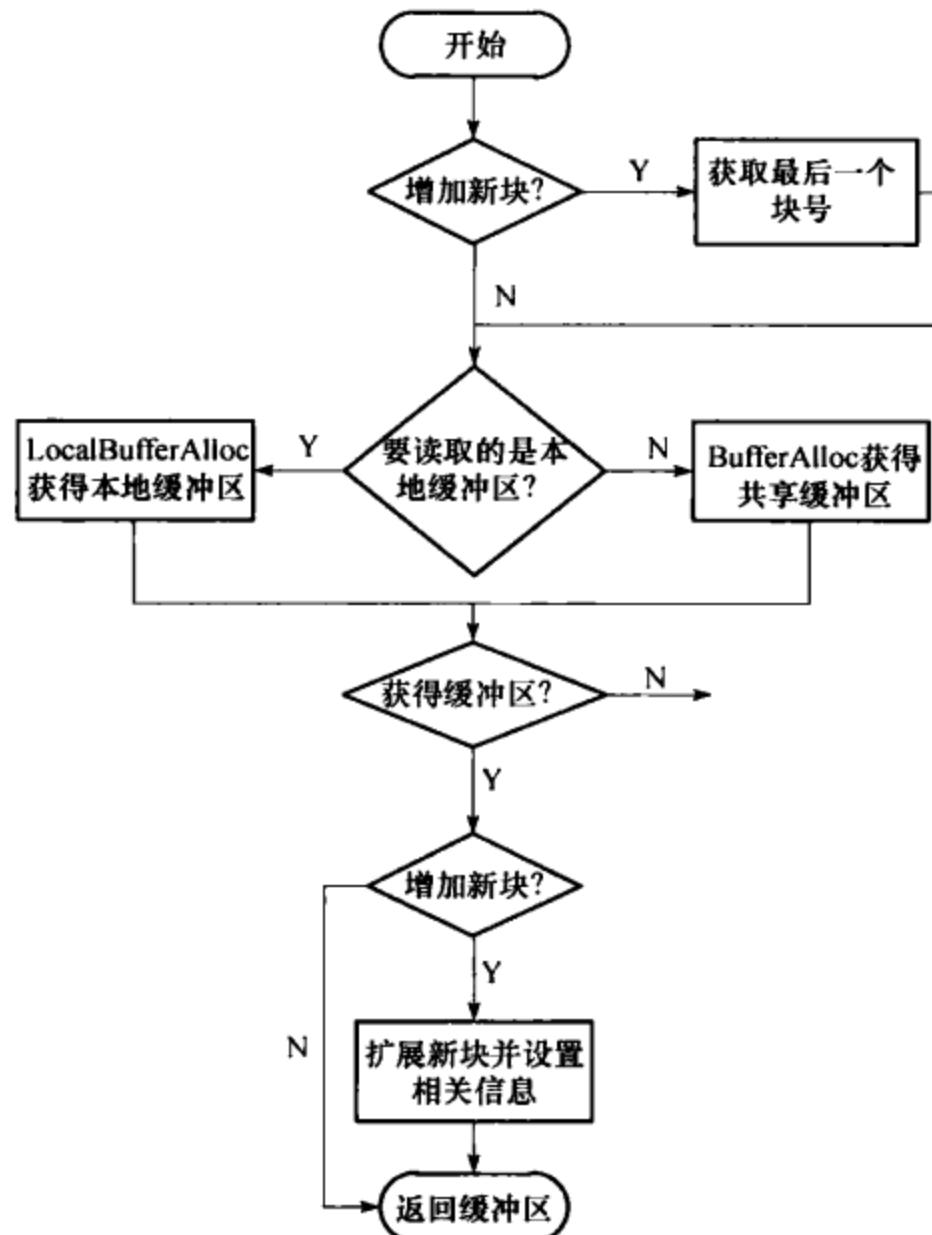


图 3-20 ReadBuffer_common 函数的流程图

行以下操作：

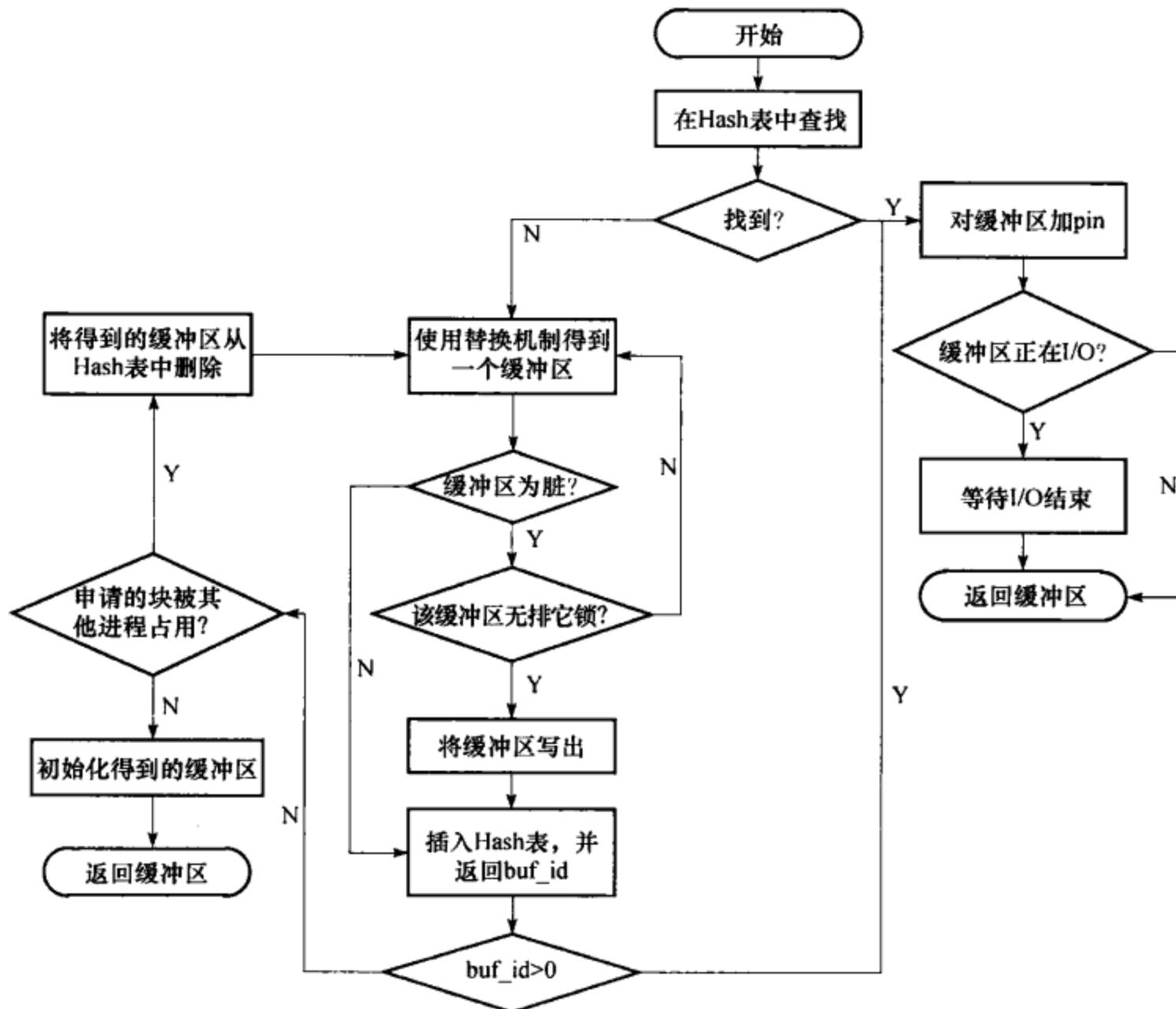


图 3-21 BufferAlloc 函数的流程

1) 通过缓冲块替换策略获取一个缓冲区，若为脏且不允许重复使用，则重新获取缓冲区，否则进行下一步。

2) 若获得的缓冲区标记为脏，则先将该缓冲区的内容刷回磁盘，此后该缓冲区就已经可以用于新请求的文件块。但是在真正装入之前需要将缓冲区插入到 Hash 表中（因为缓冲区的 BufferTag 发生了变化），这个插入过程由函数 BufTableInsert 完成，BufTableInsert 的返回值可以分成两种情况：

① 返回值为 -1，表明插入成功，直接执行步骤 3。

② 返回一个大于 0 的值，则表明有其他进程已经将装入文件块的缓冲区加入到 Hash 表中，返回值就代表了该缓冲区的索引号。这种情况说明我们请求的文件块已经被其他进程载入到缓冲池中（但是放在另外一个缓冲区里），也就是说我们没有必要再重新载入该文件块，直接使用其他进程载入的这个缓冲区即可。此时，我们需要放弃计划使用的这个缓冲区（刚刚申请到的缓冲区），并返回 BufTableInsert 指示给我们的缓冲区。

3) 即使我们可以按照原计划在新申请到的缓冲区中载入所需的文件块，仍然要考虑一种可能

的冲突：在一个进程对申请的缓冲区进行初始化时，另一个进程也申请到该缓冲区，这时需要从 Hash 表中删除前者申请到的缓冲区，并利用缓冲区替换策略重新获取缓冲区（即回到步骤 1 重新进行申请）。

这里的介绍中只分析了对通过缓冲区替换策略得到的缓冲区的后续处理，并没有对缓冲区的替换策略进行说明，下面将专门用一个小节介绍共享缓冲区的替换策略。

3. 共享缓冲区替换策略

在共享缓冲池中，初始化定义的缓冲区个数是有限的（由宏 NBuffers 定义，默认为 1000 个），并且这个值在初始化分配后将不会再被改变。因此在不断的操作过程中，可能出现缓冲区被用光的局面，这时候就需要替换一些最近未使用的缓冲区，以加载新请求的文件块。

对于缓冲区的替换策略，在 PostgreSQL 8.4.1 中有两种方案：一般的缓冲区替换策略以及缓冲环替换策略。

（1）一般的缓冲区替换策略

首先在缓冲池中维持一个 FreeList 链表，FreeList 是一个单向链表。FreeList 中的缓冲区通过其描述符中的 freeNext 字段链接起来，在 BufferStrategyControl 结构中记录了 FreeList 第一个和最后一个元素。当某缓冲区 refcount 变为 0 时（变成空闲），将其加入到 FreeList 链尾；当需要一个空闲缓冲区用于替换时，则从链首取得。若 FreeList 中某缓冲区被使用，则将其从 FreeList 中删除。

BufferStrategyControl（数据结构 3.25）主要用于对 FreeList 的控制。

数据结构 3.25 BufferStrategyControl

```
typedef struct
{
    int nextVictimBuffer;      //在 clock sweep 算法中, 用于指向下一个 Buffer
    int firstFreeBuffer;       //第一个空闲的缓冲块 ID
    int lastFreeBuffer;        //最后一个空闲的缓冲块 ID
}BufferStrategyControl;
```

当在 FreeList 中无法找到合适的缓冲区时，通过 nextVictimBuffer 对所有缓冲区进行扫描直到找到空闲缓冲区，这个扫描过程使用了一个简单的 clock-sweep 算法。clock-sweep 算法的主要流程如下：

- 1) 初始化 tryCounter = NBuffers。
- 2) 根据 nextVictimBuffer 字段找到相应的缓冲区，初始值为 0。
- 3) 将 nextVictimBuffer 加 1，如果当前 nextVictimBuffer 指向池中最后一个缓冲区，设置 nextVictimBuffer 为 0。
- 4) 如果步骤 2 中得到的缓冲区的 refcount 为 0：
 - ①若 usage_count 不为 0，则置 usage_count 减 1，并重置 trycounter 为 NBuffers。
 - ②否则获取这个缓冲区并返回。
- 5) 如果步骤 2 中得到的缓冲区的 refcount 不为 0，则将 tryCounter 减 1，如果此时 tryCounter 等于 0，抛出错误。
- 6) 返回步骤 2。

Clock-sweep 算法实际上运行了一个死循环遍历缓冲池来检测池中是否存在空闲缓冲区，使用一个计数器 tryCounter 来防止循环无止境地运行下去。当检测到一个引用计数为 0 但最近被进程使用过（即 usage_count 不为 0）的缓冲区时，并不立即返回该缓冲区，而是将其 usage_count 减 1，接着重置计数器并继续循环遍历。如果在若干次循环之后某一个缓冲区的引用计数和使用计数都变为 0，说明该缓冲区在最近的一段时间内极少有进程使用它（如果在此期间有较多次使用，usage_count 不会变为 0 反而可能增大），那么这个缓冲区就可以被替换。当然也有可能达到计数限制，这时会抛出错误，但此类情况发生的概率极低。通过这种替换策略能够有效地获取到一个最近未被使用的缓冲区。

（2）缓冲区替换策略

缓冲环策略是缓冲区策略的一种特殊情况，即对于某些请求而言，加载到内存中的磁盘块在使用过后将不会再次被使用，但其又将使用到大量的缓冲区。因此会造成内存中充斥着大量仅需要使用一次的缓冲区，这与设置缓冲区以减少 I/O 次数的初衷是违背的。对于这样的情况，PostgreSQL 采用了缓冲环策略，即对于这类请求，为其分配固定数量的缓冲区，一旦用完环中分配的缓冲区后，就从环的头部开始重复使用缓冲区。

在缓冲环策略下，PostgreSQL 会将用过的缓冲区组织成一个链表，首尾相连形成一个环状结构，称之为缓冲环。例如，对于一次连续的只读扫描，在扫描过后，驻留在内存中的缓冲区将不会再被使用到，因此使用 256KB 的缓冲池空间就足够。而对于一次 VACUUM 操作，同样的也使用 256KB 的缓冲池空间，但是清理过程中产生的脏页并不是从环中移除，而是刷新 WAL（一次性刷新整个缓冲区）以允许缓冲区被再次使用。在 PostgreSQL 8.3 之前，一次 VACUUM 操作后的缓冲区是被送到 FreeList 中的，这会导致过度地刷新 WAL，而一次刷新 256KB 的缓冲区将更有效率。

数据结构 3.26 BufferAccessStrategyData

```

typedef struct BufferAccessStrategyData
{
    BufferAccessStrategyType btype;          // 缓冲环控制策略类型
    int                    ring_size;        // 环的大小
    int                    current;          // 最近加入到环中的 Buffer
    bool                  current_was_in_ring; // 最近通过 StrategyGetBuffer 获得的 Buffer 是否是直接在环中取的
    Buffer                Buffers[1];        // 数组，用于存储加入到环中的缓冲区的索引号
} BufferAccessStrategyData;

```

缓冲环替换策略主要依靠数据结构 BufferAccessStrategy（BufferAccessStrategy 是数据结构 3.26 的指针）来控制。BufferAccessStrategy 用于缓冲区策略的控制，目前仅有缓冲环策略一种实现，在未来可能有更多的实现。

BufferAccessStrategyData 的 btype 字段用于选择不同的缓冲环策略，每一种取值适用于一类缓冲环策略，在不同的缓冲环策略下缓冲环的大小有所区别。在 PostgreSQL 8.4.1 中，定义了以下几种取值：

- 1) BAS_NORMAL：正常的随机存取。

- 2) BAS_BULKREAD: 只读扫描。
- 3) BAS_BULKWRITE: 多个块的写操作, 例如 CREATE TABLE AS SELECT 命令。
- 4) BAS_VACUUM: VACUUM 操作。

对于缓冲环策略, 通过 GetBufferFromRing 函数来获取缓冲区, 该函数只有一个类型为 BufferAccessStrategy 的参数 strategy。GetBufferFromRing 替换缓冲区的流程如下:

1) 首先将 strategy 中的 current 指针指向 strategy 的 Buffers 字段的下一个元素 (代表可能的下一个缓冲区), 如果当前指向的是 Buffers 的最后一个元素, 则将 current 置为 0 (指向 Buffers 的第一个元素)。

2) 检查 current 指针指向的元素, 如果其中记录的值为 InvalidBuffer (即值为 0), 表明环还未充满, 这个位置还没有记录一个缓冲区。这种情况下设置 strategy 的 current_was_in_ring 字段为 false 之后返回空值。GetBufferFromRing 的上层调用函数 (StrategyGetBuffer) 在检测到返回值为空之后会采用一般的替换策略取得一个空闲缓冲区, 并通过 AddBufferToRing 将该缓冲区加入到缓冲环中。

3) 如果 current 指针指向的元素中记录的是一个有效的缓冲区索引号, 则检查该缓冲区的 ref_count 和 usage_count。如果两者都为 0 则可以直接将其返回; 如果 ref_count 为 0 且 usage_count 为 1, 说明该缓冲区现在没有被使用但最近被使用过一次, 而且上一次使用它的很可能是当前进程本身, 因此也可以把这个缓冲区替换出来返回; 如果以上情况都不满足, 则表明该缓冲区仍在被其他进程使用中或最近被其他进程使用过。这时需采用和步骤 2 类似的方法, 由上层调用函数采用一般的替换策略取得空闲缓冲区。

两种策略的使用是根据不同的操作来选择的, 不同的操作会封装不同的 BufferAccessStrategy 来调用 ReadBuffer_common (有一个参数是 BufferAccessStrategy 类型) 获取缓冲区。在需要进行缓冲区替换时, ReadBuffer_common 则会通过 BufferAlloc 调用 StrategyGetBuffer, StrategyGetBuffer 会根据传入的 BufferAccessStrategy 执行替换缓冲区的查找, 如果 BufferAccessStrategy 的值设为空, 将会采用一般的缓冲区替换策略; 否则采用缓冲环的替换策略。

4. 本地缓冲池管理

本地缓冲池是每个进程所特有的, 在进程初始化时自行创建。本地缓冲池对于其他进程是不可见的。它只在创建临时表或操作的数据对其他进程不可见的特殊情况下使用。

本地缓冲池的初始化调用函数 InitLocalBuffer 来完成。与共享缓冲池不同的是, 初始化过程中并不预先创建缓冲区, 而是在使用时创建。本地缓冲池的初始化主要包括以下两步:

- 1) 分配内存并初始化缓冲区描述符 (和共享缓冲区的描述符结构一致) 数组, 默认大小为 1000。注意, 这里不分配缓冲区, 仅分配存放描述符的数组。
- 2) 创建 Hash 表用于查询本地缓冲区, 使用 hash_create 函数创建本地 Hash 表, 指定 Hash 函数为 tag_hash, 该 Hash 表用全局变量 LocalBufHash 保存。

本地缓冲区的获取主要调用 LocalBufferAlloc 函数来完成。该函数根据指定的文件块信息创建或者找到缓冲区。对本地缓冲区的操作由于仅在本进程的内存区域进行, 因此不需要进行锁的管理。其流程如图 3-22 所示。

需要注意的是, 共享缓冲区在创建时就分配好了固定大小的内存, 而本地缓冲区则是在使用的时候才实际分配内存。共享缓冲区与本地缓冲区的区别如表 3-4 所示。

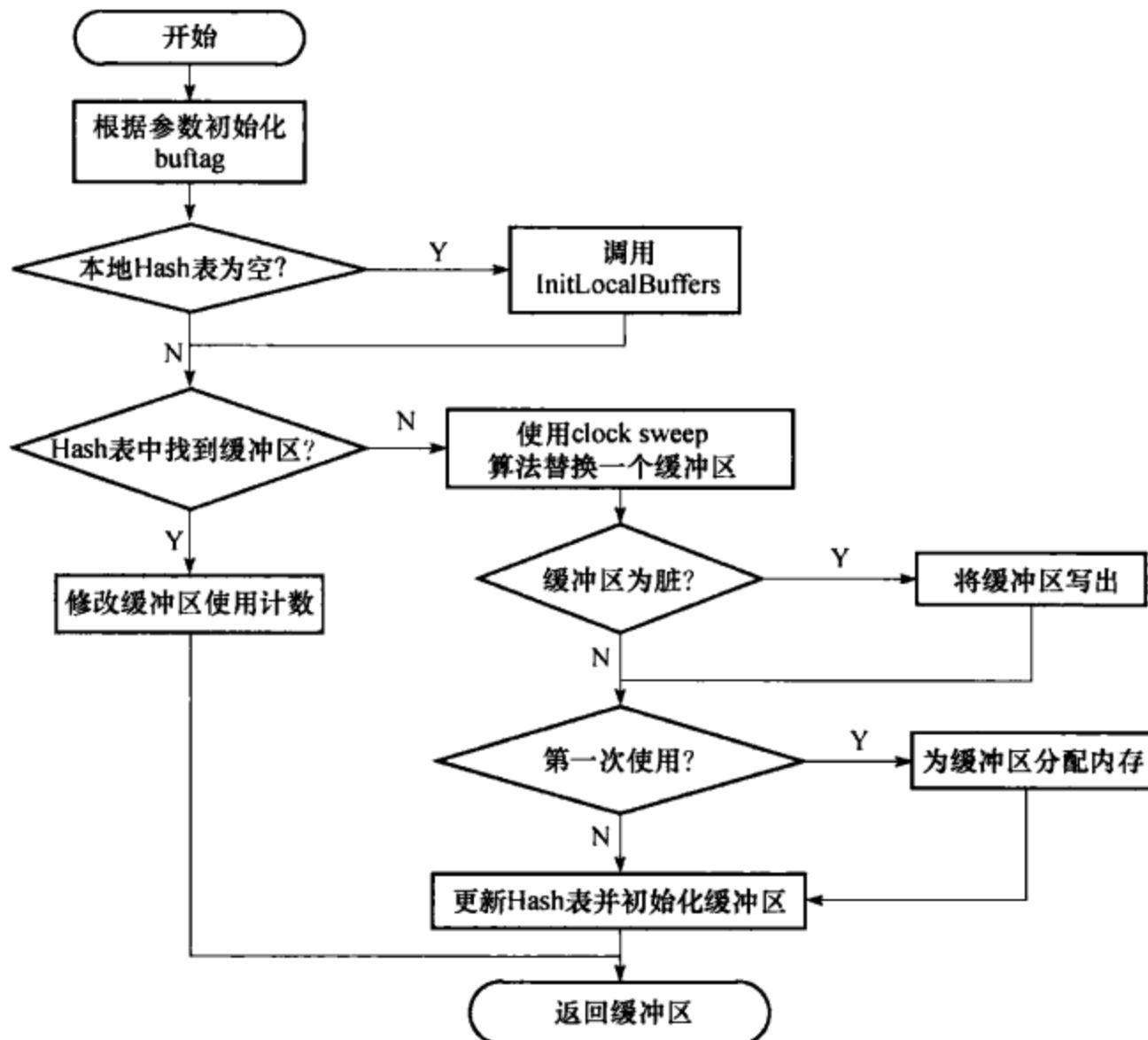


图 3-22 LocalBufferAlloc 函数流程图

表 3-4 共享缓冲区与本地缓冲区的区别

	共享缓冲区	本地缓冲区
所处内存空间	共享内存	本地内存
初始化	Postmaster 启动时初始化	进程启动时初始化
	初始化时分配实际空间	使用时分配实际空间
可见性	对其他进程可见	对其他进程不可见
缓冲区索引号	大于 0	小于 0

3.3.4 IPC

IPC（进程间通信，Inter-Process Communication）是指至少两个或者两个以上的进程交换数据或者信号的技术或者方法。由于进程是操作系统分配资源的最小单位，每个进程都有自己独立的一套系统资源，不同进程之间的资源是相互隔离的。为了使不同进程之间能够互相访问资源并且进行协同工作，操作系统提供了 IPC 机制。广义上的 IPC 不仅可以使同一计算机上的多个进程进行通信，而且可以使不同计算机上的多个进程进行通信，但 PostgreSQL 中的 IPC 只考虑同一计算机上的进程间通信。

IPC 有多种实现方法，包括文件、Socket、共享内存等。PostgreSQL 中的 IPC 主要采用共享内存的方式来实现，即在系统中开辟一片所有进程都可以读写的内存空间，并约定好进程读写这片内存的时机和方式，这样进程之间就可以通过这片共享的内存来交换数据。在共享内存的基础上，PostgreSQL 的 IPC 机制还提供了以下功能：

- 1) 进程和 Postmaster 的通信机制。
- 2) 统一管理进程的相关变量和函数。
- 3) 提供了 SI Message 机制，即无效消息传递机制。
- 4) 有关清除的函数。

1. 共享内存管理

在 PostgreSQL 初始化过程中，系统会分配一块内存区域，它对 PostgreSQL 系统中的所有后端进程而言是可见的，该内存区域称为共享内存。

在初始化的过程中，系统为共享内存创建了一个名为“shmemp index”的 Hash 索引。当试图为一个模块分配共享内存时，会调用函数 ShmemInitStruct。该函数首先根据模块名在 Hash 索引中查找，如果找不到则再调用 ShmemAlloc 函数在内存中为其分配一块区域。

在 ipci.c 文件中提供了对共享内存和信号量进行初始化的接口，而具体的实现则在 shmem.c 文件中。函数 CreateSharedMemoryAndSemaphores 负责对共享内存和信号量进行初始化，由 Postmaster 或其子进程调用。Postmaster 调用该函数时会初始化共享内存和信号量，其他进程（如 Postgres）调用时不进行初始化工作，仅仅获得已创建的共享变量指针和信号量指针。

初始化流程如下：

- 1) 计算共享内存总共需要的大小。
- 2) 分配共享内存空间。
- 3) 初始化共享内存头指针。
- 4) 创建信号量并注册清理函数
- 5) 构建共享内存 Hash 索引。
- 6) 初始化各个模块，调用 ShmemInitStruct 函数从已分配的共享内存分配空间。

在初始化共享内存过程中，CreateSharedMemoryAndSemaphores 在共享内存中创建一个 Hash 表。当其他模块需要在共享内存中开辟一个空间时，就会调用 Shmem.c 文件中的函数在共享内存中获得一个区域，并且在 Hash 表中获得一个索引键。

PostgreSQL 的共享内存机制是利用了操作系统的共享内存的编程技术，在此基础之上利用共享内存来支持进程间的通信。相应的共享内存操作可以参考 Linux 系统下的共享内存管理，本节的重点在于介绍基于共享内存的 SI Message 机制。

2. SI Message

SI Message 的作用在前面的“Cache 同步”部分中已经介绍过，它主要用于不同进程的 Cache 进行同步操作。

为了实现 SI Message 的这一功能，PostgreSQL 在共享内存中开辟了 shmInvalBuffer 记录系统中所发出的所有无效消息以及所有进程处理无消息的进度。shmInvalBuffer 是一个全局变量，其数据类型为 SISeg（数据结构 3.27）。

数据结构 3.27 SISeg

```

typedef struct SISeg
{
    int minMsgNum;
    int maxMsgNum;
    int nextThreshold;
    int lastBackend;
    int maxBackends;
    slock_t msgnumLock;
    SharedInvalidationMessage Buffer[MAXNUMMESSAGES];
    ProcState procState[1];
}SISeg;

```

在 `shmInvalBuffer` 中，无效消息存储在由 `Buffer` 字段指定的定长数组中（其长度 `MAXNUMMESSAGES` 预定义为 4096），该数组中每一个元素存储一个无效消息，也可以称该数组为无效消息队列。无效消息队列实际是一个环状结构，最初数组为空时，新来的无效消息从前向后依次存放在数组的元素中，当数组被放满之后，新的无效消息将回到 `Buffer` 数组的头部开始插入。`minMsgNum` 字段记录 `Buffer` 中还未被所有进程处理的无效消息编号中的最小值，`maxMsgNum` 字段记录下一个可以用于存放新无效消息的数组元素下标。实际上，`minMsgNum` 指出了 `Buffer` 中还没有被所有进程处理的无效消息的下界，而 `maxMsgNum` 则指出了上界，即编号比 `minMsgNum` 小的无效消息是已经被所有进程处理完的，而编号大于等于 `maxMsgNum` 的无效消息是还没有产生的，而两者之间的无效消息则是至少还有一个进程没有对其进行处理。因此在无效消息队列构成的环中，除了 `minMsgNum` 和 `maxMsgNum` 之间的位置之外，其他位置都可以用来存放新增加的无效消息。

`PostgreSQL` 在 `shmInvalBuffer` 中用一个 `ProcState` 数组（`procState` 字段）来存储正在读取无效消息的进程的读取进度，该数组的大小与系统允许的最大进程数 `MaxBackends` 有关，在默认情况下这个数组的大小为 100（系统的默认最大进程数为 100，可在 `postgresql.conf` 中修改）。`ProcState` 的结构如数据结构 3.28 所示。

数据结构 3.28 ProcState

```

typedef struct ProcState
{
    pid_t procPid;           //进程 PID,如果 ProcState 还没有用于记录一个进程的状态则
                            //其 procPid 为 0
    int nextMsgNum;          //进程下一个要读的无效消息编号(在 Buffer 数组中的下标),
                            //procPid 为 0 或者 resetState 为真时无效
    bool resetState;         //该进程是否需要重置状态
    bool signaled;           //该进程是否接受到 catchup 信号
    LocalTransactionId nextLXID;
}ProcState;

```

在 ProcState 结构中记录了 PID 为 procPid 的进程读取无效消息的状态，其中 nextMsgNum 的值介于 shmInvalBuffer 的 minMsgNum 值和 maxMsgNum 值之间。

如图 3-23 所示，minMsgnum 和 MaxMsgnum 就像两个指针，它们区分出了哪些无效消息已经被所有的进程读取以及哪些消息还在等待某些进程读取。在 minMsgnum 之前的消息已被所有进程读完；maxMsgnum 之后的区域尚未使用；两者之间的消息是还没有被所有进程读完的。当有进程调用函数 SendSharedInvalidMessage 将其产生的无效消息添加到 shmInvalBuffer 中时，maxMsgnum 就开始向后移动。SendSharedInvalidMessage 中将调用 SIIInsertDataEntries 来完成无效消息的插入。

在向 SI Message 队列中插入无效消息时，可能出现可用空间不够的情况（此时队列中全是没有完全被读取完毕的无效消息），需要清空一部分未处理无效消息，这个操作称为清理无效消息队列，只有当当前消息数与将要插入消息数之和超过 shmInvalBuffer 中 nextThreshold 时才会进行清理操作。这时，那些还没有处理完 SI Message 队列中无效消息的进程将收到清理通知，然后这些进程将抛弃其 Cache 中的所有元组（相当于重新载入 Cache 的内容）。

显然，让所有进程重载 Cache 会导致较高的 I/O 次数。为了减少重载 Cache 的次数，PostgreSQL 会在无效消息队列中设置两个界限值 lowbound 和 minsig，其计算方式如下：

- $\text{lowbound} = \text{maxMsgNum} - \text{MAXNUMMESSAGES} + \text{minFree}$ ，其中 minFree 为需要释放的队列空间的最小值（minFree 指出了需要在无效消息队列中清理出多少个空位用于容纳新的无效消息）。
- $\text{minsig} = \text{maxMsgNum} - \text{MAXNUMMESSAGES}/2$ ，这里给出的是 minsig 的初始值，在进程重载过程中 minsig 会进行调整。

可以看到，lowbound 实际上给出了此次清理过程中必须要释放的空间的位置，这是一个强制性的限制，nextMsgNum 值低于 lowbound 的进程都将其 resetState 字段置为真，这些进程将会自动进行重载 Cache 的工作。对于那些 nextMsgNum 值介于 lowbound 和 minsig 之间的进程，虽然它们并不影响本次清理，但是为了尽量避免经常进行清理操作，会要求这些进程加快处理无效消息的进度（Catch Up）。清理操作会找出这些进程中进度最慢的一个，向它发送 SIGUSR1 信号。该进程接收到 SIGUSR1 后会一次性处理完所有的无效消息，然后继续向下一个进度最慢的进程发送 SIGUSR1 让它也加快处理进度。

清理无效消息队列的工作由函数 SICleanupQueue 实现，该函数的 minFree 参数给出了这一次清理操作至少需要释放出的空间大小。该函数的流程如下：

- 1) 计算 lowbound 和 minsig 的值。
- 2) 对每一个进程的 ProcState 结构进行检查，将 nextMsgNum 低于 lowbound 的进程 resetState 字段设置为 true，并在 nextMsgNum 介于 lowbound 和 minsig 之间的进程中找出进度最慢的一个。
- 3) 重新计算 nextThreshold 参数。

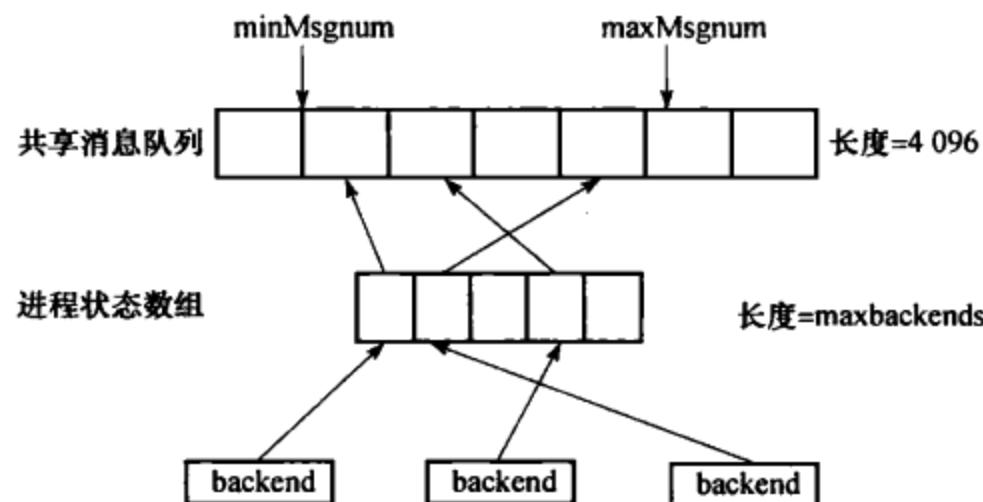


图 3-23 SI Message 机制

4) 向步骤 2 中找到的进度最慢的进程发送 SIGUSR1 信号。

Postgres 进程通过函数 `CatchupInterruptHandler` 来处理 SIGUSR1 信号，该函数最终将调用 `ReceiveSharedInvalidMessages` 来处理所有未处理的无效消息，最后调用 `SICleanupQueue`（`minFree` 参数为 0）向下一个进度最慢的进程发送 SIGUSR1 信号。

每个进程在需要刷新其 Cache 时也会调用 `ReceiveSharedInvalidMessages` 函数用于读取并处理无效消息，函数参数为两个函数指针：

- 1) `intvalFunction`: 用于处理一条无效消息。
- 2) `resetFunction`: 将该后台进程的 Cache 元组全部抛弃。

对于 `resetState` 设置为真的进程，函数 `ReceiveSharedInvalidMessages` 会调用 `resetFunction` 抛弃其所有的 Cache 元组。否则，`ReceiveSharedInvalidMessages` 将从消息队列中读取每条无效消息并调用 `intvalFunction` 对消息进行处理。如果该进程是根据 SIGUSR1 信号调用该函数，那么还将调用 `SICleanupQueue` 函数将这个信号传给比它进度慢的进程。

3. 其他

在 `PMsignal.c` 中，包含后台进程向 Postmaster 发送信号的相关函数。在实现中，后台进程是这样通知 Postmaster 的：

- 1) 首先在共享内存中开辟一个数组 `PMSignalFlags` (`PMsignal.c`)，数组中的每一位对应于一个信号。
- 2) 然后如果后台进程希望向 Postmaster 发送一个信号，那么后台首先将信号在数组 `PMSignalFlags` 中相应的元素置 1 (逻辑真)，然后调用 `kill` 函数向 Postmaster 发送 SIGUSR1 信号。
- 3) 当 Postmaster 收到 SIGUSR1 后首先检测共享存储中 `PMSignalFlags`，确认具体的信号是什么。同时将信号在数组 `PMSignalFlags` 中相应的元素置 0 (逻辑假) 然后作出相应反应。

每一个后台进程都有一个结构 PGPROC 存储在共享内存中。`Procarray.c` 在共享内存中分配 `ProcArrayStruct` 类型的数组 `procArray`，统一管理这些 PGPROC 结构。PGPROC 结构中包含很多的信息，`Procarray.c` 中的函数主要处理 PGRROC 中的 pid、databaseId、roleId、xmin、xid、subxids 等字段。这些函数的功能或是统计事务的信息，或是通过 databaseId 统计有多少个 pid (也就是多少个后台进程) 与指定数据库相连接等统计信息。

IPC 负责的清除工作有两个方面：一个是与共享内存相关的清除，另一个是与各个后台进程相关的清除工作。与共享内存相关的清除并不是将共享内存丢弃，而是重新设置共享内存。清除工作的流程可以描述如下：首先在申请资源的时候，系统会同时为该资源注册一个清除函数，当要求做清除操作时，系统将会调用对应的清除函数。

3.4 表操作与元组操作

上层模块通过表和元组的操作来调用存储模块的诸多功能，这些操作是存储模块与上层其他模块交互的接口，在文件夹 `backend/access/heap` 以及文件 `heaptuple.c` 中对它们进行了详细的定义。接下来将介绍关系表的操作以及对元组的操作。

3.4.1 表操作

表的操作由 `heapam.c` 提供接口。有两种操作表的方式，一种以表名进行操作；另一种以表的 OID 进行操作。而实际上，以表名进行操作的函数，只不过是先通过表的名称获取 OID，然后再调

用以 OID 进行操作的函数（要注意可能存在下述特殊情况，即某个表被删除后又创建了一个同样表名的表，这样在获取 OID 前需要处理所有的无效信息保证根据表名得到的 OID 是最新的）。在 heapam.c 中实现了关系表的打开、关闭、删除、扫描等操作。

1. 打开表

表的打开并不是打开具体的物理文件，仅仅是返回该表的 RelationData 结构体，主要有两种打开方式：

1) 根据表的 OID 打开：调用函数 relation_open，该函数将根据表的 OID 从 RelCache 中找到相应的关系描述符 RelationData，并将引用计数加 1。如果是第一次打开，那么会在 RelCache 中创建一个新的 RelationData（从系统表中获取与表相关的信息填充 RelationData 结构），并将引用计数置为 1。然后根据参数中的锁类型（lock mode）对其加锁，并且返回 RelationData。

2) 根据表名打开：调用函数 relation_openrv，该函数与 relation_open 类似，但它会首先获取到表的 OID，然后调用 relation_open 得到 RelationData 结构体。

可以看到，每当一个表被打开时，都会为这个关系表创建一个 RelationData 结构。RelationData 是关系描述符（relation descriptor），它记录了该表的全部相关信息，前面的章节里已经多次提到该函数。对于表的其他操作都会以表的 RelationData 结构作为重要的参数。

2. 扫描表

PostgreSQL 中对表的扫描实际上有两类方法，第一类是不依赖索引的扫描（顺序扫描），第二类是以索引为基础的扫描（索引扫描）。在本章中我们只介绍顺序扫描的实现，而索引扫描将在第 4 章中介绍。

为了从表中获取我们想要的数据，在没有索引的情况下需要对表文件进行顺序扫描。从表获取一个指定元组的顺序扫描过程需要经过如图 3-24 所示的几个层次。

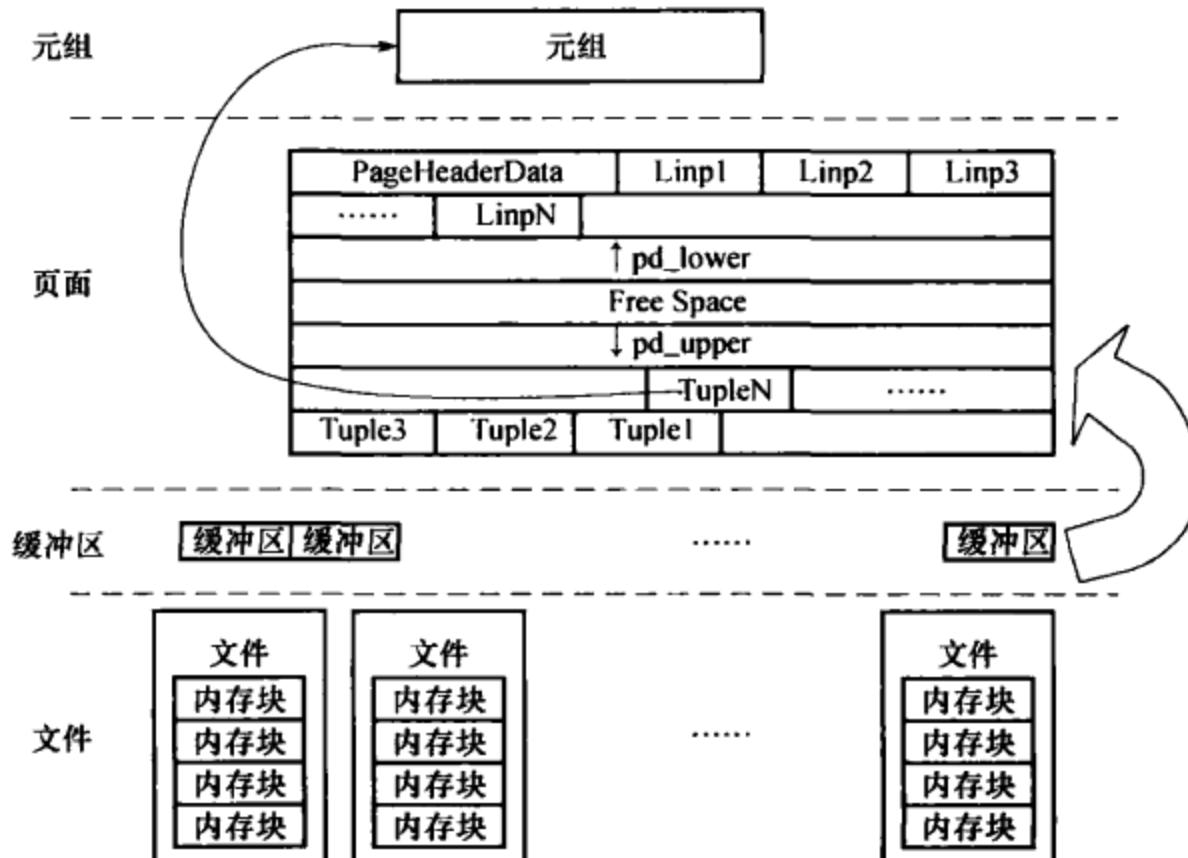


图 3-24 顺序扫描表获取元组

首先，我们将文件块逐一加载到缓冲区中，然后扫描每个缓冲区中的每一个元组，以找到满足查询需求的元组。

在对一个表进行扫描的时候，会利用结构体 HeapScanDescData（数据结构 3.29）来保存表的基本信息以及当前的扫描状态，我们称之为扫描描述符。

数据结构 3.29 HeapScanDescData

```
typedef struct HeapScanDescData
{
    /* 第一部分,设定扫描参数 */
    Relation      rs_rd;                      // 表描述信息
    Snapshot      rs_snapshot;                 // 快照,记录了一些事务相关的信息
    int           rs_nkeys;                   // 扫描键的个数
    ScanKey       rs_key;                     // 扫描键描述结构体
    bool          rs_bitmapscan;              // 是否为位图扫描
    bool          rs_pageatatime;             // 是否允许 page-at-a-time (IsMVCCSnapshot)
                                              // 是否允许使用缓冲区控制策略
    bool          rs_allow_strat;
    bool          rs_allow_sync;                // 是否允许使用 syncscan

    /* 第二部分,初始化该结构体 (initscan) */
    BlockNumber   rs_nblocks;                 // 需要读的表块数
    BlockNumber   rs_startblock;               // 扫描起始块号
    BufferAccessStrategy rs_strategy;         // 缓冲区控制策略
    bool          rs_syncscan;                // 是否报告本次扫描的扫描位置?

    /* 当前扫描状态 */
    bool          rs_initied;                 // 是否已初始化扫描
    HeapTupleData rs_ctup;                   // 当前扫描到的元组
    BlockNumber   rs_cblock;                  // 当前扫描的块号
    Buffer        rs_cbuf;                   // 当前扫描的缓冲区
    ItemPointerData rs_mctid;                // 记录当前扫描的元组项指针

    /* 用于 page-at-a-time 模式和位图扫描 */
    int           rs_cindex;                  // 当前元组在 rs_vistuples 数组中的位置
    int           rs_mindex;                  // 用于记录某个扫描的位置
    Int           rs_ntuples;                 // 当前块内可见的元组数
    OffsetNumberrs_vistuples[MaxHeapTuplesPerPage]; // 存储当前块内可见的元组偏移数组
} HeapScanDescData;
```

PostgreSQL 中对于顺序扫描的实现也有两种策略：基本策略和同步扫描（syncscan）策略。其中，数据结构 HeapScanDescData 中的 rs_allow_sync、rs_startblock 和 rs_syncscan 字段主要是用于同步扫描策略。在本节先对基本的顺序扫描进行介绍，同步扫描将在下一节介绍。

从 HeapScanDescData 的结构体中，我们可以看到对该结构体的填充分为两种：在扫描前初始化

以及在扫描过程中填充扫描状态。

在初始化扫描时将调用函数 `heap_beginscan_internal` 对该结构体进行设置，该函数具有以下几个主要参数：

- `relation`: 将要进行扫描的表的 `RelationData` 结构。
- `nkeys`: 扫描键的个数。
- `key`: 扫描键。
- `allow_strat`: 是否允许使用缓冲区控制策略。
- `allow_sync`: 是否使用同步扫描策略。

`heap_beginscan_internal` 函数的处理流程如下：

1) 首先根据传递给函数的参数对 `HeapScanDescData` 的相关扫描参数进行设置，包括扫描的关系表、快照、扫描键、`rs_allow_strat`、`rs_allow_sync` 等信息。

2) 调用函数 `initscan` 对该结构体第二部分字段进行初始化：

- ① 调用函数 `RelationGetNumberOfBlocks` 获取该表的文件块个数。
- ② 如果该表并非临时表，块的个数大于缓冲池的四分之一，且 `rs_allow_strat` 为 `true` 的话，则设置其缓冲区控制策略类型为 `BAS_BULKREAD`（见缓冲池管理部分）。

在初始化扫描描述符完毕后，就可以调用函数 `heap_getnext` 进行表扫描，该函数每调用一次都返回表中的下一个满足条件的元组。扫描描述符中的 `rs_pageatatime` 字段如果为真，表明当前扫描模式为“page at-a-time”模式，则调用函数 `heapgettup_pagemode` 进行扫描，该函数将会只扫描块中的可见元组；否则调用函数 `heapgettup` 扫描块中的所有元组进行对比。这两个函数的实现逻辑基本类似，但前者无需对缓冲区加 `content` 锁，此外两者在调用函数 `heapgetpage` 时会做不同的处理。两个函数都支持从前往后和从后往前两种扫描方式，这两种方式的实现基本类似，只是读取文件块的顺序不同，因此这里仅介绍从前往后扫描的方式。这两个函数的总体流程如图 3-25 所示。

函数的大致流程如下：

- 1) 设置当前需要读的块号。
- 2) 调用函数 `heapgetpage`：
 - ① 将指定块号的文件块读入内存。
 - ② 若 `scan->rs_pageatatime` 为 `false`，则直接返回；否则，继续执行以读取可见元组。
 - ③ 调用 `heap_page_prune_opt` 清理块内无效 HOT 链。
 - ④ 遍历该文件块，将所有可见的元组偏移量记录到扫描描述符结构体的 `rs_vistuples` 数组中。
 - ⑤ 设置扫描描述符结构体中其他一些相关变量。

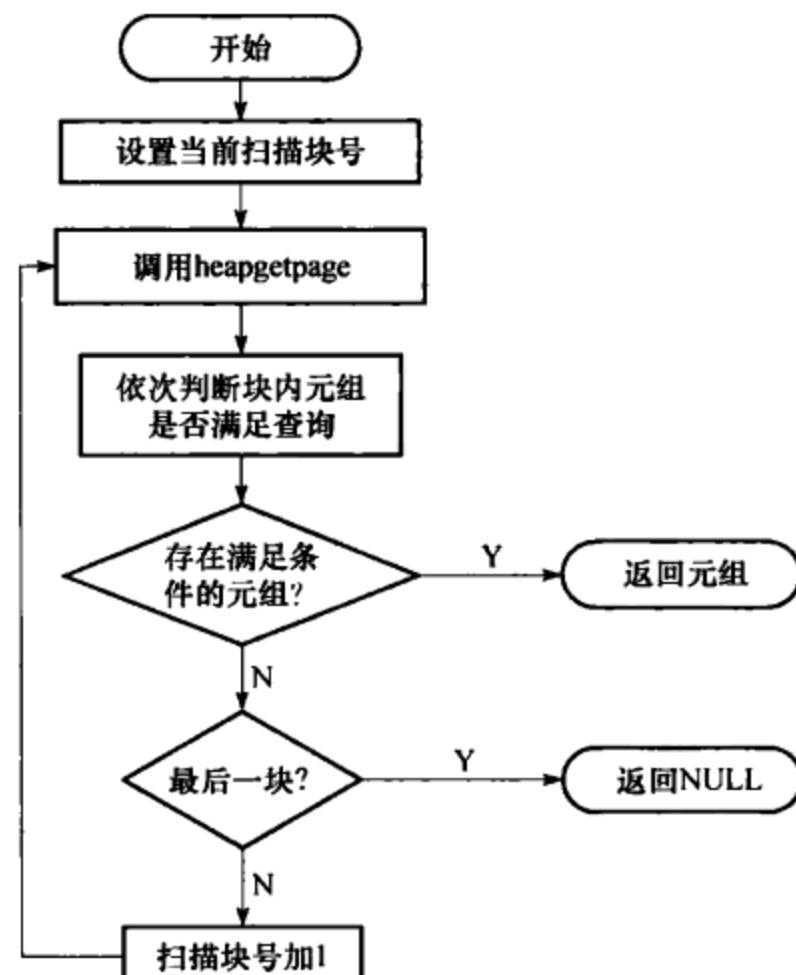


图 3-25 表的扫描过程

3) 判断块内元组是否满足请求。在函数 `heapgettup_pagemode` 中仅扫描数组 `rs_vistuples` 中记录的元组，而在 `heapgettup` 中则依次扫描块内的每一个元组；若找到满足请求的元组则返回，否则继续。

4) 若已到达最后一块，表明该文件中没有满足条件的元组，返回。

5) 否则将页号加 1，并返回到步骤 2。

3. 同步扫描

对于那些没有建立索引的表进行扫描时顺序扫描是我们唯一的选择。如果表中的大部分元组都需要作为结果返回时，顺序扫描的效果最好。但是，在一些情况下顺序扫描的效果也并不好，下面用一个例子进行说明。

如图 3-26 所示，表 T 共有 $N + 1$ 个文件块，同时有三个扫描 A、B、C 对表 T 进行顺序扫描。三个扫描都是从第 0 块开始扫描，目前的扫描进度是 A 最快、B 其次、C 最慢。从理论上来说，A 扫描过的文件块都应该缓存在缓冲池中，这样比 A 进度慢的 B 和 C 应该都能够从缓冲池中找到自己需要的文件块，从而避免对于物理文件的直接访问。但是，在实际运行中会有大量的并发进程同时访问数据库，它们访问到的数据也需要放入缓冲池中，但缓冲池的大小是有限的，因此缓冲池中的旧数据会被新数据替换。在这样的条件下，因为 A 的访问而放入缓冲池的文件块极有可能在 B 需要访问的时候已经被替换出去，那么 B 就需要读取物理文件，然后再将自己读入的文件块放入缓冲池。同理，当 C

要访问 B 访问过的文件块时也极可能出现同样的状况。所以，在这种情况下，虽然在同一个表上进行的多个扫描都是在访问同样的文件块，但它们之间并不能共享缓冲数据，这无疑会大大增加 I/O 的时间开销。

针对上面的这种情况，PostgreSQL 中增加了一种新的顺序扫描策略：同步扫描。同步扫描的思想非常简单。同样结合图 3-26 中的例子来说，就是让 A 仍然从第 0 块开始扫描，B 和 C 开始扫描的位置并不是从第 0 块开始，而是从 A 的当前扫描位置开始扫描，这样由于后续扫描始终紧跟着 A 进行文件块的访问，那么 A 缓存的文件块被替换出去的可能性就极小，因此能够有效地降低 I/O 开销。当然，由于 B 和 C 并不是从第 0 块开始扫描，所以当 B 和 C 扫描到第 N 块时需要回转到第 0 块，将第 0 块到其扫描起始位置之间的文件块继续扫描完，这时虽然 B 和 C 不能再共享 A 的缓冲数据（A 已经结束扫描），但是 B 和 C 之间仍然能够共享。

在实现上，为了能够让后续的扫描能够获取到第一个扫描的当前扫描位置（当前块号），Post-

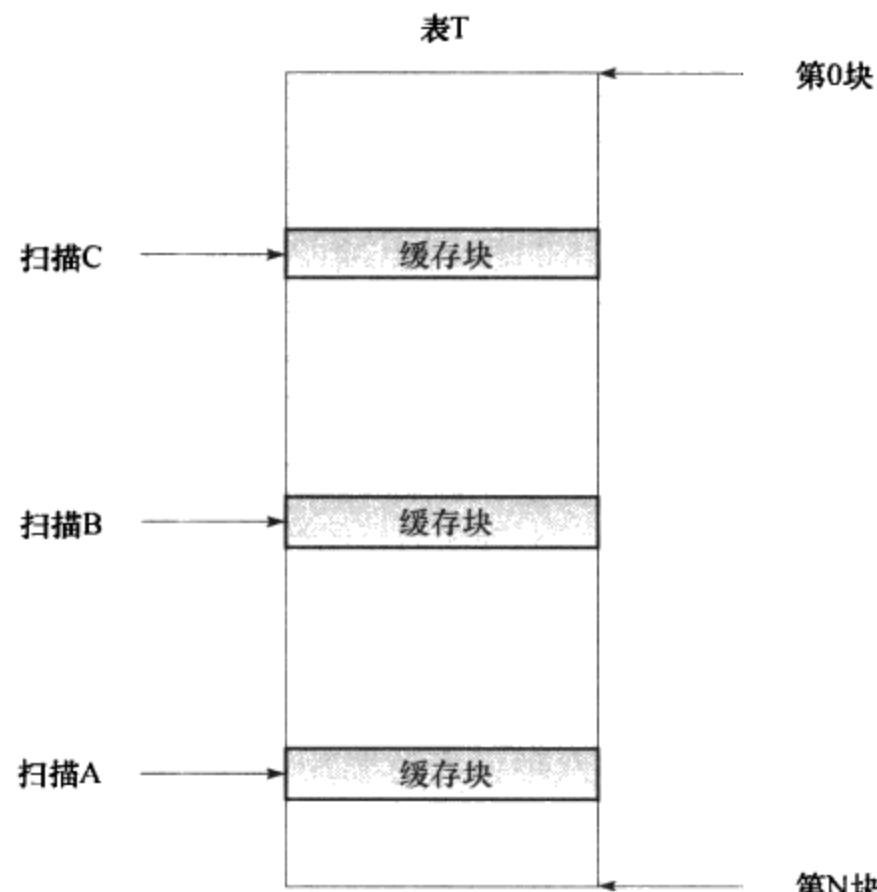


图 3-26 基本顺序扫描的例子

greSQL 在共享内存中放置了一个 `ss_scan_locations_t` 结构（数据结构 3.30）来存放对于某个表的当前扫描位置，每个进程中都有一个全局指针 `scan_locations` 指向这个结构。

`ss_scan_locations_t` 结构中管理了一个 `ss_lru_item_t` 类型（数据结构 3.31）的数组 `items`，数组中的每一个元素对应于一个表，并且记录了该表目前扫描到的块号。`Items` 数组元素中的 `prev` 和 `next` 字段将这些数组元素链接成了一个双向链表，而 `ss_scan_locations_t` 结构的 `head` 和 `tail` 字段则分别指向了这个链表的头部和尾部。

数据结构 3.30 `ss_scan_locations_t`

```
typedef struct RelationData
{
    ss_lru_item_t*head;
    ss_lru_item_t*tail;
    ss_lru_item_t items[1];
}ss_scan_locations_t;
```

数据结构 3.31 `ss_lru_item_t`

```
typedef struct RelationData
{
    struct ss_lru_item_t*prev;
    struct ss_lru_item_t*next;
    ss_scan_location_t location;           //记录了表的标识符,以及当前扫描的块号
}ss_scan_locations_t;
```

为了能够在 `scan_locations` 中记录和获取表的当前扫描块，PostgreSQL 提供了函数 `ss_report_location` 和 `ss_get_location` 分别用于报告某个表的当前扫描块和读取某个表的当前扫描位置。每当扫描到新的一块时，都会由 `heapgettup` 或者 `heapgettup_pagemode` 调用 `ss_report_location` 来报告所扫描表的当前扫描块号。而在初始化扫描时，`initscan` 会根据 `ss_get_location` 返回的结果来设置 `HeapScanDescData` 的 `rs_startblock` 字段，如果 `ss_get_location` 找不到扫描表在 `scan_locations` 中对应的记录将返回 0（表示当前初始化的扫描是第一个在该表上的扫描）。

`HeapScanDescData` 函数中的 `rs_allow_sync` 字段用来控制是否在当前扫描中使用同步扫描，而 `rs_syncscan` 字段则用来控制是否报告当前扫描的扫描位置。目前，PostgreSQL 中只有在建立 B-tree、Hash 和 GiST 三种索引时才可能用到同步扫描。之所以说“可能”，是因为虽然在建立这三种索引时通过调用 `heap_beginscan_internal` 设置了 `HeapScanDescData` 的 `rs_allow_sync` 字段为真，但是在 `initscan` 中还会根据当前扫描所需要扫描的文件块总数是否超过缓冲区总数的四分之一来重新设置该字段，如果超过则保持原设置不变，否则将 `rs_allow_sync` 字段和 `rs_syncscan` 字段都置为假。

4. 其他表操作

在创建表之前先判断数据库目录是否存在，如果不存在，则相应地创建目录。确保表的数据库目录存在后，磁盘管理器调用函数 `mdcreate` 立即申请创建一个文件，这个文件在创建的时候是一个空文件，只有当事务提交之后才会写入相应的内容。

关闭表就是一个内存空间释放和将修改过的数据写回磁盘的过程。写回磁盘时，根据缓冲区的脏标记决定是否需要写回磁盘。

表的删除由函数 `smgrdounlink` 实现。首先它会释放物理内存，将缓冲池中所有该表的缓冲区

全部移除（包括脏块，因为表将被删除，所以脏块也没有回写的必要）并关闭表。然后将根据表的物理路径删除表文件。然而，这个函数通常不会被直接调用。当用户想删除一个表的时候，PostgreSQL 只是会将它加入到 PendingDeletes 这个链表中（其每个节点的类型为数据结构 3.32），在所有的事务提交之后，SMGR 管理器将会调用 smgrDoPendingDeletes 函数将所有需要删除的表删除。

数据结构 3.32 PendingRelDelete

```
typedef struct PendingRelDelete
{
    RelFileNode           relnode;      //要删除的表的物理信息
    bool                  isTemp;       //要删除的表是否为临时表
    bool                  atCommit;     //用于标识事务提交时,是要删除表还是要新建表
    int                   nestLevel;   //插入该节点的事务嵌套层次
    struct PendingRelDelete *next;     //下一个要删除的表,用于链接 PendingDeletes 链表
} PendingRelDelete;
```

注意，在新创建一个关系表时，该表也会被记录到 PendingDeletes 链表中，使用标志 atcommit 来标志：atcommit 设置为 true 标记被删除的表，表示事务提交则删除；atcommit 设置为 false 标记新创建的表，表示事务失败则删除。

PendingRelDelete 的 nestLevel 字段标识插入该数据结构的事务的嵌套层次。当一个子事务提交时，将调用函数 AtSubCommit_smgr，把 PendingDeletes 链表中比当前事务 nestLevel 值更大的节点（都是它的子事务）的 nestLevel 值进行减 1 操作，这相当于把子事务中删除或创建文件的节点合并到其父事务中，当子事务从最低层逐层提交后，最终所有的删除或创建文件的节点都将合并到最顶层事务中；当子事务中断时，则调用函数 AtSubAbort_smgr，对 PendingDeletes 链表中 nestLevel 大于当前事务 nestLevel 的节点进行回滚，将已创建的文件删除，将标记为删除的文件的节点从链表中删除，即放弃对该文件的删除。

smgrDoPendingDeletes 是用于处理 PendingDeletes 的主要函数，它用布尔型参数 isCommit 来标识事务是否执行成功。在成功提交事务（函数 CommitTransaction）或事务失败（AbortTransaction 和 AtSubAbort_smgr）时，均会调用到这个函数来对链表中的记录集中处理。

3.4.2 元组操作

对元组的操作包括插入、删除和更新三种基本操作，这三种操作都是把元组当作一个整体进行处理。除些之外，在 heaptuple.c 这个文件中还实现了元组内部结构的相关操作，包括元组的构造、修改、分解、复制、释放等操作。

一个完整的元组信息将对应一个 HeapTupleData 结构和一个 TupleDesc 结构，在 HeapTupleData 中还包含一个前面介绍过的 HeapTupleHeaderData 结构。

TupleDesc 是关系结构 RelationData 的一部分，也称为元组描述符，它记录了与该元组相关的全部属性模式信息。通过元组描述符可以读取磁盘中存储的无格式数据，并根据元组描述符构造出元

组的各个属性值，元组描述符的结构见数据结构 3.33。

其中：

- natts 表示元组的属性个数。
- attrs 数组表示元组每个属性的相关信息，从 pg_attribute 系统表读取，每个元素表示一个属性。
- constr 数组用于表示元组的约束条件，从 pg_constraint 系统表中读取，每个元素表示一个约束条件。
- tdtypeid 是元组的复合类型 ID。PostgreSQL 会自动为每个表建立一个行类型，以表示该表的行结构。
- tdtypmod 是元组模式。当 tdtypeid 有值时，tdtypmod 置位 -1。当 tdtypeid 为默认值 RECORDOID 时，tdtypmod 可以为 -1 或其在 typecache（建立在 syscache 之上）中的位置。
- tdhasoid 表示元组是否有 OID。
- tdrefcount 表示该元组描述符的引用计数。只有在引用计数为 0 时，该元组描述符才能被删除。

HeapTupleData 是元组在内存中的拷贝，它是磁盘格式的元组读入内存后的存在方式，HeapTupleData 的结构如数据结构 3.34 所示。

数据结构 3.33 tupleDesc

```
typedef struct tupleDesc
{
    int natts;
    Form_pg_attribute *attrs;
    TupleConstr *constr;
    Oid tdtypeid;
    int32 tdtypmod;
    bool tdhasoid;
    int tdrefcount;
} *TupleDesc;
```

数据结构 3.34 HeapTupleData

```
typedef struct HeapTupleData
{
    uint32 t_len;           // 表示该元组的长度
    ItemPointerData t_self; // 记录块号以及元组在块内的偏移量
    Oid t_tableOid;        // 记录元组所在表的 OID
    HeapTupleHeader t_data; // 用于记录元组头信息
} HeapTupleData;
```

从上述三个数据结构可以看到，在数据结构中并没有出现存储元组实际数据的属性，这是因为 PostgreSQL 通过编程技巧，巧妙地将元组的实际数据存放在 HeapTupleHeaderData 结构后面的空间中。

1. 插入元组

在插入元组之前，我们首先要根据元组内数据和描述符等信息初始化 HeapTuple 结构，函数 heap_form_tuple 实现了这一功能。函数原型如下：

```
Heap Tuple heap_form_tuple(TupleDesc tupleDescriptor, Datum *values, bool *isnull)
```

其中，values 参数是将要插入的元组的各属性值数组，isnull 数组用于标识哪些属性为空值。heap_form_tuple 根据 values 和 isnull 数组调用函数 heap_compute_data_size 计算形成元组所需要的内

存大小，然后为元组分配足够的空间。在进行必要的元组头部设置后，调用函数 `heap_fill_tuple` 向元组中填充实际数据。

当完成了元组数据在内存中的构成后，下一步就可以准备向表中插入之组了。插入元组的函数接口为 `heap_insert`，其流程如图 3-27 所示。

- 1) 首先我们会为新插入的元组（`tup`）调用 `newoid` 函数为其分配一个 OID。
- 2) 初始化 `tup`，包括设置 `t_xmin` 和 `t_cmin` 为当前事务 ID 和当前命令 ID、将 `t xmax` 置为无效、设置 `tableOid`（包含此元组的表的 OID）。

3) 找到属于该表且空闲空间（`freespace`）大于 `newtup` 的文件块，将其载入缓冲区以用来插入 `tup`（调用函数 `RelationGetBufferForTuple`）。

4) 有了插入的元组 `tup` 和存放元组的缓冲区后，就会调用 `RelationPutHeapTuple` 函数将新元组插入至选中的缓冲区。

5) 向事务日志（`XLog`）写入一条 `XLog`（见第 7 章）。

6) 当完成上述过程后，将缓冲区解锁并释放，并返回插入元组的 OID。

2. 删除元组

在 PostgreSQL 中，使用标记删除的方式来删除元组，这对于多版本并发控制（Multi-version Concurrency Control，MVCC）是有好处的，其 `Undo` 和 `Redo` 速度是相当高速的，因为只需重新设置标记即可。被标记删除的磁盘空间会通过运行 `VACUUM`（清理数据库命令，通常每天运行一次）收回。

删除元组主要调用函数 `heap_delete` 来实现，其主要流程如下：

- 1) 根据要删除的元组 `tid` 得到相关的缓冲区，并对其加排他锁（Exclusive Lock）。
- 2) 调用 `HeapTupleSatisfiesUpdate` 函数检查元组对当前事务的可见性。如果元组对当前事务是不可见的（`HeapTupleSatisfiesUpdate` 函数返回 `HeapTupleInvisible`），那么对缓冲区解锁并释放，再返回错误信息。
- 3) 如果元组正在被本事务修改（`HeapTupleSatisfiesUpdate` 函数返回 `HeapTupleSelfUpdated`）或已经修改（`HeapTupleSatisfiesUpdate` 函数返回 `HeapTupleUpdated`），则将元组的 `ctid` 字段指向被修改后的元组物理位置，并对缓冲区解锁、释放，再分别返回 `HeapTupleSelfUpdated` 和 `HeapTupleUpdated` 信息。
- 4) 如果元组正在被其他事务修改（`HeapTupleSatisfiesUpdate` 函数返回 `HeapTupleBeingUpdated`），那么将等待该事务结束再检测。如果事务可以修改（`HeapTupleSatisfiesUpdate` 函数返回 `Heap-`

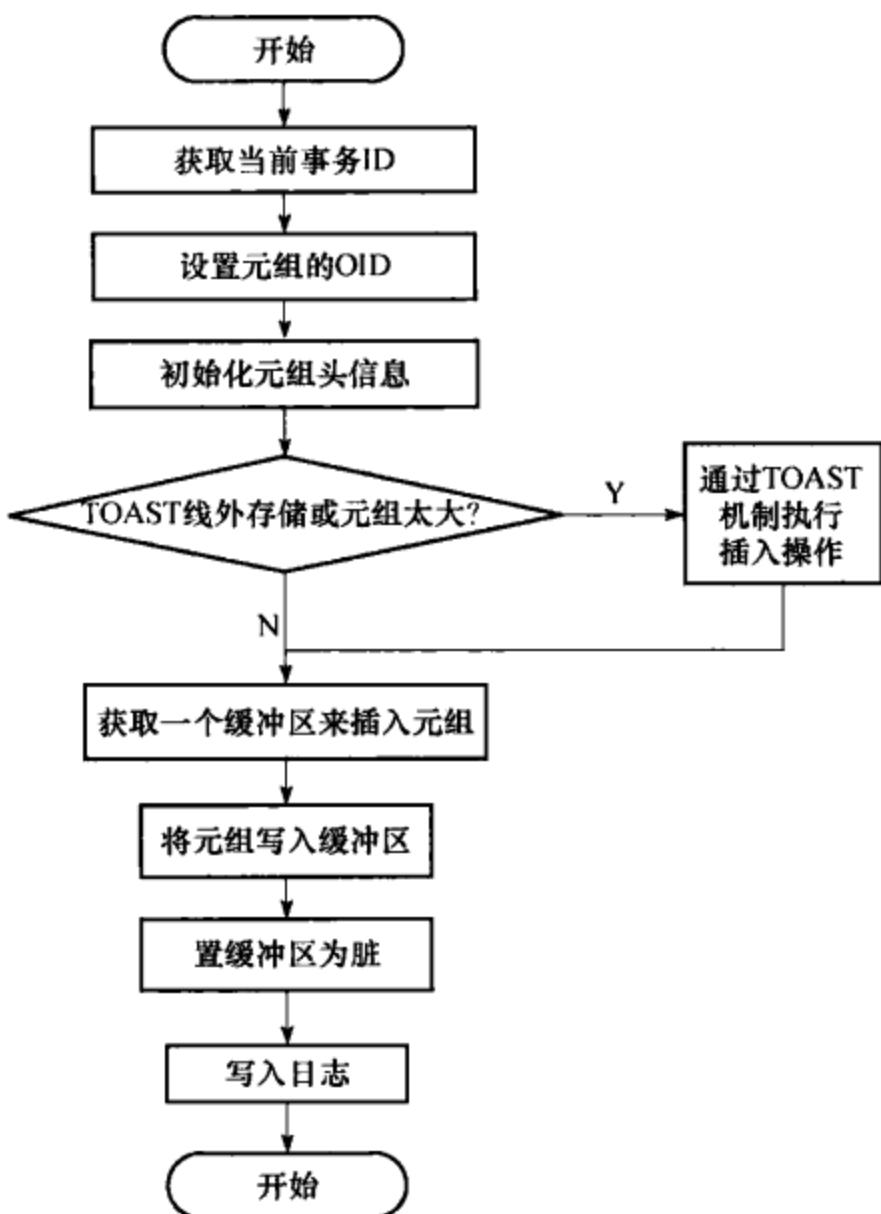


图 3-27 插入元组流程

`TupleMayBeUpdated`), 那么 `heap_delete` 会继续向下执行。

5) 这时就进入临界区域, 设置 `t_xmax` 和 `t_cmax` 为当前事务 ID 和当前命令 ID。到此为止该元组已经被标记删除, 或者可以说该元组已经被删除了。

6) 记录 XLog。

7) 如果此元组存在线外数据, 即经过 TOAST 的数据, 那么还需要将其 TOAST 表中对应的数据删除。调用函数 `toast_delete` 完成相关工作。

8) 如果是系统表元组, 则发送无效消息。

9) 设置 FSM 表中该元组所处文件块的空闲空间值。

3. 更新元组

元组的更新操作实际上是删除和插入操作的结合, 即先标记删除旧元组, 再插入新元组。元组的更新由函数 `heap_update` 实现。

值得注意的是, PostgreSQL 中进行删除和更新操作时, 被删除或修改的元组并不会从物理文件中删除, 而是在事务标记中被标记为无效。因此, 当进行过大量的删除和更新操作之后, 数据库数据文件中由于有大量的无效元组, 其尺寸会变得异常庞大, 此时需要对数据库进行一定的清理操作, 这就需要用到下一节要介绍的 VACUUM 机制。

3.5 VACUUM 机制

当修改元组属性、执行重建索引等操作时, 特别是当进行多次更新、删除操作时, 磁盘上会出现很多无效的元组, 占据了很多磁盘空间并且会导致系统性能下降, 这时就需要进行 VACUUM 操作清除掉这些“垃圾”。

很明显, 那些经常更新或者删除元组的表需要比那些较少更新的表清理得更频繁一些。PostgreSQL 开启了一个辅助进程 AutoVacuum 来执行自动清理工作。在创建表时可以加上 `WITH (auto-vacuum_enabled)` 来允许 AutoVacuum 进程自动清理该表。用户还可以在系统配置文件中定义一些值, 用于“指挥” AutoVacuum 进程在表中的无效数据到达某个数量时执行清理操作。

3.5.1 VACUUM 操作

VACUUM 操作的具体实现定义在文件 Vacuum.c 中。根据用户输入的命令, VACUUM 分为两种: 一种为 Full VACUUM, 它会对表进行完全清理; 一种为 Lazy VACUUM, 它仅标记无效数据空间为可用。VACUUM 的命令格式如下:

- `VACUUM[FULL][VERBOSE][table]`
- `VACUUM[FULL][VERBOSE]ANALYZE[table[(column[,...])]]`

其中, FULL 选项表示完全清理, 不加 FULL 则默认为 Lazy VACUUM。VERBOSE 选项将在清理完毕后打印一份清理报告, ANALYZE 选项则用于更新优化器的统计信息。

在介绍 VACUUM 的实现之前, 首先介绍两个在 Lazy VACUUM 和 Full VACUUM 中皆会用到的重要函数。

第一个函数是 `PageRepairFragmentation`, 当将一个文件块内无效的元组标记为可用后, 该函数将

把块内所有的空闲碎片移动到块的空闲区域，从而实现了块内的碎片整理。

第二个函数是 `heap_page_prune`，该函数用于清理单个文件块的 HOT 链，并进行块内碎片的整理。碎片整理调用 `PageRepairFragmentation` 函数来完成。

VACUUM 操作主要由函数 `vacuum` 实现，其主要流程如图 3-28 所示。

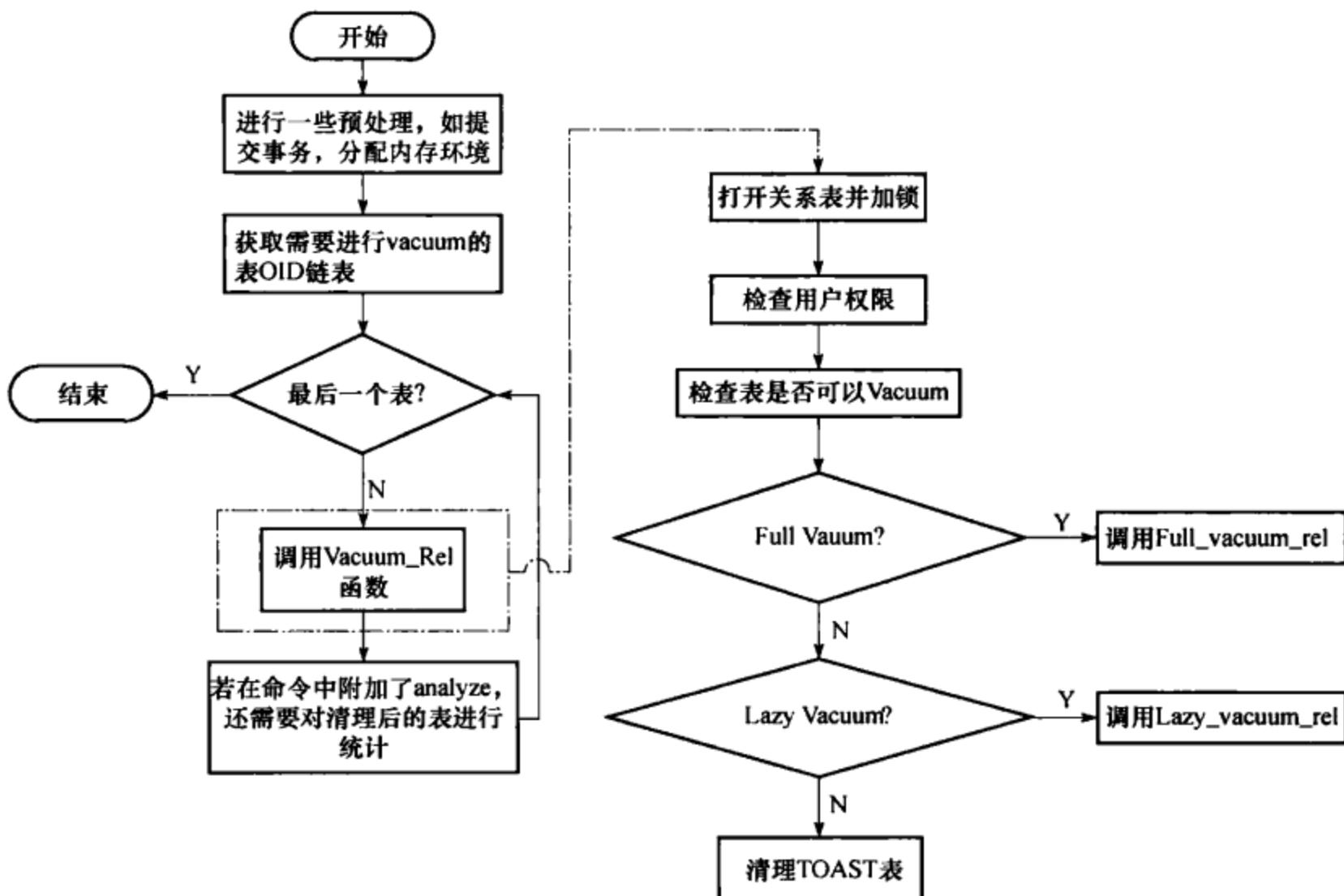


图 3-28 `vacuum` 函数及 `vacuum_rel` 函数流程

由图 3-28 可见，通过在命令中指定 VACUUM 的类型，在 `vacuum_rel` 函数中会使用不同类型的 VACUUM，然后调用相应的接口函数。

`Lazy` VACUUM 调用函数 `Lazy_vacuum_rel` 来实现，该操作会在表和索引中标记无效数据为可用。它并不试图立即回收这些无效数据使用的空间（除非位于表末尾并且很容易锁定此表）。因此，表文件不会缩小，并且任何文件中没有使用的空间都不会返回给操作系统。这种 VACUUM 操作可以和通常的数据库操作并发执行。

`Full` VACUUM 主要在函数 `Full_vacuum_rel` 中进行实现，这个形式的 VACUUM 使用一种更加激进的方式来回收无效元组占据的空间。通过 VACUUM FULL 回收的空间都立即返回给操作系统。但是，在进行 VACUUM FULL 操作的时候要求加一个排他锁。因此，经常使用 VACUUM FULL 会对并发数据库查询有不利的影响。

3.5.2 Lazy VACUUM

对一个表进行 `Lazy VACUUM` 通过 `lazy_vacuum_rel` 函数实现，该函数在执行过程中将对单个表

进行清理，并清理其索引、更新页面数和元组数的统计值，主要流程如图 3-29 所示。

`lazy_scan_heap` 函数是执行 Lazy VACUUM 的核心函数，该函数将首先扫描表，找到无效的元组和具有空闲空间的页面，然后计算表的有效元组数目，最后执行表和索引的清理操作。

注意，这里在扫描关系表的时候用到了前面介绍过的 VM 表，`lazy_scan_heap` 函数的执行流程如图 3-30 所示。

在该流程图中，调用了以下几个比较重要的函数：

1) `lazy_vacuum_heap`: 该函数将对无效元组链上的元组按文件块为单位进行逐块清理，清理函数为 `lazy_vacuum_page`，并调用函数 `RecordPageWithFreeSpace` 对文件块的空闲空间值进行更新；

2) `lazy_vacuum_page`: 负责在文件块内，检查每一个元组指针（Linp 指针，见 3.2.1 节），如果该指针被标记为 LP_DEAD 则将其标记改为 LP_UNUSED，然后对文件块内的碎片进行整理，将空闲元组空间移动到文件块的空闲区域。

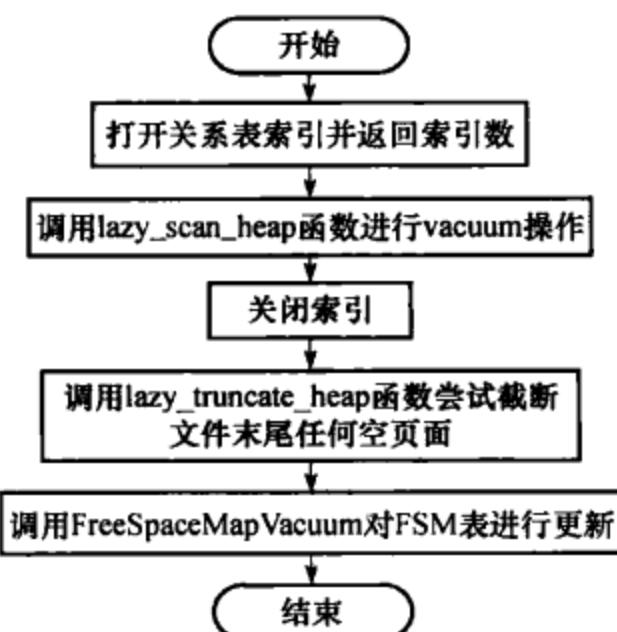


图 3-29 `lazy_vacuum_rel` 函数流程

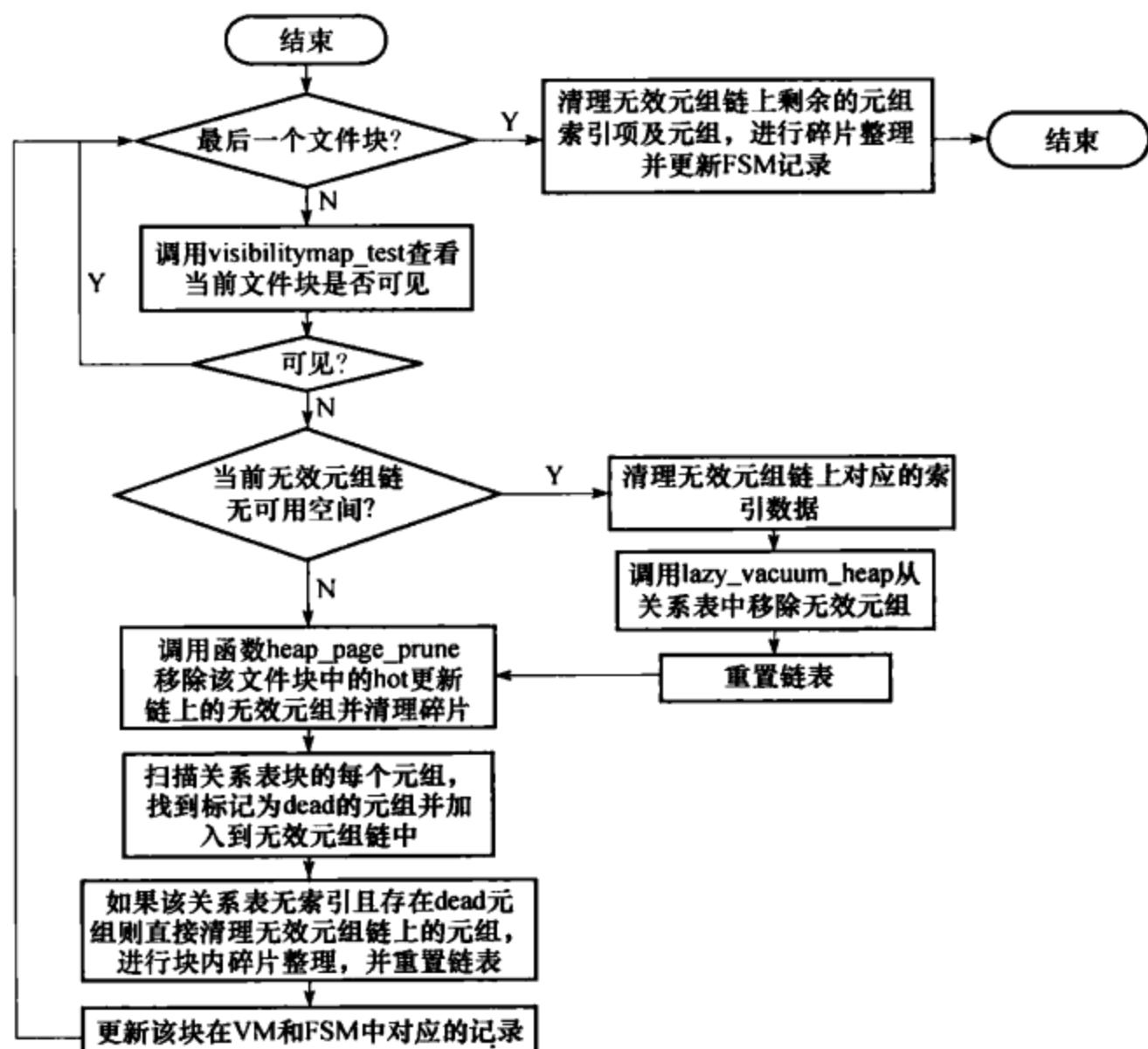


图 3-30 `lazy_scan_heap` 函数

3) heap_page_prune 函数：该函数用于清理单个文件块的 HOT 链，并进行块内碎片的整理。

3.5.3 Full VACUUM

与 Lazy VACUUM 不同，Full VACUUM 除了对文件块进行整理外，还实现了更加复杂的跨块移动元组的操作，其接口函数为 full_vacuum_rel。该函数用于对单个关系表或索引进行清理操作，相比较 Lazy VACUUM 而言，其实现更加复杂。

在介绍具体的函数之前，首先介绍用于控制 Full VACUUM 的两个链表，即 fraged_pages 链表和 vacuum_pages 链表，前者用于记录具有空闲空间的文件块，用作作为元组跨块移动的目标块（即可以把别的文件块中的元组移动到这个链表中的文件块里）；后者则记录了所有需要清理的文件块。这两个链表的内容有可能重复，两者节点的结构体相同，如数据结构 3.35 所示。

数据结构 3.35 VacPageListData

```
typedef struct VacPageListData
{
    BlockNumber    empty_end_pages;           //最后一个空文件块号
    int            num_pages;                 //在 pagedesc 数组中的 vacpage 数目
    int            num_allocated_pages;        //在 pagedesc 中分配的文件数
    VacPage        *pagedesc;                 //需要进行 VACUUM 的文件块数组
} VacPageListData;
```

每个 VacPage 结构体用于跟踪每一个存在可用空闲空间的文件块，其结构如数据结构 3.36 所示。

数据结构 3.36 VacPage

```
typedef struct VacPageData
{
    BlockNumber    blkno;                   //文件块块号
    Size           free;                    //在该块中的空闲空间大小
    uint16         offsets_used;           //空闲的 linp(项指针)数目
    uint16         offsets_free;           //空闲或将被释放的 linp 数目
    OffsetNumber   offsets[1];             //数组, 用于记录 offsets_free 个 linp 的数组
} VacPageData;
typedef VacPageData* VacPage;
```

full_vacuum_rel 函数的主要流程如图 3-31 所示。

通过分析源代码，我们发现，在执行完全清理的过程中，并没有使用 VM 表来加快清理过程（VM 表中记录有需要清理的块）。这是因为在 FULL VACUUM 过程中，需要扫描每一个文件块来找到有足够的空闲空间的文件块，VM 文件在这里没有太大的作用。

fraged_pages 链表和 vacuum_pages 链表的生成由 scan_heap 函数来完成。该函数对表进行扫描，标记无效元组并计算每个表块的空闲空间，将有足够的空间（一般为 BLOCKSZ/10）可用于装

填元组的文件块加入到 `fraged_pages` 链表中，将需要进行清理操作的文件块加入到 `vacuum_pages` 链表中。

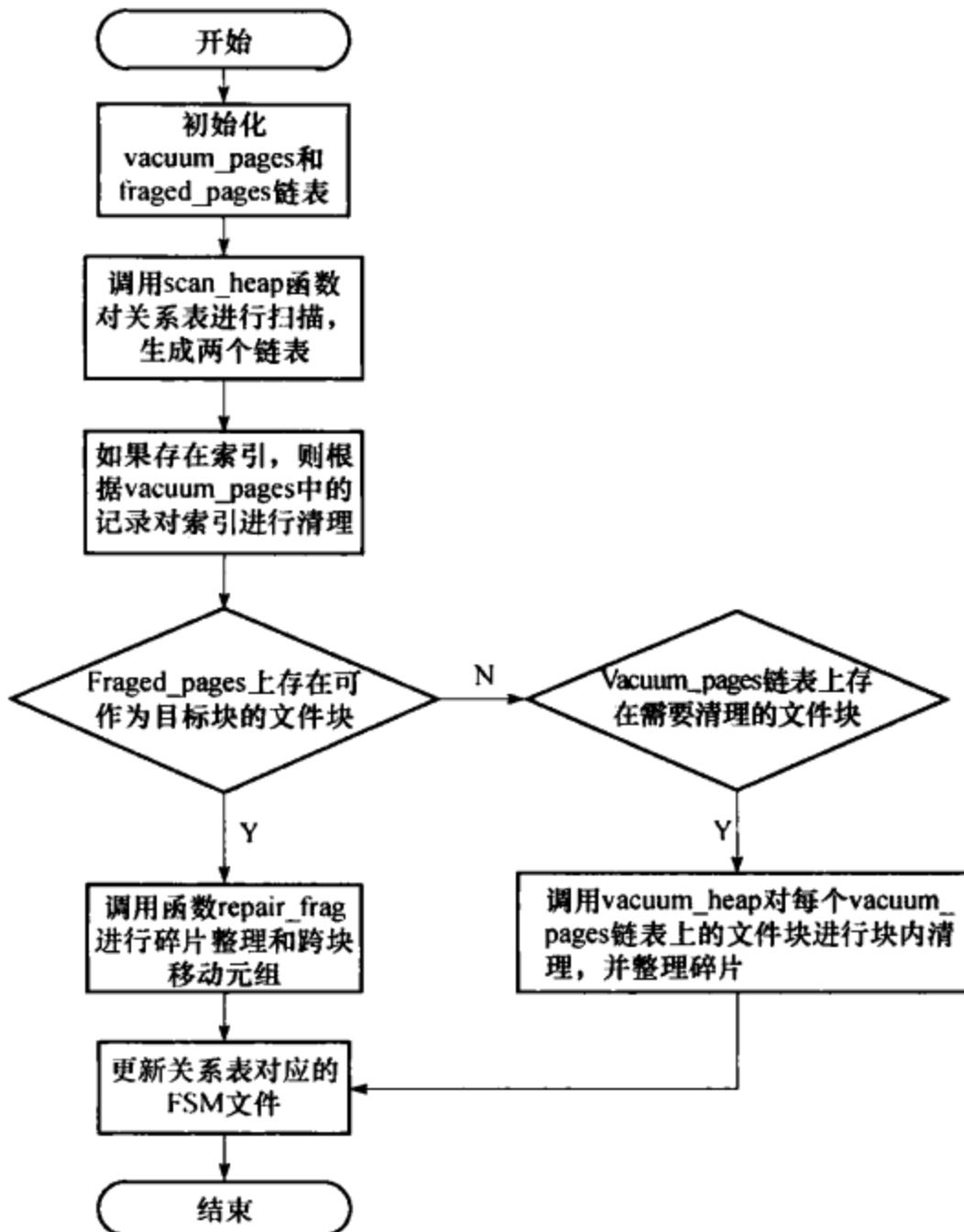


图 3-31 `full_vacuum_rel` 函数流程

在对表进行扫描的过程中，`scan_heap` 和 Lazy VACUUM 操作一样，也调用了函数 `heap_page_prune` 来对每个关系表块的 HOT 更新链进行清理和块内碎片整理，然后再对文件块中的每个元组进行扫描，记录块中未使用的项指针。若存在标记为 `dead` 的元组，则也将其项指针记录下来；在扫描文件块结束后，将文件块复制到一个临时块中进行碎片整理并计算可用的空闲空间（包括未使用和将被释放的），从而决定是否将这个块加入到 `fraged_pages` 链表中。

`repair_flag` 是碎片整理函数，该函数利用 `fraged_pages` 链表和 `vacuum_pages` 链表来进行整理，不仅包括对块内进行整理，还包括跨块移动元组，其实现繁琐复杂。大体流程如下所示：

1) 从 `vacuum_pages` 上最后一个非空文件块开始向前扫描。

- ①如果当前块存在无效数据，需要调用 `vacuum_page` 函数对文件块进行清理和块内碎片整理。
- ②逐个扫描块内元组，并将其向 `frag_pages` 链表上的文件块中移动元组，目标块总是在当前块的前面。元组跨块移动函数为 `move_plain_tuple` 和 `move_chain_tuple`，前者用于移动单个元组，后者

则用于移动最近更新产生的元组 HOT 链。

③当没有可作为移动目标的文件块时，跳出循环。

2) 调用 `vacuum_page` 函数处理 `vacuum_pages` 链表上剩余的文件块。

3) 由于跨块移动，文件尾可能产生空文件块，因此需要对文件进行截断，去除空块。

总的来说，Lazy VACUUM 开销较小，但是空间回收的效果也较差；而 Full VACUUM 开销较大，但是回收空间的效果很好。在实际的应用中，需要根据实际情况来选择。

3.6 ResourceOwner 资源跟踪

在 PostgreSQL 中，每个事务（包括其子事务）在执行时，需要跟踪其占用的内存资源。例如，需要跟踪一个缓冲区的 pin 锁或者一个表的锁资源，以保证在事务结束或者失败的时候能够被及时释放。PostgreSQL 中使用 ResourceOwner 对象来对资源集中进行跟踪，我们也称 ResourceOwner 为资源跟踪器。ResourceOwner 记录了大多数事务过程中使用到的资源，主要包括：

- 父节点、子节点的 ResourceOwner 指针，用于构成资源跟踪器之间的树状结构。
- 占用的共享缓冲区数组和持有的缓冲区 pin 锁个数。
- Cache 的引用次数以及占用的 Cache 中存储的数据链表，包括 CatCache、RelCache 以及缓存查询计划的 PlanCache。
- TupleDesc 引用次数以及占用的 TupleDesc 数组。
- Snapshot（见 7.10.3 节）引用次数以及占用的 Snapshot 数组。

ResourceOwner 对象与 MemoryContext 类似，相互之间会构成一个树结构（如图 3-32 所示），每一个 Portal（见 6.2 节）、事务及其子事务都有自己的 ResourceOwner。为了便于管理 ResourceOwner 树，在内存中保持了三个全局的 ResourceOwner 结构变量：

- CurrentResourceOwner：记录当前使用的 ResourceOwner。
- CurTransactionResourceOwner：记录当前事务的 ResourceOwner。
- TopTransactionResourceOwner：记录顶层事务的 ResourceOwner。

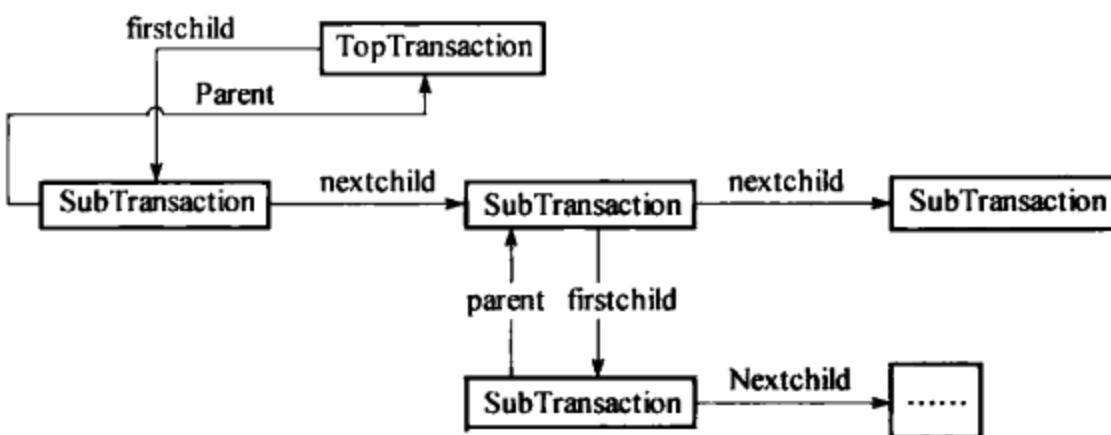


图 3-32 事务过程中的 ResourceOwner 树

每一个事务的资源跟踪器都通过其 `parent` 指针指向父事务的资源跟踪器，如果该事务本身就是顶层事务，则 `parent` 指针为空。同一事务的子事务的资源跟踪器之间通过 `nextchild` 指针链接在一起，而父事务的资源跟踪器的 `firstchild` 指针则指向其第一个子事务的资源跟踪器。通过 `parent`、`firstchild`、`nextchild` 三个指针，资源管理器之间就构成了一种树状结构。

在 PostgreSQL 中建立了多个 ResourceOwner 树，根节点分别命名为 TopTransacion、background writer、buildflatfiles、Wal Writer，它们有着不同的作用。每一个 ResourceOwner 的节点内存都在 TopMemoryContext 中进行分配。

无论对于哪种资源，都为其定义了三种操作：在 ResourceOwner 中分配内存、向 ResourceOwner 中增加资源跟踪、从 ResourceOwner 中移除资源跟踪。

这里以缓冲区资源管理为例加以说明。

(1) 在 ResourceOwner 中分配内存（函数 ResourceOwnerEnlargeBuffers）

当需要读一个缓冲区时，首先要确保 ResourceOwner 中有足够的空间（buffers 数组）来存储缓冲区指针。如果 buffers 原来没有分配内存，将会在 TopMemoryContext 中为 buffers 分配长度为 16 的空间；如果 buffers 中的空间不够，将会把 buffers 的空间增加一倍。

(2) 增加对缓冲区的跟踪（函数 ResourceOwnerRememberBuffer）

当获取一个缓冲区并对其加 pin 后，需要将该缓冲区加入到 ResourceOwner 的 buffers 数组中，并增加 ResourceOwner 对缓冲区的跟踪计数（把 nbuffers 加 1）。

(3) 移除对缓冲区的跟踪（函数 ResourceOwnerForgetBuffer）

当涉及对缓冲区减 pin 操作时，则将该缓冲区从资源跟踪器的 buffers 数组中移除，并减少资源跟踪器对缓冲区的跟踪计数（把 nbuffers 减 1）。

对于其他资源而言，基本上也可以为不同的资源定义上述三种操作。当事务完成并提交后，就从 TopTransactionResourceOwner 开始检查所有占用的资源一次性进行释放。

对于资源跟踪器本身，主要有两种操作：创建资源跟踪器和释放资源跟踪器所跟踪的资源。

1) 资源跟踪器的创建由函数 ResourceOwnerCreate 实现。该函数在 TopMemoryContext 中为要创建的资源跟踪器分配内存，将 parent 指针指向其父节点，将 nextchild 指针指向父节点当前的第一个子节点，然后将父节点的 firstchild 指针指向自身。这里要注意，创建资源跟踪器时并没有为用于记录所跟踪资源的数组（如 buffers）等分配内存，这些数组的内存将在实际添加跟踪资源时分配。

2) 资源跟踪器中跟踪资源的释放工作由函数 ResourceOwnerRelease 完成。释放资源时会递归地调用该函数来释放子节点中跟踪的资源。对于每一个叶子节点来说，释放资源的工作分为三个阶段：

①释放共享缓冲区和 RelCache (RESOURCE_RELEASE_BEFORE_LOCKS)：在提交时进行这一阶段释放时，理论上是不应该存在需要释放的共享缓冲区和 RelCache 的，如果这种情况下还有未释放的类似资源，该函数会打印出资源泄漏警告信息。

②释放锁 (RESOURCE_RELEASE_LOCKS)：如果当前资源跟踪器属于一个顶层事务，则释放所有的非会话锁（会话锁是跨事务存在的）；否则将该跟踪器跟踪的锁转移给它的父节点。

③释放其他资源 (RESOURCE_RELEASE_AFTER_LOCKS)：释放掉除以上两类资源之外的其他资源。和第一阶段相似，在提交阶段发现有资源没有释放时也会打印资源泄漏的警告。

从资源跟踪器对资源的释放过程可以看到，虽然 PostgreSQL 中期望事务会主动释放其不再使用的资源，但是并不会因为事务忘记释放资源（通常是在对 PostgreSQL 源代码的修改过程中忘记了主动释放资源）而发生泄漏。资源跟踪器会跟踪所有在事务执行过程中用到的资源，如果事务忘记了释放资源，在提交时会通过释放资源跟踪器来释放这些资源。

3.7 小结

本章首先从整体结构上介绍了 PostgreSQL 的存储管理架构，然后分别从外存和内存两个方面进行了详细的描述，并介绍了建立在存储管理上的表和元组操作，最后介绍了 PostgreSQL 的数据清理机制和资源跟踪器。存储管理模块是整个系统的基础部分，对其他模块提供了访问数据的接口，尤其是内存管理模块与其他模块的交互最多。读者在阅读其他章节时，将会对 PostgreSQL 内存管理模块的使用有更深入的了解。

习题

习题 3.1 建立一个表，采用任意一种编程语言，随机产生 N 条元组，并在任意一段时间内，模拟对该表进行随机操作，执行两种 VACUUM 操作，分析两种操作对数据库表文件的影响，并撰写实验报告。

习题 3.2 建立一个表，并向其中插入一个文本类型大对象，文本内容不超过 2KB，利用 GDB，找到相应的大对象处理函数，跟踪其插入流程，并撰写实验报告。

第 4 章



索引是指按表中某些关键属性或表达式建立元组的逻辑顺序，它是由一系列表元组的标识号组成的一个列表。使用索引可快速访问表中的特定元组，被索引的表称为基表。索引并不改变表中元组的物理顺序，索引技术会将对于元组的逻辑排序保存在索引文件中。基表文件中的元组被修改或删除时，索引文件会自动更新以保证能够准确地找到新的数据。任何一个成熟的数据库都离不开索引的支持，PostgreSQL 当然也不例外。

在 PostgreSQL 8.4.1 中目前支持的索引有：B-Tree 索引、Hash 索引、GiST 索引和 GIN 索引，本章将分别对这些索引进行分析。

4.1 概述

假设有下面这样一个表：

```
CREATE TABLE student(
    id integer,
    name varchar
);
```

需要大量使用类似下面这样的语句进行查询：

```
SELECT name FROM student WHERE id = 007;
```

通常，数据库系统需要一行一行地扫描整个 student 表以寻找所有匹配的元组。如果表 student 的规模很大，但是满足 WHERE 条件的只有少数几个（可能是零个或一个），那么这种顺序扫描的性能就比较差了。如果让数据库系统在 id 属性上维护一个索引用于快速定位匹配的元组，那么数据库系统只需要在搜索树中查找少数的几层就可以找到匹配元组，这将大大提高数据查询的性能。类似的，在数据库中进行更新、删除操作时也需要先找到要操作的元组，利用索引同样可以提升这些操作的性能。

PostgreSQL 8.4.1 中一共提供了 5 种索引方式：唯一索引、主键索引、多属性索引、部分索引、

表达式索引，以及 4 种索引类型：B-Tree、Hash、GiST、GIN。每种索引类型都分别适合某些特定的查询类型，因为它们用了不同的索引结构。CREATE INDEX 命令可以用来创建索引，默认情况下该命令将创建一个 B-Tree 索引。

PostgreSQL 里的所有索引都是“从属索引”，也就是说，索引在物理上与它描述的表文件分离。作为一种数据库对象，每个索引都在 pg_class 表里面有记录，一个索引的内部结构与该索引的访问方法（索引类型）相关。PostgreSQL 中所有索引访问方法都通过页面来组织索引的内部结构，这样可以使用存储管理器提供的接口来访问索引。所有现有的索引访问方法都使用 3.2.1 节描述的标准的页面布局。

索引从本质上来说就是一些数据的键值与元组标识符（TID）之间的映射，这些标识符确定了该索引键值在表中对应的元组。

4.1.1 索引方式

PostgreSQL 中共有 5 种索引方式，即唯一索引、主键索引、多属性索引、部分索引和表达式索引。下面将分别加以介绍：

1) 唯一索引：如果索引声明为唯一索引，那么就不允许出现多个索引值相同的元组。唯一索引可以用于强迫索引属性数值的唯一性，或者是多个属性组合值的唯一性。唯一索引通过在创建索引命令中加上 UNIQUE 关键字来创建。一个多字段唯一索引认为只有两个元组的所有被索引属性都相同的时候才是相同的，这种重复元组才被拒绝。目前，只有 B-Tree 可以创建唯一索引。

2) 主键索引：如果一个表上定义了一个主键，那么 PostgreSQL 会自动在主键属性上创建唯一索引来实现主键约束。可以说，主键索引是唯一索引的特殊类型。

3) 多属性索引：如果一个索引定义在多于一个的属性上，就称其为多属性索引，它多用于组合查询。目前，PostgreSQL 中的 B-Tree、GiST 和 GIN 支持多属性索引，最多可在 32 个属性上创建索引[⊖]。虽然 PostgreSQL 提供了多属性索引的功能，但大多数情况下在单个属性上的索引就足够了。除非表的查询模式非常固定，否则超过三个属性的索引几乎没太多用处。当对一个表创建多属性索引时，对于表中的一个元组，会依次去读取出该元组被索引属性的值，使用这些值一起作为该元组的索引键值。多属性索引中不仅可以使用表中的属性，也可以是使用函数或表达式计算得到的值。

4) 部分索引：建立在一个表的子集上的索引，该子集由一个条件表达式定义（表达式即部分索引的谓词），该索引只包含表中那些满足这个谓词的元组。下面的语句在 student 表中，对 id 在“1”到“255”中的元组的 name 属性创建索引，这种索引就是部分索引。

```
CREATE INDEX stu_name_idx ON student(name) WHERE (id > 1 AND id < 255);
```

在 4.1.3 节讲解 pg_index 系统表时，会看到如果是部分索引，则会将其表达式保存在 indpred 属性中。使用部分索引，能够减小索引的规模，提高索引的查询效率。

5) 表达式索引：索引并非一定要建立在一个表的属性上，还可以建立在一个函数或者从表中一个或多个属性计算出来的标量表达式上。例如，可以在 student 表的 name 字段上通过小写函数来

[⊖] 见 PostgreSQL 8.4.1 文档。

创建表达式索引并进行查询：

```
创建:CREATE INDEX stu_low_name_idx ON student(lower(name));
查找:SELECT* FROM student WHERE lower(name) = 'jack';
```

表达式索引只有在查询时使用与创建时相同的表达式才会起作用。与部分索引类似，表达式索引的表达式是保存在 pg_index 系统表的 indexprs 属性中。在创建索引过程中，会根据表达式计算出实际索引的键值，这会导致插入或更新元组时效率变慢（在更新索引时会使用表达式进行重新计算）。

部分索引和表达式索引可能会让人觉得有些混淆。部分索引的特点是通过表达式来限制被索引元组的数量；而表达式索引则是用表达式来计算被索引元组的键值。两者的表达式的作用对象和用途都是不同的，在使用时一定要注意。

4.1.2 索引类型

PostgreSQL 8.4.1 中共有 4 种索引类型，下面先简单介绍这 4 种索引类型。

1) B-Tree：B-Tree 索引使用一种类似于 B+ 树的结构来存储数据的键值，通过这种结构能够快速地查找索引。B-Tree 索引适合支持比较查询及范围查询。在一个建立了 B-Tree 索引的属性涉及使用操作符（>、=、< 操作符）进行比较的时候，PostgreSQL 的查询优化器会考虑使用 B-Tree 索引进行查找。

2) Hash：Hash 索引使用 Hash 函数对索引的关键字进行散列。Hash 索引只能处理简单的等于比较。当一个建立了 Hash 索引的属性涉及使用“=” 操作符进行比较的时候，查询优化器会考虑使用 Hash 索引。

3) GiST：GiST (Generalized Search Tree) 意为通用搜索树。严格来说，GiST 索引不是一种独立的索引类型，而是一种架构或者索引模板，可以在这种架构（模板）上实现不同的索引策略。因此，可以使用 GiST 索引的操作符类型高度依赖于索引策略（操作符类）。在以前版本的 PostgreSQL 中还有 R-Tree 索引，但后来随着 GiST 的出现，取消了 R-Tree 索引，因为使用 GiST 索引的架构可以很容易地实现 R-Tree 索引。

4) GIN：GIN (Generalized Inverted Index) 索引是倒排索引，它可以处理包含多个键的值（比如数组）。与 GiST 类似，GIN 支持用户定义的索引策略，对于不同的索引策略，可以使用的操作符也是不同的。在后面将介绍的 Tsearch2 全文搜索，既可以通过 GiST 来实现，也可以通过 GIN 实现，两种实现各有特点。

4.1.3 索引相关系统表

为了管理各种索引类型，PostgreSQL 定义了相关系统表，这些系统表记录了索引相关的信息，下面将针对索引相关的几个重要的系统表进行分析。

每种索引类型都在系统表 pg_am (access method) 里面用一个元组来记录。pg_am 表中的每一个元组包括了该种索引类型提供的访问函数，这些函数是引用自 pg_proc 系统表中注册的函数（该系统表存储了关于函数的信息，每一个元组表示一个函数）。另外，pg_am 中的元组还记录了索引类型的一些特性，比如它是否支持多属性索引等。

目前 pg_am 中共有 4 个元组，分别对应 B-Tree、Hash、GiST 和 GIN 索引，pg_am 中包含的属性如表 4-1 所示。

表 4-1 pg_am 系统表

名称	类型	引用	说 明
amname	name		索引类型名称
amstrategies	int2		该索引类型支持的操作符类的数目
amsupport	int2		用于支持该种索引的程序数目
amcanorder	bool		该索引类型是否支持有序扫描
amcanbackward	bool		该索引类型是否支持逆向扫描
amcanunique	bool		该索引类型是否支持唯一索引
amcanmulticol	bool		该索引类型是否支持多字段索引
amoptionalkey	bool		该索引类型是否支持在第一个索引列上的没有条件的扫描
amindexnulls	bool		该索引类型是否支持空索引项
amsearchnulls	bool		该索引类型是否支持 IS NULL 查找
amstorage	bool		是否允许索引存储的数据类型与列的数据类型不同
amclusterable	bool		该索引类型是否支持聚簇
amkeytype	oid	pg_type.oid	索引存储的数据类型，无固定类型为 0
aminsert	regproc	pg_proc.oid	“插入索引元组” 所对应的函数
ambeginscan	regproc	pg_proc.oid	“开始新的扫描” 对应的函数
amgettuple	regproc	pg_proc.oid	“获取下一个有效元组”的函数，没有则为 0
amgetitemmap	regproc	pg_proc.oid	“获取所有有效元组”的函数，没有则为 0
amrescan	regproc	pg_proc.oid	“重新开始扫描”的函数
amendscan	regproc	pg_proc.oid	“结束扫描”的函数
ammarkpos	regproc	pg_proc.oid	“标记当前扫描位置”的函数
amrestrpos	regproc	pg_proc.oid	“恢复已标记的扫描位置”的函数
ambuild	regproc	pg_proc.oid	“创建索引”的函数
ambulkdelete	regproc	pg_proc.oid	“批量删除索引”的函数
amvacuumcleanup	regproc	pg_proc.oid	VACUUM 后的清理函数
amcostestimate	regproc	pg_proc.oid	估计一个索引扫描代价的函数
amoptions	regproc	pg_proc.oid	“分析确认该索引的 reoptions”的函数

对其中特别的几个字段予以说明：

- amsupport：该属性表示支持该类型的索引需要的程序数目，这里的“程序”也可以理解为“函数”。以 GiST 索引为例，在 4.4.2 节中将会分析实现该种索引需要实现的 7 个函数，那么 GiST 索引对应的“amsupport”属性的值就为 7。
- amoptionalkey：该属性表示索引类型是否允许在第一个被索引属性上进行无条件扫描。如果该索引类型的 amcanmulticol 属性值为假，那么 amoptionalkey 的值实际上表示该索引类型是否允许不带条件的扫描。
- amindexnulls：该属性表示索引类型是否支持对索引键为空的元组创建索引项。由于某些操作符非常严格，它们无法对于空值给出正确的判断，因此含有空值的索引项不可能被作为结果返回（因此也不会存储在索引中）。若设置了 amoptionalkey 为真，则必须索引空元组，因为查询编译模块可能会决定在没有扫描键值的时候使用索引，而无条件的索引扫描应该

返回所有的元组（包括含有空值的元组）。另外，当 amcanmulticol 为真时，amindexnulls 也应该为真。例如一个索引建立在 a、b 两个属性上，如果使用“a = 4”进行索引扫描，那么那些“a = 4”且 b 为空值的元组也应该被返回，因此必须允许索引为含空值的元组创建索引项。

- amoptions：该属性存储了一个函数的 OID，这个函数用于验证该索引的“reloptions”。“reloptions”是 pg_class 系统表中的一个属性，定义了数据库对象特定的选项。比如，创建 B-Tree 索引时，指定了 fillfactor = 80，则该索引对应的 pg_class 元组中的 reloptions 属性就会存储“fillfactor = 80”。而 B-Tree 索引类型对应的 pg_am 元组中的 amoptions 属性指定的函数就可以分析出填充因子为 80。

pg_am 系统表的最后 13 个属性记录了每种索引提供给其他模块调用的 13 个接口函数，在 4.1.4 节会详细介绍它们。

对于每一个创建的索引，会在 pg_class 系统表中添加一个元组，同时还会在 pg_index 系统表中添加一个元组。pg_index 用于记录与索引有关的信息，其包含的属性如表 4-2 所示。

表 4-2 pg_index 系统表

名称	类型	引用	说 明
indexrelid	oid	pg_class.oid	该索引的 OID
indrelid	oid	pg_class.oid	创建该索引的基表的 OID
indnatts	int2		索引中的属性数
indisunique	bool		是否是唯一索引
indisprimary	bool		索引的属性是否为表的主键
indisclustered	bool		若为真，则该表最后在该索引上建了聚簇
indisvalid	bool		该索引是否可以用于查询
indcheckxmin	bool		若为真，则在查询时需检查事务的 xmin 判断索引是否可用
indisready	bool		该索引是否可以进行插入/更新操作
indkey	int2vector	pg_attribute.attnum	该索引是在基表中哪些属性上创建的
indclass	oidvector	pg_opclass.oid	索引中每个属性对应的操作符 oid
indoption	int2vector		索引中每个属性的标志位
indexprs	text		索引条件的表达式树
indpred	text		部分索引的表达式树

其中几个重要的属性说明如下：

- indexrelid：表示该索引在 pg_class 里的 OID。
- indrelid：表示该索引依赖的基表的 OID。
- indisvalid：若为真，那么该索引可以用于查询；若为假，表示该索引可能不完整，需要在 INSERT/UPDATE 操作时进行更新。
- indcheckxmin：若为真，则在查询时，需要比较该元组的 xmin 值与查询事务的 xmin 值的大小。若前者比后者小，才能使用该索引进行查询；否则不能使用该索引。
- indkey：该属性为一个数组，记录了这个索引是在基表中的哪些属性上建立的。数组中的 0 值表示对应的索引属性是在表属性上的一个表达式。

- indoption：索引中的每个属性的标志位信息，具体含义由索引类型确定，是新增属性。目前对于所有未排序的索引该值都为 0，已排序的有 DESC 等值。
- indexprs：保存创建索引条件的表达式树。对于每一个索引属性，若索引条件不是简单的引用，则存在表达式树。该属性存储的表达式树主要用来在对索引进行插入更新操作时计算键值。
- indpred：若索引为部分索引，则该属性保存部分索引的表达式树，在 indexprs 中则不再保存。当对部分索引进行查找、插入更新时，需要使用该属性判断新增的元组是否需要添加到索引中。

每一种索引类型并不直接设定该类型的索引所要操作的数据类型（即对什么类型的数据创建索引），而是由操作符类系统表 pg_opclass（operator class）进行管理，该表表明了索引方法在操作特定数据类型的时候需要使用的操作集合。其包含的属性如表 4-3 所示。

表 4-3 pg_opclass 系统表

名称	类型	引用	说 明
opmethod	oid	pg_am.oid	该操作符类所服务的索引类型
opname	name		该操作符类的名称
opnamespace	oid	pg_namespace.oid	该操作符类的名字空间
opowner	oid	pg_authid.oid	操作符类的属主
opcfamily	oid	pg_opfamily.oid	包含该操作符类的操作符集合（operator family）
opcintype	oid	pg_type.oid	该操作符的输入数据类型
opcdefault	bool		如果操作符类是 opcintype 的缺省，则为真
opckeytype	oid	pg_type.oid	索引数据的类型，如果和 opcintype 相同则为零

其中，opcdefault 可以为同一个数据类型和索引类型定义多个操作符类。比如，一个 B-Tree 索引需要对所有的数据类型进行排序。对一个复数数据类型来说，可以定义对复数绝对值排序的操作符类，也可以定义根据复数的实部排序的操作符类等，通常其中的一个操作符类会被设定为缺省的。被设定为缺省的操作符类对应的 pg_opclass 元组的 opcdefault 属性将被设置为真。

在 pg_opclass 中，每一个操作符类都引用了系统表 pg_opfamily 中的一个元组，表明该操作符类的操作符集合。每个 pg_opfamily 元组定义了一个操作符的集合，在后面介绍 pg_amop 和 pg_amproc 系统表时，会看到对于每一个操作符或函数过程，它都会有一个 amopfamily 的属性表示该操作符所属的操作符集合，不同的操作符的 amopfamily 可能是相同的，这表明这些操作符属于同一个集合。对于属于同一个集合的操作符，它们是“兼容的”（还与索引类型有关），即允许对不同的数据类型进行操作。pg_opfamily 包含的属性如表 4-4 所示。

表 4-4 pg_opfamily 系统表

名称	类型	引用	说 明
opfmethod	oid	pg_am.oid	该操作符集合所服务的索引类型
opfname	name		操作符集合的名称
opfnamespace	oid	pg_namespace.oid	该操作符集合的名字空间
opfowner	oid	pg_authid.oid	操作符集合的属主

观察 pg_opclass 和 pg_opfamily 可以发现，这两个系统表都没有指定索引方法对各种数据类型的具体操作函数，这些信息保存在相关的系统表 pg_amop 和 pg_amproc 中。pg_amop 存储每个索引操作符集合（pg_opfamily）与具体的操作符（pg_operator）的关联信息，若 pg_operator 中的一个操作符属于 pg_opfamily 中的一个集合，则在 pg_amop 表中就有一条元组来记录这个对应关系。同一个操作符可能属于多个操作符集合。pg_amop 包含的属性如表 4-5 所示。

表 4-5 pg_amop 系统表

名称	类型	引用	说 明
amopfamily	oid	pg_opfamily.oid	该操作符所属的操作符集合
amoplefttype	oid	pg_type.oid	该操作符的左输入数据类型
amoprighttype	oid	pg_type.oid	该操作符的右输入数据类型
amopstrategy	int2		操作符策略号
amopopr	oid	pg_operator.oid	该操作符的 oid
amopmethod	oid	pg_am.oid	该操作符所服务的索引类型

与 pg_amop 类似，pg_amproc 存储了每个索引操作符集合（pg_opfamily 表）与其相关联的支持过程或函数（pg_proc）的关联信息，如果 pg_proc 中的一个函数属于 pg_opfamily 中的一个操作符集合，则在 pg_amproc 中就有一个元组来记录这个对应关系。pg_amproc 包含的属性如表 4-6 所示。

表 4-6 pg_amproc 系统表

名称	类型	引用	说 明
amprocfamily	oid	pg_opfamily.oid	该过程（函数）属于的操作符集合
amproclefttype	oid	pg_type.oid	该过程（函数）的左输入数据类型
amprocrighttype	oid	pg_type.oid	该过程（函数）的右输入数据类型
amprocnum	int2		支持该过程（函数）的编号
amproc	regproc	pg_proc.oid	该过程（函数）的 OID

4.1.4 索引的操作函数

在 PostgreSQL 中，每一种索引类型都在 pg_am 中注册了操作函数，不同的索引类型的操作函数数目并不相同，最多可以有 13 个操作函数。为了方便后面介绍，我们用存储每一种操作函数的 pg_am 元组的属性名称来指代对应的操作函数。

每一种索引类型可以实现 ambuild、aminsert、ambulkdelete、amvacuumcleanup、amcostestimate、amoptions、ambeginscan、amgettuple、amgetbitmap、amrescan、amendscan、ammakpos 和 amrestrpos 共 13 个函数，会有相应的 13 个函数来实现这些接口。注意，这 13 个接口并不要求每一种索引类型全部实现，每一种索引类型可以选择地实现其中的某些接口。当其他模块需要使用索引时，将会根据索引的类型找到 pg_am 中对应的元组，然后在该元组中找到适当的操作函数来完成所需的操作。

下面将分别介绍这些函数的功能。

1) ambuild 函数用于创建一个新索引。在调用 ambuild 之前，索引文件已经在物理上创建好了，但是是空的，还需要使用索引元组进行填充。ambuild 的工作就是生成索引元组并将它们填充到索引文件中。通常，ambuild 函数会调用 IndexBuildHeapScan 函数扫描基表以获取表中的元组并以之为

基础构建索引元组。但 ambuild 是否会将索引元组填充到索引文件中需要根据具体的索引类型来确定，比如 B-Tree 索引会在其 ambuild 函数中完成填充工作。

2) aminsert 函数向现有索引中插入一个新的索引元组。如果该索引要求是唯一索引，则在插入过程中会检查该索引元组是否已经存在。函数返回值标识插入是否成功。

3) ambulkdelete 函数从索引中删除元组。该函数所作的删除可能是标记删除，也可能是物理上的删除。ambulkdelete 是一个“大批删除”的操作，通常都会扫描整个索引，检查每个元组是否需要被删除。一般调用该函数时，会给它传递一个回调函数。回调函数负责判断索引元组是否需要删除。ambulkdelete 函数可能会调用 amvacuumcleanup 函数来执行实际删除操作。

4) amvacuumcleanup 函数会在一个 VACUUM 操作（一个或多个 ambulkdelete 调用）之后调用，主要做一些额外的清理工作，比如对 FSM 文件进行清理。该函数通常用于批量删除中。

5) amcostestimate 函数用于估算一个索引扫描的代价。在最简单的情况下，该函数实现的所有功能都可以通过调用优化器里面的标准过程完成。该函数存在的目的是允许索引访问方法提供与索引类型相关的信息，这样可以改进标准的代价估计。

6) amooptions 函数用于分析和验证一个索引的 reloptions 数组，仅当一个索引存在非空 reloptions 数组时才会被调用。reloptions 是一个 text 类型的数组，每个元素的类型为：“name = value”（其中 name 是参数名，value 是参数值），可以参见 pg_am 系统表的 amooptions 字段。amooptions 的作用就是通过分析这些元素来验证索引的各种参数。

创建索引的目的当然是支持那些包含可以使用索引的带 WHERE 条件的查询操作（扫描），之后这些函数是一个索引访问方法提供的与扫描有关的函数。

7) ambeginscan 函数开始一个新的扫描。该函数主要功能是构造索引扫描描述符结构 IndexScanDescData，然后用该描述符执行扫描。索引访问方法必须通过调用 ambeginscan 函数来创建这个结构，然后调用 amrescan 函数来执行扫描。在大多数情况下，ambeginscan 本身除了调用 RelationGetIndexScan 函数之外几乎不干别的事情。

IndexScanDescData 结构（数据结构 4.1）是索引通用的扫描描述符，记录了使用索引进行扫描的过程中的状态信息。

数据结构 4.1 IndexScanDescData

```
typedef struct IndexScanDescData
{
    Relation      heapRelation;           //当前扫描的基表的描述符,或为 NULL
    Relation      indexRelation;          //索引描述符,
    Snapshot      xs_snapshot;            //用于可见性判断的快照
    int           numberofKeys;           //扫描的关键字的个数
    ScanKey       keyData;                //有关扫描关键字信息的数组
    bool          kill_prior_tuple;       //上一个返回的元组是否无效
    bool          ignore_killed_tuples;  //是否返回被 killed 的元组
    void          *opaque;                //进行扫描时,记录具扫描状态,由具体索引类型指定
    HeapTupleData xs_ctup;               //扫描成功后得到的元组
    Buffer        xs_cbuf;                //扫描成功后得到的元组所在的缓冲区
}
```

```

    bool          xs_recheck;           //若为真,说明扫描关键字必须进行核查
    //在 index_getnext 中,HOT 链中的一些状态
    bool          xs_hot_dead;         //若为真,说明 HOT 链中所有的成员都 dead
    OffsetNumber  xs_next_hot;        //HOT(Heap Only Tuples)链中的下一个成员的偏移量
    TransactionId xs_prev_xmax;      //HOT 链中前一个 HOT 成员的 XMAX 值
} IndexScanDescData;

```

8) amgettuple 函数在给出的扫描描述符里获取下一个元组,向给出的方向移动(在索引里向前或者向后)。如果获取成功会将该元组的 TID 值返回。

9) amgetbitmap 函数在给出的扫描描述符里获取所有可用的元组。返回的元组通过一个 bitmap(位图)的方式来表示。

10) amrescan 函数用于重启一个扫描,通过该函数可以使用一个新的扫描关键字。实际上,函数 RelationGetIndexScan 中也是调用 amrescan 来完成扫描描述符的创建的。因此 amrescan 既可用于索引扫描的初始化扫描,也用于重复扫描。

11) amendscan 函数结束扫描并释放资源,扫描过程中使用的任何 lock 锁和 pin 锁都应该释放。

12) ammarkpos 函数标记当前扫描位置,事实上该函数会将当前的扫描位置记载在扫描描述符的 opaque 字段中。

13) amrestrpos 函数把扫描恢复到最近标记的位置。

为了给上层模块提供一个使用索引的统一接口,PostgreSQL 在 indexam.c 和 index.c 文件中提供了一套用于操纵索引的函数,这些函数和前面介绍的 13 个接口函数一一对应(见表 4-7)。

表 4-7 索引上层操作函数和下层接口函数的对照

上层操作函数	下层接口函数	上层操作函数	下层接口函数
index_insert	amininsert	index_restrpos	amrestrpos
index_beginscan	ambeginscan	index_build	ambuild
index_getnext	amgettuple	index_bulk_delete	ambulkdelete
index_getbitmap	amgetbitmap	index_vacuum_cleanup	amvacuumcleanup
index_rescan	amrescan	index_costestimate	amcostestimate
index_endscan	amendscan	index_reoptions	amoptions
index_markpos	ammarkpos		

这些 index_* 函数都是通过系统表信息取得并调用对应的下层接口函数来完成各自的工作。由于下层接口函数只负责对索引文件本身的操作,因此在上层操作函数中除了调用下层之外,还需要做一些额外的工作,例如 index_build 函数中还需要对索引在系统表中的一些状态信息进行更新。

通过这一套统一的操作函数,上层模块可以很方便地使用索引,而每一种索引的实现则交由索引各自实现,下面我们将对 PostgreSQL 中的四种索引类型分别进行介绍。

4.2 B-Tree 索引

在 PostgreSQL 中的 B-Tree 索引是在 Lehman 和 Yao 的论文“Efficient Locking for Concurrent Operations on B-Trees[⊖]”的基础上,结合 PostgreSQL 的应用要求来实现的,其结构类似于 B+ 树。

⊖ <http://portal.acm.org/citation.cfm?id=319663>

首先简单回顾一下 B + 树。B + 树为 B 树的变形，它的结构特点如下：

- B + 树的所有关键字均在叶子节点中出现，并按关键字排序以顺序链的方式链接，同时，叶子节点还保存了指向相应记录的指针。
- 所有非叶子节点可看成是索引部分，并不指向实际的存储位置；非叶子节点中仅仅包含其子树节点中最小的关键字。

在 B + 树上进行随机查找、插入和删除的过程基本上与 B 树类似。只是在查找时，如果非叶子节点上的关键字等于查找值，查找过程也不终止，而是继续向下直到叶子节点。因此，对于 B + 树来说，不管查找成功与否，每次查找都是经过了一条从根到叶子节点的路径。一棵 B + 树的结构如图 4-1 所示。

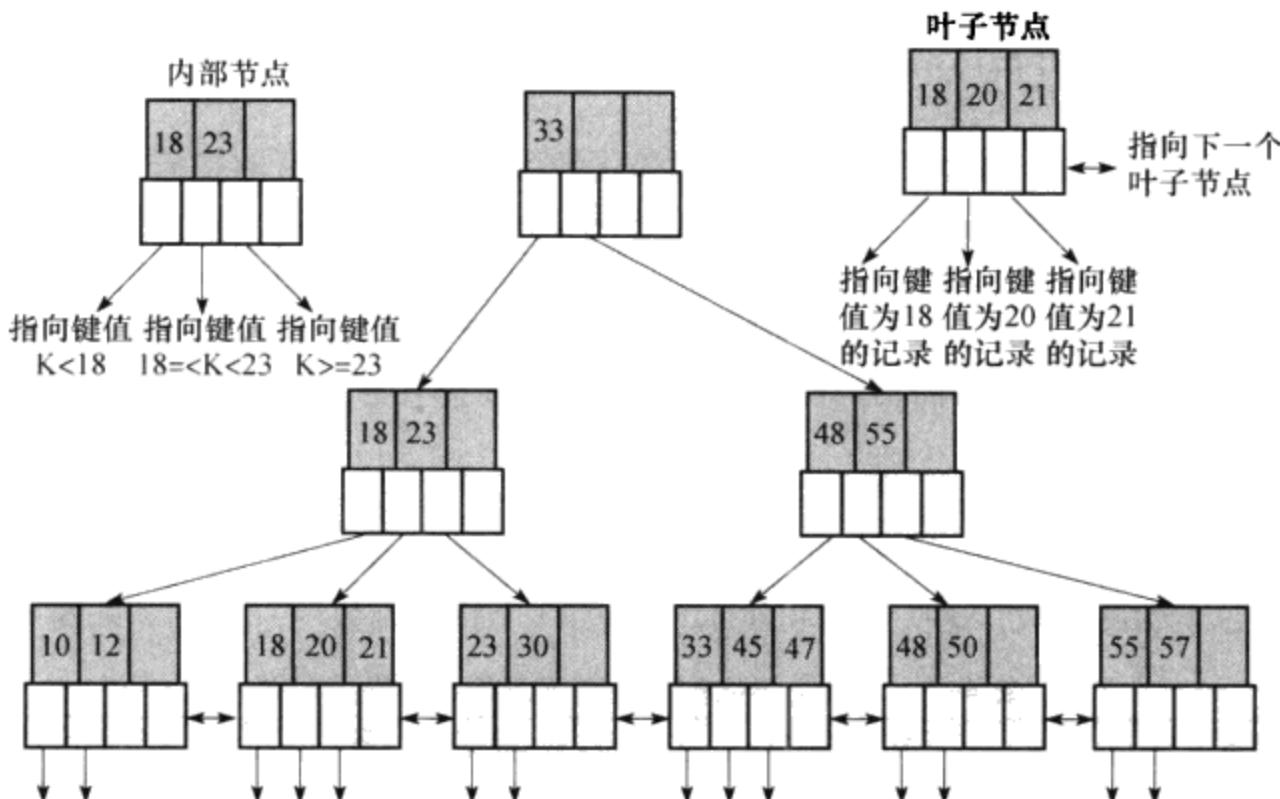


图 4-1 B + 树结构图

在 Lehman 和 Yao 的论文中，修改了 B 树的结构，成为 B^{link}-Tree，不管是内部节点还是叶子节点，都有一个指针指向兄弟节点，其节点的结构如图 4-2 所示。

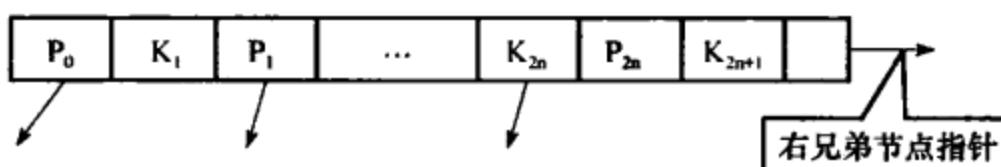


图 4-2 B^{link}-Tree 节点

在图 4-2 中，对于指针 P_i ，其指向的子节点包含的键值的范围为 $(K_{i-1}, K_i]$ 。B^{link}-Tree 对除每层最右节点外的其他节点，都设置了一个“High-key”（节点中的 K_{2n+1} ），该值为该节点及该节点的所有子节点中的最大值，“High-key”并不作为索引结构中的一个元组，而是标记了一个最大范围。根据 Lehman 和 Yao 的论文，设置“High-key”在对节点进行并发的插入、查找时会用到，在后面分析对 B-Tree 索引的操作时会进行介绍。

4.2.1 B-Tree 索引的组织结构

在 PostgreSQL 数据库中，索引也是按照图 3-4 中的页面结构进行存储，B-Tree 索引的页面组织

结构如图 4-3 所示。

图 4-3 中的 $itup_1$ 、 $itup_2$ 、 $itup_3$ 等都是索引元组，它们是有序的。在页面的“Page header”结构之后，是 $linp_1$ 、 $linp_2$ …，它们存储了各个索引元组在页面内的实际位置，通过 $linp$ 可以快速访问到索引元组。

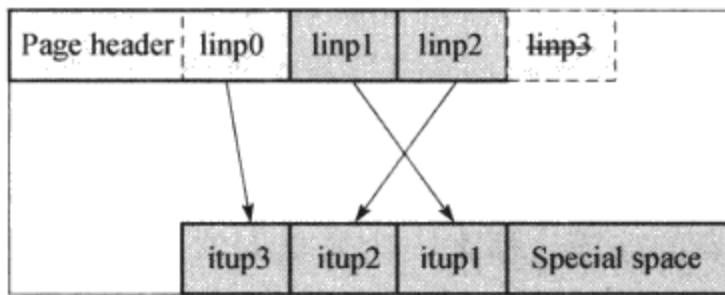
根据 B^{link} -Tree 的要求，在每层非最右节点中，需要一个最大关键字（High-key），High-key 并非一个真正的关键字，它只不过是给出该节点中的索引关键字的范围，在节点中的索引关键字都应小于（等于）其 High-key。所以，如果要插入的关键字大于 High-key，就需要向右边的节点移动来寻找合适的插入节点。当构建一个索引页面完成后，需要设置该值。图 4-3 可以理解为在索引创建过程中的结构，当填充完成一个页面后，会进行调整。最右节点和非最右节点的结构是不一样的。每一个页面的“Page header”结构中保存了 $linp_0$ （如图 4-3，图中的 $linp_0$ 实际是包含在 page header 中的，为了描述方便，这里把它单独画出来），在填充页面过程中， $linp_0$ 并没有赋值。当页面填充完成后，根据节点类型进行以下两种不同的操作（这里假设一个页面填充了 3 个元组即不再允许插入新的元组，此时页面中有 $linp_1$ 、 $linp_2$ 、 $linp_3$ 分别指向 $itup_1$ 、 $itup_2$ 、 $itup_3$ ）：

1) 若该节点不为本层最右节点：

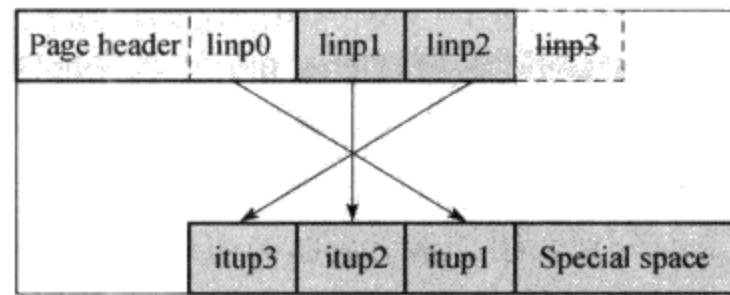
- 首先将 $itup_3$ （节点中值最大的索引元组）复制到该节点的右兄弟节点中，然后将 $linp_0$ 指向 $itup_3$ （页面中的 High-key）。
- 然后去掉 $linp_3$ 。也就是说，使用 $linp_0$ 来指向页面中的 High-key，由于 High-key ($linp_0$) 只是作为一个索引节点中键值的范围，并不是指向实际元组 ($itup_3$)，所以去掉指向 $itup_3$ 的链接 $linp_3$ 。

2) 若该节点是本层的最右节点：由于最右节点不需要 High-key，所以 $linp_0$ 不需要用于保存 High-key，则将所有的 $linp$ 都递减一个位置， $linp_3$ 同样不再使用。

这两种情况下会采取不同操作，在 4.2.2 节讲解索引创建时会用到。对于这两种情况，操作完成后的页面组织如图 4-4 所示。



a)



b)

图 4-4 完成操作后的 B-Tree 索引页面结构图

a) 非最右节点；b) 最右节点

按照图 4-2 的要求，在每个节点都有一个指针指向其右侧的兄弟节点，而 PostgreSQL 在实现时，使用了两个指针，分别指向左右兄弟节点。这两个指针是由页面尾部的一块称为 Special 的特殊区域保存的，其中存放了一个由 BTPageOpaqueData 结构表示的数据，该结构记录了该节点在树结构中的左右兄弟节点的指针以及页面类型等信息。BTPageOpaqueData 如数据结构 4.2 所示。

数据结构 4.2 BTPageOpaqueData

```

typedef struct BTPageOpaqueData
{
    BlockNumber          btpo_prev;      //前一页块号,用于索引反向扫描
    BlockNumber          btpo_next;     //后一页块号,用于索引正向扫描
    union
    {
        uint32           level;         //页面在索引树中的层次,0 表示叶子层
        TransactionId    xact;         //删除页面的事务 ID,用于判断该页面是否可以被重新分配使用
    }btpo;
    uint16              btpo_flags;    //页面类型
    BTCycleId          btpo_cycleid; //页面对应的最新的 Vacuum cycle ID
}BTPageOpaqueData;

```

其中成员变量 `btpo_flags` 用于记录有关页面类型的标志，包括以下几种标志：

- `BTP_LEAF`: 叶子页面，没有该标志则表示非叶子页面。
- `BTP_ROOT`: 根页面（根页面没有父节点）。
- `BTP_DELETED`: 页面已经从树中删除。
- `BTP_META`: 元页面。
- `BTP_HALF_DEAD`: 空页面，但是还保留在树中。
- `BTP_SPLIT_END`: 在一次页面的分裂中，待分裂的最右一个页面。
- `BTP_HAS_GARBAGE`: 页面中含有 `LP_DEAD` 元组。当对索引页面中某些元组进行了删除后，该索引页面并没有立即从物理上删除这些元组，这些元组仍然保留在索引页面中，只是对这些元组进行了标记，同时索引页面中其他有效的元组保持不变。

若页面发生分裂，`btpo_cycleid` 字段记录了当前页面最新的“Vacuum cycle ID”，该值可用于确定哪些页面需要进行 VACUUM。比如，现在有页面“`A ←→ B`”，当页面 `A` 发生分裂后，得到“`A ←→ C ←→ B`”，此时页面 `A` 和 `C` 中的元组发生了更改，其 `btpo_cycleid` 字段会重新分配值，而与 `B` 的值不同，当运行 VACUUM 时，由于 `B` 页面中的元组并未改变，因此不需要进行回收，这就可以通过比较页面的 `btpo_cycleid` 字段来实现。

上面分别介绍了索引页面的内部结构和兄弟节点的连接，那么一个完整的索引组织结构如图 4-5 所示。可以看到，左侧叶子节点的 `itup3`（该节点的“High-key”）和其右侧兄弟节点的 `itup1`（该节点的最小值）都指向相同的堆元组，这正是在上面介绍“High-key”时的操作决定的。

在图 4-5 中，虚线上方的表示索引结构，虚线下方的为表元组。在叶子节点层，索引元组的指针指向表元组。每一个索引节点对应一个索引页面，内部节点（包括根节点）与叶子节点的内部结构是一致的，不同的是内部节点指向下一层的指针是指向索引节点，而叶子节点是指向物理存储的某个位置（也就是实际存放元组的位置）。

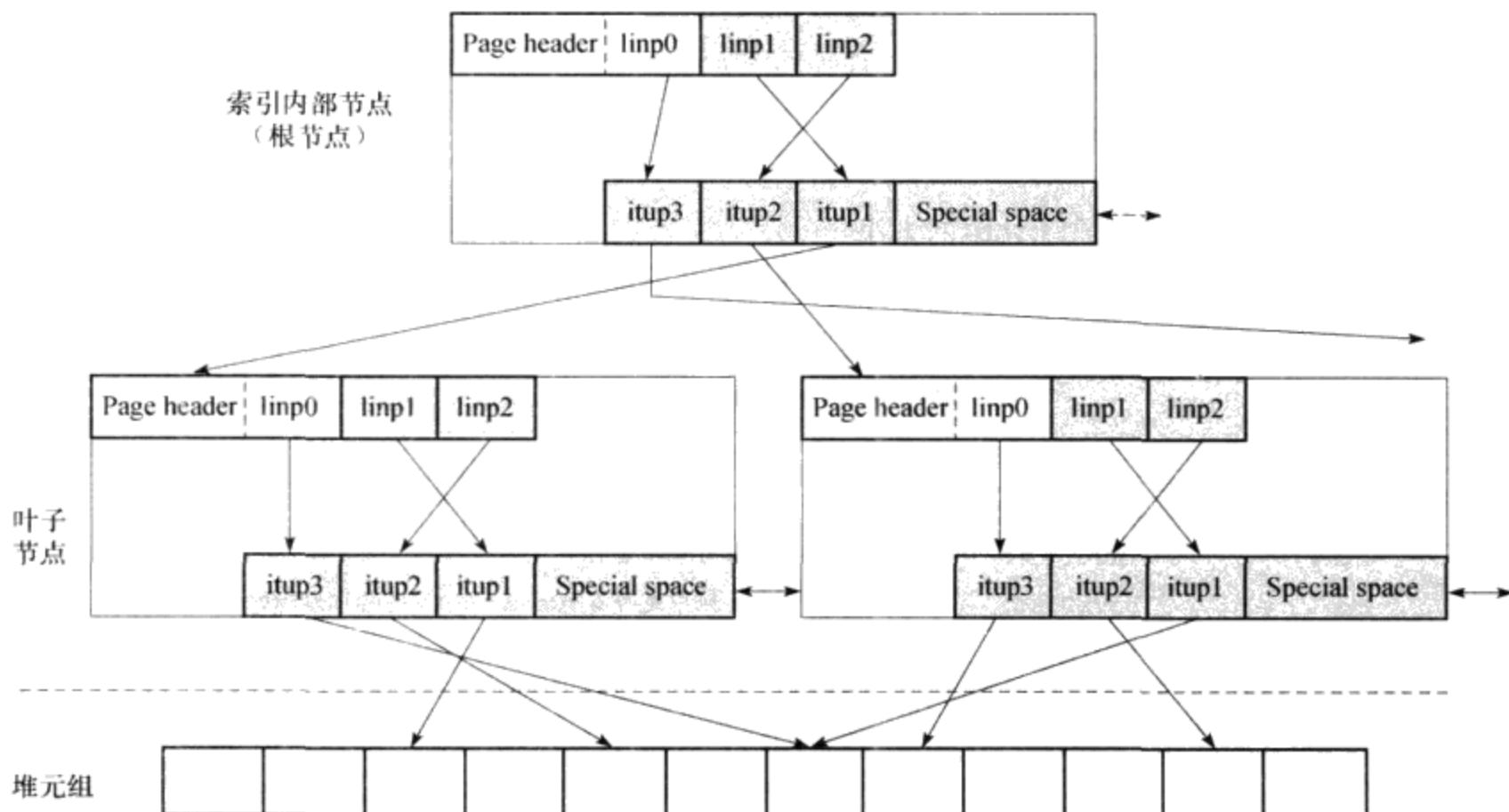


图 4-5 B-Tree 索引的组织结构

4.2.2 B-Tree 索引的操作

按照 pg_am 系统表的规定，B-Tree 索引所对应的统一接口调用的函数为：btbuild、btinsert、btbeginscan、btgettupple、btgetbitmap、btrescan、btendscan、btmarkpos、btrestrpos、btbulkdelete、btvacuumcleanup、btcostestimate 和 btoptions。这些函数都在 src\backend\access\nbtree 目录中实现，nbtree 目录中的其他函数都直接或间接为以上 13 个函数服务，下面将选取其中一些重要的函数进行详细分析，对其他的函数只作功能上的简单说明。

1. 索引的创建

PostgreSQL 系统创建 B-Tree 索引时，首先将对每一个需要索引的表元组生成对应的索引元组，然后调用 tuplesort 函数对所有的索引元组进行排序，最后创建索引。

索引元组是一个索引结构的基本单位，由 IndexTupleData 表示（数据结构 4.3）。B-Tree 索引和 Hash 索引的索引元组结构都是一样的，都是由 IndexTupleData 进行存储，该结构包括了该索引元组所指向的表元组以及一些索引元组的信息。

数据结构 4.3 IndexTupleData

```

typedef struct IndexTupleData
{
    ItemPointerData t_tid;          // 标记该索引元组指向的表元组的 ID
    unsigned short t_info;          // 关于索引元组的一些信息
} IndexTupleData;
typedef IndexTupleData* IndexTuple;

```

其中 t_info 内各位表示的信息如下：

- 第 15 位：该元组是否为 null。
- 第 14 位：是否有可变长度的属性。
- 第 13 位：未使用。
- 第 12 位至第 0 位：索引元组的大小。

在将表元组封装成索引元组（索引项）的过程中，会生成一个 BTBuildState 结构用于保存索引元组，其结构如数据结构 4.4 所示。

数据结构 4.4 BTBuildState

```
typedef struct
{
    bool          isUnique;      // 该索引是否为唯一索引
    bool          haveDead;     // 扫描的表元组中是否有“dead tuple”(标记删除的元组)
    Relation      heapRel;      // 索引的基表
    BTSpool       *spool;       // 保存生成的索引元组
    BTSpool       *spool2;      // 保存“dead tuple”对应的索引元组
    double        indtuples;   // 总共得到的索引元组个数(spool 和 spool2 之和)
}BTBuildState;
```

BTBuildState 结构保存了两个 BTSpool 类型：spool 和 spool2。spool2 只是在创建唯一索引时才需要使用。由于从 PostgreSQL 8.3 开始使用 HOT 链，在扫描表元组时，获取到的元组可能是“dead tuple”（实际更新后的表元组可以通过 HOT 链找到），对于这样的表元组，会将封装得到的索引元组放在 spool2 里。由于 spool2 里面的索引元组都不是最新的，那么对于唯一索引，就不需要检查这些元组是否唯一，可以直接插入到索引结构中。

每个 B-Tree 索引都有一个元页（metadataPage），它主要说明该 B-Tree 的版本、根节点位置（块号）以及根节点在树中的层次等信息，元页始终位于 B-Tree 索引的第一页（编号为 0）。

由于元页相关的信息只有在索引创建完成后才能够获得，因此元页会预留，一直等到索引创建完成后，才会生成元页并赋值。

元页用 BTMetaPageData 结构表示，该结构主要是用来保存 B-Tree 的根节点（root）以及有效根节点（fastroot）的相关信息，主要包括根节点、有效根节点所在的磁盘块号以及有效根节点所在的层次。其结构如数据结构 4.5 所示。

数据结构 4.5 BTMetaPageData

```
typedef struct BTMetaPageData
{
    uint32        btm_magic;      // btree 的 magic 号
    uint32        btm_version;   // 版本号信息
    BlockNumber   btm_root;      // 根节点对应的块号
    uint32        btm_level;     // 根节点在树中的层号
    BlockNumber   btm_fastroot;  // fastroot 对应的块号
}
```

```

    uint32          btm_fastlevel;           // fastroot 对应的层号
} BTMetaPageData;

```

其中：

- magic 号标识是 B-Tree 的一个编号，这个值在所有的 B-Tree 里都是一样的，用于完整性检查和辅助调试，确定该结构确实是 B-Tree 的元页。
- fastroot 是有效根节点。由于大量的删除操作可能导致根节点下面出现单节点层（所谓单节点是指从根节点到子节点之间其实只有唯一的一条路径），用 fastroot 记录索引中最底层的单节点层（叶子层为最底层），对 B-Tree 的操作（查找、插入等）均从 fastroot 开始，这样可提升效率。下面以一个简单的 B-Tree 结构来说明 root 和 fastroot，如图 4-6 所示。root 为实际的根节点，但 root 只有一个子节点。所以对该树进行查找时，可直接从 fastroot 节点开始查找，从而加快速度。

创建一个 B-Tree 索引，首先会生成一个 BTWriteState 结构，它用于记录整个索引创建过程中的信息，如数据结构 4.6 所示。

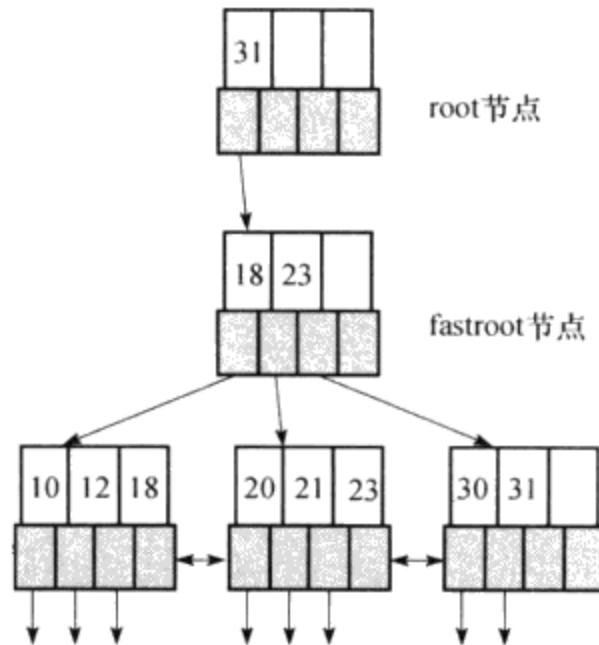


图 4-6 root 和 fastroot 节点

数据结构 4.6 BTWriteState

```

typedef struct BTWriteState
{
    Relation      index;           // 正在创建的索引关系
    bool         btws_use_wal;     // 是否开启了 WAL，若开启则在每次更新 BTWriteState 时会
                                  // 写 LOG
    BlockNumber   btws_pages_allocated; // 索引下一次申请 page 时的块号
    BlockNumber   btws_pages_written; // 已经写入文件的页面块号
    Page         btws_zeropage;     // 需要填充“0”值的页面
} BTWriteState;

```

其中一些字段说明如下：

- btws_pages_allocated**：在创建索引的过程中，随着不断添加索引元组，索引关系所占用的页面也会不断增加。根据分配的先后顺序，索引结构的页面编号（块号）依次递增，btws_pages_allocated 字段记录的正是分配给下一次申请的 page 的块号。
- btws_pages_written**：在创建索引时，若一个页面（节点）已填充满，后面的元组则会填充到新申请的页面中，这时就会把填满的页面写入文件中，该页面中的内容就不再改变，btws_pages_written 字段正是记录了当前已写入文件的块号。
- btws_zeropage**：若当前需要写入的页面块号大于 btws_pages_written，则位于二者中间的页面需要填充为“0”。可能的一种情况是：由于元页都位于索引结构中的第一页（编号为 0），但元页是

在整个索引结构创建完成后才写入的（预留），所以最开始分配给索引的块号是 1，当去写块号为 1 的页面到文件中时，会发现块号为 0 的页面还没有写入，就会对 0 号页面填充“0”值。在创建 B-Tree 索引时，对于树结构中的每一层都会生成一个结构 BTPageState，如数据结构 4.7 所示。

数据结构 4.7 BTPageState

```
typedef struct BTPageState
{
    Page                btps_page;           // 当前进行插入操作的页面
    BlockNumber         btps_blkno;          // 面对应的磁盘块号
    IndexTuple          btps_minkey;         // 页面中最小元组的拷贝
    OffsetNumber        btps_lastoff;        // 页面中最后插入索引元组的位移
    uint32              btps_level;          // 页面所在树结构中的层号，0 表示叶子层
    Size                btps_full;           // 页面允许的最小空闲空间，与 fillfactor 有关
    struct BTPageState * btps_next;          // 指向其父节点的对应的结构
} BTPageState;
```

在 BTPageState 结构中，使用 btps_next 指针来指向父节点的目的是：当该节点中关键字变化导致需要调整父节点中的关键字时，可通过该指针快速定位到父节点。当所有索引元组都插入到索引结构中后，需要调整每一层最右节点，就会用到这个指针，在后面讲解函数_bt_uppershutdown 时会分析。

在创建索引的过程中，对于每一个层次的所有页面，只有一个 BTPageState 结构。当一个页面填充满后，会申请一个新的页面，这时 BTPageState 结构（page 信息、minkey 等）随即更新为新页面中的信息。

上面介绍了在创建索引过程中使用到的主要数据结构，下面将使用这些结构来完成对索引的创建。

首先将待索引的表元组封装为索引元组，并对索引元组进行排序。然后将排好序的索引元组填充到当前叶子节点中，若当前叶子节点填充满，则新申请一个叶子节点作为当前叶子节点的右兄弟节点，然后在父节点中添加指向新叶子节点的指针（没有父节点则创建父节点），接下来把新申请的节点作为当前叶子节点。重复这个过程，直到添加完成所有的索引元组为止。

创建索引的入口函数是 btbuild，其流程如图 4-7 所示。

btbuild 函数首先调用 IndexBuildHeapScan 函数对表进行扫描，将每一个表元组都封装成索引元组，在扫描的过程中判断表元组是否为“dead tuple”，然后将得到的索引元组插入到不同的 spool 中。

执行完对表元组的扫描后，即调用_bt_leafbuild 函数。该函数首先定义创建索引结构需要的 BTWriteState 结构，然后调用排序函数对 spool（和 spool2）中的索引元组进行排序。若为唯一索引，在排序时会进行检查是否有重复的元组。排序完成后，会调用_bt_load 函数顺序读取出 spool 中的元组（已经排序）。对每个索引元组调用_bt_buildadd 函数（将索引元组添加到索引结构中，会处理分裂等情况）。若 spool2 不为空，则在添加索引元组到索引结构中时会使用归并的方法：按照索引元组的偏序关系，从小到大把 spool 和 spool2 中的索引元组添加到索引结构中。当所有索引元组添加完成后，调用_bt_uppershutdown 函数完善索引结构，并写入元页的信息（在前面讲到，元页信息是在索引结构创建完成后才写入的）。总的来说，_bt_leafbuild 函数扫描有序的索引元组，并构建出索引的树结构，其流程如图 4-8 所示。

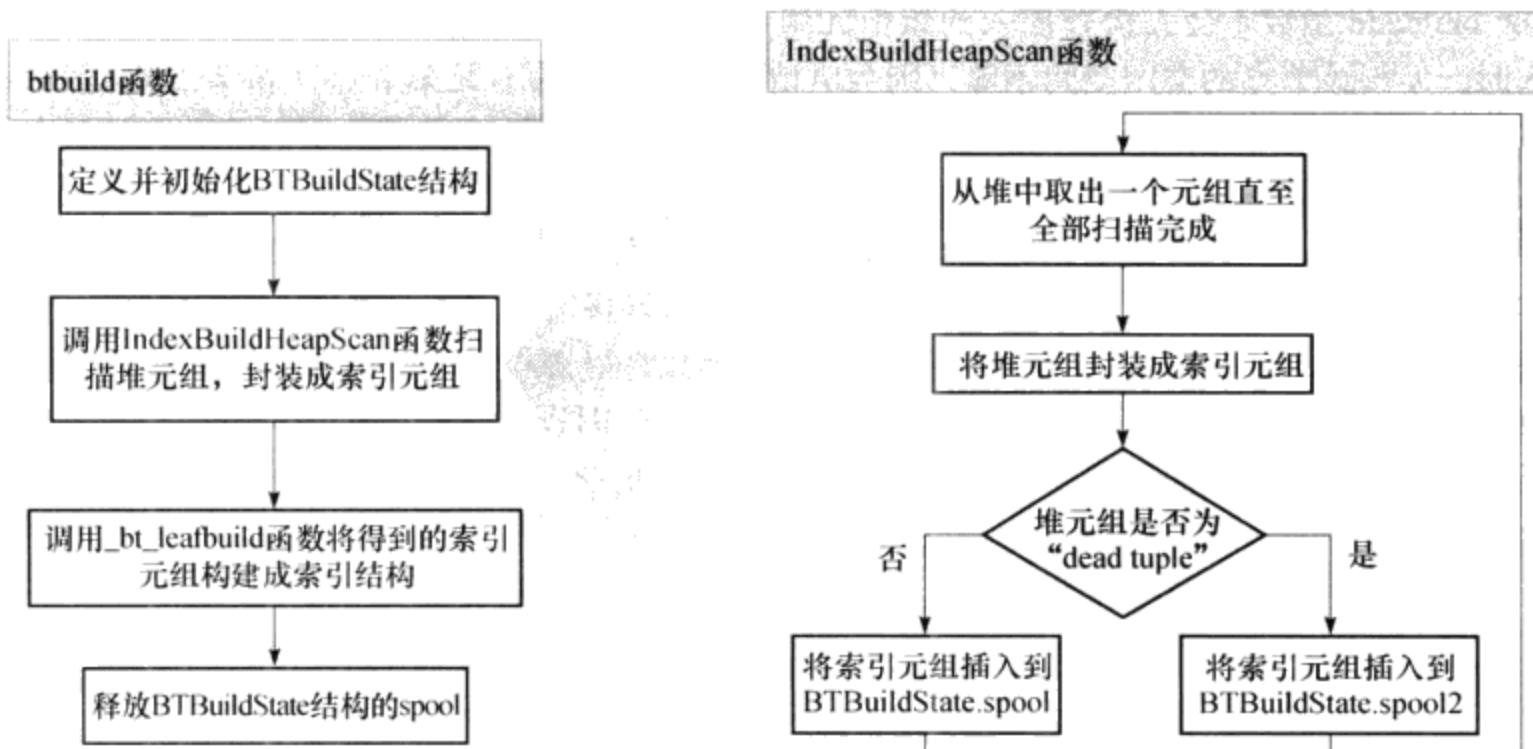


图 4-7 B-Tree 索引创建流程

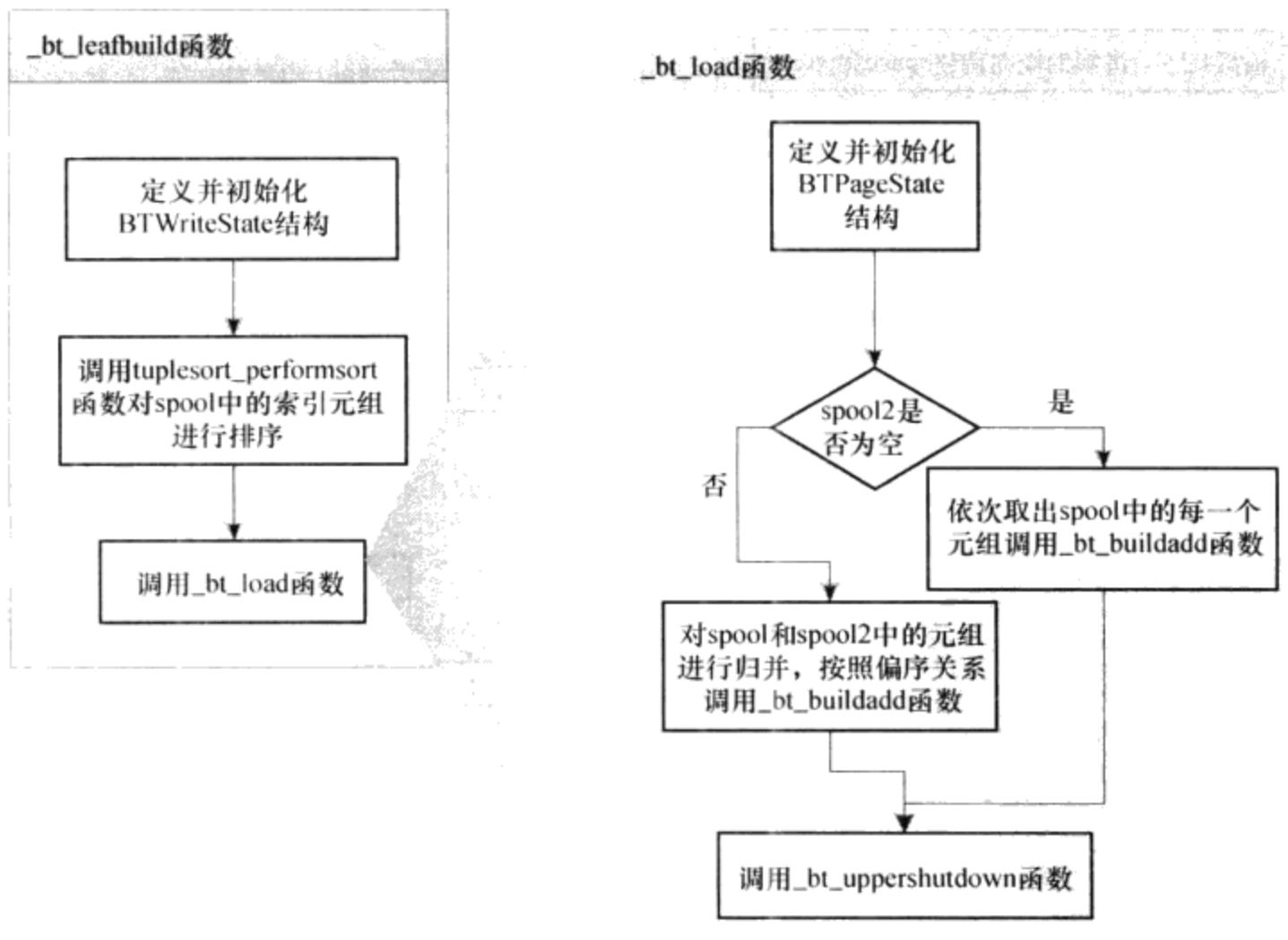


图 4-8 生成索引结构 (_bt_leafbuild 函数) 流程图

在调用_bt_buildadd 函数时，spool (spool2) 中所有的索引元组都已经有序，则依次取出 spool 中的索引元组，将其添加到索引结构中，_bt_buildadd 函数正是完成将一个索引元组插入到索引结构的工作。

`_bt_buildadd` 函数不会去检查插入的索引元组与节点中已有元组的偏序关系，而是直接插入。若插入时，发现该节点的空闲空间不够，则会申请一个页面作为右兄弟节点（新节点），然后设置旧节点的“High-key”。接着通过当前层次的 `BTPageState. btps_next` 查找其父节点（若没有父节点则创建一个），将旧节点的 `min key`（最小键值）插入到父节点中（还是调用 `_bt_buildadd` 函数，这种递归调用实现了向上层调整索引结构），然后设置新旧节点的兄弟关系（左右指针链接）。由于旧节点已完成填充不会再修改，因此调用 `_bt_bwritepage` 函数将旧节点的信息写入索引文件，同时修改索引的 `BTWriteState` 结构。最后调用 `_bt_sortaddtup` 函数将待插入的索引元组插入节点（若一开始检查节点有足够的空间则直接跳到这一步），并更新该层的 `BTPageState` 结构（若申请了新页面，则以后插入节点都是对新页面进行操作）。`_bt_buildadd` 函数的执行流程如图 4-9 所示。

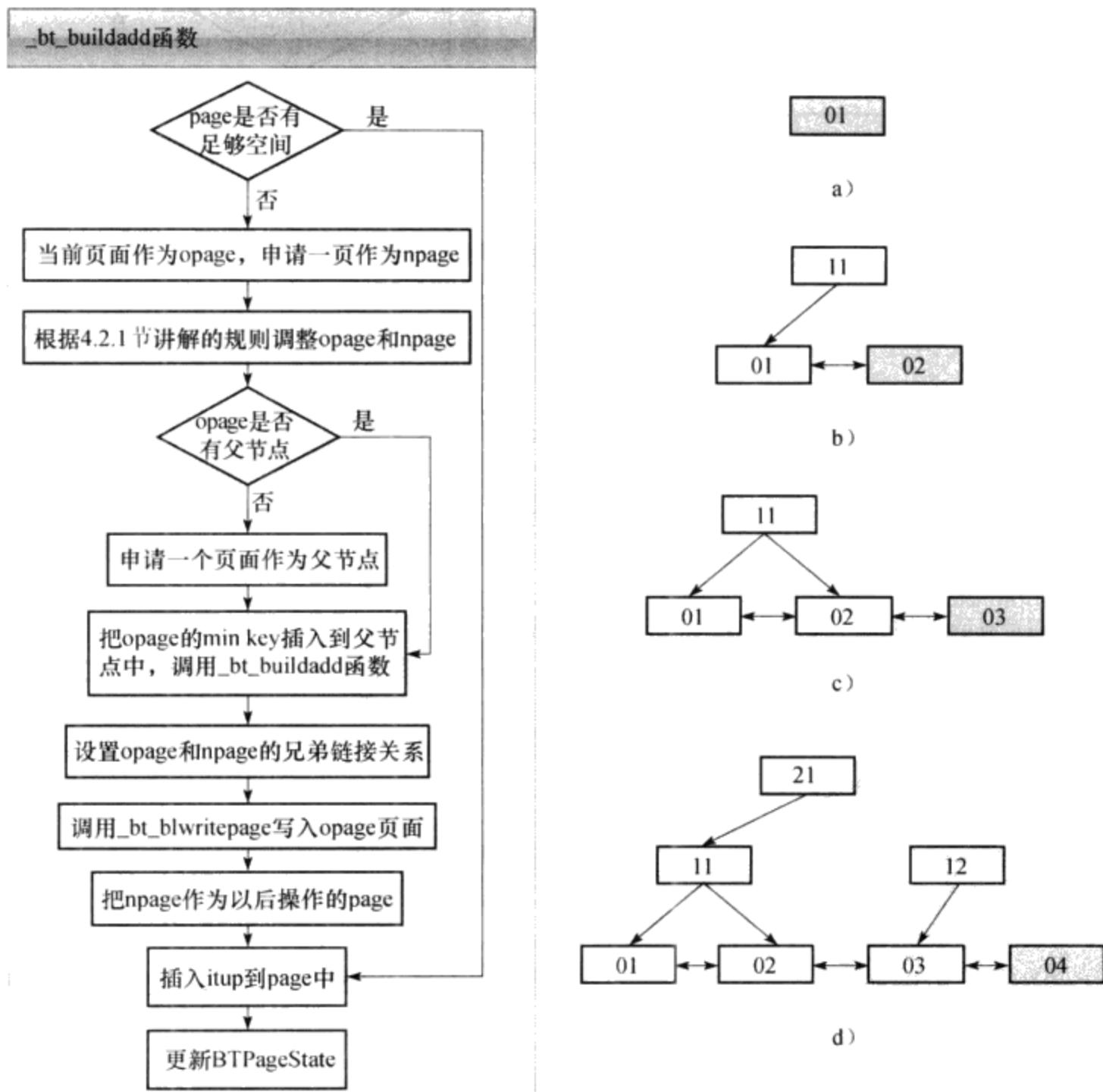


图 4-9 创建索引时插入索引元组流程及实例

- a) 当前插入节点为 01；b) 01 节点填充满后，申请节点 02 及父节点 11，当前插入节点为 02；
- c) 02 节点插满后，构建父节点与 02 的链接，然后申请节点 03，当前插入节点为 03；
- d) 03 节点插满后，申请父节点 02，然后申请新节点 04，当前插入节点为 04

在图 4-9 中，给出了一个插入过程中索引结构变化的例子，这里假设一个节点里最多存放两个元组。图 4-9a 到图 4-9d 反映了索引结构的变化。图中数字的第一位表示层数，第二位表示该节点是该层第几个节点。可以看出，当申请一个新的最右节点时，它与左兄弟节点的链接关系会立即构建，但与父节点的链接关系并没有设定。当所有索引节点插入完成后，每一层的最右节点的链接关系是由函数_bt_uppershutdown 来完成的。

当读取完成 spool 结构中的索引元组后，得到了初步的索引结构，但结构中每层最右节点与父节点的指针链接还没有构建，这需要调用函数_bt_uppershutdown 来完成。之前讲解 BTPageState 结构时提到，该结构内有一个指针指向一层的 BTPageState 结构，可以把这个结构理解为一个“栈”，在调整最右节点时，依次出栈，设置子节点和父节点的链接关系，直到“栈”为空则调整完毕，索引结构构建完成。_bt_uppershutdown 函数的流程如图 4-10 所示。

通过上面讲解的函数调用流程，函数_bt_uppershutdown 执行完成后，就得到完整的索引结构，也设置了索引结构的元页信息，最后再释放内存，整个索引即创建完毕。

2. 插入索引元组

若对一个表中的某个属性创建了索引，当表中有新的元组插入时，索引也需要进行相应的更新，也就是将新插入的元组封装成索引元组并插入到索引中，这个过程由 btinsert 函数完成。

btinsert 函数首先将表元组封装为索引元组，然后沿着 B-Tree 往下找到合适的插入节点。找到正确的节点后，若该索引关系是唯一索引，则会在节点中验证待插入的元组是否已存在，若存在则报错结束；如果索引不是唯一索引或者没有重复元组，则将索引元组插入到索引中。

在从根节点往下查找合适的节点过程中，会使用结构 BTStackData 来保存查找过程中的父节点，即保存查找路径。当插入后需要分裂叶子节点时，可以根据栈中存储的对应关系找到所有的父节点，并根据情况依次对父节点进行调整。BTStackData 结构的定义如数据结构 4.8 所示。

数据结构 4.8 BTStackData

```
typedef struct BTStackData
{
    BlockNumber      bts_blkno;    // 存储该元组对应的值的物理块号
    OffsetNumber     bts_offset;   // 在块中的偏移量
    IndexTupleData   bts_btentry;  // 整个索引元组
    struct BTStackData *bts_parent; // 在 B 树下降扫描时记录其父节点的信息，在分裂叶子节点的时候会
                                    // 按照该路径回溯
} BTStackData;
```

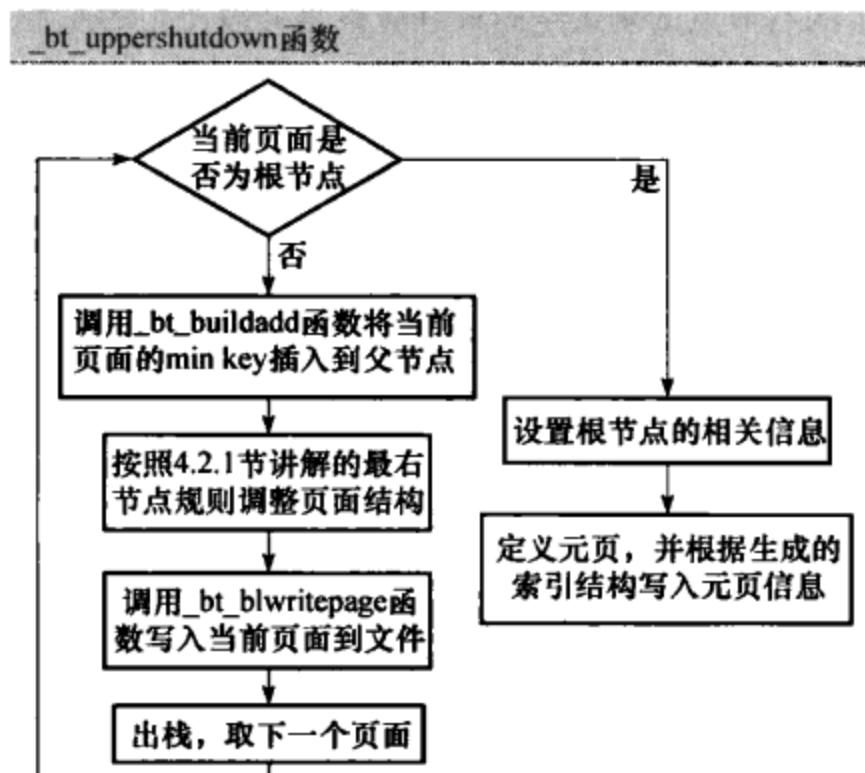


图 4-10 _bt_uppershutdown 函数流程图

在插入过程中还应考虑插入元组后节点的分裂问题。如果该节点需要进行分裂，则需要在该节点中找到合适的分裂位置。在 PostgreSQL 中规定，如果该节点为该层中的最右节点，那么分裂后产生的右节点的剩余空间最好应该是分裂产生的左节点剩余空间的 2 倍；否则分裂产生的两个节点的剩余空间最好应该相等。

当需要分裂节点时，会生成 FindSplitData 结构，它用来记录寻找节点分裂位置时的相关信息，特别是节点的最佳分裂位置以及待插入的新元组的相关信息。其结构的定义如数据结构 4.9 所示。

数据结构 4.9 FindSplitData

```
typedef struct FindSplitData
{
    Size          newitemsz;           // 插入新元组的大小
    int           fillfactor;         // 当分裂最右节点时需要
    bool          is_leaf;            // 是否是分裂叶子节点
    bool          is_rightmost;       // 是否是分裂该层中最右节点
    OffsetNumber  newitemoff;         // 新元组的插入位置
    int           leftspace;          // 左节点的剩余空间
    int           rightspace;         // 右节点的剩余空间
    int           olddataitemstotal;   // 旧元组所占空间
    bool          have_split;         // 是否找到合法的分裂
    // 下面的字段当 have_spilt 为真时有效
    bool          newitemonleft;       // 新元组是否插入到分裂左节点中
    OffsetNumber  firstright;        // 最优分裂点位置
    int           best_delta;         // 当前最好的分裂参考值
}FindSplitData;
```

说明：

- **best_delta**：是根据页面的空闲空间除以 16 得到的一个参考值，用于确定分裂位置。
- **fillfactor**：一个索引的填充因子，它是一个百分比，表示创建索引时每个索引页的数据填充率。对于 B-Tree 来说，在创建索引时节点将按照此百分比填充数据，在右侧（最大的键值）扩展索引时同样也按照此百分比填充数据。如果后来某个页被完全填满，那么该页将被分裂，从而导致索引性能退化。B-Tree 默认的填充因子是 90，但是有效的取值范围是 10 ~ 100。对于静态的不会发生改变的表，最佳值 100 可以让索引的物理体积最小，但是对于不断增长的表，较小的填充因子更合适，因为这将尽可能减少对页的分裂。其他索引方法对填充因子的理解与此类似，但是其默认值各不相同。如果有条件周期性地重建索引，那么建议使用较大的填充因子以减少索引的物理体积。用户可以在创建索引时指定 fillfactor 的值。

在实际构建 B-Tree 索引时，每个节点往往并没有按照该节点能够容纳关键字的个数完全填满，而是保留了一定的“空位”。这样虽然会带来一定空间的浪费，但可以避免在插入元组时过于频繁地分裂节点。

用户在创建索引时，可通过指定一个索引的填充因子（fillfactor），来设定创建索引时每个索引页面的数据填充率。若用户没有设定，则B-Tree索引默认的规定如下：

- 在叶子节点中加入的索引元组的总大小如果超过页面大小的90%就视为充满。
- 在内部节点中加入索引元组的总大小如果超过页面大小的70%就视为充满。
- 分裂节点时，如果该节点为该层的最右节点，必须保证分裂后的左节点的空闲空间为 $(100 - \text{fillfactor})\%$ ，否则应保持左右节点空闲空间相等。

如下面的语句将在表films上的title属性上创建B-Tree索引，但不使用默认的fillfactor，设定为70：

```
CREATE INDEX stu_id_idx ON student(id) WITH(fillfactor=70);
```

在代码实现时，封装了多个函数来实现插入元组的功能，实现这些操作的人口函数是_btinsert函数。该函数首先将表元组封装成索引元组，然后调用_bt_doinsert函数将索引元组插入到索引，_bt_doinsert的流程如图4-11所示。

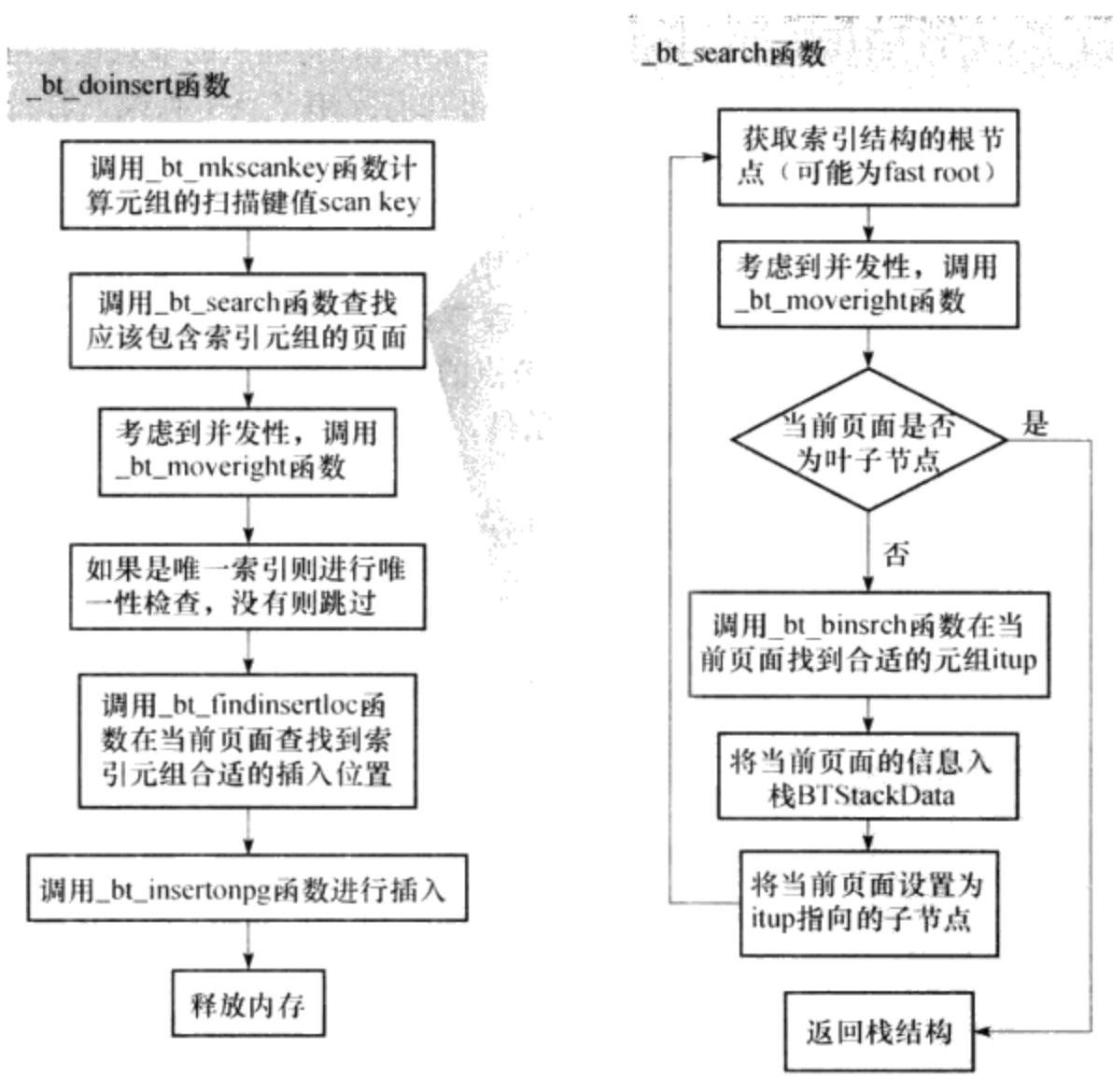


图4-11 插入索引元组到索引的流程

在图4-11中，首先调用_bt_search函数根据索引结构查找该索引元组应该存放的节点。这些节点都放在BTStackData栈中，这主要是考虑到后面插入索引元组时可能会导致节点分裂的情况。_bt_search函数调用了_bt_binsrch函数来实现当前页面的查找，根据查找节点的不同类型，_bt_binsrch函

数的返回值可能为：

- 若当前节点为叶子节点，则返回当前节点中第一个满足“key >= scankey”的位置。
- 若当前节点为内部节点，则返回当前节点中最后一个满足“key < scankey”的位置。

若为内部节点，则根据查找到的位置的元组找到其指向的子节点，然后循环向下查找。

从图 4-11 中可以看到，_bt_doinsert 函数和_bt_search 函数都调用了_bt_moveright 函数，该函数是根据 Lehman 和 Yao 的论文中的要求而实现的。在高并发的情况下进行当前查找操作的同时，可能有另一个操作正在分裂当前页面，那么当前获得的页面在分裂操作完成之后，可能不再是正确的页面（假如不是正确的，那么正确的页面只可能是在右兄弟页面），那么这个时候就要往右查找到正确的页面，_bt_moveright 就是实现这个功能的，具体信息可以参考原论文。

函数_bt_search 执行完成后，即得到了查找路径的栈结构及用于插入的叶子节点。若是唯一索引，则在待插入的叶子节点中进行唯一性检查，若检查到重复，则报错并结束。否则接下来调用_bt_findinsertloc 函数，该函数实现在当前节点中查找到待插入元组合适的位置。

_bt_findinsertloc 函数查找的节点可能不会局限于当前的节点，还会到右兄弟节点中查找。如果当前节点中的最大元组值和右兄弟节点中的最小节点值相等（即待插入的节点也是可以插入到右兄弟节点的），且当前节点中没有足够的剩余空间，则会在右兄弟节点中查找，若右兄弟节点还是没有足够的空间，则结束查找并在后面具体进行插入操作时进行分裂操作。

查找到待插入的节点后，即调用_bt_insertonpg 函数执行插入操作。该函数首先判断当前页面是否有足够的剩余空间，如果有则直接插入。如果没有则会调用_bt_findsplitloc 函数遍历节点中的每个位置，检查其作为最佳分裂点的可行性，该查找过程有几个参数：

- leftfree：分裂点为当前偏移量时，分裂后左节点的空闲空间（rightfree 含义与此类似）。
- Delta：用于判断当前偏移量是否满足分裂条件，best_delta 见数据结构 4.8 中 FindSplitData 的定义。

找到合适的分裂点后，即调用函数_bt_split 进行分裂操作。该函数首先申请一个页面作为新的节点（右节点），然后将页面中的元组根据分裂点元组确定是放入新节点还是旧节点（左节点）中，最后将待插入的元组插入到节点中。

该函数执行完成后，得到一个新的节点，需要将新节点作为一个元组插入到其父节点中，插入到父节点中元组的键值是左节点的“High-key”，也就是右节点的最小值。_bt_insert_parent 函数首先封装得到需要插入父节点的索引元组，然后调用_bt_insertonpg 函数将元组插入到父节点中。在查找父节点时，就会使用到上面讲到的 BTStackData 栈结构。

综上所述，_bt_insertonpg 的执行流程如图 4-12 所示。

3. 扫描索引

在利用索引进行查找等相关工作时需要使用扫描函数。使用 B-Tree 索引进行范围查询时，首先就要对索引进行扫描。在 PostgreSQL 中，与扫描相关的函数主要有以下几个：

- btgettuple 函数：得到扫描中的下一个满足条件的索引元组。如果已经初始化了扫描信息（也就是说之前执行了扫描），只需要得到上一次的扫描位置，然后调用_bt_next 函数根据扫描方向获取下一个满足扫描条件的索引元组；如果没有初始化，就调用_bt_first 函数开始新的扫描。该函数首先会预处理扫描信息，然后调用前面分析的_bt_search 函数来实现扫描。最后将获取到的索引元组的 TID 返回。

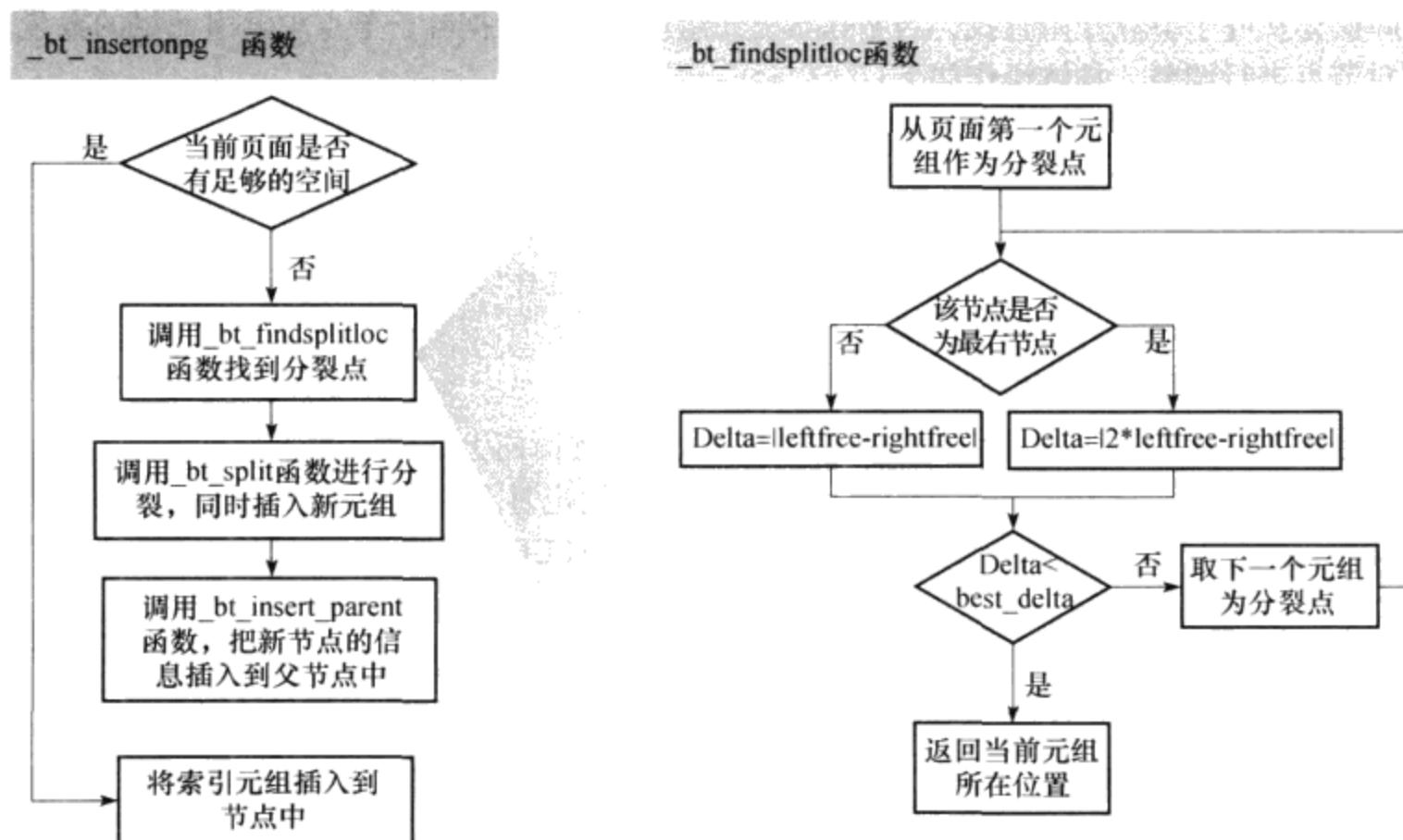


图 4-12 插入索引元组的流程

- **btbeginscan** 函数：开始索引扫描。该函数根据需要扫描的索引的相关信息生成一个 **IndexScanDesc** 结构，该结构为指向 **IndexScanDescData** 的指针。**IndexScanDescData** 结构中保存了扫描索引的相关信息，如扫描的键值、查找到的位置等，最后该函数返回 **IndexScanDesc** 结构。
- **btrescan** 函数：在某些情况下可能需要重新扫描索引，这时候就会调用 **btrescan** 函数，它的功能是重新开始一个索引扫描过程。
- **btendscan** 函数：此函数与 **btbeginscan** 函数成对出现，主要功能是释放进行一个索引扫描所占用的系统资源。
- **btmarkpos** 函数：当一个正在进行的索引扫描由于某种原因需要停止时，就需要调用 **btmarkpos** 函数，它将保存当前扫描位置的相关变量。
- **btrestrpos** 函数：与上面的 **btmarkpos** 相对应，用 **btmarkpos** 中所存储的最后扫描位置信息，导入到当前扫描位置信息变量中，从而恢复到上次的扫描位置以后再开始扫描，这样可以节省扫描时间。

4. 删除索引元组

在 PostgreSQL 中，删除 B-Tree 索引元组的函数主要有两个：一个是查找需要删除的节点信息的函数 **btvacuumcleanup**，该函数寻找可以删除的页面；另一个是进行批量删除的函数 **btbulkdelete**，该函数批量删除指向一个表元组集合所对应的所有索引元组。

在 PostgreSQL 中，若删除了基表中的一个元组，系统并不会立即删除该表元组对应的索引元组，而是在 VACUUM (清空) 的时候删除。

删除操作从叶子节点开始查找，当删除索引元组后，需要根据叶子节点中索引元组的数量来决定是否对节点进行调整，具体操作如下：

```

if(删除后, 节点仍有最小数目的关键字)
    then 不做任何操作;
else if(同父相邻兄弟节点关键字个数大于最小限值)
    then 从同父兄弟节点调剂关键字, 并修改父节点中“键 - 指针”对;
else if(同父相邻兄弟节点关键字个数刚好等于最小限值)
    then 与兄弟节点合并, 并修改父节点中“键 - 指针”对。

```

如果删除父节点关键字导致其数目低于最小限值，则逐渐向上层做上述修改。

在 B-Tree 索引中，由函数 `btbulkdelete` 通过扫描索引确定基本表中哪些元组被删除，而后逻辑上删除（标记删除）那些元组对应的索引元组，最后由 `VACUUM` 在物理上删除。

4.3 Hash 索引

在实际的数据库系统中，除了 B-Tree 外，还有多种数据结构可做索引，Hash 表就是其中的一种。通过 Hash 表可以把键值“分配”到各个桶中。使用这种结构的索引称为 Hash 索引，其核心是构建 Hash 表。

在 Hash 表中，有一个 Hash 函数 h ，它以查找键为参数，并计算出一个介于 0 到 $B-1$ 的数值，其中 B 是桶的数目。如果记录的查找键为 K ，那么该记录将被放在编号为 $h(K)$ 的桶中。一般有两种方法实现 Hash 表：

1) 静态 Hash 表：桶数目 B 不变。

2) 动态 Hash 表：允许 B 改变，使 B 近似于记录总数除以块中能容纳的记录数所得到的商。也就是说，每个桶大约有一个存储块。动态 Hash 表又可分为两种：

①可扩展 Hash 表：它在简单的静态 Hash 表结构上进行了扩充，其特点如下：

- 为桶增加了一个间接层，即用一个指向数据块的指针数组来表示桶，而不是用数据块本身组成的数组来表示桶。
- 指针数组能增长，它的长度总是 2 的幂，因而数组每增长一次，桶的数目就翻倍。但并非每个桶都有数据块；如果某些桶的所有记录都可以放在一个块中，那么这些桶可能共享一个块。
- Hash 函数 h 为每个键计算出一个 K 位二进制序列，该 K 值足够大，比如 32。但是，无论何时桶的数目都使用从序列第一位开始的若干位来表示，位的数值等于 K 。比如，使用了 r 位，则这时桶中将有 2^r 个项。

②线性 Hash 表。它的桶的增长较为缓慢，具有如下特点：

- 桶数目 B 的选择总是使存储块的平均记录数与存储块所能容纳的记录总数保持一个固定的比例，例如 80%。假设每个磁盘块能存放两条记录， r 为当前 Hash 表中总的记录数， B 为当前的桶数目，那么 $r \leq 1.6B$ 。
- 由于存储块并不是总能分裂，所以允许有溢出块。
- 用来做桶数组项序号的二进制位数是 $\lceil \log_2 B \rceil$ ，其中 B 是当前的桶数。这些位总是从 Hash 函数得到的位序列的右（低位）端开始取。

- 假定 Hash 函数值的 i 位正在用来给桶数组项编号，且有一个键值为 K 的记录想要插入到编号为 $a_1 a_2 \dots a_i$ 的桶中，即 $a_1 a_2 \dots a_i$ 是 $h(K)$ 的后 i 位。那么， $a_1 a_2 \dots a_i$ 当作二进制数，设它为 m 。如果 $m < B$ ，那么编号为 m 的桶存在并把记录存入该桶中。如果 $B \leq m < 2^i$ ，那么桶 m 还不存在，因此把记录存入桶 $(m - 2^{i-1})$ 中，也就是把 a_i 改为 0 时对应的桶。

在 PostgreSQL 系统中使用的是动态 Hash 表，桶的个数随着分裂次数的增加而增加，但每次分裂后，所有增加的桶并不会立即使用，而只是使用需要存入元组的桶，其他新增加的桶都被保留。这种方式能够有效地减少重新分配元组到桶中的工作量，提高效率。下面将从 Hash 索引的页面组织结构和实现流程两方面详细分析。

4.3.1 Hash 索引的组织结构

在 PostgreSQL 8.4.1 的实现中包含两种类型的 Hash 表，一个是用做索引的外存 Hash 表，另一个是用于内部数据查找的内存 Hash 表（在第 3 章中已多次提到）。本节主要对前者进行分析，其源代码在 `src/access/hash` 目录中，它采用的是线性 Hash 表的一个实现。

在 Hash 表中有四种不同类型的页面，分别为元页、桶页、溢出页、位图页，这 4 种页面的结构如图 4-13 所示。

Hash 索引的页面结构也是按照图 3-4 所示的结构来组织的，对于每种页面末尾的“Special Space”，Hash 索引填充的是 HashPageOpaqueData 结构。其定义如数据结构 4.10 所示。

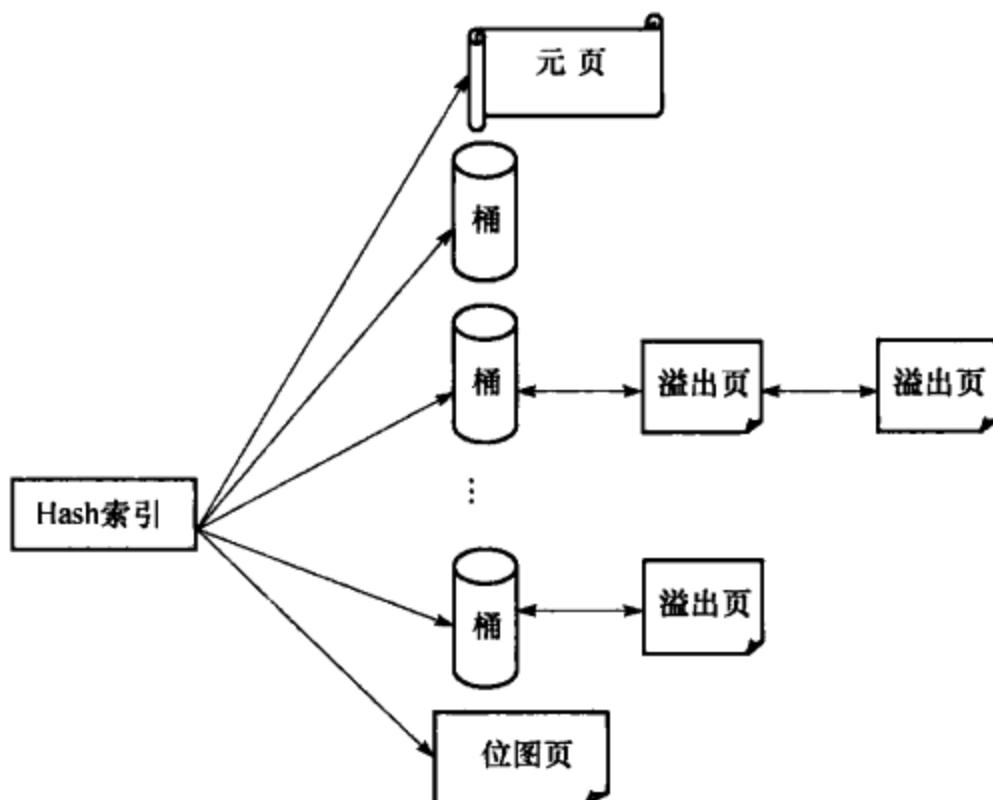


图 4-13 Hash 索引 4 种页面组织结构图

数据结构 4.10 HashPageOpaqueData

```

typedef struct HashPageOpaqueData
{
    BlockNumber hasho_prevblkno;      //前一页(桶或溢出块)的块号
    BlockNumber hasho_nextblkno;      //后一页(桶或溢出块)的块号
    Bucket      hasho_bucket;        //该页属于某个桶页的桶号
    uint16       hasho_flag;          //标识该页的类型,是桶页、溢出页还是位图页等
    uint16       hasho_page_id;        //用于标识该页是 Hash 索引的 ID
} HashPageOpaqueData;

```

下面将分别详细介绍这四种页面。

1. 元页

每个 Hash 索引都有一个元页，它是该索引中的 0 号页（也就是第 1 页，这里的编号是 C 语言的编号习惯），其 HashPageOpaqueData.hasho_bucket 为 -1，表示它不属于任何桶。字段 hasho_prevblkno 和 hasho_nextblkno 均置为 InvalidBlockNumber。

元页中记录了 Hash 的版本号、Hash 索引记录的索引元组数目、桶的信息、位图等相关信息。通过元页可以了解该 Hash 索引在总体上的分配和使用情况，在索引元组的插入、溢出页的分配回收以及 Hash 表的扩展等过程中，都需要使用元页。其结构及各数据成员的意义如数据结构 4.11 所示。

数据结构 4.11 HashMetaPageData

```

typedef struct HashMetaPageData
{
    uint32          hashm_magic;                      //Hash 表的 magic 号,类似 B-Tree 索引的
                                                       //magic 号
    uint32          hashm_version;                    //PostgreSQL 的 Hash 版本号
    double          hashm_ntuples;                   //表中存储的元组的数目
    uint16          hashm_ffactor;                  //装填因子,与索引元组个数/桶数的比值进
                                                       //行比较,确定是否分裂以增加桶的个数
    uint16          hashm_bsize;                     //桶的大小(以字节计)
    uint16          hashm_bmssize;                  //位图的大小(以字节计),必须为 2 的整数
                                                       //次幂
    uint16          hashm_bmshift;                 //log2(bitmap 数组的大小)
    uint32          hashm_maxbucket;                //可用的最大桶号
    uint32          hashm_highmask;                 //桶号的高 16 位的掩码
    uint32          hashm_lowmask;                  //桶号的低 16 位的掩码
    uint32          hashm_ovflpoint;               //当前的分裂次数,记录表成倍增长的情况,
                                                       //分裂一次加 1
    uint32          hashm_firstfree;                //第一个可能的空闲溢出页号
    uint32          hashm_nmmaps;                  //位图的个数
    RegProcedure    hashm_procid;                  //Hash 函数的 OID
    uint32          hashm_spares[HASH_MAX_SPLITPOINTS]; //记录了从第一个分裂点开始到当前的分裂
                                                       //点总共分配的溢出页数目
    BlockNumber     hashm_mapp[HASH_MAX_BITMAPS]; //记录了位图的块号
} HashMetaPageData;

```

其中的字段说明如下：

- hashm_bsize 字段实际保存的是一个桶页中用于存放索引元组的空间大小（按字节），由于使用一个磁盘块作为一个桶页，该值为“磁盘块所占空间（8k） – PageHeaderData 所占空间 – HashPageOpaqueData 所占空间”。
- hashm_highmask 和 hashm_lowmask 用于根据索引元组的 Hash 值计算其应放入的桶号。

- `hashm_firstfree` 表示当前可能空闲的最小溢出页号，注意是可能的，并不一定就是空闲的。在查找空闲的溢出页时，还要对该页进行核查。
- `hashm_spares` [`HASH_MAX_SPLITPOINTS`] 保存了截止到各个分裂点时整个 Hash 索引分配的溢出页总数。当 `splitpoint` 增加后（即产生了一次新的分裂），在之前保存在这个数组中的值就不能再改变，即不管在该分裂点之前分配的溢出页是否回收都不会修改 `hashm_spares` 数组中的值。`HASH_MAX_SPLITPOINTS` 的最大值为 32，即最多可分裂 32 次。
- `hashm_mapp` [`HASH_MAX_BITMAPS`] 记录了各个位图的块号，`HASH_MAX_BITMAPS` 的最大值为 128，即最多可分配 128 块位图。通过这个数组，能够快速找到各个位图的块号。

从上面的结构定义可以看出，元页记录了 Hash 索引的基本信息，每次对 Hash 索引进行读写操作时都要将元页读入内存，并对其进行更新。Hash 索引初始化的时候会将已知的关于即将要建立的索引的信息都写到元页中去。

2. 桶页

Hash 表由多个桶组成，每个桶由一个或多个页组成，每个桶的第一页称为桶页，其他页称为溢出页，桶页随着桶的建立而建立。桶页结构中的 `HashPageOpaqueData.hasho_bucket` 为该桶的标识，若该桶有溢出页，字段 `hasho_nextblkno` 指向该溢出页，用于构建桶页的链结构，实现快速查找。

桶页是成组分配的，每次分配的数目都是 2 的幂次。0 号和 1 号桶页是在 Hash 表初始化的时候分配的，当 `splitpoint` 为 1 时（进行第一次分裂），将会同时分配 2 号和 3 号桶页。当需要第 5 个桶时，也就是 `splitpoint` 增加到 2 时，同时分配 4 ~ 7 号桶页。当 `splitpoint` 增加到 3 时，同时分配 8 ~ 15 号桶页，以此类推。而每次分配的桶页在磁盘上都是连续存储的，通过这种方式，可以快速了解到各个桶是在哪一次分裂时产生的及其磁盘块号。

在 4.3.2 节，会讲到当 `splitpoint` 增加 1 的时候，并不是立即使用所增加的桶页，而是根据需要分裂一个桶页中的元组到两个桶页（原来的桶页 + 新增加的桶）中。例如，当 `splitpoint` 增加到 2 时，并不会立即使用 4 ~ 7 号桶页，可能先会使用 4 号桶，5、6 和 7 号桶现在并没有使用但保留着，当 4 号桶填满需要增加溢出页时，并不会将溢出页分配在 4 号桶页之后，而是 7 号桶页的后面，因为 5、6 和 7 号桶的磁盘已经在 4 号桶之后连续分配。这样一种分配方式可以保证前面提到的“每次分裂产生的桶页，都是连续存储”的原则，这也为下面将要介绍的计算桶所对应的磁盘块号的方法提供了基础。

在每一个 `splitpoint`，首先分配（预留） $2^{\text{splitpoint}}$ 个桶页，然后在这些桶页的后面，再根据需求分配一定数目的溢出页，这个数目将由数组 `hashm_spares` 记录。当 `splitpoint` 增加 1 的时候，再分配 $2^{\text{splitpoint}}$ 个桶页，以此类推，这样的实现机制使得很容易求得页面在磁盘中的块号。公式如下：

$$\text{blocknum} = \begin{cases} \text{pagenum} + \text{hashm_spares}[\lceil \log_2(\text{pagenum} + 1) \rceil - 1] + 1 & \text{pagenum} \neq 0 \\ 1 & \text{pagenum} = 0 \end{cases}$$

可以由桶页的页号 `pagenum` 求出其在在磁盘的块号 `blocknum`。

3. 溢出页

如果某个元组在它所属的桶中放不下的时候，就需要给该桶增加一个溢出页。溢出页和桶页之间是通过双向链链接起来的，即每一页都记录了前一页的块号（`prevblkno`）和后一页的块号（`nextblkno`），这个双向链保存在 `HashPageOpaqueData` 结构中。当溢出页分配给某个桶页使用时，溢出页

的 HashPageOpaqueData.hasho_bucket 将设置为桶页的块号。

前面说过，在每次扩展 Hash 表后，都会给该次扩展预留桶的编号（不管是否实际已经分配）。所以，分配的溢出页的磁盘号都是在预留的桶之后的。同时，在元页中使用数组 hashm_spares 记录了在各次桶的扩展时分配的溢出页数，从而能够快速地计算出各个桶页的块号。当 splitpoint 的值由 k 增加到 $k+1$ 时，hashm_spares [k] 的值就不再允许更改，以保证能够有效地找到每次扩展的桶。也就是说，溢出页一旦分配，便一直存在，即使是被回收也只是标记它为空闲，并没有释放其物理空间。

当该桶存放的记录已经达到饱和时，以后插入到该桶中的记录会添加到该桶的溢出页中，而溢出页通过双向指针与桶页链接，可灵活实现查找、删除等操作。

4. 位图

位图用于管理 Hash 索引的溢出页和位图页本身的使用情况。如果某个溢出页上的元组都被移除或删除，就要将此溢出页回收，但并不把它还给操作系统，而是继续由 PostgreSQL 管理，以便下次需要溢出页的时候使用。因此，需要一种机制来记录每一个溢出页是否可用，这时便用到了位图。所有位图页的 HashPageOpaqueData.hasho_bucket 都为 -1，HashPageOpaqueData.hasho_prevblkno 和 hasho_nextblkno 均置为 InvalidBlockNumber（无效），通过元页中的 hashm_mapp 数组可以找到位图对应的块号。

在位图里，对于每一个溢出页和位图页都有一个比特位标识其使用情况，0 表示该页可用，1 表示不可用，可以把这些 0、1 位看作一个数组。当新生成一个位图页时，由于每一位并没有实际对应的溢出页，所以都置为 1。位图本身也在位图中有记录，可以把位图看作是一个特殊的溢出页。

位图的页面格式与图 4-5 的格式稍有不同，位图中不包含任何元组，页头为 PageHeaderData，页尾有 Special Space，这和一般的页面相同。位图的页头和页尾中间的部分就是位数组，每一位对应一个溢出页或位图。位图中所含位的数目必须是 2 的整数幂。位图的格式如图 4-14 所示。

那么，如何通过位图页中的 0、1 标记来找到其对应的溢出页的块号呢？首先可以根据元页中的 hashm_mapp 数组找个各个位图页的块号，从而找到位图页。在这里还要使用到上面提到的 hashm_spares 数组，该数组记录每一次分裂点时分配的溢出页个数。每分配一个溢出页，就会在位图页中有一个比特位与其对应，随着溢出页的不断增加，位图页中的比特位也依次被使用。也就是说，通过位图页中某一个比特位在该页（数组）中的编号，可以得出该位对应的溢出页是在哪一次桶的分裂点之后分配的，再结合 hashm_spares 数组（保存每次桶的分裂点之后分配的溢出页数目）即可找到该位对应的溢出页号。可以看出，正是利用位图页中的位序号和溢出页页号一一对应及溢出页一经分配即永远使用的特性，巧妙地实现了快速查找空闲溢出页的功能，其查找流程如图 4-15 所示。

下面举例说明 Hash 表在进行不断的插入过程中各种页面分配情况。

1) 初始化 Hash 表。初始状态下，Hash 表中的 0 号块是元页，1 号和 2 号块分别分配给 0 号桶和 1 号桶，并且在桶的后面分配 0 号位图页（块号为 3）。相关字段的值如下：

PageHeaderData			1	1	0	1	0
1	0	0	1	1	1	0	0
.....							
1	1	0	1	1	1	1	0
0	1	1	1	0	Special Space		

图 4-14 位图格式

```

maxbucket = 1;           //最大的桶号,下同
splitpoint = 1;          //目前的分裂次数,下同
hashm_spares[1] = 1;     //1个位图,0个溢出页,当前分配的溢出页总数,下同

```

页面	元页 0	位图 0	桶 0	桶 1
块号	0	3	1	2

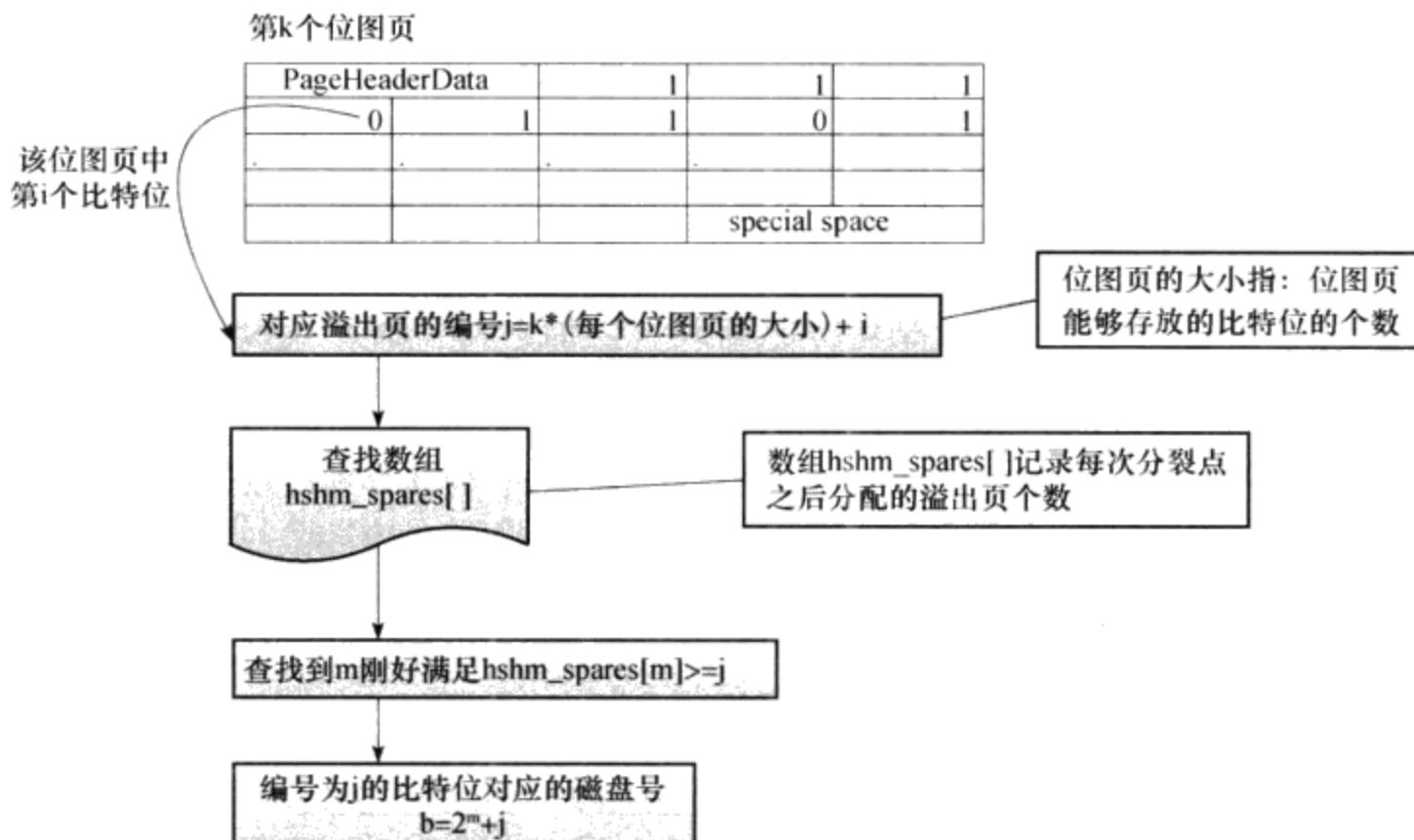


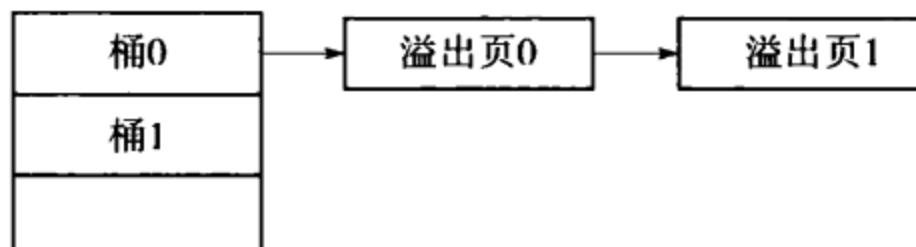
图 4-15 根据位图查找对应的溢出页流程图

2) 增加两个溢出页。为 0 号桶增加的两个溢出页只能在 0 号位图的磁盘块之后分配，因此其磁盘块号分别为 4 号和 5 号。由于目前还未进行分裂，因此最大桶号和分裂次数保持不变，但分配的溢出页总数发生了变化。

```

maxbucket = 1;
splitpoint = 1;
hashm_spares[1] = 3      //1个位图,2个溢出页

```

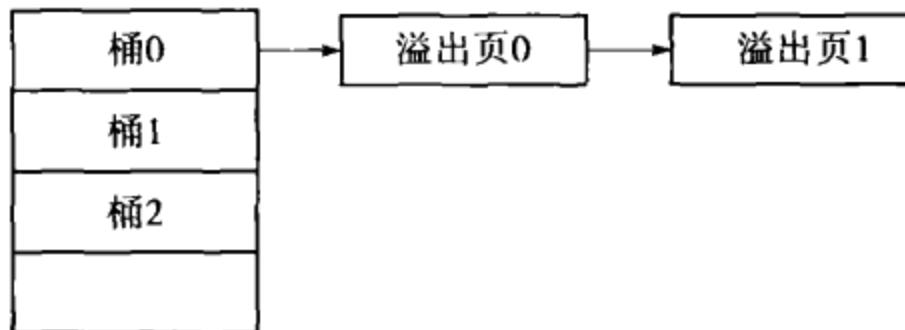


页面	元页 0	位图 0	桶 0	桶 1	溢出页 0	溢出页 1
块号	0	3	1	2	4	5

3) 增加一个桶。增加一个 2 号桶需要进行分裂，这时实际增加了 2 号桶和 3 号桶，只是 3 号

桶并未投入使用，因此最大桶号只能设置为 2，分裂次数自然也要加 1。分配的 2、3 号桶的块号分别为 6 号和 7 号。

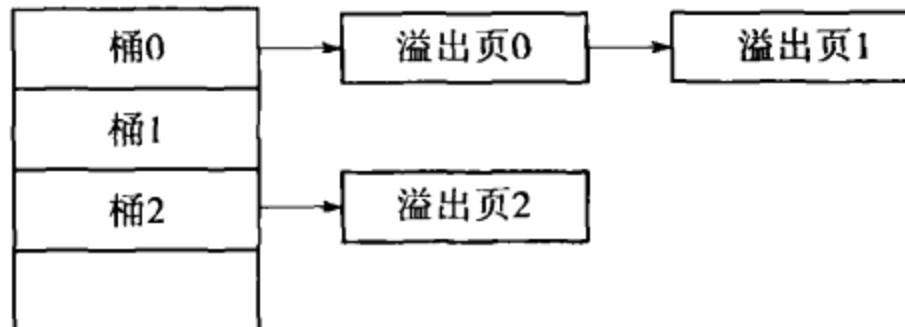
```
maxbucket = 2;
splitpoint = 2;
hashm_spares[2] = 3 //1 个位图,2 个溢出页
```



页面	元页 0	位图 0	桶 0	桶 1	桶 2	桶 3	溢出页 0	溢出页 1
块号	0	3	1	2	6	7	4	5

4) 为 2 号桶增加一个溢出页。虽然目前只使用了 2 号桶，但 3 号桶已经分配保留，所以为 2 号桶增加的溢出页的块号只能在 3 号桶之后分配，因此新溢出页块号为 8。

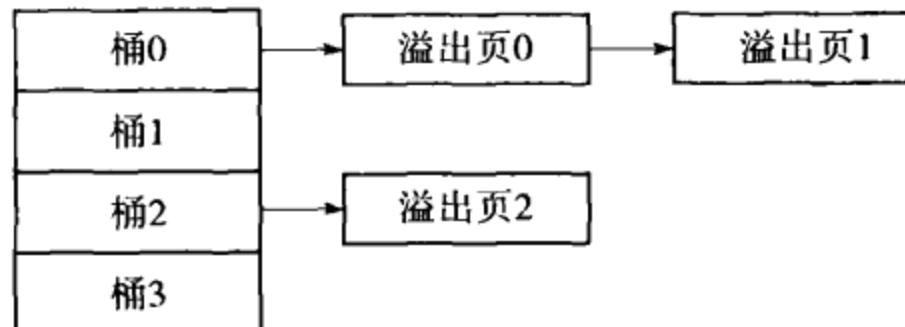
```
maxbucket = 2;
splitpoint = 2;
hashm_spares[2] = 4 (1 个位图,3 个溢出页)
```



页面	元页 0	位图 0	桶 0	桶 1	桶 2	桶 3	溢出页 0	溢出页 1	溢出页 2
块号	0	3	1	2	6	7	4	5	8

5) 启用 3 号桶。此时最大桶号改变为 3，但是由于 3 号桶是上次分裂时分配的，因此分裂次数不变。

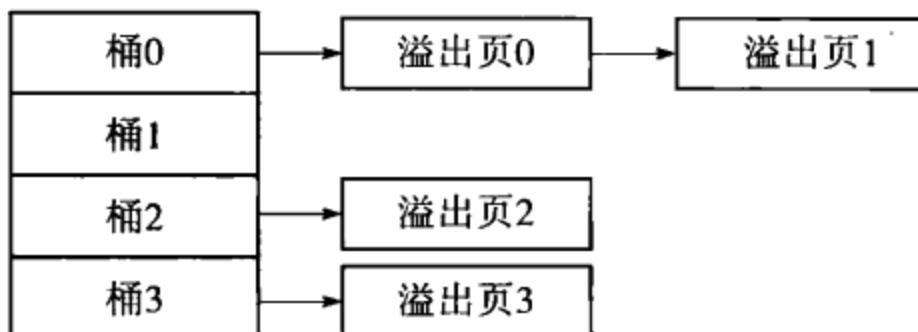
```
maxbucket = 3;
splitpoint = 2;
hashm_spares[2] = 4 //1 个位图,3 个溢出页
```



页面	元页0	位图0	桶0	桶1	桶2	桶3	溢出页0	溢出页1	溢出页2
块号	0	3	1	2	6	7	4	5	8

6) 为3号桶增加一个溢出页。

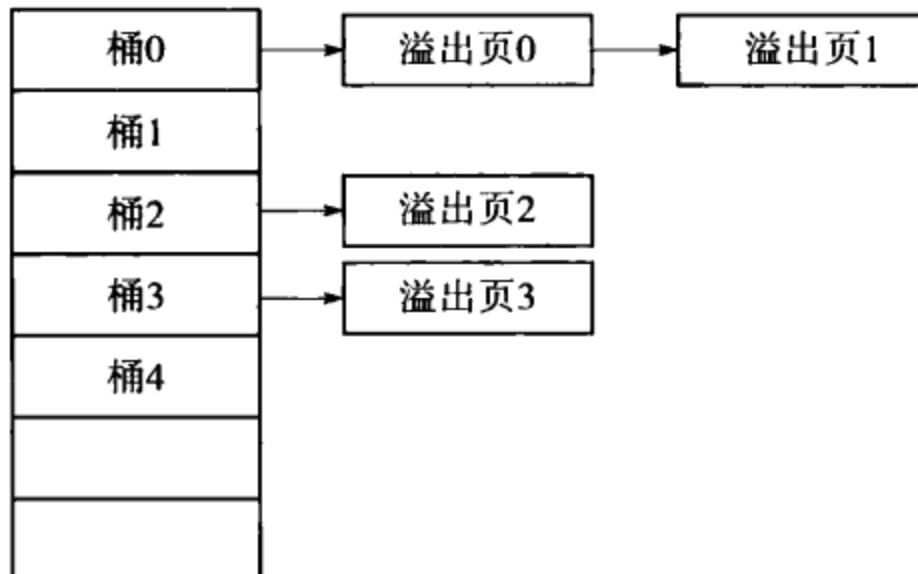
```
maxbucket = 3;
splitpoint = 2;
hashm_spares[2] = 5 //1个位图,4个溢出页
```



页面	元页0	位图0	桶0	桶1	桶2	桶3	溢出页0	溢出页1	溢出页2	溢出页3
块号	0	3	1	2	6	7	4	5	8	9

7) 增加一个桶。当需要新桶时，又要进行一次分裂，因此分裂次数加1，并且分配4~7号桶。但只使用4号桶，5~7号桶只分配位置但不使用，因此最大桶号为4，并且4~7号桶只能在上一次分配的溢出页（3号溢出页）之后分配（9号块）。

```
maxbucket = 4;
splitpoint = 3;
hashm_spares[3] = 5 //1个位图,4个溢出页
```



页面	元页0	位图0	桶0	桶1	桶2	桶3	桶4	桶5	桶6	桶7	溢出页0	溢出页1	溢出页2	溢出页3
块号	0	3	1	2	6	7	10	11	12	13	4	5	8	9

4.3.2 Hash索引的实现

实现Hash索引主要涉及以下几项工作：Hash表的构建、索引元组的插入、溢出页的分配与回

收、以及 Hash 表的扩展。下面将针对这 4 个方面选取其中重要的函数进行详细分析，对其他的函数则只作简单说明。读者若需要查看相关函数的代码，可参考 src\backend\access\hash 文件夹。

1. Hash 表的构建

正常情况下，构建一个 Hash 索引时，首先要初始化该索引的元页、桶以及位图页。之后，调用扫描函数对待索引的表进行扫描，生成对应的索引元组，最后将这些索引元组插入到 Hash 表中。

PostgreSQL 8.4.1 使用全局变量来记录所创建的 Hash 索引，它提供了一个外部函数 hashbuild 来建立 Hash 索引，函数执行的流程如图 4-16 所示。

下面将对 hashbuild 运行时所调用的几个子函数进行分析。

1) 函数 _hash_metapinit，该函数初始化 hash 表，其中包含一个元页（即图 4-7 中的 metapage）、N 个桶以及一个位图，最后返回初始化桶的数目 N，其中 N 为 2 的 k 次方。其执行流程如图 4-17 所示。

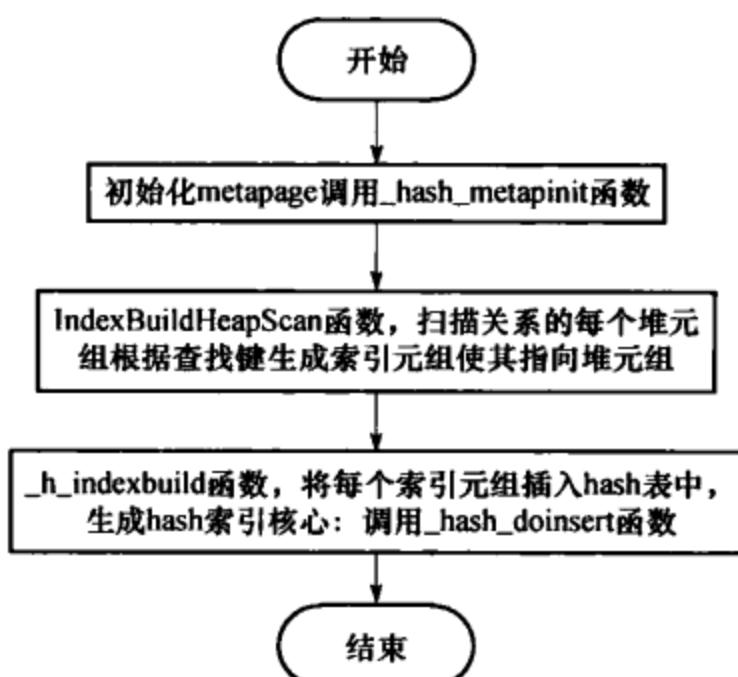


图 4-16 hashbuild 函数流程

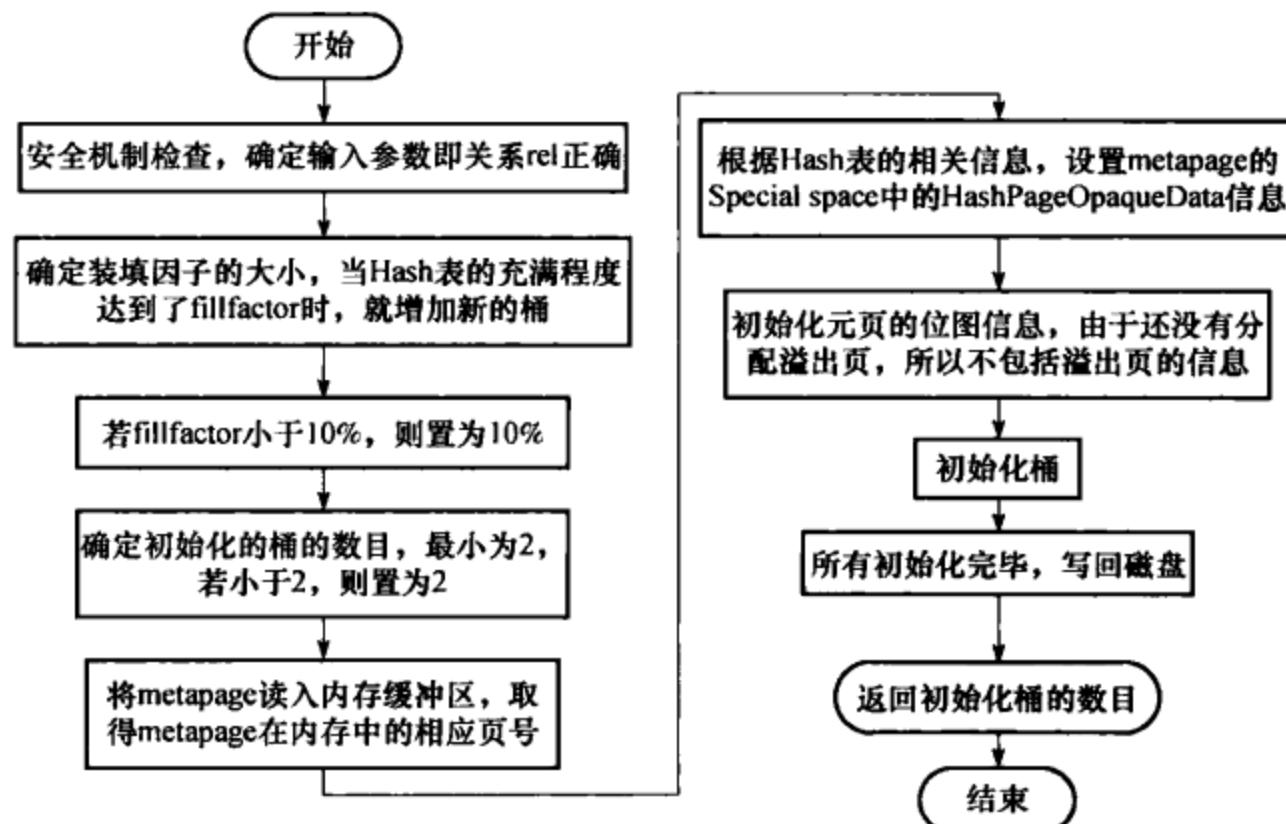


图 4-17 函数 _hash_metapinit 执行流程

与其他程序相比，本程序并没有严格的锁的机制。在初始化一个 Hash 表时，此 Hash 表还并没有建成，所以不会有程序对此 Hash 表进行访问，当然也不会有死锁的情况发生，因此本程序并没有严格的加锁与释放锁机制。

2) 函数 IndexBuildHeapScan，该函数扫描整个基表，根据查找键对其中的每一个元组都生成一个 Hash 索引元组，并插入到 Hash 表中。在其中还将调用函数 hashbuildCallback。

3) 函数 hashbuildCallback: 该函数根据表中的元组生成 Hash 索引元组，并插入到 Hash 表中，其执行流程如图 4-18 所示。

2. 元组的插入

完成了元页的初始化，根据查找键生成索引元组后，即可调用函数将该索引元组插入到 Hash 表中，这是创建 Hash 索引关键的一步。系统首先会读取元页中的相关信息，计算出该索引元组应插入的桶号，若该桶有足够的空间，则将元组插入到该桶中；否则申请溢出页并插入到溢出页中。完成插入后，系统会再次根据元页的信息来判断是否需要增加新的桶。

函数 _hash_doinsert 的执行流程如图 4-19 所示。

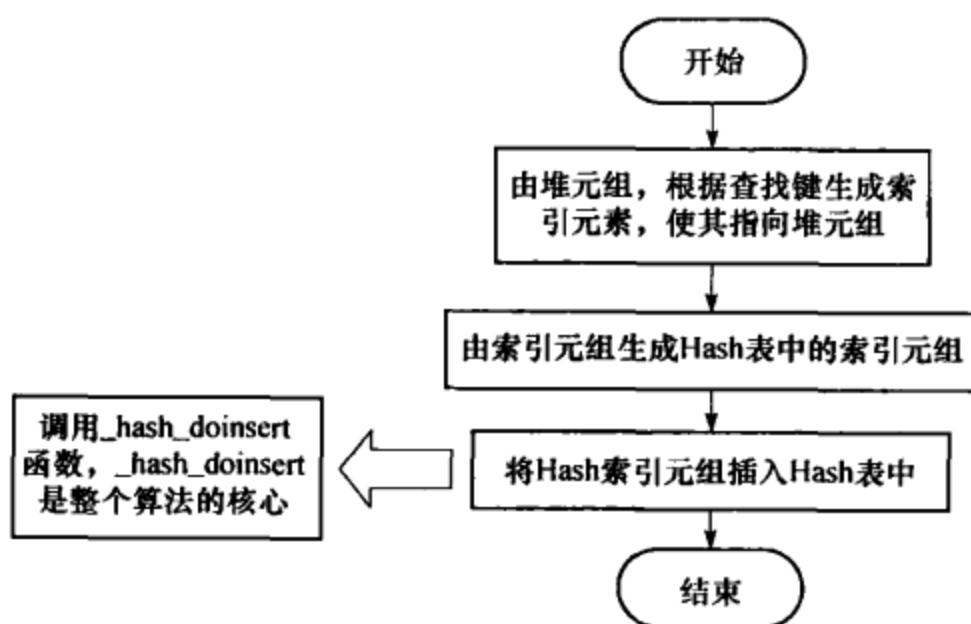


图 4-18 函数 hashbuildCallback 实现流程

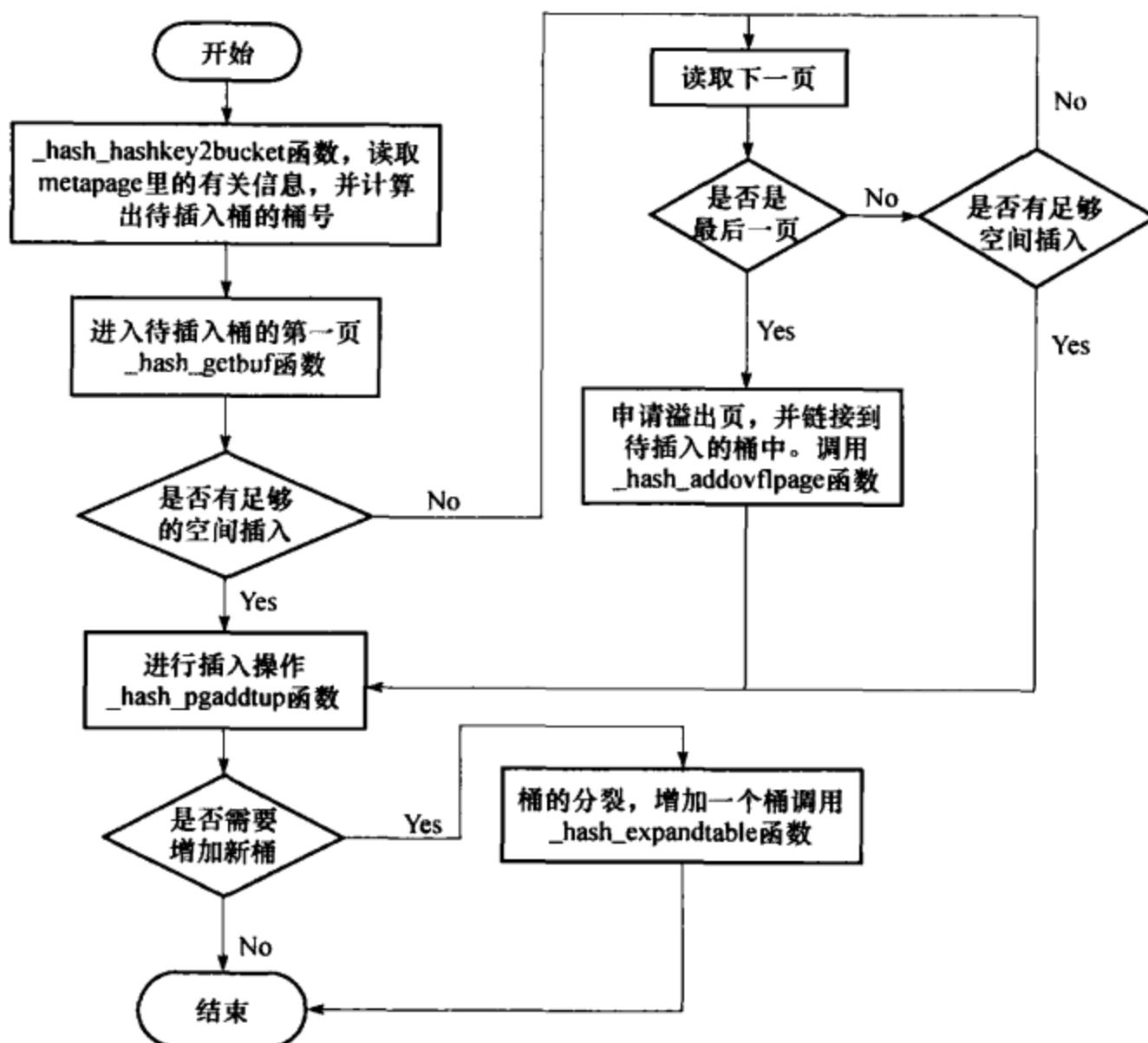


图 4-19 Hash 元组的插入流程图

因为插入操作并不会改变桶中已经存在的元组的位置，所以插入操作可以和扫描操作并发进行。

3. 溢出页的分配与回收

在索引元组的插入过程中，若待插入的桶中没有空间，就需要创建一个溢出页，并把它链接到该桶上，然后将索引元组插入到分配的溢出页中。当对 Hash 表进行扩展之后，原来存放在溢出页中的索引元组可能就会移到新增加的桶中，这时就需要对溢出页进行回收。下面将针对溢出页的分配和回收过程，讲解其中重要的函数和流程。

1) 函数 `_hash_addovflpage`: 当要插入的元组在所属桶链上找不到足够的空间的时候，就要增加一个溢出页，由此函数完成这项操作。其中，要插入元组所属的桶链的当前最后一页在 `buf` 这个缓冲区里。执行流程如图 4-20 所示。

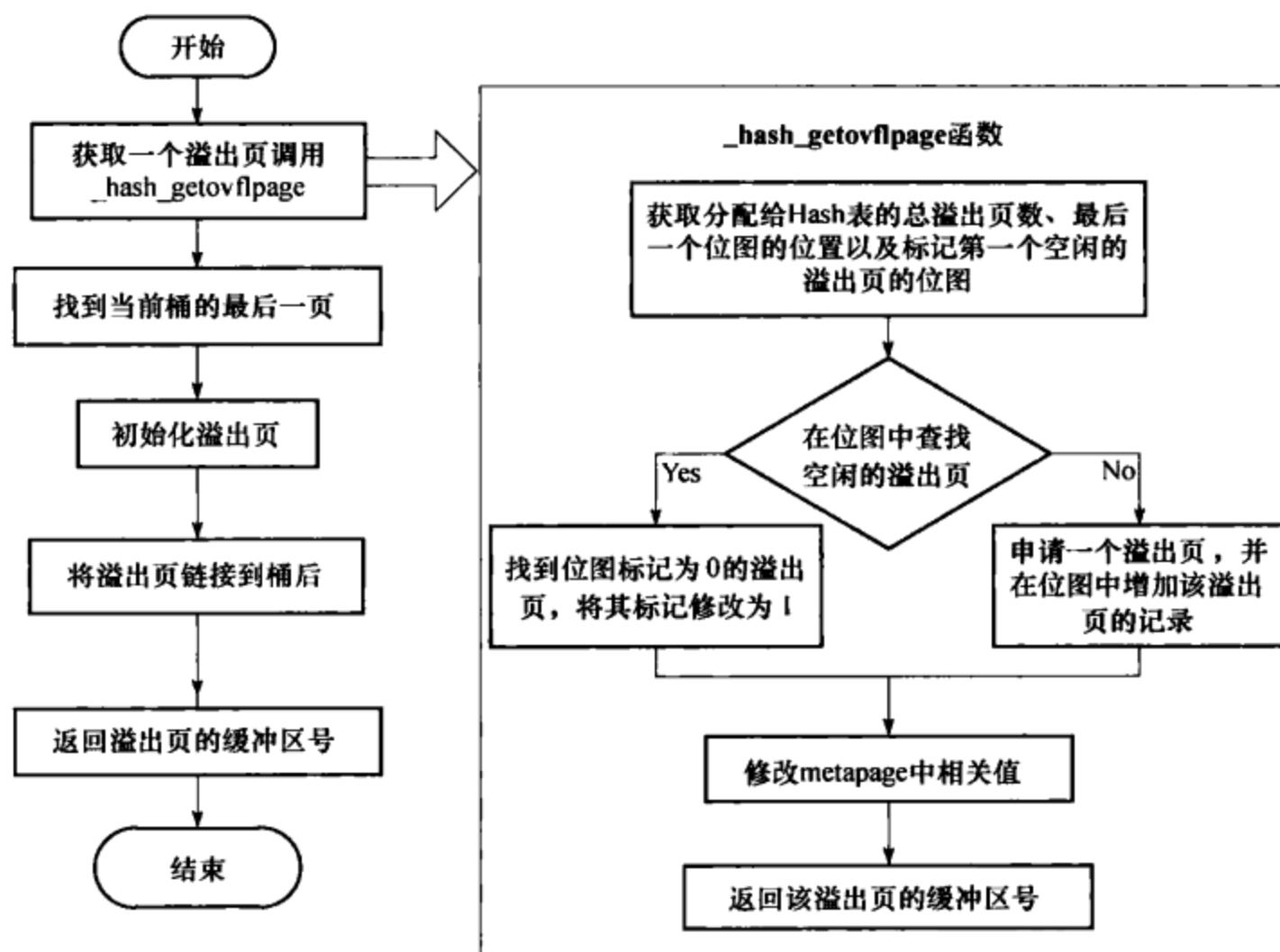


图 4-20 `_hash_addovflpage` 函数执行溢出页的分配

2) 函数 `_hash_getovflpage`, 该函数获取一个溢出页，并返回它在磁盘中的块号。基本执行过程如下：寻找溢出页时，首先通过位图查找可用的溢出页，若找到则将其分配，并修改位图中相应的位；若位图中没有可用的溢出页，即分配给该 Hash 表的溢出页都在使用，则从磁盘上申请一个块作为溢出页，并且在位图中增加该溢出页的记录。

- 在函数执行的过程中，调用者应拥有插入元组所属桶链的共享 “`lmgr lock`”。如果有某个进程调用函数 `void_hash_squeezebucket` 试图压缩桶中的元组，使桶中的元组更加紧凑，于是，在某些情况下此进程释放了 `buf` 这个溢出页（这种情况下 `buf` 是溢出页，而不是桶页），这

将会导致溢出页 buf 不在桶链中，最终使得程序出错。

- 程序中并没有释放返回的溢出页上的写锁，这是为了保证所返回的溢出页没有被插入任何元组。
- 此函数的调用者只需对 buf 加 pin，而不需要对其上锁，这保证了此函数可以由多个进程并发执行。
- 此函数可以由多个进程并发执行，因此，在调用函数 _hash_getovflpage 以得到一个空的溢出页这个过程中，可能有某个进程已经往桶链中插入了一些溢出页，这就使得进入程序时桶链的最后一页不一定是待插入的溢出页的前一页，所以在程序中有寻找当前桶链的最后一页这一操作。可以看到每次进入 for 循环时都要给 buf 上写锁，并且退出 for 循环时，也不释放 buf 上的写锁，直到将溢出页链接到桶链上以后，才释放写锁，这些机制将保证溢出页链接到的一定是当前桶链的最后一页。

3) 函数 _hash_freeovflpage：该函数将断开桶链，释放某个溢出页，并在位图中标识其为可用。参数 ovflbuf 为此溢出页所在的缓冲区号，进入函数前该缓冲区置有写锁，退出函数时释放该写锁。函数的返回值是桶链中此溢出页的下一页的块号，若下一页不存在，则返回 InvalidBlockNumber。此外，在整个函数的执行过程中，调用者始终拥有此溢出页所在的桶上的排他 lmgr 锁。执行流程如图 4-21 所示。

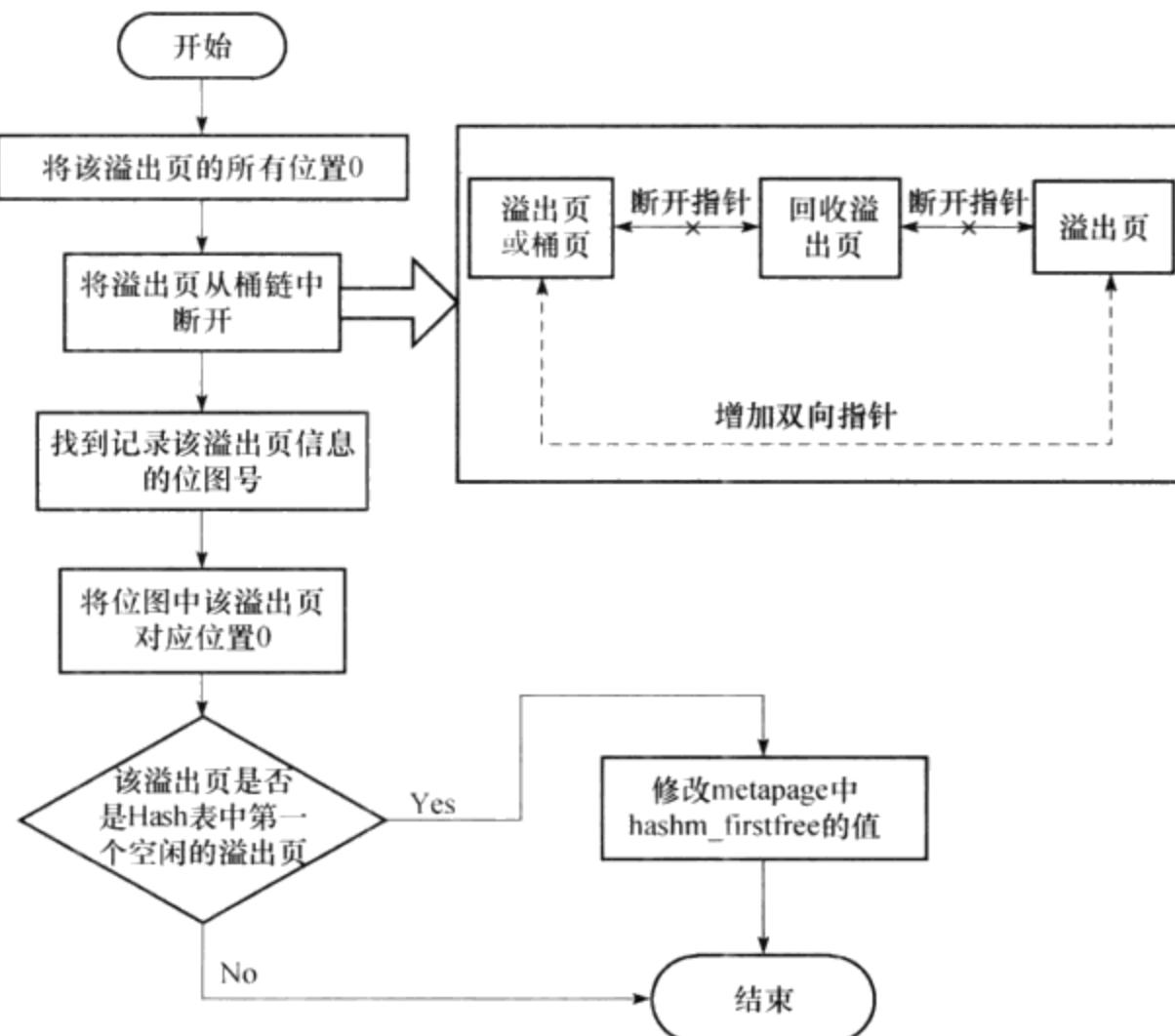


图 4-21 溢出页回收流程

假如有两个进程同时调用此函数，即它们都在同一个桶链上释放溢出页，这将有可能造成

进行断开桶链操作时出错，导致桶链的前后链接不一致，故多个进程在同一个桶链上释放溢出页是不能并发执行的，这就是为什么要求调用此函数前必须拥有此溢出页所在的桶上的排他 lmgr 锁的缘故。

4. Hash 表的扩展

每次插入后，都用当前的记录总数 r 和当前的桶数目 n 相除计算 r/n ，若比率太大，就对 Hash 表进行扩展，即增加一个桶到 Hash 表中，新增加的桶和发生插入的桶之间没有任何必然的联系。如果新加入的桶号的二进制表示为 $1a_2a_3\cdots\cdots a_i$ ，那么就试图分裂桶号为 $0a_2a_3\cdots\cdots a_i$ 的桶中的元组。如图 4-22 所示。

以下对 PostgreSQL 8.4.1 中与表扩展相关的函数进行分析。

1) 函数 `_hash_expandtable`: 该函数试图增加一个新桶来扩展 Hash 表。其执行流程如图 4-23 所示。

- 在修改完 metapage 的信息之后，进行桶的分裂之前，程序释放了 metapage 上的“lmgr lock”，这使得桶的分裂操作可以和其他操作并发执行，特别是多个对桶的分裂操作可以并发执行。
- 如果未能取得分裂桶上的“lmgr lock”，则放弃增加新桶来扩展 Hash 表的操作。因为如果继续等待分裂桶上的“lmgr lock”，而某个拥有分裂桶上的“lmgr lock”的进程碰巧也在等待本进程拥有的某个锁，这将会导致死锁。

2) 函数 `_hash_splitbucket`: 该函数将桶号为 obucket 的旧桶中的一部分元组分裂到桶号为 nbucket 的新桶里。执行流程如图 4-24 所示。

3) 函数: `_hash_squeezebucket`: 该函数压缩桶中的元组，使桶中的元组更加紧凑。图 4-25 给出了该函数的执行过程。其中要注意：

- 在压缩桶中元组时，从桶的最后一页开始顺序向前进行，对扫描到的每一个元组都把它移动到桶的前面。具体来说，有两类页面，分别是 rpage 和 wpage。开始时，rpage 是桶链的最后一页，wpage 是桶链的第一页，从 rpage 中读取每一个元组，并一次把它们写到 wpage 中，并将它们在 rpage 中的记录删除。若 rpage 中的所有元组都已扫描完毕，则令 rpage 的前一页为 rpage，并对其继续扫描；若 wpage 中已经没有足够的空间进行插入，则令 wpage 的下一页为 wpage，并将元组插入其中。这样，rpage 从桶链的最后一页开始向前移，wpage 从桶链的第一页开始向后移，直到 rpage 与 wpage 是同一个桶页为止，程序结束。
- 当把位于桶链里后面的元组都移到前面后，可能会出现一些不含任何元组的溢出页，把这些溢出页回收，并在位图中标识其可用，这是调用了 `_hash_freeovflpage` 函数完成的。
- 程序结束时，桶链上的所有页都不空，除非一开始整个桶链就是空的。
- 程序的调用者必须持有整个桶链的“lmgr lock”，以防止某些并发执行的进程访问本桶链，造成程序出错。

`_hash_squeezebucket` 的流程如图 4-26 所示。

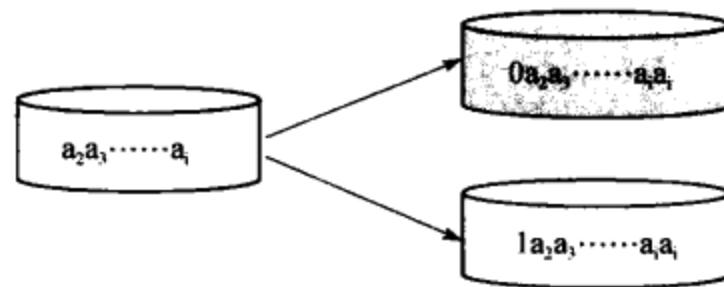


图 4-22 Hash 表扩展示意图

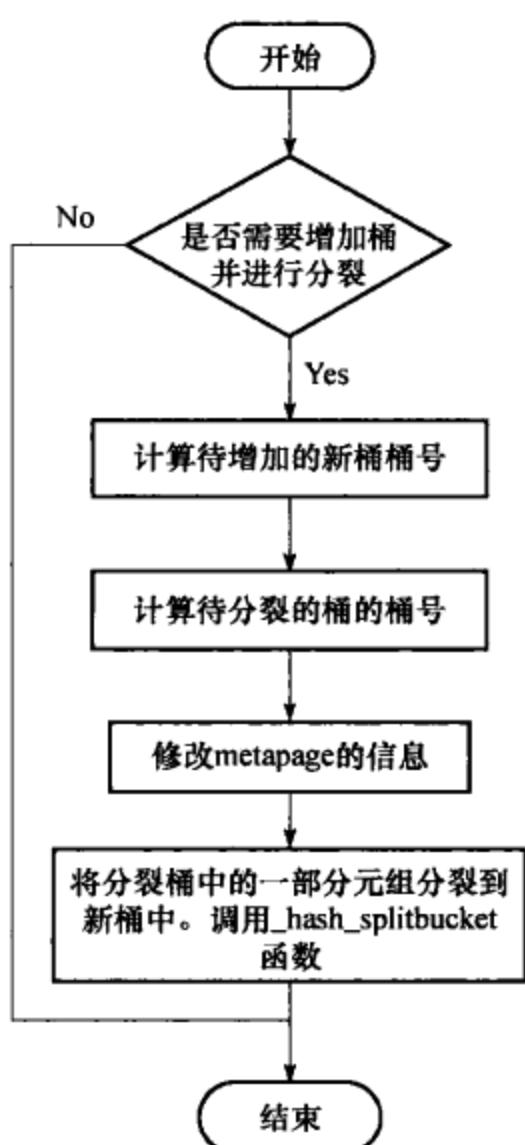


图 4-23 Hash 表扩展流程

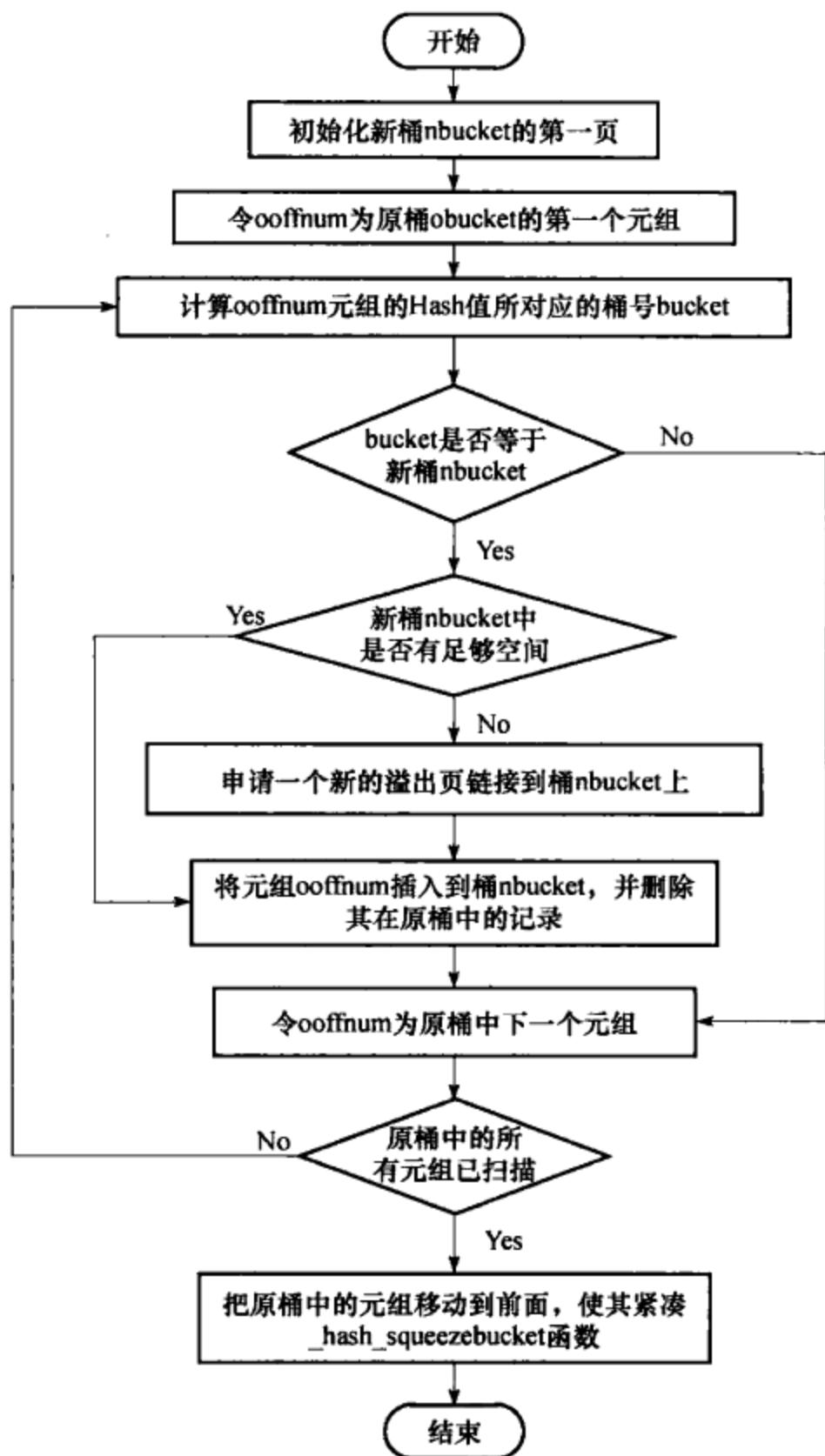


图 4-24 桶的分裂流程图

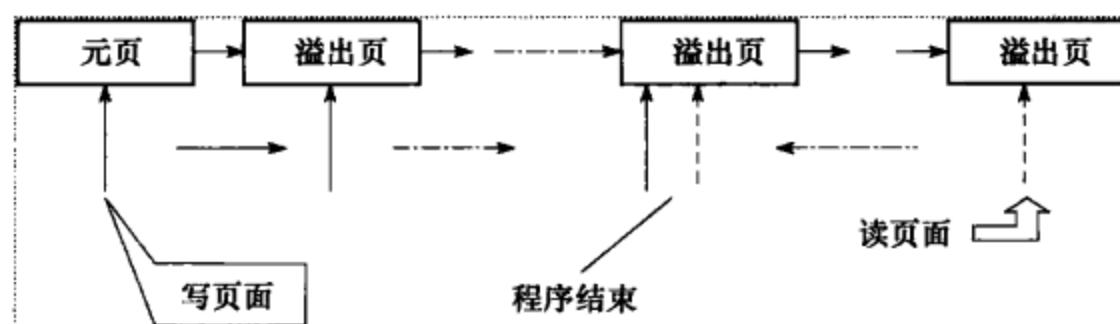
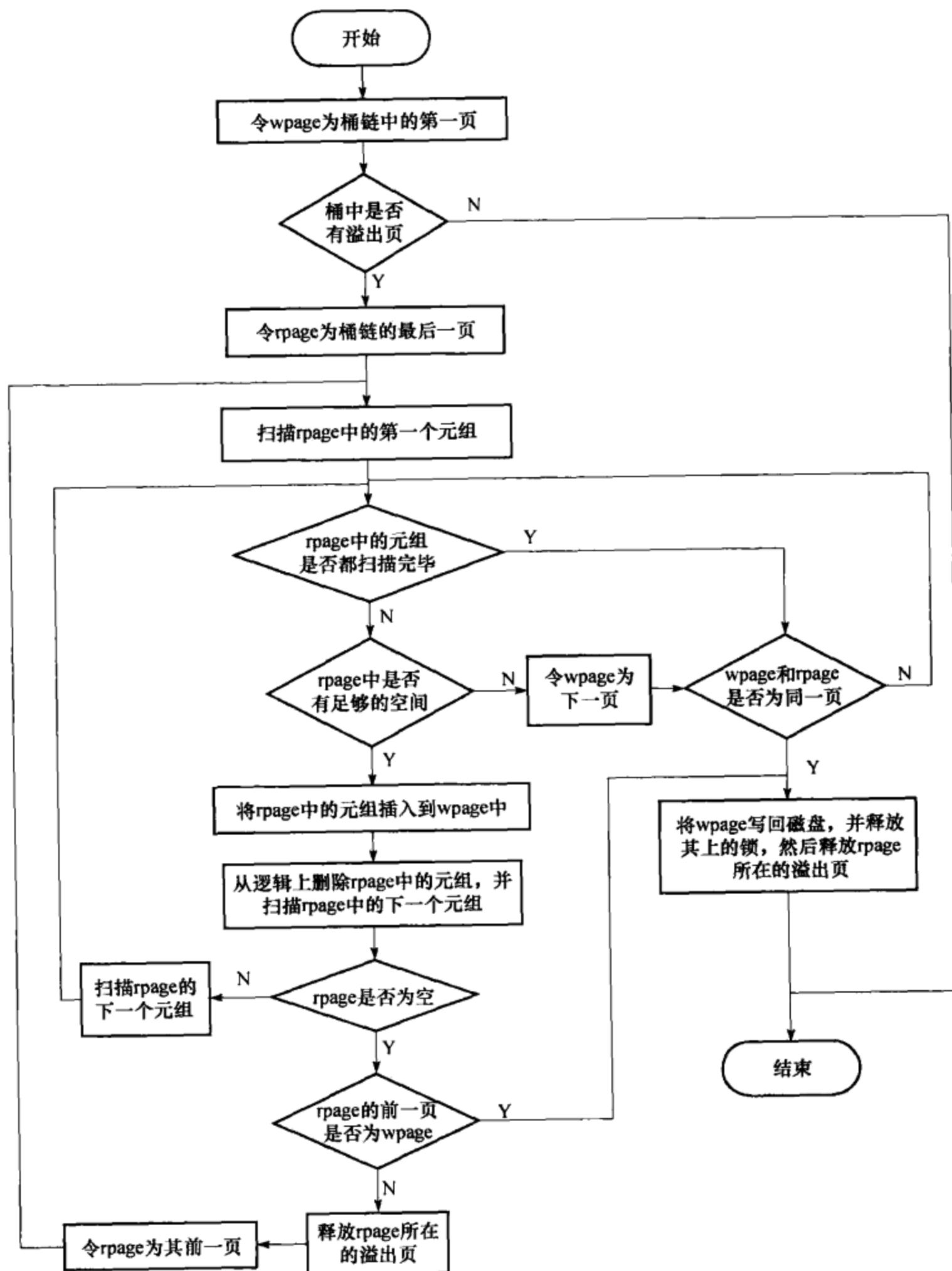


图 4-25 _hash_squeezebucket 执行流程

图 4-26 函数 `function_hash_squeezebucket` 实现流程

如果 rpage 的前一页就是 wpage，那么 wpage 上的写锁将会阻止函数 `_hash_freeovflpage` 释放 rpage 所在的溢出页，因为当对 wpage 加写锁的时候，进程是不能修改 wpage 所在页的指向桶链中下一页所在块的指针的（详见对函数 `_hash_freeovflpage` 的代码分析），所以在这里必须分情况讨论。

4.4 GiST 索引

GiST (Generalized Search Tree, 通用搜索树) 是一种平衡的、树状结构的访问方法。它在系统中相当于一个基础的模板，几乎可以使用它实现任意索引模式。B-trees、R-trees 和许多其他的索引模式都可以用 GiST 实现。它可以建立一种可扩展的索引结构，包括数据类型和查询谓词的扩展。

GiST 允许用户（并非数据库专家）开发自己的数据类型，并通过相应的访问方法来在该数据类型上使用 GiST 索引。通常，实现一种新的索引访问方法意味着大量艰苦的工作。因为必须理解数据库的内部工作机制，比如锁的机制和预写式日志。GiST 接口提供了一个高层的抽象，只要求访问方法的实现者实现被访问数据类型的语义。GiST 层本身会处理并发、日志和搜索树结构的任务。

GiST 的代码分布在 `\src\backend\access\gist` 目录下，包括了 GiST 索引的创建、查找、删除等代码。

4.4.1 GiST 的扩展性

PostgreSQL 支持可以扩展的 B-Tree 等标准搜索树，但不要把 GiST 的扩展性和其他标准搜索树的扩展性混淆在一起，比如它们所能处理的数据类型等方面。例如，B-Tree 索引支持对多种数据类型创建索引，但只支持范围查询。

这就是说，可以用 PostgreSQL 在多种数据类型上建立 B-Tree。但是 B-Tree 只支持范围谓词 (`<`, `=`, `>`)。所以，如果用 PostgreSQL 的 B-Tree 索引一个图像集，那么就只能发出类似“图像 x 和图像 y 相等吗”或“图像 x 是不是比图像 y 小”、“图像 x 是否大于图像 y”这样的查询。这些查询根据在这个环境下定义的“等于”、“小于”、“大于”的含义可能有意义，也有可能没有意义。

而使用一个基于 GiST 的索引，可以创建一些方法来发出和数据类型所处领域相关的问题，比如“找出所有马的图像”或者“找出所有曝光过度的图像”。

GiST 可以建立一种可扩展的索引结构，包括数据类型和查询谓词的扩展。这种结构支持研究人员快速为新的数据类型开发索引方法，其特点是在扩展数据类型的同时对谓词进行相应的扩展。

举例来说，颜色无法绝对排序，但是可以定义 `(red redthan blue)`、`(blue redthan green)` 这样的谓词，在扩展了数据类型的同时增加了谓词 `redthan`，当然数据类型可以是一组数据（如扩展为 R-Tree）。

4.4.2 GiST 索引的组织结构

GiST 是一棵平衡树，除根节点的子树数目在 $2 \sim M$ 之间外，每个节点的子树数目在 $k \cdot M \sim M$ 之间，常量 k 称作该树的最小填充因子，满足 $2/M \leq k \leq 1/2$ ， M 为一个节点可以容纳索引项的最大数目。

索引项形式为 (p, ptr) ，其中 p 是搜索的谓词。在叶子节点中， ptr 为指向数据库中某一元组的指针；而在非叶子节点中， ptr 为指向其子树节点的指针。

一个典型的 GiST 结构如图 4-27 所示。

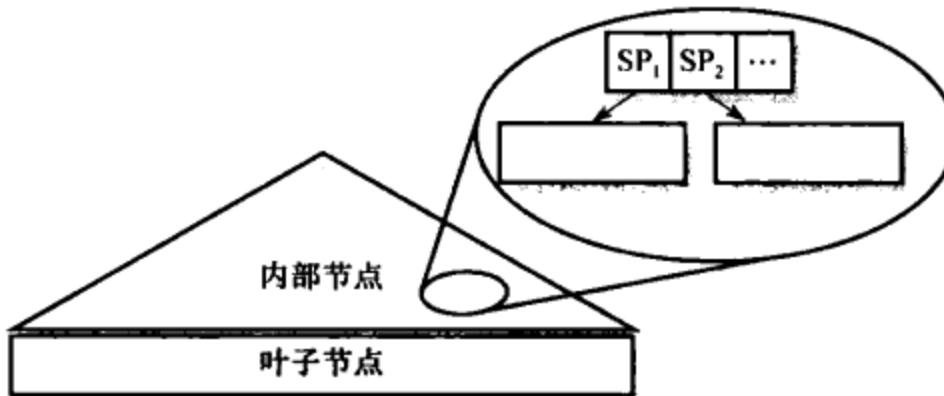


图 4-27 基本的 GiST 结构^①

在图 4-27 中， SP_1 、 SP_2 （subtree predicates）是指用于分隔数据的谓词。可以看出，GiST 的结构和 B-Tree 索引的结构有一定的相似性。

GiST 内置实现了索引项查询、插入和删除等算法。用户通过定义索引项并提供与索引项管理有关的方法，便可以实现某一特定的索引结构。这些方法包括：

- 1) Consistent (E, q)：对于给定的索引项 $E (p, \text{ptr})$ 和查询谓词 q ，判断索引项 E 是否与查询谓词 q 匹配，若肯定不能匹配，则返回假；否则，返回真。
- 2) Union (P)：对于给定集合 P 中索引项 $(p_1, \text{ptr}_1), \dots, (p_n, \text{ptr}_n)$ ，返回谓词 r ，满足 $(p_1 \vee \dots \vee p_n) \rightarrow r$ ， r 代表了 ptr_1 到 ptr_n 所指向的所有记录。
- 3) Same (E_1, E_2)：如果两个条目相同，返回真，否则返回假。
- 4) Penalty (E_1, E_2)：给定两个索引项 $E_1 (p_1, \text{ptr}_1)$ 和 $E_2 (p_2, \text{ptr}_2)$ ，返回值为将 E_2 插入到以 E_1 为根的子树中时的惩罚值。这用来辅助 Split 和 Insert 算法，有利于将一个记录对应的索引项插入到最适合的位置（该插入引起的惩罚值最小）。
- 5) PickSplit (P)：对于包含 $M + 1$ 个索引项 (p, ptr) 的集合 P 而言，将 P 划分为两个集合 P_1 和 P_2 ，每个集合至少包含 $k \cdot M$ 个索引项。该方法确定了节点分裂的原则。
- 6) Compress (E)：对于给定的索引项 (p, ptr) ，返回 (π, ptr) ， π 为 p 的压缩形式。
- 7) Decompress (E)：对于给定的压缩表示索引项 (π, ptr) ， $\pi = \text{Compress} (p)$ ，返回 (r, ptr) ， r 满足 $p \rightarrow r$ 。（可能为有损压缩，并不要求 $p \leftarrow r$ ）。

^① KORNACKER M. Access Methods for the Next-Generation Database Systems. University of California, Berkeley, 2001.

4.4.3 GiST 索引的实现

由于 GiST 已内置实现了索引项的创建、查找和删除等操作，而这些操作都依赖于 4.4.2 节中介绍的 7 个方法，因此用户只需要实现这 7 种方法。而索引的创建、查找等 PostgreSQL 自动会进行扩展。下面将对数据库如何使用这 7 种方法来实现对数据库的操作进行分析。

1. GiST 索引创建

Postgres 中 GiST 索引的创建由函数 `gistbuild` 完成，实现流程如图 4-28 所示。

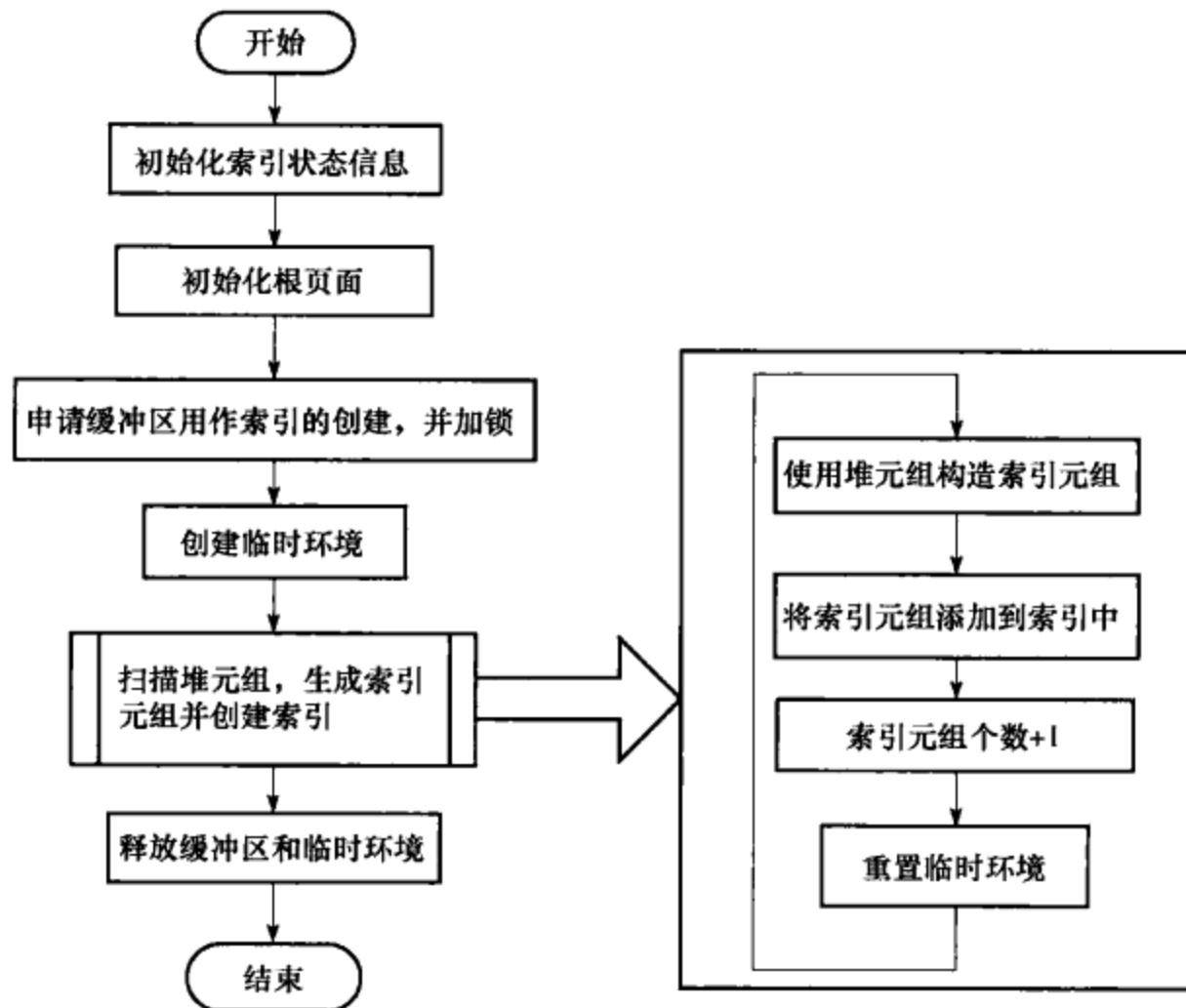


图 4-28 `gistbuild` 函数的实现流程

在索引创建过程中，索引元组的插入在函数 `gistdoinsert` 中完成，该函数将一个索引元组插入到索引结构中。其实现过程是先从搜索树的根节点开始遍历，找到插入代价最小（由 `Penalty` 方法实现）的叶子节点进行插入。若叶子节点已满，插入新索引项会导致叶子节点的分裂，也可能造成分裂的向上传播。分裂时将调用用户定义的 `PickSplit` 方法来决定新老节点中索引项的布局。而在向上更新描述谓词时，将会调用 `Union` 方法来确定父节点相应索引项中的描述谓词。其执行流程如图 4-29 所示。

在插入索引元组的过程中，会生成一个类型为 `GISTInsertStack` 的栈结构，用于记录当前插入节点及其父节点的相关信息。通过该结构能够找到当前插入节点的父节点，当插入节点发生分裂需要插入新节点到父节点中，或更新父节点的信息时，通过该栈保存的信息能够依次向上层进行调整。其结构如数据结构 4.12 所示。

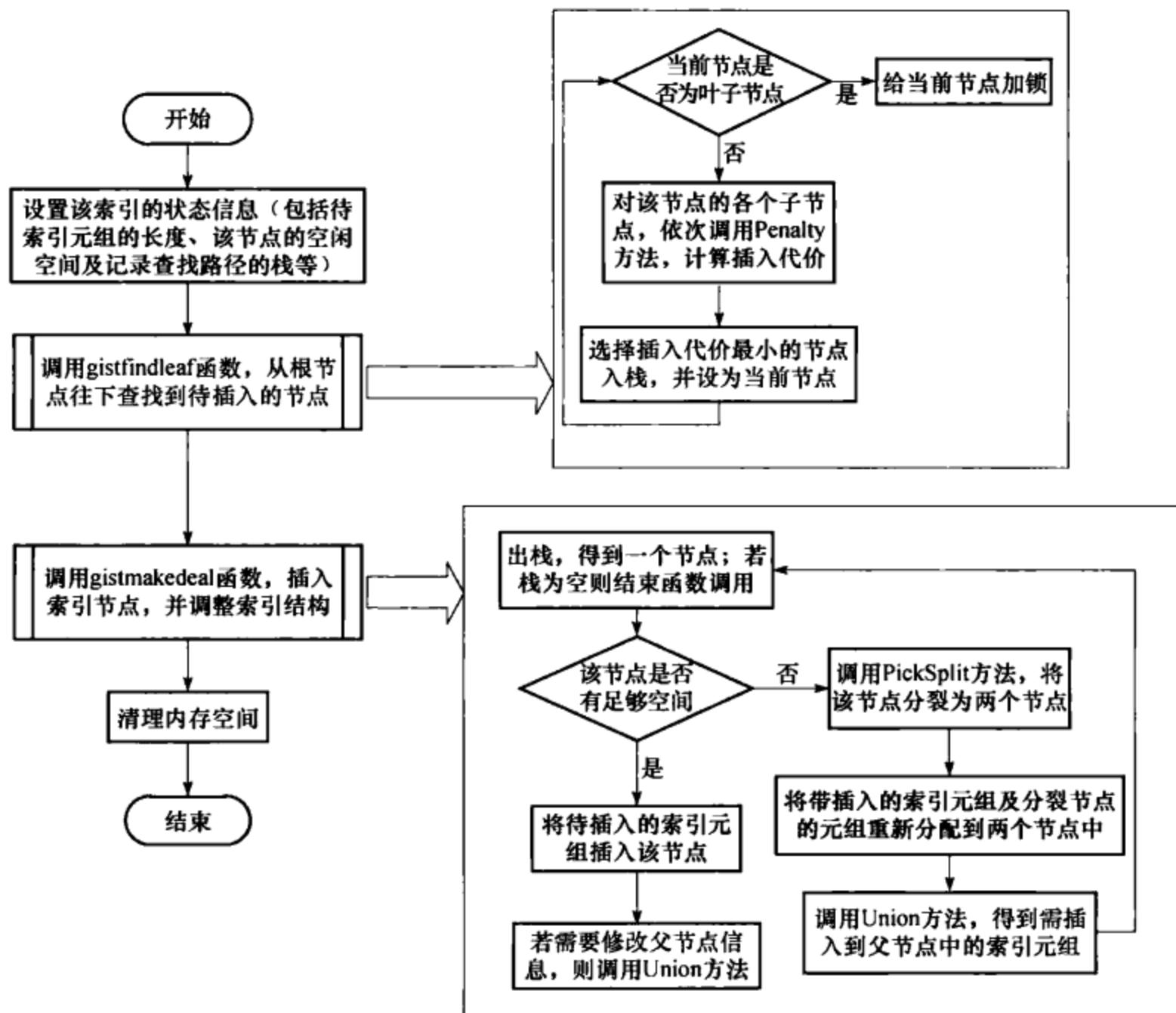


图 4-29 gistdoinsert 函数执行流程图

数据结构 4.12 GISTInsertStack

```

typedef struct GISTInsertStack
{
    BlockNumber          blkno;           // 当前页面的块号
    Buffer               buffer;          // 当前页面所在的缓冲区
    Page                 page;            // 当前页面的页面号
    GistNSN              lsn;             // 用于确定页面是否更新或分裂
    OffsetNumber         childoffnum;     // 子节点的偏移
    struct GISTInsertStack *parent;        // 指向父节点
    struct GISTInsertStack *child;         // 指向子节点
    struct GISTInsertStack *next;          // 指向下一个入栈的节点，用于路径查找
}GISTInsertStack;

```

2. GiST 索引查询

在 PostgreSQL 的代码中，GiST 索引的查询流程与 B-Tree 索引类似，从根节点开始按深度优先原则自上而下进行检索：

- 若当前节点 R 是内部节点，检查 R 上每个索引项 E 是否与检索谓词 q 相符合，对于满足 Consistent (E, q) 的索引项，递归向下检索以 E.ptr 为根的子树。
- 若 R 是叶子节点，检查 R 上每个索引项 E 是否满足 Consistent (E, q)；对于满足 Consistent (E, q) 的索引项，则通过 E.ptr 取得相应记录与 q 进行准确匹配；将匹配成功的记录放入结果集中。

GiST 索引查询主要通过函数 `gistnext` 来实现，该函数通过从上往下搜索索引结构，使用上面的方法找到结果。在扫描的过程中会生成一个类型为 `GISTSearchStack` 栈结构，用于保存扫描过程中内部节点里面满足前述的 `Consistent` 方法的节点：如果扫描的是内部节点且找到满足 `Consistent` 方法的元组，则将该元组指向的下一层的节点入栈。其数据结构见数据结构 4.13 所示。

数据结构 4.13 GISTSearchStack 结构的定义

```
typedef struct GISTSearchStack
{
    struct GISTSearchStack *next;           // 指向下一个元素
    BlockNumber            block;           // 本节点对应的块号
    GistNSN                lsn;             // 确认该节点是否发生更改
    GistNSN                parentlsn;        // 确认该节点是否发生分裂
}GISTSearchStack;
```

执行完 `gistnext` 函数后，即得到结果集。`gistnext` 函数的执行流程如图 4-30 所示。

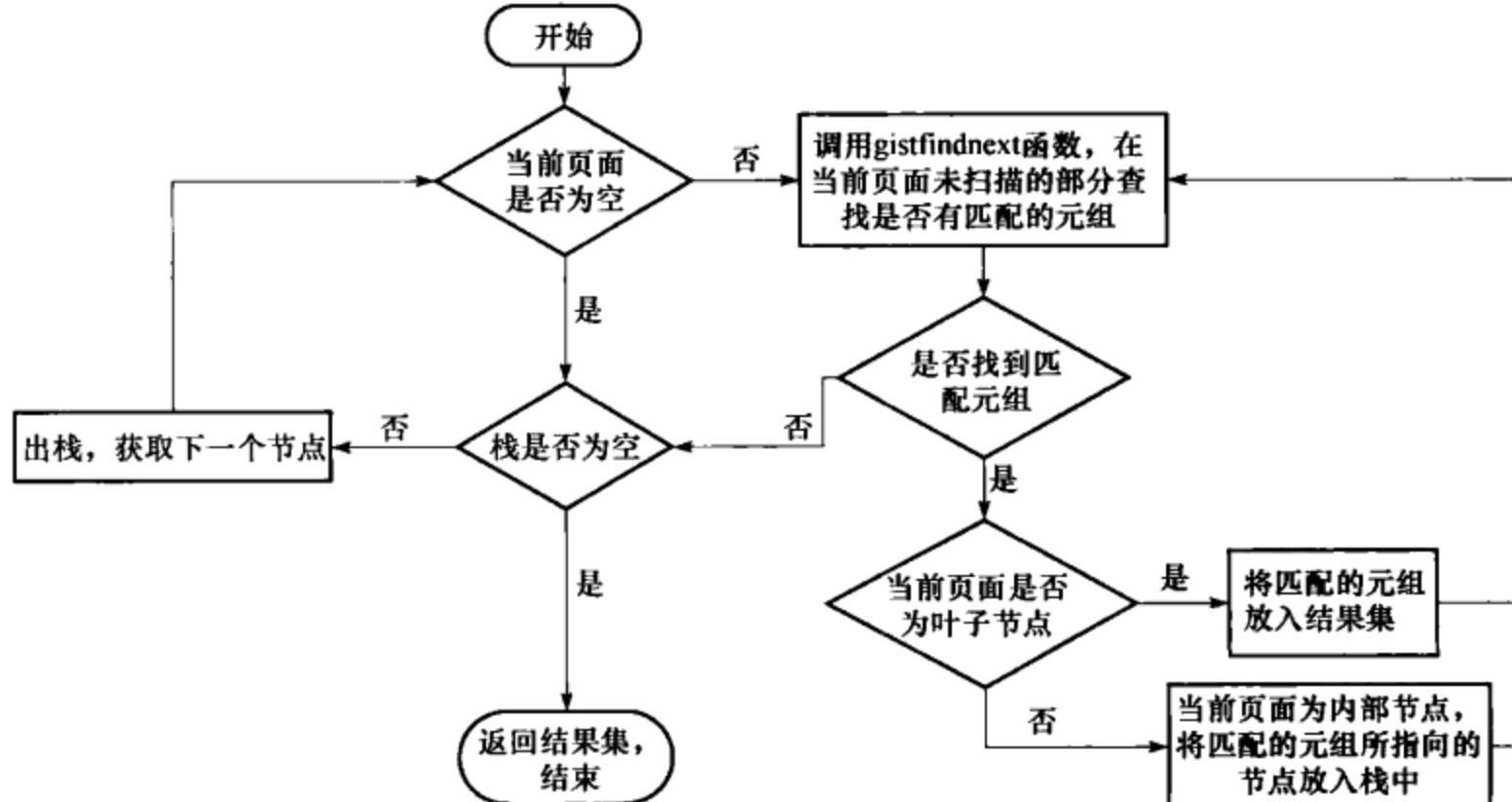


图 4-30 `gistnext` 函数执行流程图

在图 4-30 中，调用了 `gistfindnext` 函数在当前页面查找匹配的元组，最终通过调用 4.4.2 节中介绍的 `Consistent` 方法来确定是否匹配。

3. GiST 索引删除

与 B-Tree 索引类似，当表元组被删除时，GiST 索引中与之对应的索引元组不是立即被删除，而是由 VACUUM 操作来批量完成无效索引元组的清除。若是普通的 VACUUM (Lazy Vacuum) 则只更新 FSM (“空闲空间映射表”，参见 3.2.4 节)，实际并不删除索引结构中索引元组；如果是 Full VACUUM，则需要修改索引页面中的元组信息。

在删除过程中，首先找到叶子节点中需删除的索引项，然后从叶子节点往上回溯更新索引，若删除索引项后，存在空节点，则删除该节点。这项工作主要由 `gistvacuumcleanup` 函数完成。

4.4.4 GiST 索引实例

下面以 GiST 的 R 树扩展为例给出简单说明。设待索引的数据是二维多边形，用最小外接矩形表示。关键字 $(X^{ul}, Y^{ul}, X^{lr}, Y^{lr})$ 中 (X^{ul}, Y^{ul}) 为外接矩形的左上角， (X^{lr}, Y^{lr}) 为外接矩形的右下角。

假设在此支持的操作有 Contains (包含)、Overlap (部分重叠) 和 Equal (相等)。具体解释如下：

- Contains $((x^1ul, y^1ul, x^1lr, y^1lr), (x^2ul, y^2ul, x^2lr, y^2lr))$

如果 $(x^1lr \geq x^2lr) \wedge (x^1ul \leq x^2ul) \wedge (y^1lr \leq y^2lr) \wedge (y^1ul \geq y^2ul)$ ，则返回真，否则返回假。如图 4-31 所示。

- Overlap $((x^1ul, y^1ul, x^1lr, y^1lr), (x^2ul, y^2ul, x^2lr, y^2lr))$

如果 $(x^1ul \leq x^2lr) \wedge (x^2ul \leq x^1lr) \wedge (y^1lr \leq y^2ul) \wedge (y^2lr \leq y^1ul)$ ，则返回真，表示两者重叠，否则返回假。如图 4-32 所示。

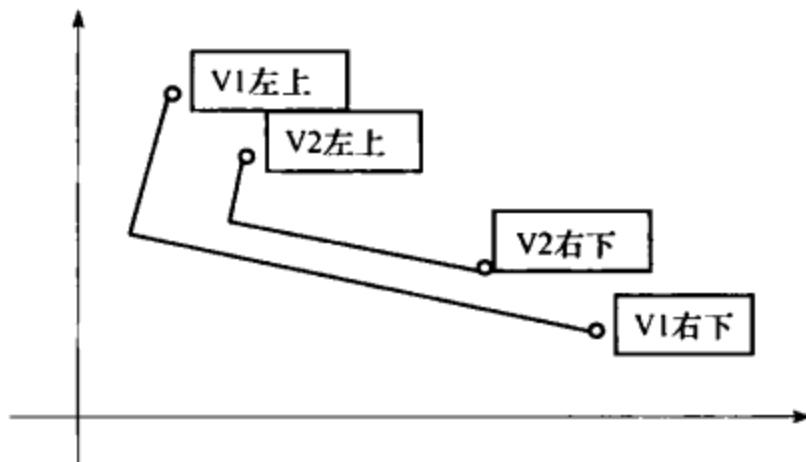


图 4-31 GiST 索引实例图一

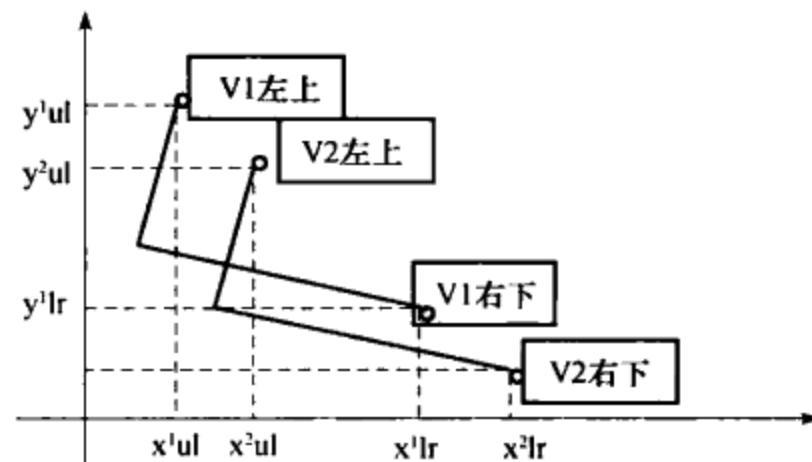


图 4-32 GiST 索引实例图二

- Equal $((x^1ul, y^1ul, x^1lr, y^1lr), (x^2ul, y^2ul, x^2lr, y^2lr))$

如果 $(x^1lr = x^2lr) \wedge (x^1ul = x^2ul) \wedge (y^1lr = y^2lr) \wedge (y^1ul = y^2ul)$ 则返回真，表示两者相等，否则返回假。

对于以上定义的谓词，GiST 的 6 种方法可扩展如下：

(1) Consistent (E, q)

给定一个索引项 $E = (P, \text{ptr})$ ，其中 $P = \text{Contains}$ 为 $((X^1ul, Y^1ul, X^1lr, Y^1lr), v)$ ， q 是 `Contains`、`Overlap` 或者 `Equal`，参数为 $(X^2ul, Y^2ul, X^2lr, Y^2lr)$ 。对于这些查询，如果 $\text{Overlap}((X^1ul, Y^1ul, X^1lr, Y^1lr), (X^2ul, Y^2ul, X^2lr, Y^2lr))$ 为假，则返回假。否则返回真。

(2) `Union (E1 = ((X1ul, Y1ul, X1lr, Y1lr), ptr), …, En = ((Xnul, Ynul, Xnlr, Ynlr), ptr))` 方法返回 $(\text{MIN} (X^1\text{ul}, \dots, X^n\text{ul}), \text{MAX} (Y^1\text{ul}, \dots, Y^n\text{ul}), \text{MAX} (X^1\text{lr}, \dots, X^n\text{lr}), \text{MIN} (Y^1\text{lr}, \dots, Y^n\text{lr}))$ 。

(3) `Compress (E = (P, ptr))`

形成一个多边形的最小包围矩形。假定多边形以线段集的形式给出, $li = (X^i_1, Y^i_1, X^i_2, Y^i_2)$, 令 $\pi = (\text{MIN} (X^i\text{ul}), \text{MAX} (Y^i\text{ul}), \text{MAX} (X^i\text{lr}), \text{MIN} (Y^i\text{lr}))$, 返回 (π, ptr) 。

(4) `Decompress (E = (X1ul, Y1ul, X1lr, Y1lr), ptr))`

可简单实现为识别函数, 即直接返回 E。

(5) `Penalty (E1, E2)`

给定 $E1 = (p1, \text{ptr1})$, $E2 = (p2, \text{ptr2})$, 通常表示将 E2 插入到以 E1 为根的子树时从 E1.P 到 `Union ((E1, E2))` 的某种增量, 可以计算 $q = \text{Union} (E1, E2)$, 并返回 `area (q) - area (E1.P)`。

(6) `PickSplit (P)`

对于包含 $M + 1$ 个索引项 (p, ptr) 的集合 P 而言, 把 P 划分为两个索引项集合 P1、P2, 每个集合包含的索引项数目最少为 $k \cdot M$ 个。

PostgreSQL 8.4.1 内建了对 `box` (矩形)、`polygon` (多边形) 和 `circle` (圆形) 等数据类型的 GiST 索引的支持, 可以在 `postgresql-8.4.1\src\backend\access\gist\gistproc.c` 文件中查看到相关代码。这些数据类型可以直接创建 GiST 索引, 其他数据类型如果需要创建 GiST 紴引, 则需要用户手动添加。

在 `postgresql-8.4.1\contrib` 文件夹下提供了很多扩展的模块, 其中也有关于 GiST 的模块, 如 `cube` 文件夹下是用于多维立方体的 GiST 紴引; `ltree` 文件夹下是用于树状结构的 GiST 紴引等。编译这些模块, 就可以使数据库支持该种数据类型, 同时可以在该类型上创建 GiST 紴引。下面就以多维立方体 `cube` 文件夹下的扩展模块为例, 讲解如何将其添加到数据库系统中。

查看 `cube` 文件夹下的文件可以看出, 这些文件已经编写好了增加一种数据类型、相应的操作符、该数据类型的支持函数、以及 GiST 紴引对该数据类型的支持。那么, 只需要将这些信息编译进数据库即可。主要包括以下几个步骤:

1) 执行 `make` 编译 C 代码得到 `cube.so` 链接库。这一步主要完成的工作是将该种数据类型所需要的支持函数、GiST 紴引需要实现的 7 种函数方法编译到链接库里, 供数据库调用。

2) 执行 `make install`, 这时系统就会自动将刚才得到的链接库添加到数据库中, 然后执行新生成 `cube.sql`, 添加相关命令到数据库系统中。

当执行完成 `make` 命令之后, 可以看到在当前目录下生成了一个文件 `cube.sql`, 该文件是修改了 `cube.sql.in` 得到的 (将 `make` 得到的 `cube.so` 链接库的目录写入了)。若读者感兴趣, 也可以自己手动将 `cube.sql.in` 中的命令添加到数据库系统中, 只需在执行 `make` 之后执行以下命令:

3) 修改 `cube.sql.in` 文件。该文件的主要功能是向数据库系统中添加数据类型 (`pg_type` 系统表)、添加该数据类型及索引的操作符 (`pg_opclass`、`pg_amop`、`pg_operator` 系统表) 以及支持函数 (`pa_amproc` 系统表), 由于这些函数都要链接到刚才编译得到的 `cube.so`, 所以需要修改 `cube.sql.in` 文件里 `cube.so` 文件的路径。

4) 登录一个数据库, 然后执行 `cube.sql.in` 语句, PostgreSQL 会自动执行 `cube.sql.in` 里所有的语句, 将相关信息添加到数据库系统中。

通过以上四个步骤就完成了向数据库系统里添加 cube 数据类型、操作符、相关函数及其对 GiST 索引类型的支持，下面就可以像使用数据库自带的数据类型一样使用 cube 类型。如首先创建一个表：

```
CREATE TABLE test_cube (
    name varchar,
    cub cube);
```

插入一些数据后，即可在该表的 cub 字段上创建 GiST 索引，并使用该索引进行查找：

```
CREATE INDEX test_cube_ix ON test_cube USING gist (cub);
SELECT * FROM test_cube WHERE cub && '(3000,1000), (0,0)';
```

如果想了解具体添加了哪些信息到数据库中，可以查看 cube\cube.sql.in 文件。

4.5 GIN 索引

GIN (Generalized Inverted Index, 通用倒排索引) 是一个存储对 (key, posting list) 集合的索引结构，其中“key”是一个键值，而“posting list”是一组出现过“key”的位置。如 (“hello”, “14: 17, 23: 1, ……”) 中，“hello”表示一个键值，而“14: 17, 23: 1, ……”则表示“hello”这个键值出现的位置。每一个位置又由“元组 ID: 位置”来表示，位置“14: 17”说明“hello”在第 14 号元组的被索引属性中第 17 个位置出现。每一个被索引的属性值都可能包含多个键值，因此同一个元组的 ID 可能会出现在多个“posting list”中。通过这种索引结构可以快速地查找到包含指定关键字的元组，因此 GIN 索引特别适合于支持全文搜索。而开发 GIN 索引的主要目的就是为了让 PostgreSQL 能支持功能高度可扩展的全文搜索。

4.5.1 GIN 索引的扩展性

GIN 索引具有很好的可扩展性，允许在开发自定义数据类型时由该数据类型的领域专家（而不是数据库专家）设计适当的访问方法，这些访问方法只需要考虑对于数据类型本身语义的处理，而 GIN 索引自身可以处理并发操作、记录日志、搜索树结构等操作。

定义一个 GIN 访问方法所要做的就是实现五个用户定义的方法，这些方法定义了键值、键值与键值之间的关系、被索引值、能够使用索引的查询及部分匹配。简而言之，GIN 结合了扩展性、普遍性、代码重用、清晰的接口等特点。

一个 GIN 索引需要实现的五个方法如下：

1) compare 方法：比较两个键值 a 和 b，然后返回一个整数值，返回负值表示 a 小于 b，返回 0 表示 a 等于 b，返回正值表示 a 大于 b。其原型如下：

```
int compare(Datum a, Datum b)
```

2) extractValue 方法：根据参数 inputValue 生成一个键值数组，并返回其指针，键值数组中元素的个数存放在另一个参数 nkeys 中。其原型如下：

```
Datum* extractValue(Datum inputValue, uint32 *nkeys)
```

3) extractQuery 方法：根据一个查询（由参数 query 给定）生成一个用于查询的键值数组，并返

回其指针。extractQuery 通过参数 n 指定的操作符策略号来决定 query 的数据类型以及需要提取的键值，返回键值数组的长度存放在 nkeys 参数中。如果 query 中不包含键值，则 extractQuery 将根据操作符的语义在 nkeys 中存储 0 或 -1。nkeys 值为 0 表示索引中所有值都能满足 query，因此将执行完全索引扫描，例如当 query 为空字符串时就会发生此种情况。当 nkeys 值为 -1 时表明索引中没有键值能匹配查询，因此可以跳过索引扫描。当索引支持部分匹配时，输出参数 pmatch 用于记录返回的键值数组中每一个键值是否要求部分匹配。extractQuery 将分配一个长度为 nkeys 的布尔型数组，将这个数组的指针通过 pmatch 输出，数组中的每一个元素都记录了键值数组中对应位置键值的部分匹配情况，如果为真就表示该键值要求部分匹配。输出参数 extra_data 用来向 consistent 和 comparePartial 方法传递用户定义需要的数据，extra_data 指向一个长度为 nkeys 的指针数组，extractQuery 可以在每一个数组元组所指向的内存空间存放任何数据。如果 extra_data 被设为非空，则 extra_data 指向的整个数组都将被传递给 consistent 方法，而其中适当的元素将被传递给 comparePartial 方法。extractQuery 的原型如下：

```
Datum *extractQuery(Datum query, int32 *nkeys,
    StrategyNumber n, bool **pmatch,
    Pointer **extra_data)
```

4) consistent 方法：该方法用于检查索引值是否满足查询。如果查询与索引满足策略号为 n 的操作符则返回真，如果返回真且输出参数 recheck 也设置为真则说明查询和索引可能满足，还需要进一步检查。check 数组的长度必须与先前由 extractQuery 为该查询返回的键值数量相同。如果索引值包含相应的查询，那么 check 数组中的每一个元素都是真。也就是说，如果 check [i] 为真，那么 extractQuery 返回的键值数组的第 i 个键值存在于索引值当中。如果 GIN 索引的比较是精确的，recheck 将被设置为假，否则 recheck 将被设置为真，通过索引找到的基表元组还需要进行是否满足操作符的检查。图 4-33 为一个 consistent 方法实现的实例，图中查询条件中的“&”符号表示“与”，而“|”符号则表示“或”。

```
bool consistent(bool check[], StrategyNumber n,
    Datum query, int32 nkeys,
    Pointer extra_data[], bool *recheck)
```

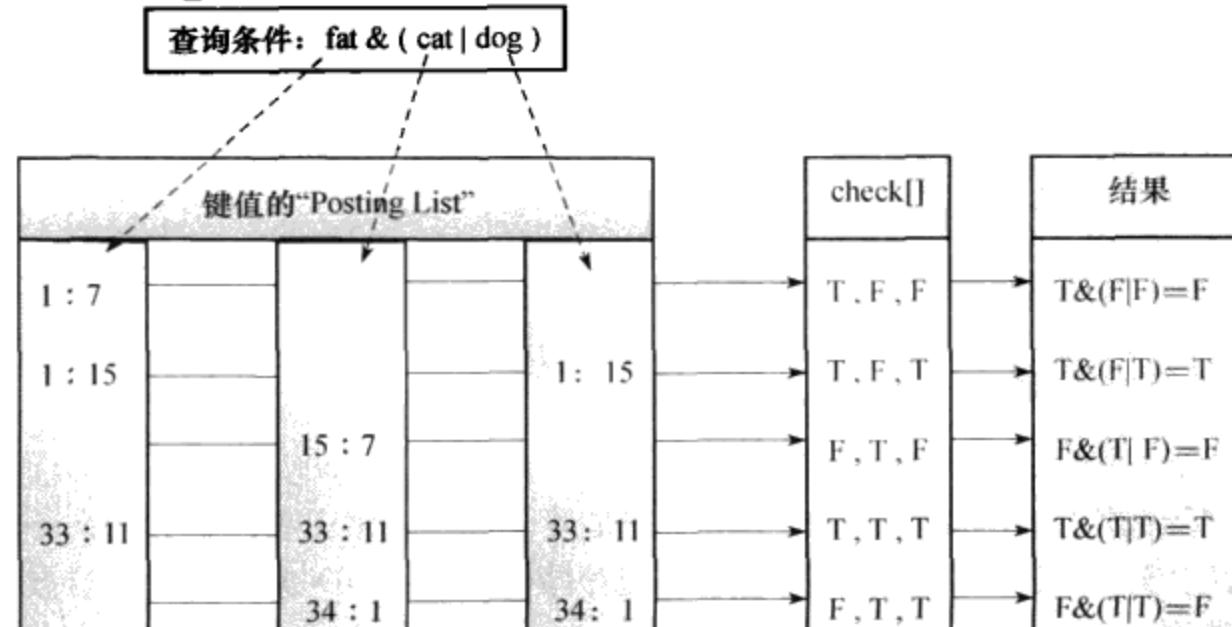


图 4-33 consistent 方法实例

5) comparePartial 方法：将部分匹配的查询与索引值进行比较，返回值为负值表示两者不匹配，但继续扫描索引；返回值为 0 表示两者匹配；返回值为正值表示停止扫描。n 是操作符的策略号；extra_data 是由函数 extractQuery 传递过来的数组。这个函数是在 PostgreSQL 8.4.1 中新增的函数。该方法的原型如下：

```
int comparePartial(Datum partial_key, Datum key,
    StrategyNumber n, Pointer extra_data)
```

如果要在 PostgreSQL 中添加一种新的数据类型并且让 GIN 索引能够支持该数据类型，则需要完成以下步骤：

1) 添加数据类型：

① 实现新数据类型的输入输出函数，并通过“CREATE OR REPLACE FUNCTION”语句将其注册到数据库内部，这一步操作将会在 pg_proc 系统表中为每一个函数增加一个元组。

② 利用上一步注册的输入输出函数，通过“CREATE TYPE”语句创建数据类型，这一步会在 pg_type 系统表中为新数据类型增加一个元组，并通过元组中的属性值与上一步创建的输入输出函数关联起来。

③ 为新数据类型实现并注册各种操作符所需的函数，然后通过“CREATE OPERATOR”语句为新数据类型创建操作符。

2) 为新数据类型实现 GIN 索引所需要的 5 种支持函数（compare、extractValue 等），并通过“CREATE OR REPLACE FUNCTION”语句将这些函数注册到数据库内部。

3) 用“CREATE OPERATOR CLASS”语句为新的数据类型创建一个操作符类，该 SQL 语句需指定 GIN 索引所需要的 5 个支持函数。

完成上述步骤之后，就可以用以下 SQL 语句在新数据类型上创建 GIN 索引了：

```
CREATE INDEX index_name ON table_name(column_name) USING GIN;
```

其中，index_name 是要创建的索引名称，table_name 是要创建 GIN 索引的基表名称，column_name 是要创建 GIN 索引的基表属性名。

4.5.2 GIN 索引的组织结构

GIN 索引包括以下四个部分：

- Entry：是 GIN 索引中的一个元素（可以认为是一个词位），Entry 也可以理解为 GIN 索引中的一个键值。
- Entry Tree：在一些 Entry 上构建的 B-Tree。
- Posting Tree：在一个 Entry 出现的物理位置上构建的 B-Tree。
- Posting List：一个 Entry 出现的物理位置的列表。

GIN 索引包含一个构建在 Entry（键值）上的 B - Tree 紴引（Entry Tree），“Entry Tree”的内部节点与普通 B - Tree 納引的内部节点一样，不同的是叶子节点中索引项的指针指向不同，普通 B - Tree 的叶子节点中索引项指针指向基表元组的物理存储位置，而 GIN 索引的叶子节点中索引项的指针指向位置分两种情况：

- 1) 如果该叶子节点中索引项内包含的键值（Entry）所出现的物理位置大于宏 TOAST_INDEX_

TARGET 所指定的值，则在这些物理位置上构建“Posting Tree”，然后将索引项指针指向“Posting Tree”的根节点。

2) 否则，叶子节点中索引项指针直接指向“Posting List”。

也就是说，如果某个 Entry 出现的位置较多（超过了 TOAST_INDEX_TARGET 的值），则在其出现的位置（也就是“Posting List”）上再创建一个 B-Tree 结构，以加快查找的速度。

一个 GIN 索引的组织结构示例如图 4-34 所示。

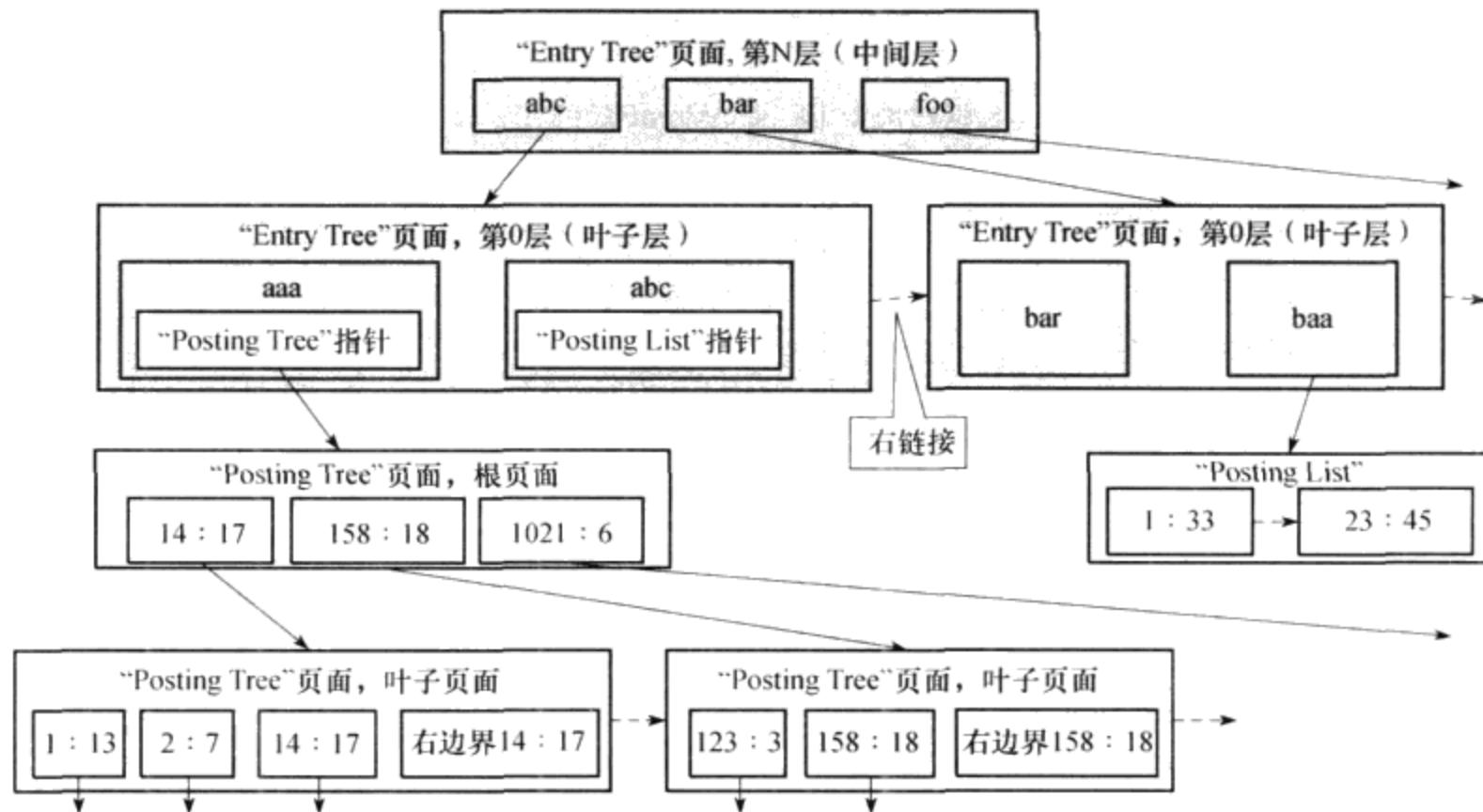


图 4-34 GIN 索引组织结示例

在实现上，“Entry Tree”的非叶子节点部分和 B-Tree 索引类似，都是以一个页面来组织索引元组。但“Entry Tree”叶子节点和普通的 B-Tree 索引页面结构有些区别：

- 若叶子节点中索引元组（IndexTuple，见数据结构 4.3）的位置信息采用“Posting List”形式存储，则在该索引元组键值之后的空间连续存放“Posting List”，这样可快速获取到该键值所有的位置信息。由于是连续存储（不需要访问其他节点），因此索引元组的 t_tid 字段的 ip_blkid 部分用于存储“Posting List”的长度（即该键值出现了多少次），而 ip_posid 部分用于记录该索引元组的大小（不包括“Posting List”的大小）。
- 若叶子节点中索引元组的位置信息采用“Posting Tree”形式存储，则该索引元组的 t_tid 字段的 ip_blkid 部分记录了“Posting Tree”根节点所在的磁盘块号，而 t_tid 字段的 ip_posid 部分则设置为 GIN_TREE_POSTING，用于与“Posting List”区分开来。

通过这种方式，巧妙地使用索引元组中的 t_tid 字段，使得对“Posting List”和“Posting Tree”的操作得到统一，根据 t_tid 字段中的 ip_posid 部分的取值可以区分“Posting List”和“Posting Tree”，然后根据两种不同的类型调用不同的方法来读取位置信息。PostgreSQL 提供了一系列的宏定义用于支持这些操作。

GIN 索引中共设置了 5 种页面类型，如表 4-8 所示。

表 4-8 GIN 索引页面类型

GIN 页面类型	描述	GIN 页面类型	描述
GIN_META	GIN 索引的元页	GIN_DELETE	被标识删除的页面
GIN_DATA	存放“Posting Tree”的页面	GIN_LIST	待建 GIN 索引页面
GIN_LEAF	“EntryTree”的叶子节点页面	GIN_LIST_FULLROW	被填满的 GIN_LIST 页面

其中，GIN 索引的元页主要用来管理索引创建过程中待索引的页面信息，其定义如数据结构 4.14 所示。

数据结构 4.14 GinMetaPageData

```
typedef struct GinMetaPageData
{
    BlockNumber head;           //待建索引页面链头
    BlockNumber tail;          //待建索引页面链尾
    uint32      tailFreeSize;   //尾页面空余空间大小
    BlockNumber nPendingPages;  //待建索引页面个数
    int64       nPendingHeapTuples; //所有待建索引页面中的索引元组个数
} GinMetaPageData;
```

在创建 GIN 索引时，首先会将所有所有待索引的 Entry 及其位置信息统计出来，并将这些信息存储在 GIN_LIST 页面中，GIN_LIST 页面将构成一个双向链表（也称为“Pending List”），在元页中通过 head、tail 两个字段分别指向这个链表的头部和尾部。在存储 Entry 及其位置信息时，先申请一个 GIN_LIST 页面，将其填满之后链接在“Pending List”尾部，然后再申请下一个 GIN_LIST 页面进行填充。因此在“Pending List”中，只有尾部的页面可能是不满的，其他页面都是被充满的，在元页中用 tailFreeSize 字段记录最尾部页面的剩余空间大小。在创建 GIN 索引结构时，会依次读取出“Pending List”中页面里的信息，并插入到索引结构中。

4.5.3 GIN 索引的操作

在 PostgreSQL 8.4.1 中，GIN 索引可以支持 tsvector 和所有内置数据类型的一维数组。另外，在源代码的 contrib 目录下还有 btree-gin、hstore、intarray 和 pg_trgm 四个子目录，每一个子目录中的代码都可以用来在数据库系统中增加一种新的数据类型并让 GIN 紴能支持它，其基本实现过程和 4.5.1 节所介绍的一致。本节将对 GIN 紴的基本操作进行分析。

1. GIN 紴的创建

GIN 紴的创建过程就是从基表中依次取出基表元组，然后从基表元组构建出若干个 Entry 及其“Posting List”，然后将这些 Entry 插入到 GIN 紴结构中。在实现中，GIN 紴的创建函数并不是构建出一个 Entry 就立即将其插入到 GIN 紴中，而是建立了一个“蓄积池”，构建出的 Entry 都将先放入“蓄积池”，当“蓄积池”填充到一定程度之后才会将其中缓存的 Entry 依次插入到 GIN 紴中去。

GIN 紴的创建函数使用了一个 GinBuildState 类型（数据结构 4.15）的变量来对整个创建过程

的状态加以控制。其中，ginstate 字段存放当前使用的 GIN 索引相关的 5 个支持函数（见 4.5.1 节）信息，这些信息是与要创建 GIN 索引的基表属性的数据类型相关的。而 accum 字段则指向索引创建过程中所使用的“蓄积池”，其数据类型为 BuildAccumulator（数据结构 4.16）。

数据结构 4.15 GinBuildState

```
typedef struct
{
    GinState      ginstate;           // 当前使用的 GIN 索引的基本信息
    double        indtuples;         // 索引中索引元组的总数
    MemoryContext tmpCtx;           // 创建过程中使用的临时内存上下文
    MemoryContext funcCtx;          // 创建过程中进行函数调用时的内存上下文
    BuildAccumulator accum;         // “蓄积池”，用于缓存 Entry
} GinBuildState;
```

BuildAccumulator 中的 maxdepth 和 allocatedMemory 两个字段反映了“蓄积池”中填充 Entry 的程度，如果 maxdepth 超过 GIN_MAX_TREE_DEPTH（值为 100）或者 allocatedMemory 超过一个阈值（阈值为 $16384 * 1024$ ），则创建过程会先将“蓄积池”中的 Entry 都插入到 GIN 索引中，然后清空“蓄积池”。

数据结构 4.16 BuildAccumulator

```
typedef struct
{
    GinState      *ginstate;         // 当前使用的 GIN 索引的基本信息
    EntryAccumulator *entries;       // 指向“蓄积池”中 Entry 构成的树
    uint32        maxdepth;          // “蓄积池”中 Entry 树的最大深度
    EntryAccumulator **stack;        // 记录找到上一个插入 Entry 的路径
    uint32        stackpos;           // stack 数组的长度
    long          allocatedMemory;    // “蓄积池”中已经分配的内存大小
    uint32        length;             // “蓄积池”内已存放的 Entry 个数
    EntryAccumulator *entryallocator; // 内存分配器，用于为“蓄积池”分配内存，指向这一次分配的
                                      // 内存空间头部
} BuildAccumulator;
```

在“蓄积池”中会使用到 EntryAccumulator（数据结构 4.17）来保存找到的 Entry，每一个 Entry 都有一个 EntryAccumulator 结构，池内的 EntryAccumulator 构成一棵排序二叉树，“蓄积池”的 entries 字段指向树的根节点。

对于 EntryAccumulator 结构中的字段说明如下：

- attnum 字段：由于 GIN 支持多属性索引，attnum 记录了 Entry 出现在第几个被索引的属性中，该字段从 1 开始递增。也就是说，如果一个关键词出现在同一个元组的两个不同属性中，则会为该关键词创建两个 EntryAccumulator 添加到树中。

- length 字段：由于一个 Entry 出现的次数是无法预知的，所以 list 字段所指向的数组长度是动态变化的，当 number（频率）大于等于 length（目前分配给 list 的长度）时，即将 length 翻倍，同时 list 的长度也翻倍。
- shouldSort 字段：若该 Entry 在多个地方出现，依次将读取到的位置信息插入到 list 字段（Posting List）中，由于后面插入的位置可能会使得 list 中的位置变成无序的，就需要标记 shouldSort 字段为真，在插入 GIN 索引时会根据这个字段的值来判断是否需要进行排序。（注：在插入出现位置到“Posting List”时并不排序，而是在构建整个 GIN 索引结构时，根据该字段判断是否需要排序。）

数据结构 4.17 EntryAccumulator

```
typedef struct EntryAccumulator
{
    OffsetNumber          attnum;           // 该属性是被索引属性的第几个
    Datum                 value;            // Entry 的键值
    uint32                length;           // list 字段中数组元素的个数
    uint32                number;           // Entry 出现的次数
    ItemPointerData       *list;            // Entry 所有出现位置
    bool                  shouldSort;        // Entry 的“Posting List”是否需要排序
    struct EntryAccumulator *left;          // 左节点
    struct EntryAccumulator *right;         // 右节点
} EntryAccumulator;
```

对于同一个 Entry，即使在多个元组的同一个属性中出现，也都是保存在同一个 EntryAccumulator 中，其出现的不同位置保存在 list 结构中。但是同一个 Entry 如果出现在不同的属性中，则要用不同的 EntryAccumulator 来表示。EntryAccumulator 通过其 left 和 right 字段指向其左右子节点，EntryAccumulator 构成的二叉树是有序的，左子节点的 Entry 比父节点的 Entry “小”，父节点的 Entry 比右子节点的 Entry “小”。这里的大小可通过 compareAttEntries 函数来进行比较。

PostgreSQL 中提供了 ginInsertRecordBA 函数将属于同一个元组的同一个属性的多个 Entry 插入到“蓄积池”。ginInsertRecordBA 插入 Entry 所采用的策略如下：

- 1) 将 Entry 数组进行排序。

- 2) 首先将 Entry 数组中间元素插入到排序二叉树中，然后采用同样的策略递归处理 Entry 数组的左半部分和右半部分。

例如，要用 ginInsertRecordBA 函数插入以下的 Entry 数组（这里没有列举 Entry 出现的位置，只给出了关键词）：

Entry 数组（已排序）： abstract, binary, hash, hello, index, tree, world
--

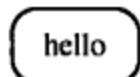
实际 ginInsertRecordBA 执行的过程如下：

- 1) 插入数组中间元素 hello，递归插入以下两个数组：

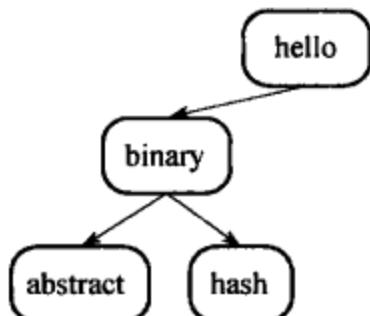
abstract, binary, hash

index, tree, world

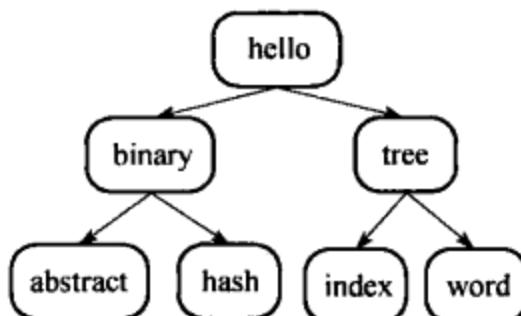
此时树中有一个节点：



2) 处理完前面一个数组，二叉树的结构为：



3) 插入第二个数组后，二叉树的结构为：



采用这种插入策略的好处是能够尽量保证插入记录后的二叉树趋于平衡，不会使得其高度太大影响查找效率。

创建 GIN 索引函数是 ginbuild 函数，其流程如下：

1) 调用 initGinState 初始化创建状态变量 buildstate (GinBuildState 类型) 的 ginstate 字段，该函数会根据当前要创建的索引的信息以及被索引属性的信息从系统表中抽取 5 个支持函数的信息用于填充 ginstate 字段的内容。

2) 初始化 GIN 索引的元页，并初始化 buildstate 的其他字段，包括置 indtuples 为 0、创建临时内存上下文和函数调用内存上下文。

3) 调用 ginInitBA 初始化“蓄积池”（即 buildstate 中的 accum 字段）。

4) 调用 IndexBuildHeapScan 对基表扫描，并将 ginBuildCallback 函数的指针传递给 IndexBuild-HeapScan。每扫描到一个基表示元组，IndexBuildHeapScan 都将其被索引属性的值解析出来然后传递给 ginBuildCallback 处理，在 ginBuildCallback 中将进行以下处理：

①对于每一个被索引属性的值，调用 ginHeapTupleBulkInsert 对其进行处理，该函数会调用 GIN 索引的 extractValue 方法将被索引属性的值解析成若干个 Entry (一个属性值中可能包含多个关键词)，并将这些获得的 Entry 插入到“蓄积池”中 (ginInsertRecordBA 函数)，该函数值将返回获得的 Entry 的个数，这个数量将被累加到 buildstate 的 indtuples 字段中。

②对“蓄积池”的填充情况进行检查，如果 allocatedMemory 或 maxdepth 字段的值超过阈值则循环调用 ginGetEntry 从“蓄积池”中取出一个 Entry，并将其用 ginEntryInsert 插入到 GIN 索引中。

③当“蓄积池”中的 Entry 都已经被插入到 GIN 索引中后，通过重置 buildstate 中的临时内存上下文来清空“蓄积池”，然后调用 ginInitBA 重新初始化“蓄积池”。

5) 当 IndexBuildHeapScan 扫描完成后，“蓄积池”中有可能还有一部分 Entry 没有被插入到 GIN 索引中 (allocatedMemory 和 maxdepth 字段的值都没有超过阈值)，同样循环调用 ginGetEntry 获取 Entry 再通过 ginEntryInsert 插入 GIN 索引中。

6) 返回统计信息，包括被索引的基表元组以及索引元组的数目。

将 Entry 及其位置信息插入到 GIN 索引结构中，是通过调用 ginEntryInsert 函数来完成的，该函数的流程如图 4-35 所示。

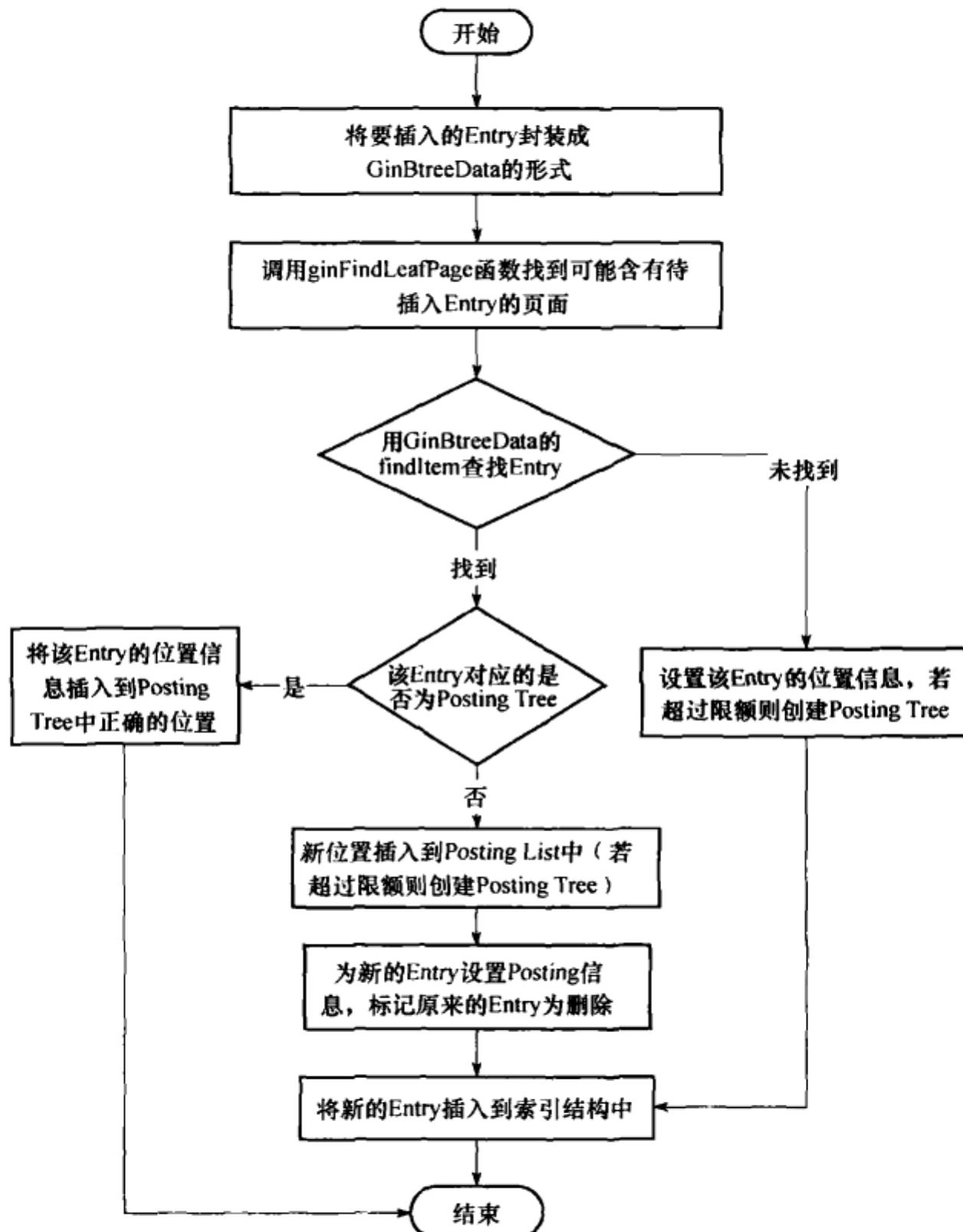


图 4-35 ginEntryInsert 插入一个 Entry 到 GIN 索引中

在 ginEntryInsert 中，待插入的 Entry 用一个数据结构 GinBtreeData 来管理，其定义如 4.18 所示。

数据结构 4.18 GinBtreeData

```

typedef struct GinBtreeData
{
    // 查找模式使用的函数指针
    BlockNumber (*findChildPage) (GinBtree, GinBtreeStack *);
    bool (*isMoveRight) (GinBtree, Page);
    bool (*findItem) (GinBtree, GinBtreeStack *);

    // 插入模式使用的函数指针
    OffsetNumber (*findChildPtr) (GinBtree, Page, BlockNumber, OffsetNumber);
    BlockNumber (*getLeftMostPage) (GinBtree, Page);
    bool (*isEnoughSpace) (GinBtree, Buffer, OffsetNumber);
    void (*placeToPage) (GinBtree, Buffer, OffsetNumber, XLogRecData **);
    Page (*splitPage) (GinBtree, Buffer, Buffer, OffsetNumber, XLogRecData **);
    void (*fillRoot) (GinBtree, Buffer, Buffer, Buffer);

    bool searchMode;           //是否为查找模式
    Relation index;           //索引表
    GinState *ginstate;       //用于记录 GIN 索引提供的 5 种扩展方法
    bool fullScan;            //是否为完全扫描
    bool isBuild;              //该 Entry 对应的位置信息上是否创建了 B - Tree
    BlockNumber rightblkno;   //该 Entry 所在页面的下一个(右)页面

    //以下是与 Entry 信息相关字段
    OffsetNumber entryAttnum; //Entry 在索引元组中对应属性编号
    Datum entryValue;         //Entry 的键值
    IndexTuple entry;          //Entry 对应的索引元组
    bool isDelete;             //Entry 是否被删除

    //以下是与位置信息相关字段
    ItemPointerData *items;   //Entry 的位置数组
    uint32 nitem;              //位置数组的长度
    uint32 curitem;             //当前使用的“位置”
    PostingItem pitem;          //保存叶子节点页面的页面号以及一个物理地址
} GinBtreeData;

```

注意，GinBtreeData 数据结构不仅在 GIN 索引的创建过程中会使用，在使用 GIN 索引进行查询、插入时都会使用到。其中的 items 字段在创建索引时用于保存该 Entry 所有的位置信息，即保存“Posting List”，然后根据位置数组的长度来判断是否需要构建“Posting Tree”。

可以看到，GinBtreeData 数据结构中定义了 9 个函数指针，这些函数指针指向的函数实现了对

GIN 索引结构的查找、插入、分裂等操作，对于每一个函数指针，既可以定义对 Entry 节点操作的函数，也可以指向对位置信息进行操作的函数。每一个函数的功能说明如下：

- 1) `findChildPage`: 在 `GinBtreeStack` 类型参数的 `buffer` 字段对应的非叶子层页面中找到 `Entry` (由 `GinBtreeData` 类型参数表示) 对应的索引元组 (由于是非叶子层页面，索引元组中的指针指向下一层的孩子页面)，并返回索引元组指向的孩子页面的块号。
- 2) `isMoveRight`: 比较 `Page` 类型参数对应的页面中最后一个索引元组与 `GinBtree` 对应的 `Entry` 的值大小，如果小于则返回真 (表示该页面应该右移)，否则返回假。
- 3) `findItem`: 在 `GinBtreeStack` 类型参数的 `buffer` 字段对应的叶子层页面中查找 `Entry` (由 `GinBtreeData` 类型参数表示) 对应的索引元组，若存在，将索引元组在页面中的偏移量赋给 `GinBtreeStack` 类型参数的 `off` 字段并返回真，否则返回假。
- 4) `findChildPtr`: 在 `Page` 类型参数对应的页面中查找孩子页面号为 `BlockNumber` 的索引元组的页面内偏移量并返回。
- 5) `getLeftMostPage`: 返回非叶子层页面中第一个索引元组指向的页面号。
- 6) `isEnoughSpace`: 判断 `Buffer` 类型参数指向的缓冲区中是否有足够空间存储 `GinBtree` 对应的索引元组，如果有足够空间则返回真，否则返回假。
- 7) `placeToPage`: 将 `GinBtree` 中的索引元组写到缓冲区中，`OffsetNumber` 类型参数为写入缓冲区的起始偏移量，并写日志。
- 8) `splitPage`: 将 `GinBtree` 中的索引元组插入到第一个参数指定的缓冲区中，并将该缓冲区分裂成两个大小几乎相当的缓冲区，分裂产生的新缓冲区放在第二个参数中，返回第一个缓冲区的页面号，最后写日志。
- 9) `fillRoot`: 将第三、四个参数指定的两个缓冲区最右端的索引元组填充到第二个参数指定的缓冲区中，并将被填充索引元组的指针分别指向第三、四个参数对应的页面号。

可以看出，这些函数中很多都使用了 `GinBtreeStack` 数据结构，该结构用于记录索引结构在一次查找过程中从根节点到叶子节点的路径，当插入一个 `Entry` 引起节点分裂后，需要使用该结构向上回溯调整。`GinBtreeStack` 结构的定义见数据结构 4.19。

数据结构 4.19 GinBtreeStack

```
typedef struct GinBtreeStack
{
    BlockNumber          blkno;           //磁盘号
    Buffer               buffer;          //对应的缓冲区
    OffsetNumber         off;             //查找到的索引元组在页面中的偏移量
    uint32               predictNumber;   //预测的当前层次页面的个数
    struct GinBtreeStack *parent;        //指向父节点的指针
} GinBtreeStack;
```

对于 GIN 索引的插入，若一次插入了很多元组，则会调用 `ginHeapTupleFastInsert` 函数进行批量插入，该函数的执行流程与索引创建流程类似，这里就不再赘述。

在大多数情况下，若更新的数据量很大，GIN 索引的“Posting List”结构可能较长，向 GIN 索

引插入的速度较慢。因此，对于向表中大量插入的操作，建议先删除 GIN 索引，在完成插入之后再重建它。

2. GIN 索引查询

GIN 索引的查询通过关键词与 Entry 的匹配操作来查找其所在的元组。在查找的过程中按照 B-Tree 结构进行搜索，调用 `consistent` 方法来判断是否找到。GIN 索引没有提供返回单个元组的函数（类似 B-Tree 索引的 `btgettupple` 函数），只提供了位图查询的方式，也就是说，对 GIN 索引的查询只能得到一个位图，其中包含了符合查询条件的元组的物理位置。

开发 GIN 索引的目的是让 PostgreSQL 支持高度可扩展的全文索引。我们常常会遇到全文索引返回海量结果的情形，在查询高频词时得到的很大的结果集其实并没有什么用，因为从磁盘读取大量记录并对其进行排序会消耗大量资源。为了易于控制这种情况，GIN 索引有一个结果集大小“软上限”的配置参数 `gin_fuzzy_search_limit`，其缺省值为 0，表示没有限制。如果设置了非零值，那么返回的结果就是从完整结果集中随机选择的一部分。

4.6 TSearch2 全文搜索

全文搜索（文本搜索）提供了一种可以检索出满足某个查询条件的自然语言文档的能力，并且还可以根据文档的相关性对文档进行排序。最常见的搜索是找出所有包含给出的查询词的文档，并且以它们符合查询的程度排序输出。

文本搜索操作符在数据库里已经存在很多年了。PostgreSQL 有 ~、~* 和 LIKE 操作符用于文本数据类型，但是它们缺乏许多现代的信息系统需要的重要功能，比如：

- 没有语言支持，不会对文本进行解析。
- 不提供检索结果的排序（ranking），在找到上千个匹配文档的时候，就不够高效了。
- 没有索引支持，所以会比较慢，因为它们必须为每个查询处理所有的文档。

从 PostgreSQL 8.3 开始提供了文本搜索模块 TSearch（Text Search），文本搜索提供了一种可以标识满足某个查询的自然语言文档的能力，并且还可以根据文档的相关性对文档进行排序。

PostgreSQL 核心系统提供的 TSearch 模块提供了对文档（在 PostgreSQL 里一个文档通常是一个表中的某个元组的一个文本属性，或是几个属性的组合）及查询条件进行解析的功能，但并没有提供对解析后的结果进行进一步处理（创建索引，以支持快速的查找）的功能。PostgreSQL 在扩展模块（contrib）里面提供了 TSearch2 来支持这些功能，TSearch2 实现了对文档创建 GIN 或者 GiST 索引的支持。本节将分析 PostgreSQL 核心系统提供的 TSearch 模块，其代码位于 `src/backend/tsearch` 目录下。

全文搜索一般有三个步骤：

- 1) 对文档（文本）进行预处理，得到处理后的结果，创建索引。
- 2) 解析查询条件，使用索引进行查询。
- 3) 得到查询结果，对结果进行处理后返回。

对应这三个过程，PostgreSQL 提供了一个数据类型 `tsvector` 用于存储预处理后的文档，还提供了一个数据类型 `tsquery` 用于查询（还为这两种类型提供了匹配操作符 @@）以及查询结果的处理，下面将对这三个步骤进行详细分析。

4.6.1 全文索引的创建

全文索引允许对文档进行预处理并且可以保存为用于快速搜索的索引。预处理包括文本解析、语义分析和词位存储。完成这三个过程后，解析后的词语信息就存放在 `TSVector` 结构中。从文本解析到词位存储这一系列过程是由函数 `to_tsvector_byid` 完成的，该函数首先调用 `parsetext` 函数对文本进行解析和语义分析，然后再调用 `make_tsvector` 将词位信息构建成 `TSVector` 结构。下面将对这三个过程依次进行分析。

1. 文本解析

文本解析通过解析器将文档解析成一个个记号（含位置信息，类型信息），该过程涉及的函数在 `wparser_def.c` 文件中。

目前 PostgreSQL 只提供一种解析器，但它足够处理大多数纯文本及 HTML 文件。PostgreSQL 中默认的记号对应表如表 4-9 所示。

表 4-9 默认解析器记号对照表

类型名称	描述	例子
asciword	单词，所有 ASCII 字母	elephant
word	单词，所有字母	manana
numword	单词，字母和数字	beta1
asciihword	连接成的词，所有 ASCII	up-to-date
hword	连接成的词，所有字母	lógico-matemática
numhword	连接成的词，字母和数字	postgresql-beta1
hword_asciipart	连接成的词的部分，所有 ASCII	postgresql-beta1 里的 postgresql
hword_part	连接成的词的部分，所有字母	lógico-matemática 里的 lógico 或 matemática
hword_numpart	连接成的词的部分，字母和数字	postgresql-beta1 里的 beta1
email	Email 地址	foo@example.com
protocol	协议头	http://
url	URL	example.com/stuff/index.html
host	主机	example.com
url_path	URL 路径	/stuff/index.html, in the context of a URL
file	文件或路径名	/usr/local/foo.txt, if not within a URL
sfloat	科学计数法	-1.234e56
float	小数表示法	-1.234
int	符号整数	-1234
uint	无符号整数	1234
version	版本号	8.3.0
tag	XML 标签	
entity	XML entity	&
blank	空白符号	(任意空白或者其他会被识别的标点)

2. 语义分析

语义分析是对解析器处理过的 token 文本序列通过参照词典的审核规范成标准的词（lexeme）信息。

词典用于删除那些不应该在搜索中出现的词（屏蔽词）并规范化一些有多重形式的词，这样同一个词的不同的衍生结果也可以被搜索到。成功规范化之后的词被称作词位（lexeme）。除了改进搜索质量，规范化和删除屏蔽词可以减少文档的尺寸，从而提高性能。下面对各个词典的使用进行举例介绍：

- Ispell：拼写词典，例如“likes”将转换为“like”。
- Simple：简单词典，例如“A NAUGHTY DOG”将转换为“naughty dog”。
- Synonym：同义词典，例如“man”和“person”是同义词。
- Thesaurus：知识词典，例如“personal computer”将转换为PC。

完成语义分析后，即得到一个全部处理后得到的单词信息，这些单词信息保存在 ParsedText 结构中。ParsedText 结构保存解析后的文本，其定义如数据结构 4.20 所示。

数据结构 4.20 ParsedText

```
typedef struct
{
    ParsedWord *words;           //存放解析得到的单词集合,动态内存空间,长度由 lenwords 决定
    int4        lenwords;         //当前允许存放单词的数目的最大值,当 curwords = lenwords 时, len-
                                //words 即加倍,words 的内存空间翻倍
    int4        curwords;         //当前解析得到的单词数目
    int4        pos;              //当前位置
} ParsedText;
```

其中的 words 指向一个数组，其中每一个元素都是 ParsedWord 类型，用于保存分析后的一个单词，其定义如数据结构 4.21 所示。

数据结构 4.21 ParsedWord

```
typedef struct
{
    uint16  len;                //单词的长度
    uint16  nvariant;           //在对一个组合词进行切分时,可能有多种切分结果,标记是第几种结果
    union
    {
        uint16  pos;             //该单词在 ParsedText. words 中的编号
        uint16  *apos;
    } pos;
    uint16  flags;              //类型标记,目前均为 TSL_PREFIX
    char    *word;               //保存单词
    uint32  alen;
} ParsedWord;
```

数据结构 4.21 中的 pos 和 apos 指针是用 union 结构来保存的。当完成文本解析后，可能会遇到相同的词出现了多次的情况，这时会将相同的词合并在一起：对于只出现一次的词，使用 pos 来保

存其出现的位置即可；对于出现多次的词，则使用 `apos` 指针来指向一个动态数组来保存所有出现的位置。`apos[0]` 为该词出现的次数，数组后面的值即为各次出现的位置。由于数组的长度是不确定的，所以使用 `alen` 字段来确定动态申请内存的空间，`alen` 初始化为 2，每当 `apos` 的长度不够时，`alen` 即翻倍，同时申请新内存将 `apos` 数组的长度翻倍。

上面介绍了用于存储语义分析后得到的单词信息的数据结构，而这整个处理流程是由函数 `parseText` 完成的，其执行如图 4-36 所示。

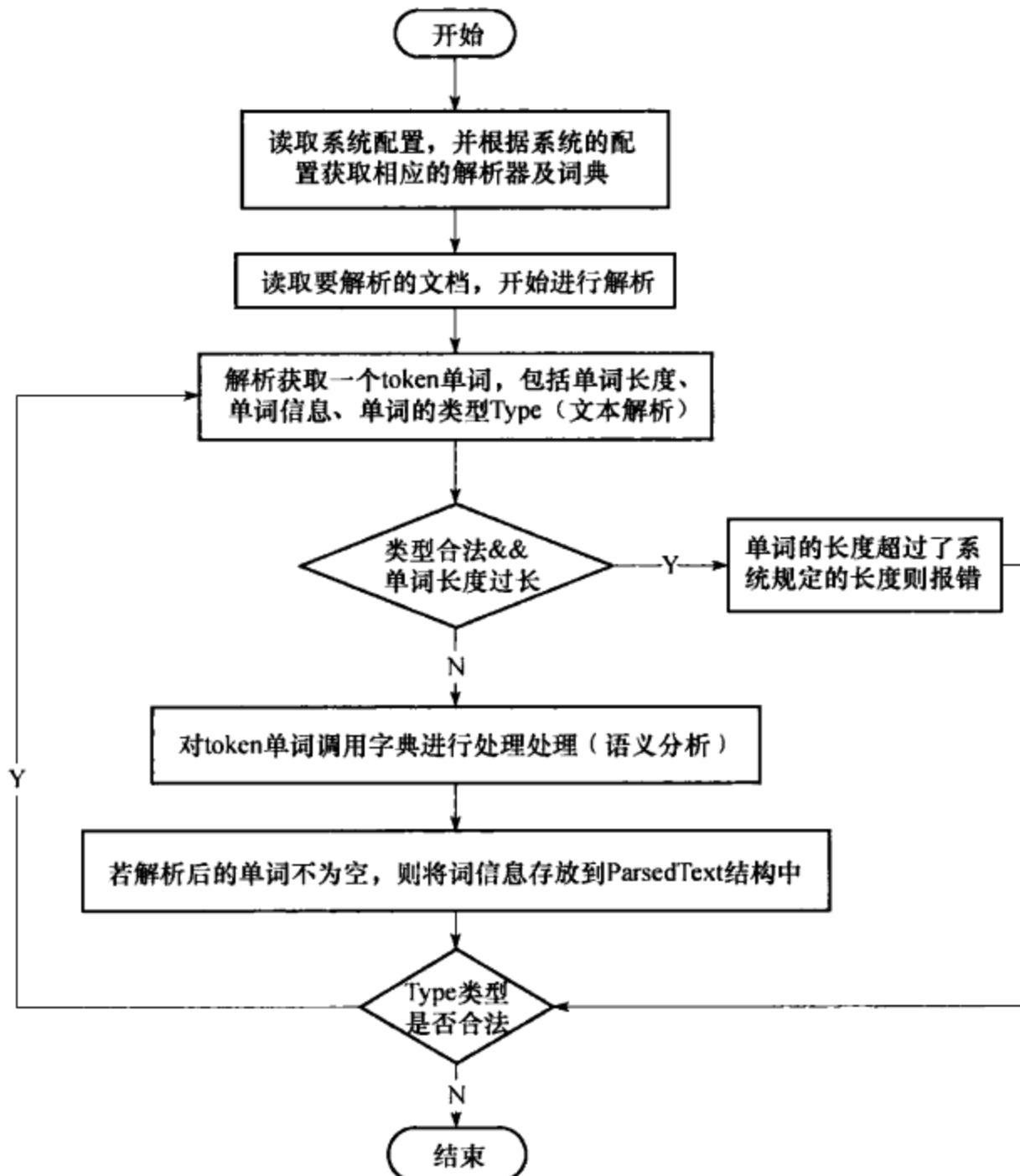


图 4-36 语义分析流程

3. 词位存储

词位存储即为归总语义分析后的单词在文档中的位置，并将其出现的次数及每个位置存储下来。

`TSVector` 是一种可搜索的数据类型，它是文档内容的一种表现形式，是出现在文档中的每个重要单词及其所有位置信息的集合。它通过一种特殊的优化结构进行组织，从而可方便快速地存取及查找。其定义如数据结构 4.22 所示。

数据结构 4.22 TSVector

```

typedef struct
{
    int32      vl_len_;           // 暂未使用
    int32      size;             // 关键字的个数 (WordEntry 数组的长度)
    WordEntry  entries [1];     // 关键词的信息
} TSVectorData;
typedef TSVectorData *TSVector;

```

TSVectorData 结构中的 WordEntry 数组用于保存所有的关键字（单词）信息，由于关键字的数目一开始并不确定，所以使用一个数组指针，该数组的实际长度根据关键字的个数在使用时进行分配。WordEntry 的定义见数据结构 4.23。

数据结构 4.23 WordEntry

```

typedef struct
{
    uint32
        haspos:1,          // 标记该单词是否多次出现
        len:11,            // 记录单词的长度, 最长 2K
        pos:20;            // 记录单词所在的位置, 最大 1M
} WordEntry;

```

上面分析了 TSVector 的数据结构，词位存储即使用语义分析得到的 ParsedText 构建 TSVector，该过程由函数 make_tsvector 完成。其执行流程如图 4-37 所示。

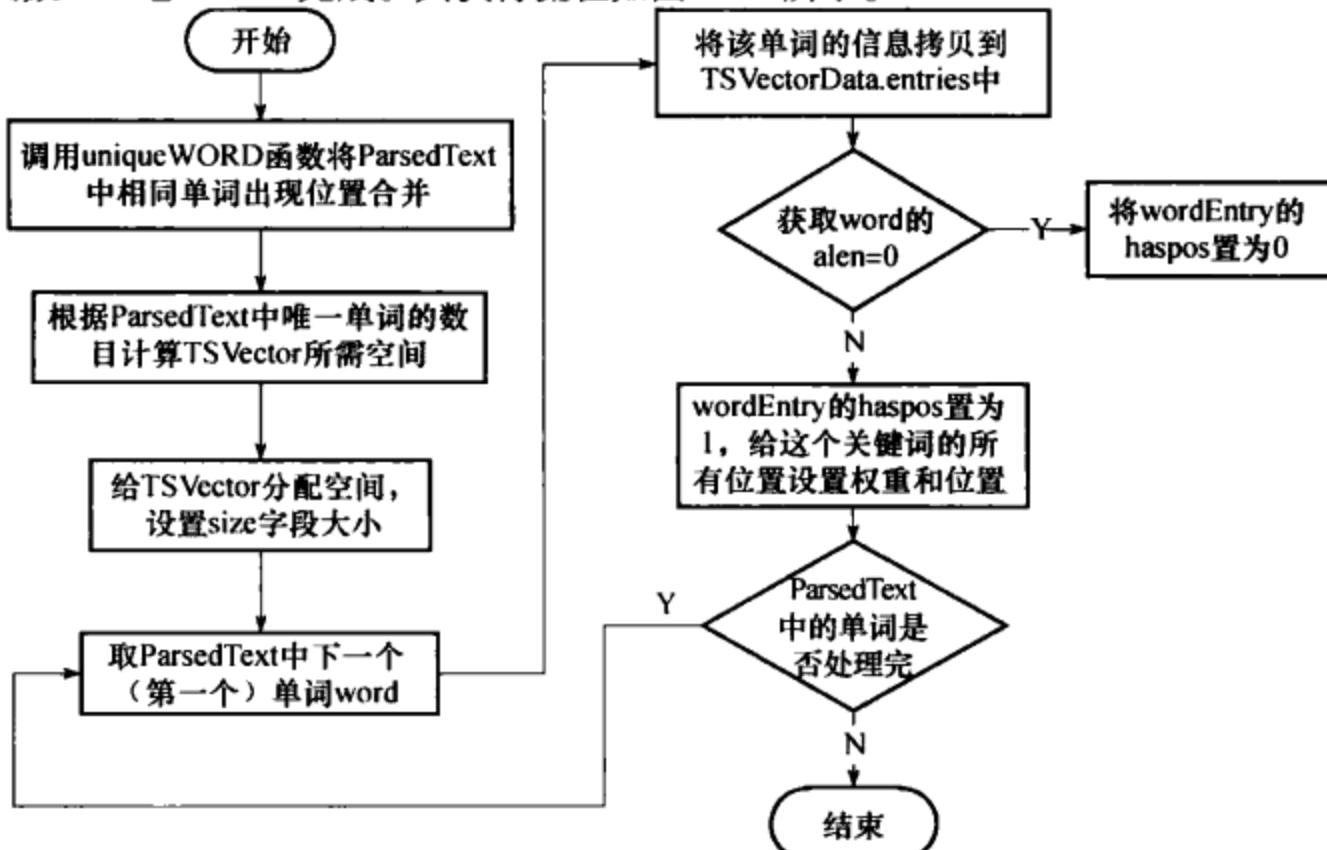


图 4-37 构建 TSVector 流程图

至此已经介绍了 PostgreSQL 内核中 TSearch 模块提供的全部功能。之前讲过，PostgreSQL 把 TSearch2 作为一个扩展模块，提供了对 Entry 创建 GIN 或者 GiST 索引的支持。通过对词位列创建 GIN 或 GiST 索引即可实现对文档的全文索引。

综上所述，全文索引的创建流程图可归纳如图 4-38 所示，其中创建 GIN (GiST) 索引部分用虚线框，表示该步骤是需要用户额外编译并安装 TSearch2 模块才具有的功能。

由于 GIN 或者 GiST 索引结构的不同，其创建、查询及更新的效率也有不同，PostgreSQL 8.4.1 官方手册上对这两种索引结构的优劣进行了比较：

- GIN 索引查询速度是 GiST 的 3 倍。
- GiST 创建索引的速度是 GIN 的 3 倍。
- GiST 索引的更新速度较 GIN 稍快。
- GIN 索引的空间比 GiST 索引大 2 至 3 倍。

前面介绍了从文本解析一直到创建全文索引的全部过程。索引创建完成后，即可利用全文索引进行查询。由于查询条件可能是自然语句，也需要对查询进行一定的处理以获取准确的查询结果。接下来将介绍如何利用这里建立的全文索引进行查询。

4.6.2 全文索引的查询

全文索引查询之前需要对检索的语句进行处理，处理过程跟全文索引的创建过程类似，要对查询语句进行文本解析及语义分析，将分析后的结果封装成 tsquery 格式，然后就可以通过对创建在 tsvector 上的索引进行匹配查询了。但与创建过程不同的是：查询处理时，各个处理后的单词需要使用布尔操作符 & (与)、| (或) 和 ! (非) 进行连接。tsquery 的格式如数据结构 4.24 所示。

数据结构 4.24 TSQueryData

```
typedef struct
{
    int32  vl_len_;           //保留，暂未使用
    int4   size;              //查询条件的单词和操作符的数量
    char   data [1];          //存储处理后的查询数据
} TSQueryData;
typedef TSQueryData *TSQuery;
```

PostgreSQL 提供了 to_tsquery 和 plainto_tsquery 两个函数用于把查询转换成 tsquery 数据类型。to_tsquery 的功能是将查询中的关键词在语义分析后转化成可与 tsvecor 进行匹配的数据类型。该函数的参数要求很严格，它们必须使用布尔操作符 “&” (与)、“|” (或) 和 “!” (非) 分隔各个单词。

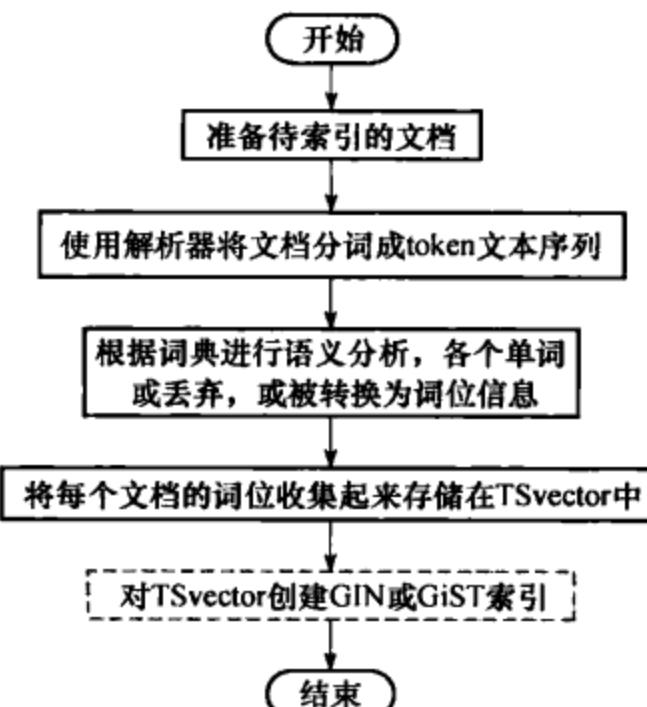


图 4-38 全文索引创建流程

`plainto_tsquery` 对查询的格式要求没有 `to_query` 那么严格，它把未格式化的文本查询转换成 `tsquery`。文本首先会像在 `to_tsvector` 里那样分析和规范化，然后用布尔操作符“&”（与）将解析后得到的单词进行连接，最后得到 `tsquery`。

`to_tsquery` 和 `plainto_tsquery` 函数的处理都是通过调用 `parse_tsquery` 函数来实现的。`parse_tsquery` 函数实现了对用户输入的查询条件的解析及语义分析过程，使用布尔操作符将得到的关键字连接起来。通过调用 `makepol` 函数，将查询语句转换成“波兰表示法”（“Polish notation”，或称为波兰记法）。对查询条件的解析过程与 4.6.1 节中的过程基本一样，通过调用 `parsetext` 函数完成。当读取到查询条件中的操作符时，只能是上面讲到的三种操作符，对于其他操作符则报错（`plainto_tsquery` 函数则不会处理用户输入条件中的操作符，它对查询条件进行解析后，全部使用“&”（与）操作符连接得到的关键词）。

当得到 `tsquery` 后，即可调用之前创建的全文索引进行查询，全文索引查询的总体流程如图 4-39 所示。

4.6.3 查询结果处理

上面分析了全文索引的创建和查询处理过程。数据库系统提供了 @@ 操作符，可以对 `TSVector`（或处理后的 `tsquery`）结构中的 `entry` 进行比较。下面给出一个使用 @@ 操作符进行全文搜索的例子。

假设现在在数据库中有一个 `messages` 表，其字段内容如表 4-10 所示。

如果要从上表中查出属性 `strMessage` 中包含“test”或者“king”的文档信息，其 `tsearch` 语句如下：

```
SELECT id, strtopic FROM Messages
WHERE to_tsvector(strMessage) @ @ to_tsquery('test | king');
```

PostgreSQL 会对 `messages` 表中的 `strMessage` 字段按 4.6.1 节中的介绍进行处理，然后对“test | king”采用 4.6.2 节中的 `to_tsquery` 函数进行处理，再调用 @@ 操作符进行匹配。最后，返回匹配成功的结果，如表 4-11 所示。

表 4-10 messages 表

Id	strTopic	strMessage
1	Testing Topic	Test message data input
2	Movie	Breakfast at Tiffany
3	Famous Author	Stephen King

表 4-11 全文查询结果

id	strtopic
1	Testing Topic
3	Famous Author

对全文搜索的结果，TSearch 还提供了一些功能对结果进行处理，包括权重设置、结果排序以及结果高亮显示，下面将对这 3 种功能进行介绍。

（1）设置文档部分的权重

权重用于标记单词来自于文档的哪个区域，比如标题和开头的摘要，这样就可以以不同的方式对待不同权重的单词。可以通过调用函数 `tsvector_setweight` 将一个 `TSVector` 中所有的关键字都设置为指定的权重。

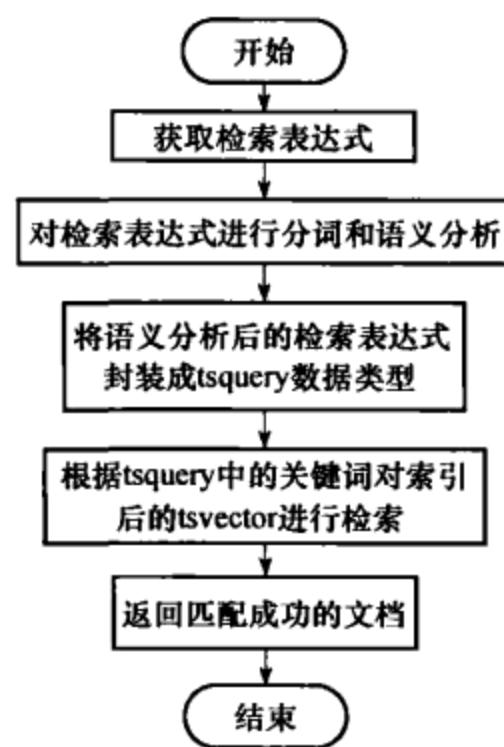


图 4-39 全文索引查询流程

(2) 对查询结果排序

排序主要是将匹配成功的搜索结果与查询表达式进行相关性比较，将最相关的结果排在前面。相关性是评判文档与特定的查询之间相关程度的一个衡量标准。如果有很多匹配的项，那么相关性高的项应该排在前面。

`ts_rank` 和 `ts_rank_cd` 函数都可以用来对查询结果进行排序。

(3) 高亮显示

呈现搜索结果时，最好是显示每个文档的一部分以及它和查询之间是如何关联的。通常，搜索引擎会显示查询关键词所在的文档片段，并且在其中将查询关键词高亮显示。PostgreSQL 也提供一个函数 `ts_headline` 实现这个功能。

`ts_headline` 接受一个文档以及对应的查询，然后返回一个文档的摘要。在摘要里面，查询是高亮显示的。`ts_headline` 使用原始的文档，而不是 `tsvector` 摘要，所以它可能比较慢，因此要小心使用。

4.7 小结

索引是提高数据库性能的常用方法，它可以令数据库服务器以更快的速度查找和检索特定的行。不过索引也增加了数据库系统的负荷，因此应该恰当地使用它们。索引的优劣不仅仅与索引的查询效率有关，同时与索引创建速度，更新速度，索引大小等因素有关。以上分析的几种索引使用环境不同，相同环境下使用的优劣也各不相同。于是，为了方便用户对特殊数据类型数据的查询，PostgreSQL 中提供了索引模板可方便用户对索引的扩展以支持自定义数据类型上的索引创建与查询。

通过使用索引，能够快速查找数据库中的数据。但在添加索引的同时，会增加数据库系统的负荷；在增删修改数据时，维护索引的一致性也需要一定的时间和空间。由于索引给查找数据带来的巨大的性能提升，因此也得到了广泛的应用。在实际的使用中，应该合理权衡使用索引带来的利弊，恰当地使用索引。

习题

习题 4.1 不同的索引有其不同的特点和长处，合理使用各种索引能为数据库性能带来很大的提升。B – Tree 索引和 Hash 索引都能处理对“=”操作符的查找，但它们也有区别。请问这两种索引类型各有什么特点？并说明 B-Tree 索引和 Hash 索引分别适合在什么情况下使用？

习题 4.2 PostgreSQL 在编译安装时，已将 GiST 和 GIN 索引进行了编译，在 `postgresql-8.4.1/src/backend/access` 文件夹下可查看到 `gist` 和 `gin` 索引的目录。作为数据库索引的扩展，在 `postgresql-8.4.1/contrib` 文件夹下，还提供了更多的扩展索引类型，如 `btree_gist` 文件夹下的程序提供了在 GiST 上模拟 B – Tree 索引，`xml2` 文件夹下的程序使得可在 PostgreSQL 中存储 XML。尝试对这些扩展模块进行编译，并测试使用其功能。在 `README` 文件中可以找到帮助信息。

习题 4.3 全文搜索作为搜索引擎中的重要的一类，在互联网中起着越来越重要的作用。使用 PostgreSQL 提供的 TSearch2 模块进行全文搜索是很简单的。PostgreSQL 系统并没有编译该模块，在 `postgresql-8.4.1/contrib/tsearch2` 文件夹下可查看其源代码。通过 4.6 节的分析，自己编译该模块，然后使用 PostgreSQL 提供的全文索引功能，编写程序对数据库中的文本进行创建全文索引，使用索引进行查找，按相关性排序，查找的匹配文本高亮显示。可参考 <http://www.michaelhinds.com/tech/pg/tsearch.html>。

第 5 章

查询编译

查询处理器是数据库管理系统中的一个部件集合，它允许用户使用 SQL 语言在较高层次上表达查询，其主要职责是将用户的各种命令转化成数据库上的操作序列并执行。查询处理分查询编译和查询执行两个阶段。查询编译的主要任务是根据用户的查询语句生成数据库中最优执行计划，在此过程中要考虑视图、规则以及表的连接路径等问题，这也是本章要介绍的主要内容。查询执行主要考虑执行计划时所采用的算法等问题，将在下一章详细介绍。

5.1 概述

当 PostgreSQL 的后台服务进程 Postgres 接收到查询语句后，首先将其传递到查询分析模块，进行词法、语法和语义分析。若是简单的命令（例如建表、创建用户、备份等）则将其分配到功能性命令处理模块；对于复杂的命令（SELECT/INSERT/DELETE/UPDATE）则要为其构建查询树（Query 结构体），然后交给查询重写模块。查询重写模块接收到查询树后，按照该查询所涉及的规则和视图对查询树进行重写，生成新的查询树。生成路径模块依据重写过的查询树，考虑关系的访问方式、连接方式和连接顺序等问题，采用动态规划算法或遗传算法，生成最优的表连接路径。最后，由最优路径生成可执行的计划，并将其传递到查询执行模块执行。PostgreSQL 中一条查询语句的完整处理流程如图 5-1 所示^①，表 5-1 是对查询执行过程中所涉及的各模块的说明。

查询优化的核心是生成路径和生成计划两个模块。由于在整个查询执行过程中，表连接操作的开销最大，因此，查询优化要处理的问题焦点在于如何计算最优的表连接路径。

^① postgresql-8.4.1\src\tools\backend\index.html

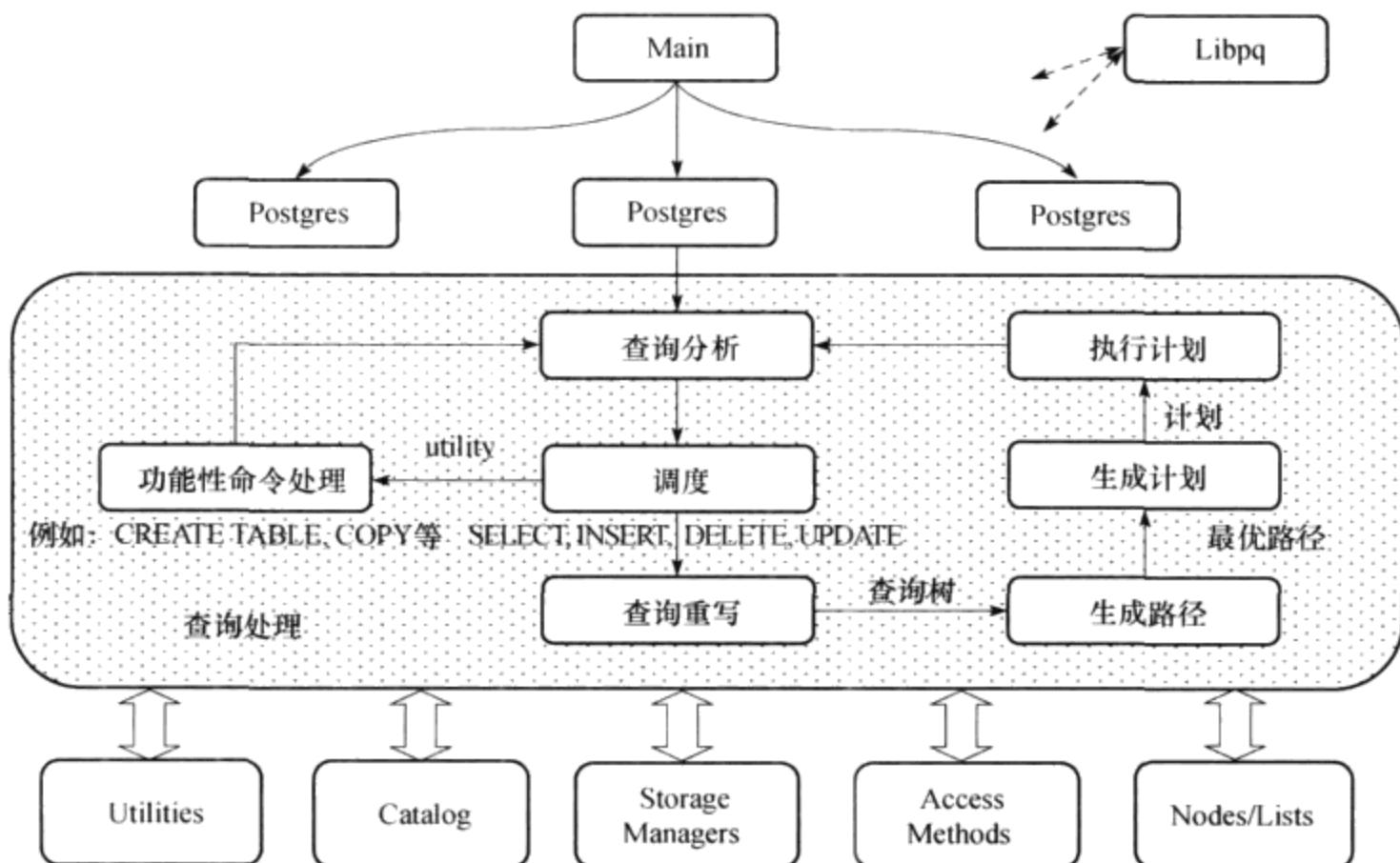


图 5-1 查询处理的完整流程

表 5-1 查询处理的各模块说明

说明阶段	模块	功能说明	源码路径
查询编译	查询分析	由 SQL 查询语句生成查询树	src/backend/parser
	查询重写	对查询树重写并生成新的查询树，以提供对规则和视图的支持	src/backend/rewrite
	查询优化	由查询树计算最优路径	src/backend/optimizer/path
		通过最优路径生成计划	src/backend/optimizer/plan
查询执行	执行计划	执行生成的计划	src/backend/executor
	调度	将请求分配到合适的处理模块	src/backend/tcop
	附件命令处理	处理建表、备份等命令	src/backend/commands

5.2 查询分析

查询分析是查询编译的第一个模块，它包括词法分析、语法分析和语义分析三个部分。它将用户输入的 SQL 命令转换为查询树（Query 结构）。其中词法分析和语法分析分别借助词法分析工具 Lex 和语法分析工具 Yacc^①来完成各自的工作。

用户输入的 SQL 命令作为字符串传递给查询分析器，对其进行词法和语法分析生成分析树，然

① <http://dinosaur.compilertools.net/>
<http://tldp.org/HOWTO/Lex-YACC-HOWTO.html>
<http://www.ibm.com/developerworks/cn/linux/sdk/lex/index.html>

后进行语义分析得到查询树。查询分析的基本流程如图 5-2 所示。图 5-3 则展示了查询分析中主要的函数调用关系。对于用户的 SQL 命令，统一由 exec_simple_query 函数处理，该函数将调用 pg_parse_query 完成词法和语法分析并产生分析树，接下来调用 pg_analyze_and_rewrite 函数逐个对分析树进行语义分析和重写：在该函数中又会调用函数 parse_analyzer 进行语义分析并创建查询树（Query 结构），函数 pg_rewrite_query 则负责对查询树进行重写。本节将介绍词法、语法和语义分析，查询重写会在 5.3 节介绍。

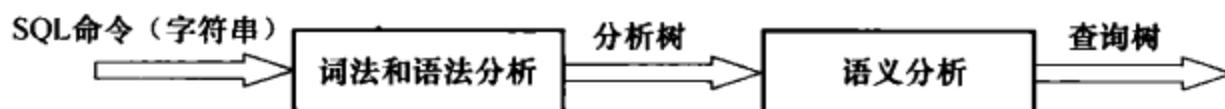


图 5-2 查询分析的执行流程

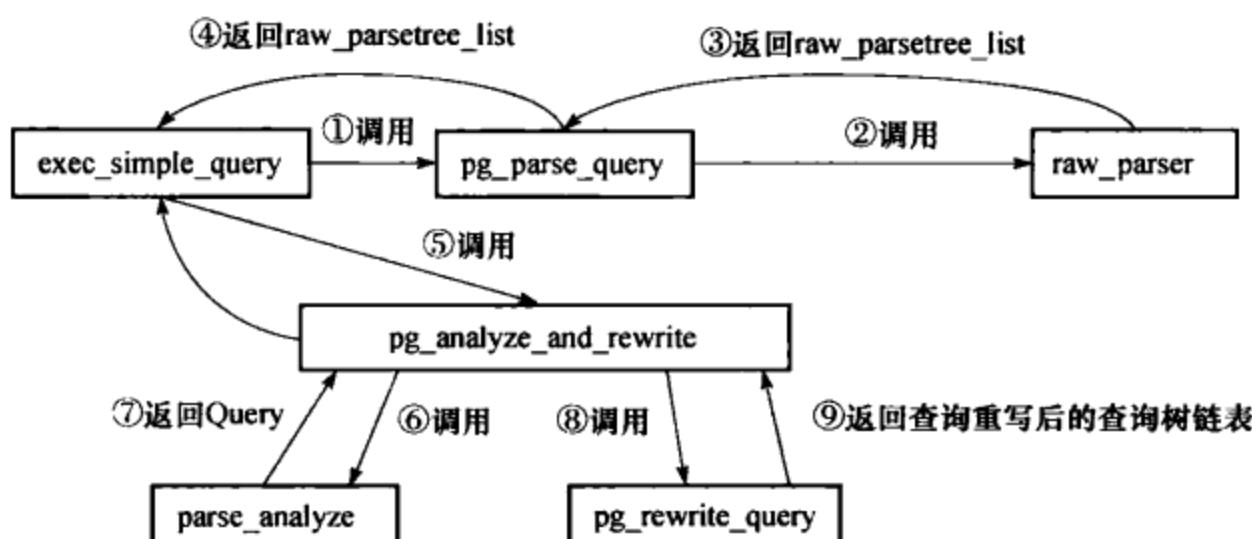


图 5-3 查询分析模块主要函数间的调用关系

查询分析的处理过程如下：

- 1) exec_simple_query 调用函数 pg_parse_query 进入词法和语法分析的主体处理过程，然后函数 pg_parse_query 调用词法和语法分析的入口函数 raw_parser 生成分析树（原始分析树链表 raw_parseetree_list）。

- 2) 函数 pg_parse_query 返回分析树给外部函数。

- 3) exec_simple_query 接着调用函数 pg_analyze_and_rewrite 进行语义分析和查询重写。首先调用函数 parse_analyze 进行语义分析并生成查询树（用 Query 结构体表示），之后会将查询树传递给函数 pg_rewrite_query 进行查询重写（在 5.3 节介绍）。

5.2.1 Lex 和 Yacc 简介

在 raw_parser 中，将通过 Lex 和 Yacc 生成的代码来进行词法和语法分析并生成分析树。Lex 和 Yacc 是词法与语法分析工具，两者相配合可以生成用于词法和语法分析的 C 语言源代码。在 raw_parser 中，就是通过调用采用 Lex 和 Yacc 预生成的函数 base_yyparse 来实现词法和语法分析工作。Lex 用来生成扫描器，其工作是识别一个一个的模式，比如数字、字符串、特殊符号等，然后将其传给 Yacc。Yacc 则用于生成语法分析器，它识别模式的组合，即语法。它们通过共同的符号表传递标识符，通过内置变量 yylval 传递表示的值。为了更好地介绍 PostgreSQL 的查询分析工作，本节首先将介绍 Lex 和 Yacc。

1. 词法分析工具 Lex

Lex 是通用的用于生成扫描器的工具，它利用正则表达式来表达模式。常见的正则表达式符号说明如下：

- “.”：匹配任何除 \n 之外的单个字符。
- “\$”：匹配行结束符。
- “^”：匹配一行的开始。
- “\”：这是转义字符，例如，正则表达式 \\$ 被用来匹配美元符号，而不是行尾。
- “[]”：匹配括号中的任何一个字符，例如 [a - z] 将匹配从 a 到 z 的任何一个字符。
- “+”：匹配 1 个或多个正好在它之前的那个字符。例如，正则表达式 9+ 匹配 9、99、999 等。
- “?”：匹配 0 或 1 个正好在它之前的那个字符。
- “|”：匹配前面的正则表达式或随后的正则表达式。例如，cow | pig | sleep 可以匹配 cow、pig 和 sleep 中的任何一个。
- “?”：匹配前面的正则表达式 0 次或一次。如 -? [0 - 9]+ 将匹配带符号的整数。
- “{}”：当括号中包含 1 个或者 2 个数字时，指示前面的模式被允许匹配多少次。例如 A {1, 3}，匹配 A 1 次到 3 次。

正则表达式示例：

- “[0 - 9]+”：匹配整数。
- “-? [0 - 9]+”：匹配带符号的整数。
- “-? ([0 - 9]*\. [0 - 9]+)”：匹配带符号的小数。
- “-? (([0 - 9]+) | ([0 - 9]*\. [0 - 9]+))”：匹配带符号的实数。
- “[eE][- +]? [0 - 9]+”：匹配用 e 表示的浮点数。
- “-? ((([0 - 9]+) | ([0 - 9]*\. [0 - 9]+))([eE][- +]? [0 - 9]+)?)”：匹配带 e 或不带 e 的任意实数。

词法分析的主要工作是利用正则表达式在给定字符序列中进行模式匹配，这些正则表达式通常写在一个后缀为 “.l” 的文件（Lex 文件）中，通过 Lex 命令可以从 Lex 文件生成一个包含有扫描器的 C 语言源代码，其他源代码可以通过调用该扫描器来实现词法分析。Lex 文件中除了正则表达式之外还有一些其他内容，下面通过一个简单的示例来说明 Lex 文件的基本语法。

Lex 文件示例

```

1.% {
2./* example.l
3.* 这是一个简单的 Lex 示例          /* 定义段 */
4.这个段里面可以包含 C 头文件以及注释
5.*/
6.#include <stdlib.h>
7.% }
8.% %
9.[ \n\t];
10.-? [0 - 9]+ {printf("num = % d\n", atoi(yytext));}
11.. ECHO;

```

```

12. % %
13. main()           /* 代码段 */
14. {
15.     yylex();
16. }

```

Lex 文件分为三段，分别是定义段、规则段和代码段，各段之间由“% %”符号分隔。

- 定义段可以包含任意的 C 语言头文件、符号说明等，其代码会被直接拷贝到生成的扫描器代码文件中。
- 规则段利用正则表达式来匹配模式，每当成功匹配一个模式，就对应其后“{}”中的代码，这些匹配规则及动作都会反映到最后生成的扫描器代码文件中。
- 代码段可以是任意 C 代码，但是其中必须要调用 Lex 提供的函数（见表 5-2）。上例中使用的 yylex 就是其中之一，它完成实际的分析工作。在代码段中，只需要调用 yylex 即可，其实现代码将由 Lex 工具根据规则段的内容来生成。代码段中的内容将会被直接拷贝到生成的扫描器代码文件中。

上面的例子中 Lex 文件生成的 C 程序可以用来识别一个整数，并将其值输出。生成扫描器代码并编译的命令如下：

```

lex example.l
cc lex.yy.c -o first-ll

```

Lex 命令从文件“example.l”生成“lex.yy.c”文件，cc 命令中的参数 -ll 用于链接 lex 库，编译生成的程序名为 first。运行 first 后输入任意字符串将会输出结果，例如输入“123”后会输出结果“num = 123”。

表 5-2 Lex 函数

函数	作用
yylex	这一函数开始分析工作。它由 Lex 自动生成
yywrap	这一函数在文件（或输入）的末尾调用。如果函数的返回值是 1，就停止解析。因此它可以用来解析多个文件。代码可以写在第三段，这样就能够解析多个文件。方法是使用 yyin 文件指针（见表 5-3）指向不同的文件，直到所有的文件都被解析。最后，yywrap 可以返回 1 来表示解析的结束
yyless	这一函数可以用来送回除了前 n 个字符外的所有读出标记
yymore	这一函数告诉 Lexer 将下一个标记附加到当前标记后

Lex 生成的扫描器在成功识别模式后，通过 Lex 变量（参见表 5-3）存储该模式对应的值，并将其返回给上层调用者（通常是 Yacc 生成的语法分析器）。

表 5-3 Lex 变量

变量	作用
yyin	FILE * 类型。它指向 lexer 正在解析的当前文件
yyout	FILE * 类型。它指向记录 lexer 输出的位置。缺省情况下，yyin 和 yyout 都指向标准输入和输出
yytext	匹配模式的文本存储在这一变量中（字符串类型）
yylen	给出匹配模式的长度
yylineno	提供当前的行数信息（某些 Lex 的实现不支持）

2. 语法分析工具 Yacc

语法分析的主要工作是从给定模式序列输入中寻找某一特定语法结构。例如，有一个简单的英文句子的语法定义“Sentence→subject verb object, Object→ noun, Subject→noun | pronoun”。这个定义表明句子（sentence）是由主语（subject）、动词（verb）和宾语（object）的序列构成，而宾语由一个名词（noun）组成，主语由名词或者代词（pronoun）组成。因此“He goes home”就是符合该语法的句子，而“He home”不符合该语法，因为后一个句子缺乏动词，不能匹配该语法的模式。Yacc 也以类似的方式来定义语法。Yacc 的工作方式和 Lex 相似，也是将语法的定义以及一些必要的 C 语言代码写在一个后缀为“.y”的文件（Yacc 文件）中，然后使用 Yacc 命令由该文件生成具有语法分析功能的 C 语言代码文件。

下面以一个简单的只识别加减法的计算器的 Yacc 文件为例，对 Yacc 文件的格式进行介绍。可以看到，Yacc 文件同样分为三个段：定义段、规则段和代码段，其中各个段由“% %”符号分隔。

- 定义段可以是 C 代码，如包含头文件和函数声明（“% { … % }”中的内容），同时也可定义 Yacc 的内部标志等。在这个计算器的例子中，% token 定义标识符，% left 指定运算法的优先级，% start 表示从标识符开始解析。
- 规则段用来表示语法规则，并在每识别一个语法规则后，根据规则后面“{|”内的代码进行相应的处理操作。其中 \$\$ 代表语法表达式的左边结构的值，而 \$1 代表语法表达式右边结构第一个标识符对应的值，依此类推。例子中的第 12 行的语法表达式就说明 expression 可以由一个 expression、“+”号和一个 NUMBER 构成，如果有一个模式能够匹配这个语法规则，则左边 expression 的值（\$\$）等于右边 expression 的值（\$1）加上 NUMBER 的值（\$3），中间跳过没有用到的 \$2 指的是“+”号。
- 代码段包括 C 代码，它将被直接拷贝到生成的 C 文件中。同样代码段中也包含一些 Yacc 函数和一些 Lex 传递给 Yacc 的变量。

简单计算器的 Yacc 文件 calculator.y

```

1.% {
2./* calculator.y*/
3.#include <stdio.h>
4.% }
5.% token NAME NUMBER
6.% left '+' '-' 
7.% start statement
8.% %
9.statement : NAME '=' expression/* 规则段 */
10.| expression{printf("=%d\n", $1);}
11.|
12.expression:expression '+' NUMBER{ $$ = $1 + $3;}
13.|           expression '-' NUMBER{ $$ = $1 - $3;}
14.|           NUMBER{ $$ = $1;}
15.|
16.% %
17.main()

```

```

18. {
19. do                                /*代码段*/
20. {
21. yyparse();
22. }
23. while(! feof(yyin));
24. }

```

从上面结构可以看出, main 函数不断地执行 yyparse 来进行语法分析。yyparse 函数识别某种语法结构并进行相应的行为处理。可以看到, 该语法结构是每当识别一个加法表达式时, 就把加法两边的数相加; 每当识别减法表达式时, 就把减号两边的数相减。这样就完成了加减功能。NAME 和 NUMBER 这些符号只是定义, 它们的值由 Lex 通过变量传送给 Yacc。注意, 我们虽然没有在例子中看到对 Lex 变量的直接使用, 但是 Yacc 在生成代码时会使用这些变量。因此, Yacc 是无法单独运行的, 要与 Lex 配合使用。计算器例子对应的 Lex 文件如下。

计算器例子中的 Lex 文件 calculator.l

```

1  % {
2  /* calculator.l */
3  #include "y.tab.h"
4  extern int yylval;
5  %
6  %
7  [0-9]+ {yylval = atoi(yytext); return NUMBER;}
8  [\t] ; /* ignore white space */
9  \n      return 0; /* logical EOF */
10 .       return yytext[0];
11 %

```

可以看到, Lex 每识别一个数字, 将以标识符 NUMBER 返回给 Yacc (Lex 文件第 7 行), 并将它的值传给 yylval。而 Yacc 根据 NUMBER 标识可以判断 yylval 中存放的是一个数字, 并从 yylval 读取其值。例如, 前面提到的 \$1 等就是从 yylval 中取值。

利用例子中的 Lex 文件和 Yacc 文件构建一个计算器程序的命令如下:

```

yacc -d calculator.y                  //生成 y.tab.c 和 y.tab.h
lex calculator.l                      //生成 lex.yy.c
cc -o calculator y.tab.c lex.yy.c -ly -ll //编译和链接 C 文件

```

注意 不同版本的 Yacc 工具其命令可能不同。例如, GNU 的 Yacc 工具是 bison。同样, 不同操作系统下编译和链接 C 文件的命令也可能并不是 cc。

最后生成的程序名为 calculator, 运行 calculator 之后可以输入一个完整的加减法公式, 该程序会识别之并计算结果返回, 例如我们可以输入 “1 + 2”, 该程序会返回结果: “=3”。

5.2.2 词法和语法分析

和前面介绍的一样, PostgreSQL 中的词法分析和语法分析也是由 Lex 和 Yacc 配合完成的。Post-

greSQL 的源代码中包含 `scan.l` 和 `gram.y` 两个文件，并且已经为两个文件预生成了 C 文件，分别是 `scan.c` 和 `gram.c`。如果对 `scan.l` 和 `gram.y` 进行了修改，在编译 PostgreSQL 的时候就会重新生成 `scan.c` 和 `gram.c`。这两个 C 文件再配合上词法和语法分析模块需要的 C 文件就构成了整个模块。词法和语法分析模块的源代码位于目录 `src/backend/parser` 下，其中主要源文件的说明见表 5-4，文件之间的关系如图 5-4 所示。`scan.l` 的主要目标是识别出 PostgreSQL 中使用的所有关键字等，这里不再详述，有兴趣的读者可以参考 PostgreSQL 文档中关于保留字的部分。本节将主要分析 SQL 的语法规则以及经过语法分析后在内存中的数据结构。

表 5-4 词法和语法分析的源文件

源文件	说 明	
<code>parser.c</code>	词法、语法分析的主入口文件，入口函数是 <code>raw_parser</code> ；对查询语句进行词法和语法分析后，返回分析树	
<code>kwlookup.c</code>	函数 <code>ScanKeywordLookup</code> ，查找单词表，返回当前标识符指向单词表中对应单词的指针	
<code>scansup.c</code>	提供词法分析时的常用函数	
	<code>scanstr</code>	处理转义字符
	<code>downcase_truncate_identifier</code>	将大写英文字符转换为小写
	<code>truncate_identifier</code>	如果标识符长度超过规定的最大标识符长度（64），则将其截断
<code>scan.l</code>	<code>scanner_isspace</code>	判断是否为空白符（空格， <code>\t</code> , <code>\n</code> , <code>\r</code> , <code>\f</code> ）
	定义词法结构，用 Lex 语言书写，用 Lex 编译后生成 <code>scan.c</code> 文件	
	定义关键字的数值编号	
	定义语法结构，用 Yacc 语言书写，用 Yacc 编译后生成 <code>gram.c</code> 文件	

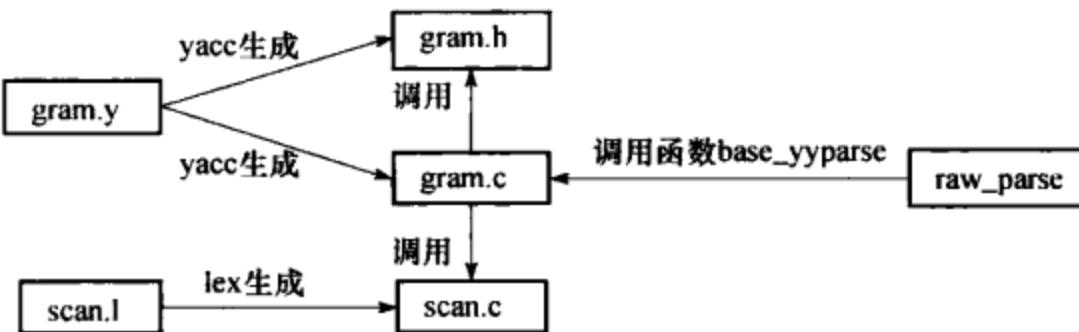


图 5-4 词法和语法分析的各文件生成和调用关系

词法和语法分析的入口函数是 `raw_parser`，该函数返回的 List 结构用于存储生成的分析树。现以 `SELECT` 语句为例讲解 PostgreSQL 中查询语句如何被解析并生成分析树。

SELECT 语句的语法

```

[ WITH [ RECURSIVE ] with_query [ , ... ] ]
SELECT [ ALL | DISTINCT [ ON ( expression [ , ... ] ) ] ]
      * | expression [ [ AS ] output_name ] [ , ... ]
      [ FROM from_item [ , ... ] ]
      [ WHERE condition ]
      [ GROUP BY expression [ , ... ] ]
      [ HAVING condition [ , ... ] ]
  
```

```
[WINDOW window_name AS (window_definition)[, ...]]
[ (UNION | INTERSECT | EXCEPT)[ ALL] select]
[ ORDER BY expression [ASC | DESC | USING operator]
    [NULLS{FIRST | LAST}] [, ...]]
[LIMIT{count | ALL}]
[OFFSET start [ROW | ROWS]]
[FETCH{FIRST | NEXT}[ count] {ROW | ROWS}ONLY]
[FOR{UPDATE | SHARE}[ OF table_name [, ...]] [NOWAIT] [...]]
```

SELECT语句在文件gram.y的定义位于第6841行到6847行：

SELECT语句在gram.y中的定义

```
6841 SelectStmt: select_no_parens          % prec UMINUS
6842     | select_with_parens           % prec UMINUS
6843 ;
6844 select_with_parens:
6845     '(' select_no_parens ')'      { $$ = $2; }
6846     | '('select_with_parens ')'   { $$ = $2; }
6847 ;
```

SELECT语句（用标识符SelectStmt表示）定义为带括号的SELECT语句（用标识符select_with_parens表示）和不带括号的SELECT语句（用标识符select_no_parens表示），如第6841行到第6843行所示。带括号的SELECT语句最终会定义为括号和不带括号的SELECT语句的序列来处理，如第6844行到第6847行所示。因此，最终要处理的是不带括号的形式，select_with_parens在文件gram.y中定义在第6860行到第6915行中。

不带括号SELECT的语法定义

```
6860 select_no_parens:
6861     simple_select{$$ = $1;}
6862     | select_clause sort_clause
6863     {
6864         insertSelectOptions((SelectStmt*)$1, $2, NIL,
6865                         NULL, NULL, NULL);
6866         $$ = $1;
6867     }
6868     | select_clause opt_sort_clause for_locking_clause opt_select_limit
6869     {
6870         insertSelectOptions((SelectStmt*)$1, $2, $3,
6871                         list_nth($4, 0), list_nth($4, 1),
6872                         NULL);
6873         $$ = $1;
6874     }
6875     | select_clause opt_sort_clause select_limit opt_for_locking_clause
6876     {
6877         insertSelectOptions((SelectStmt*)$1, $2, $4,
```

```

6878                               list_nth($3, 0), list_nth($3, 1),
6879                               NULL);
6880     $$ = $1;
6881 }
.....
6910 ;
6911
6912 select_clause:
6913     simple_select
6914     | select_with_parens
6915 ;

```

从以上代码可以看出，一条不带括号的 SELECT 语句可以定义为一条简单的 SELECT 语句（用标识符 simple_select 表示），也可以定义为在简单 SELECT 语句后接排序子句（标识符 sort_clause）、LIMIT 子句（标识符 select_limit）等构成的复杂语句（例如第 6862 行到第 6910 行）。可以看到，对于整个语句的语法分析实际上是将语句拆分成很多小的语法单元，然后分别对这些小的单元进行分析。因此，只要弄清楚各个小语法单元的做法，就可以把它们组合在一起形成整个语法树。这里我们着重分析 simple_select 的语法定义。

简单 SELECT 的语法定义

```

6940 simple_select:
6941     SELECT opt_distinct target_list
6942     into_clause from_clause where_clause
6943     group_clause having_clause window_clause
6944 {
6945     SelectStmt* n = makeNode(SelectStmt);           /* 创建 SELECT 节点 */
6946     n->distinctClause = $2;                         /* DISTINCT 子句 */
6947     n->targetList = $3;                            /* 目标属性 */
6948     n->intoClause = $4;                           /* SELECT INTO 子句 */
6949     n->fromClause = $5;                            /* FROM 子句 */
6950     n->whereClause = $6;                           /* WHERE 子句 */
6951     n->groupClause = $7;                           /* GROUP BY 子句 */
6952     n->havingClause = $8;                          /* HAVING 子句 */
6953     n->>windowClause = $9;                          /* 窗口子句 */
6954     $$ = (Node*)n;
6955 }
6956 | values_clause{ $$ = $1;}
.....
6988 ;

```

simple_select 是一条 SELECT 语句的最核心部分，从 simple_select 的语法定义可以看出，它由如下子句组成：去除行重复的 DISTINCT（标识符 opt_distinct）、目标属性（标识符 target_list）、SELECT INTO 子句（标识符 into_clause）、FROM 子句（标识符 from_clause）、WHERE 子句（标识符 where_clause）、GROUP BY 子句（标识符 group_clause）、HAVING 子句（标识符 having_clause）和

窗口子句（标识符 window_clause）。在成功匹配 simple_select 语法结构后，将创建一个 SelectStmt 结构体（数据结构 5.1），并将各子句赋值到结构体中相应的字段。

窗口函数

PostgreSQL 8.4 之后的版本中增加了窗口函数的功能。窗口函数类似于聚集函数，但使用聚集函数返回的是各个分组的结果，而窗口函数为每一行返回结果。例如：

```
SELECT depname, empno, salary, avg(salary) OVER (PARTITION BY depname) FROM empsalary;
depname | empno | salary | avg
-----+-----+-----+-----
develop | 11 | 5200 | 5020.0000000000000000
develop | 7 | 4200 | 5020.0000000000000000
develop | 9 | 4500 | 5020.0000000000000000
develop | 8 | 6000 | 5020.0000000000000000
develop | 10 | 5200 | 5020.0000000000000000
personnel | 5 | 3500 | 3700.0000000000000000
personnel | 2 | 3900 | 3700.0000000000000000
sales | 3 | 4800 | 4866.6666666666666667
sales | 1 | 5000 | 4866.6666666666666667
```

此外，simple_select 还可以定义为其他的形式（见文件 Gram.y 中第 6956 行到第 6988 行），如 VALUES 子句、关系表达式以及多个 SELECT 语句的交并差等，但这些情况最终都会转化成最基本的 simple_select 形式来处理。而对于 simple_select 来说，目标属性（标识符 target_list）、FROM 子句（标识符 from_clause）、WHERE 子句（标识符 where_clause）以及 GROUP BY 子句（标识符 group_clause）是最重要的部分，下面将结合一个具体实例对这些子句的处理加以详细介绍。

数据结构 5.1 SelectStmt

```
typedef struct SelectStmt
{
    NodeTag           type;          // 节点类型，用于标识节点的内容，在 SelectStmt
                                      // 中会取值为 T_SelectStmt
    List              *distinctClause; // DISTINCT 子句
    IntoClause        *intoClause;   // SELECT INTO / CREATE TABLE AS 子句
    List              *targetList;   // 目标属性
    List              *fromClause;   // FROM 子句
    Node              *whereClause;  // WHERE 子句
    List              *groupClause;  // GROUP BY 子句
    Node              *havingClause; // HAVING 子句
    List              *windowClause; // 窗口子句
    WithClause        *withClause;   // WITH 子句
    List              *valuesLists;   // VALUES 列表，产生“常量表”，常量表可以作为虚拟
                                      // 表出现在 FROM 中
```

```

List           * sortClause;      //ORDER BY 子句
Node          * limitOffset;    //OFFSET 子句
Node          * limitCount;     //LIMIT 子句
List           * lockingClause; //FOR UPDATE 子句
SetOperation   op;            //查询语句的集合操作，交/并/差
bool          all;            //在集合操作时是否指定了 ALL 关键字
struct SelectStmt * larg;    //左孩子结点
struct SelectStmt * rarg;    //右孩子结点

} SelectStmt;

```

表 5-5 ~ 表 5-8 给出了例子中涉及的 4 个表的模式，在这些表的基础上要完成如下查询：查询 2005 级各班高等数学的平均成绩，且只查询平均分在 80 以上的班级，并将结果按照升序排列。

表 5-5 班级信息表 class

属性名	类 型	说 明
classno	character varying (20)	班级号，主键
classname	character varying (30)	班级名
gno	character varying (20)	年级号，外键，引用 grade 表的 gno 属性

表 5-6 学生信息表 student

属性名	类 型	说 明
sno	character varying (20)	学号，主键
sname	character varying (30)	学生姓名
sex	character varying (5)	学生性别
age	integer	年龄
nation	character varying (20)	国籍
classno	character varying (20)	所在班级号，外键，引用 class 表的 classno 属性

表 5-7 课程信息表 course

属性名	类 型	说 明
cno	character varying (20)	课程号，主键
cname	character varying (30)	课程名称
credit	integer	学分
priorcourse	character varying (20)	先导课程，外键，引用 course 表的 cno 属性

表 5-8 成绩表 sc

属性名	类 型	说 明
sno	character varying (20)	学号，主键，外键，引用 student 表的 sno 属性
cno	character varying (20)	课程号，主键，外键，引用 course 表的 cno 属性
score	integer	分数

在这四个表上完成指定查询的 SELECT 语句如下：

```

1 SELECT classno,classname, AVG(score)AS avg_score
2 FROM sc, (SELECT* FROM class WHERE class.gno = '2005')AS sub

```

```

3 WHERE sc.sno IN (SELECT sno FROM student WHERE student.classno = sub.classno) AND
      sc.cno IN (SELECT course.cno FROM course WHERE course cname = '高等数学')
4 GROUP BY classno, classname
5 HAVING AVG(score) > 80.0
6 ORDER BY avg_score;

```

下面将对这个 SELECT 语句的各个子句语法结构的定义以及在 SelectStmt 中对应的数据结构进行逐一介绍。

1. DISTINCT 子句

DISTINCT 子句对应语法定义中的标识符 opt_distinct。从 opt_distinct 的语法结构可以看到，它可以匹配 DISTINCT、ALL、DISTINCT ON（表达式列表）或者为空，用来决定 SELECT 语句是否去除重复的行。

- 当匹配到 DISTINCT 时，opt_distinct 返回一个 List，该链表的第一个 ListCell 的 ptr_value 字段置为空。
- 当匹配到 DISTINCT ON 时，opt_distinct 也返回一个 List，这个 List 中包含了跟在 DISTINCT ON 之后的表达式的列表（星号或者表的属性等）。
- 当匹配到 ALL 或者空时，opt_distinct 返回 NIL，表明没有使用 DISTINCT。

DISTINCT 子句的语法定义

```

7104 opt_distinct:
7105     DISTINCT                                { $$ = list_makel(NIL); }
7106     | DISTINCT ON '(' expr_list ')'
7107     | ALL                                     { $$ = NIL; }
7108     | /* EMPTY */ /
7109 ;

```

2. 目标属性

目标属性是 SELECT 语句中所要查询的属性列表，对应着语法定义中的标识符 target_list。target_list 由若干个 target_el 组成（文件 gram.y 中第 9828 行到 9831 行）。target_el 定义为取别名的表达式、表达式以及“*”等（文件 gram.y 中第 9833 行到 9877 行）。

目标属性的语法定义

```

9828 target_list:
9829     target_el
9830     | target_list ',' target_el           { $$ = lappend($1, $3); }
9831     ;
9832
9833 target_el:    a_expr AS ColLabel
9834             {
9835                 $$ = makeNode(ResTarget);
9836                 $$ -> name = $3;
9837                 $$ -> indirection = NIL;

```

```

9838             $$->val = (Node* ) $1;
9839             $$->location = @ 1;
9840         }
9841         .....
9877     ;

```

当成功匹配一个 target_el 时，创建一个 ResTarget 结构体（数据结构 5.2），该结构体中存储了该属性的全部信息。最终 target_list 将返回一个由 ResTarget 构成的 List。

数据结构 5.2 ResTarget

```

typedef struct ResTarget
{
    NodeTag      type;
    char         *name;           //用 AS 指定的目标属性名称,没有则为空
    List         *indirection;   //通过星号、下标等引用的目标属性,没有则为 NIL
    Node         *val;            //指向各种表达式
    int          location;       //符号出现的位置
}ResTarget;

```

target_list 包括若干 ListCell 节点，每个节点中的 data 字段指向结构体 ResTarget，用来表示目标属性中的一项。图 5-5 展示了目标属性在内存中的组织结构。当目标属性中的某项涉及函数调用时，ResTarget 中的字段 val 会指向结构体 FunCall。FunCall 的字段 funcname 存储函数的名称，字段 args 指向结构体 ColumnRef 构成的链表，每一个 ColumnRef 存储了函数调用中所使用到的表的一个属性。如果没有函数调用，则结构体 ResTarget 中的字段 val 直接指向结构体 ColumnRef，存储该项目标属性所涉及的表的字段（此种情况没有在图 5-5 中展示）。

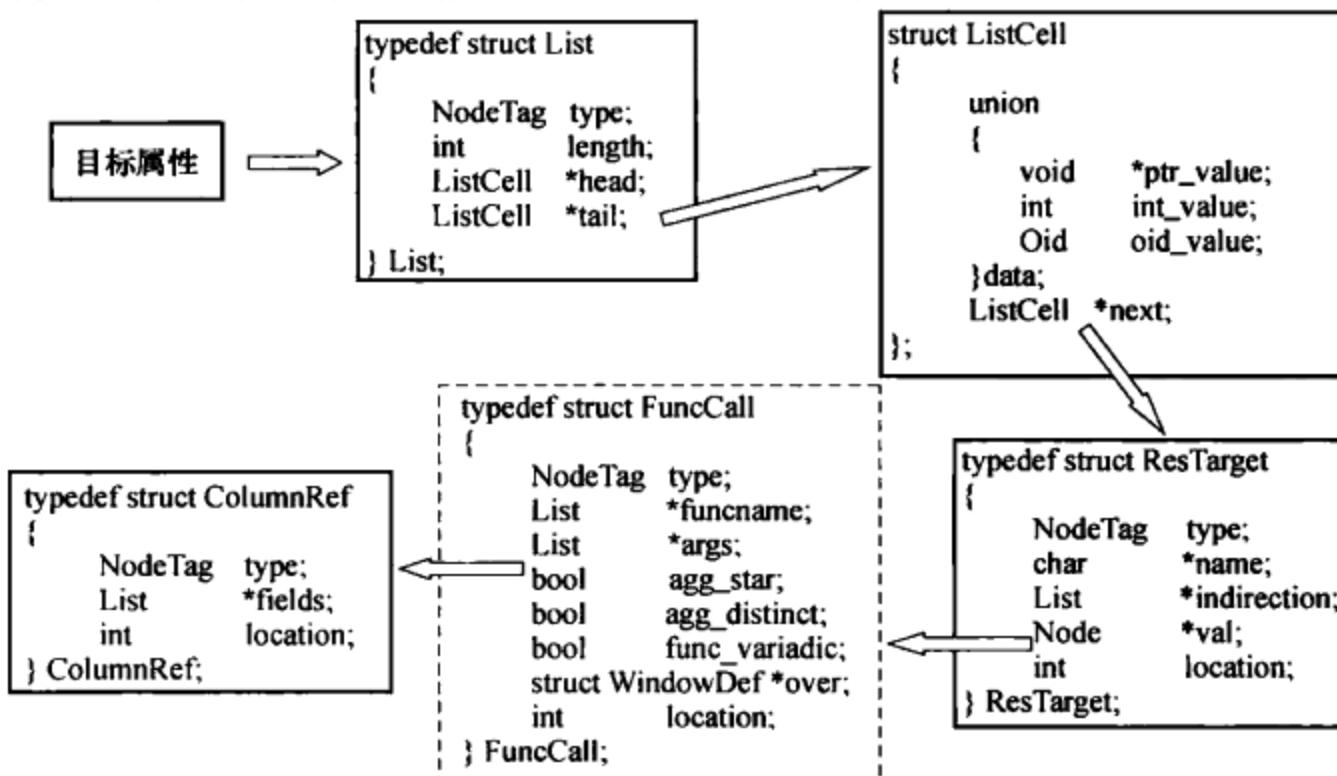


图 5-5 分析树中目标属性的数据组织结构

在省去 List 结构体、ListCell 结构体以及一些非关键字段等内容后，例 5.1 中 SELECT 语句对应分析树的目标属性在内存中的数据组织结构如图 5-6 所示。

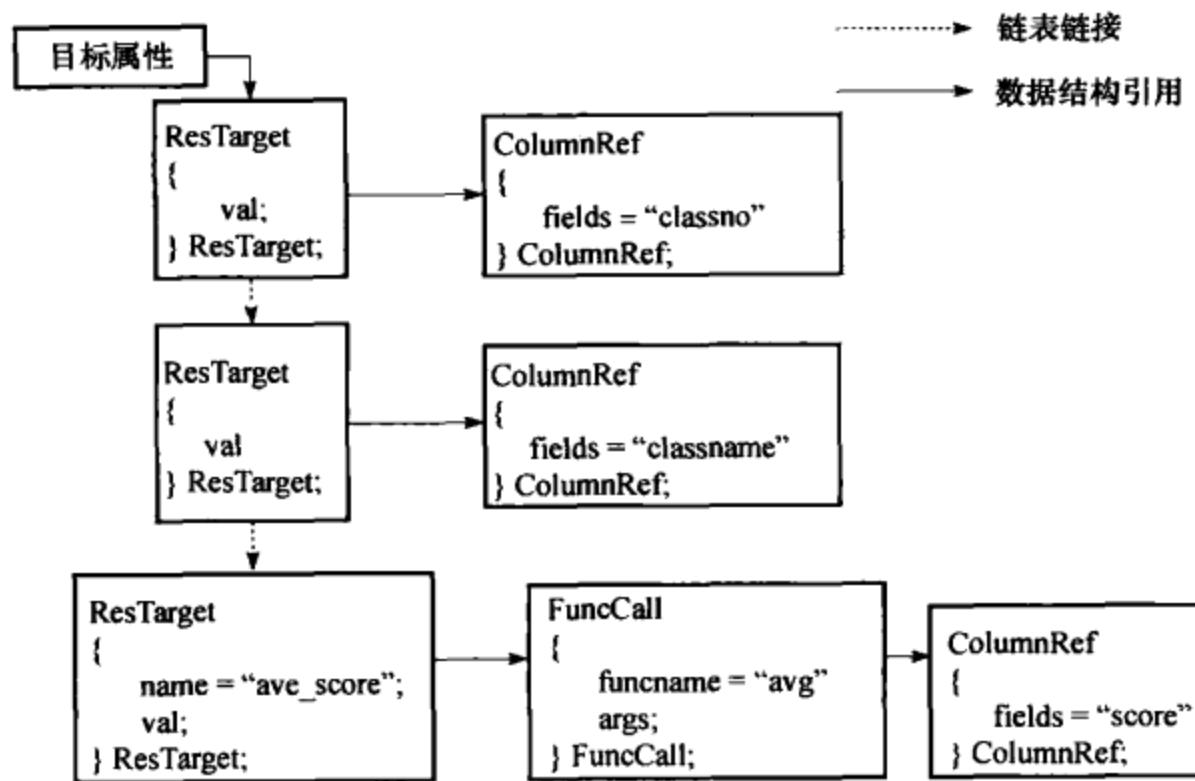


图 5-6 例 5.1 对应的目标属性内存组织结构

3. FROM 子句

文件 gram.y 中定义的标识符 from_clause 表示 SELECT 语句中的 FROM 子句，from_clause 由 FROM 关键字和 from_list 组成。而 from_list 则由若干个标识符 table_ref 组成，每一个 table_ref 表示 FROM 子句中用逗号分隔的每个子项，它表示在 FROM 子句中出现的一个表或者一个子查询。

FROM 子句的语法定义

```

7297 from_clause:
7298     FROM from_list                                     { $$ = $2; }
7299     /* EMPTY* /                                       { $$ = NIL; }
7300 ;
7301
7302 from_list:
7303     table_ref                                         { $$ = list_makel($1); }
7304     | from_list ',' table_ref                         { $$ = lappend($1, $3); }
7305 ;

```

标识符 table_ref 可以定义为关系表达式（文件 gram.y 的第 7314 行到 7317 行）、取别名的关系表达式（文件 gram.y 的第 7318 行到 7322 行）、带括号的 SELECT 语句、表连接等形式。

FROM 子句中每个子项 table_ref 的语法定义

```

7314 table_ref:      relation_expr
7315           {
7316               $$ = (Node*) $1;

```

```

7317          }
7318      | relation_expr alias_clause
7319      {
7320          $1 -> alias = $2;
7321          $$ = (Node*) $1;
7322      }
7323      | func_table
.....
7409      ;

```

由于 FROM 子句中子项（标识符 table_ref）的最简单和基本的形式是关系表达式（标识符 relation_expr）。因此，下面分析 relation_expr 的语义定义。

关系表达式的语义定义

```

7548 relation_expr:
7549     qualified_name
7550     {
7551         /* default inheritance */
7552         $$ = $1;
7553         $$->inhOpt = INH_DEFAULT;
7554         $$->alias = NULL;
7555     }
7556     | qualified_name '*' '
.....
7577     ;
.....
9903 qualified_name:
9904     relation_name
9905     {
9906         $$ = makeNode(RangeVar);
9907         $$->catalogname = NULL;
9908         $$->schemaname = NULL;
9909         $$->relname = $1;
9910         $$->location = @ 1;
9911     }
9912     | relation_name indirection
.....
9938     ;

```

关系表达式 relation_expr 定义成 qualified_name、带 ONLY 关键字的 qualified_name 等形式，最后 qualified_name 定义成 relation_name。在成功匹配最终的标识符 relation_name 后，创建一个 RangeVar 结构体（数据结构 5.3）用来存储该关系的信息。

数据结构 5.3 RangeVar

```

typedef struct RangeVar
{
    NodeTag          type;
    char             *catalogname;      //表所在的数据库名,如果为空表示在当前数据库中
    char             *schemaname;       //表所在的模式名,如果为空表示在当前模式中
    char             *relname;         //表或者序列名称
    InhOption        inhOpt;          //是否将该表上的操作递归到它的子表上
    bool             istemp;          //是否为临时表
    Alias            *alias;           //表的别名
    int              location;        //符号出现的位置
}RangeVar;

```

如图 5-7 所示, from_clause 子句在分析树中同样被组织成一个 List, 每一个 ListCell 中包含一个 RangeVar 结构 (或者其他结构)

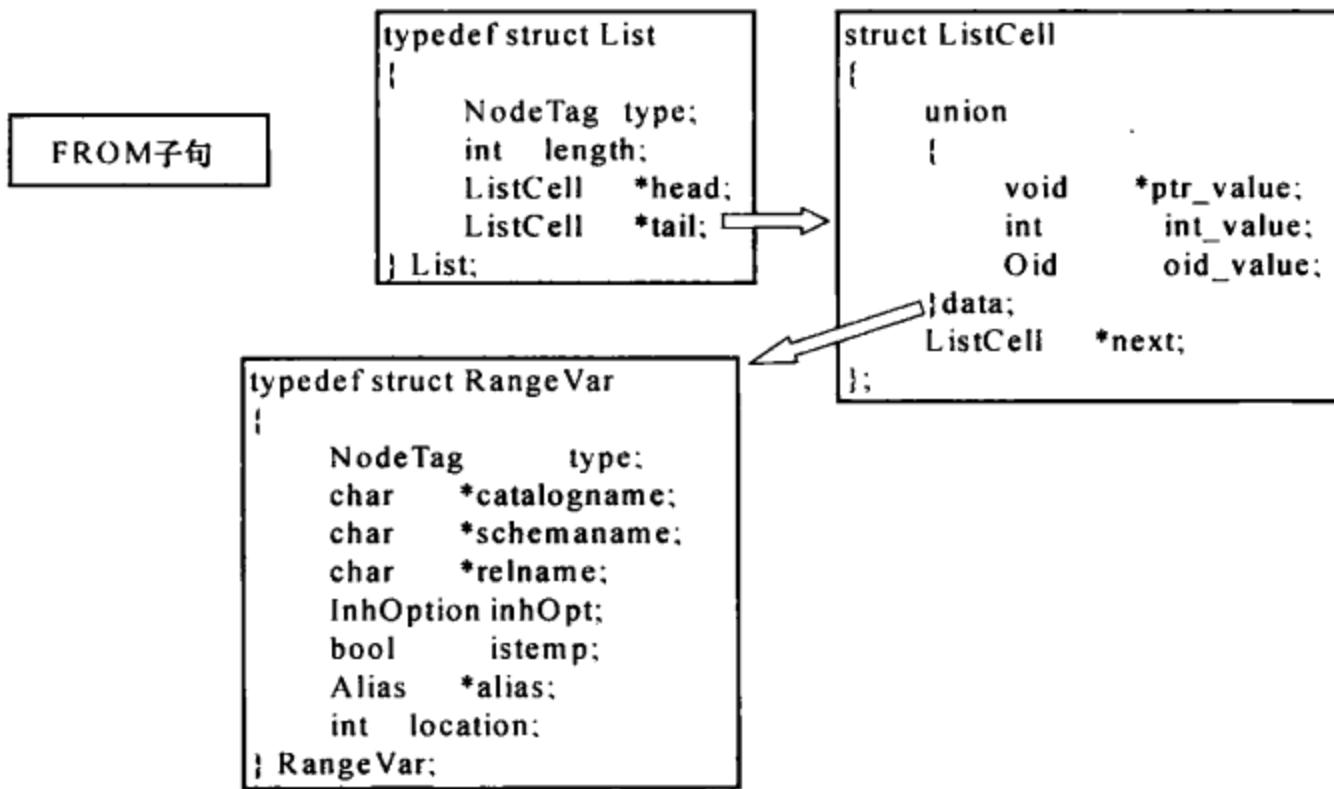


图 5-7 分析树中 FROM 子句的数据组织结构

在省去 List 结构体、ListCell 结构体以及一些非关键字段等内容后, 例 5.1 中分析树的 FROM 子句在内存中的数据组织结构如图 5-8 所示。

4. WHERE 子句

WHERE 子句中定义的是元组约束信息, 对应着语法定义中的标识符 where_clause。标识符 where_clause 定义为关键字 WHERE 和一个表达式 (标识符 a_expr)。

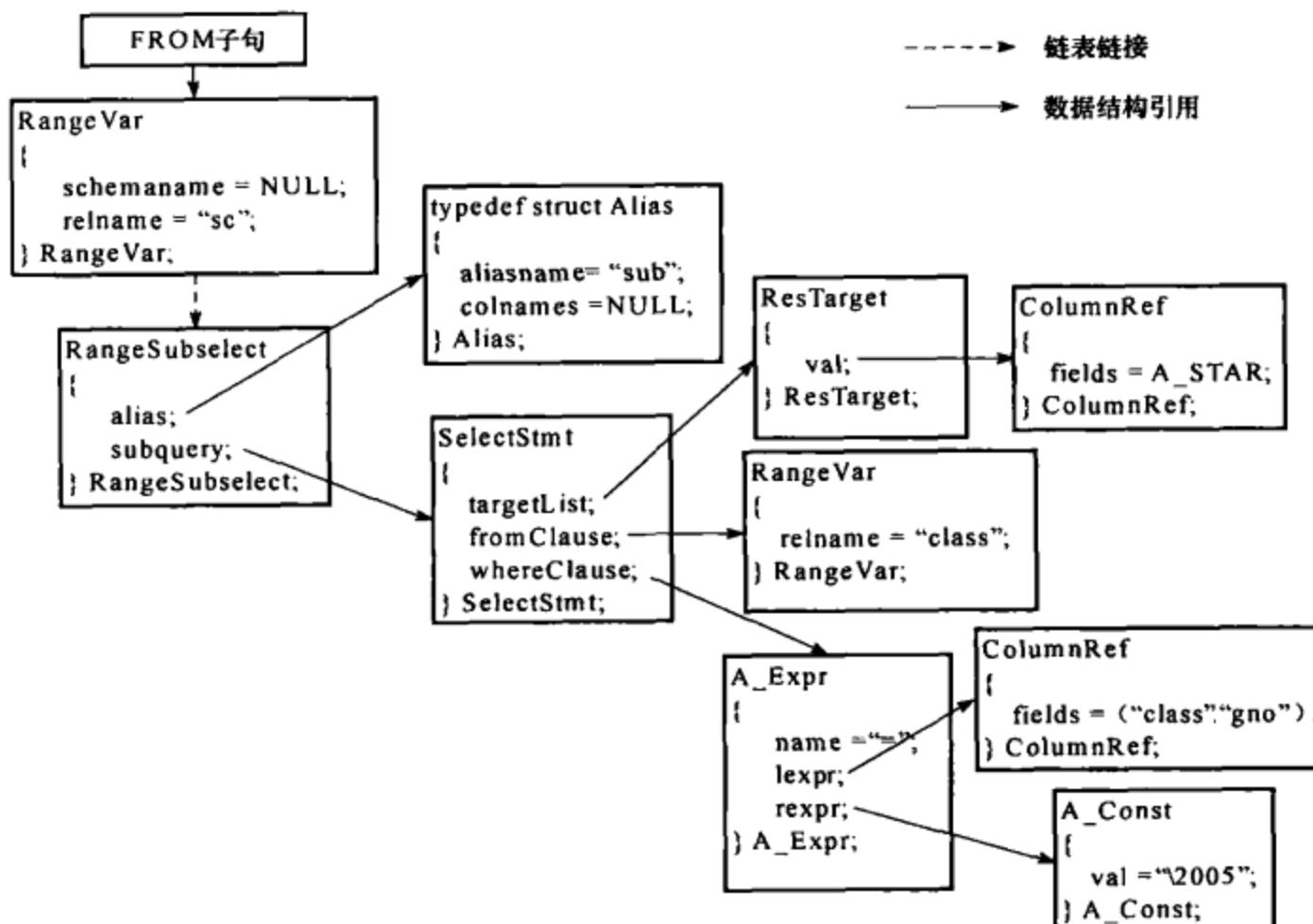


图 5-8 例 5.1 对应的 FROM 子句的内存组织结构

WHERE 子句的语法定义

```

7620 where_clause:
7621     WHERE a_expr      { $$ = $2; }
7622     /* EMPTY */      { $$ = NULL; }
7623 ;

```

由于表达式是递归定义的，因此，A_Expr 结构体（数据结构 5.4）中字段 lexpr 和 rexpr 分别代表操作符的左右两个子表达式，字段 A_Expr_kind 代表操作的类型。如果该表达式是常量或者属性等（表达式树中的叶子节点），则 lexpr 和 rexpr 都为 NULL。在 WHERE 子句中，使用到的表的属性信息用 ColumnRef 结构体来组织。

数据结构 5.4 A_Expr

```

typedef struct A_Expr
{
    NodeTag      type;
    A_Expr_Kind kind;           // 表达式类型 (AND、OR、ANY、IN 等)
    List         *name;          // 操作符名称
    Node         *lexpr;         // 左子表达式
    Node         *rexpr;         // 右子表达式
    int          location;       // 符号出现的位置
} A_Expr;

```

在省去 List 结构体、ListCell 结构体以及一些非关键字段等内容后，例 5.1 中分析树的 WHERE 子句在内存中的数据组织结构如图 5-9、5-10、5-11 所示。图 5-10 和 5-11 分别展示了图 5-9 中涉及的两个子查询在内存中的组织形式。可以看到，子查询同样是用 SelectStmt 来表示，并作为父查询的 SelectStmt 结构的一部分存在。

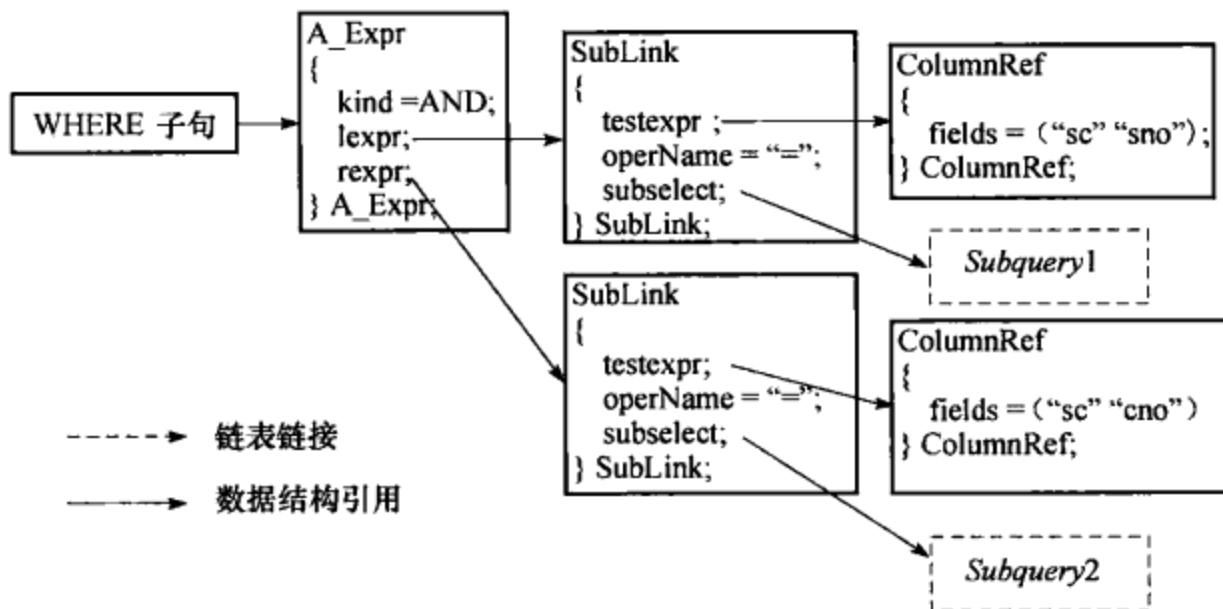


图 5-9 例 5.1 对应的 WHERE 子句的内存组织结构

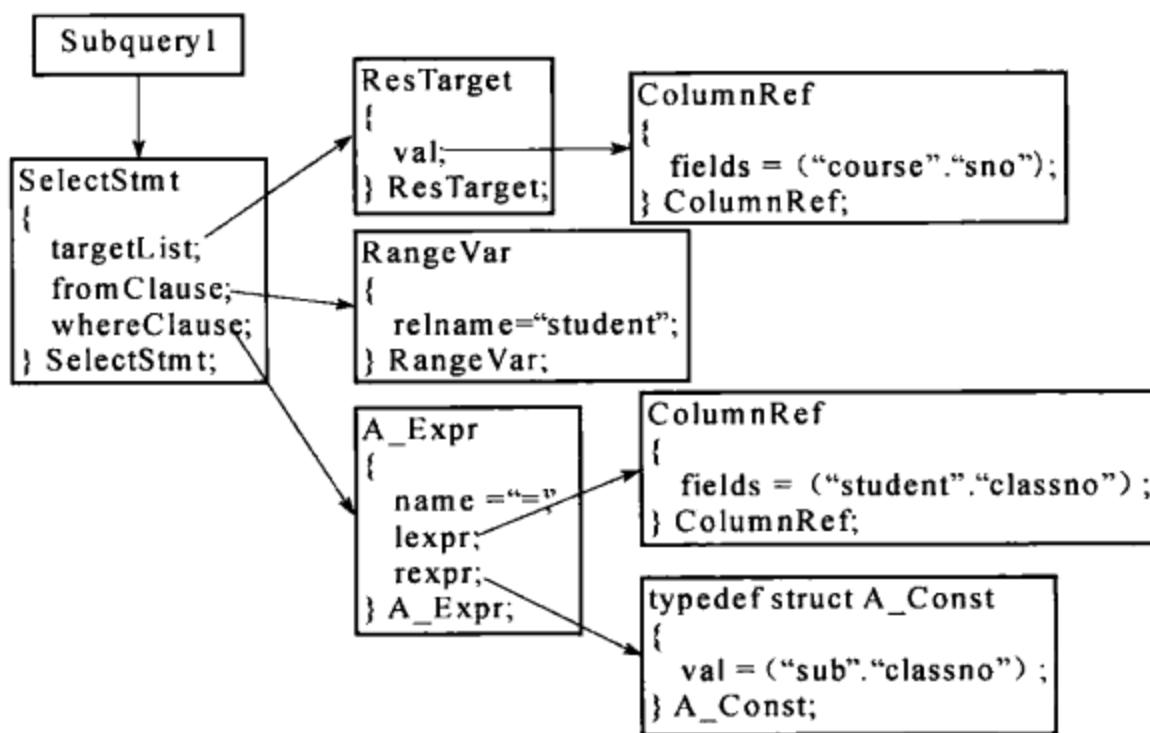


图 5-10 Subquery1 的内存组织结构

5. GROUP BY 子句

GROUP BY 子句的作用是根据所指定的属性进行分组，对应着语法定义中的标识符 group_clause。GROUP BY 子句的语法结构与 WHERE 子句非常相似，在此不再详细讨论。在省去 List 结构体、ListCell 结构体以及一些非关键字段等内容后，例 5.1 中分析树的 GROUP BY 子句在内存中的数据组织结构如图 5-12 所示。

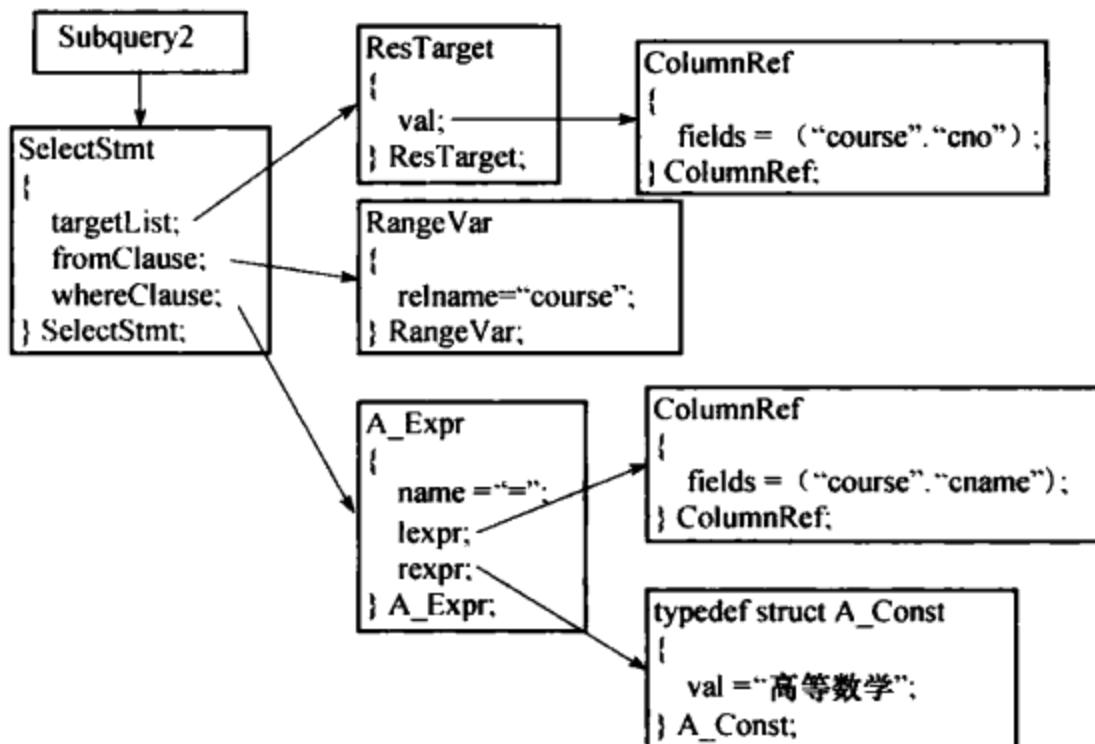


图 5-11 Subquery2 的内存组织结构

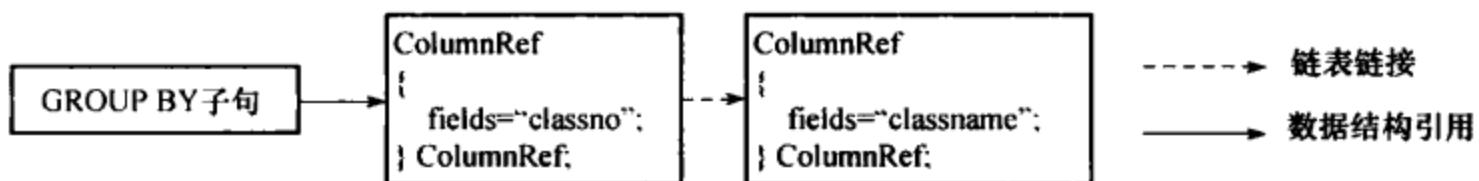


图 5-12 例 5.1 对应的 GROUP BY 子句的内存组织结构

GROUP BY 子句的语法定义

```
7620 group_clause
7621         GROUP_P BY expr_list          { $$ = $3; }
7622         /* EMPTY */                  { $$ = NIL; }
7623         ;
7624         .....
7538 expr_list:  a_expr
7539         {
7540             $$ = list_makel($1);
7541         }
7542         | expr_list','a_expr
7543         {
7544             $$ = lappend($1, $3);
7545         }
7546         ;
```

6. HAVING 子句和 ORDER BY 子句

HAVING 子句的作用是根据指定的条件对 GROUP BY 的分组进行过滤，对应于 `having_clause` 标识符。ORDER BY 子句的作用是根据指定属性对整个查询的结果进行排序，对应于 `sort_clause` 标

识符。

HAVING 子句和 ORDER BY 子句的语法定义

```

7228 having_clause:
7229         HAVING a_expr          ($$ = $2; )
7230         | /* EMPTY */        ($$ = NULL; )
7231         ;
7232         .....
7116 sort_clause:
7117         ORDER BY sortby_list    ($$ = $3; )
7118         ;
7119
7120 sortby_list:
7121         sortby                ($$ = list_make1($1); )
7122         | sortby_list', 'sortby   ($$ = lappend($1,$3); )
7123         ;
7124
7125 sortby: a_expr USING qual_all_Op opt_nulls_order
7126         {
7127             $$ = makeNode(SortBy);
7128             $$ -> node = $1;
7129             $$ -> sortby_dir = SORTBY_USING;
7130             $$ -> sortby_nulls = $4;
7131             $$ -> useOp = $3;
7132             $$ -> location = @ 3;
7133         }
7134         | a_expr opt_asc_desc opt_nulls_order
7135         {
7136             $$ = makeNode(SortBy);
7137             $$ -> node = $1;
7138             $$ -> sortby_dir = $2;
7139             $$ -> sortby_nulls = $3;
7140             $$ -> useOp = NIL;
7141             $$ -> location = -1;      /* no operator */
7142         }
7143         ;

```

在省去 List 结构体、ListCell 结构体以及一些非关键字段等内容后，例 5.1 中分析树的 HAVING 子句和 ORDER BY 子句在内存中的数据组织结构如图 5-13 和图 5-14 所示。

至此，已分析完 SELECT 查询中核心部分的语法定义以及其在内存中的数据组织结构，对于 LIMIT 子句、OFFSET 子句等可做类似分析，这里不再详述。本节没有分析带有集合操作的 SELECT

语句（交并差），因为这些复合语句最终可分解为单个 SELECT 语句来处理。对于 INSERT/DELETE/UPDATE 等也是做类似于 SELECT 的处理，最终也会生成类似于 SelectStmt 的 InsertStmt、DeleteStmt 和 UpdateStmt 等结构。词法和语法分析器会将这些分析树封装成一个 List 结构（raw_parse-tree_list）返回给 exec_simple_query，List 中的每一个 ListCell 包含一个分析树。

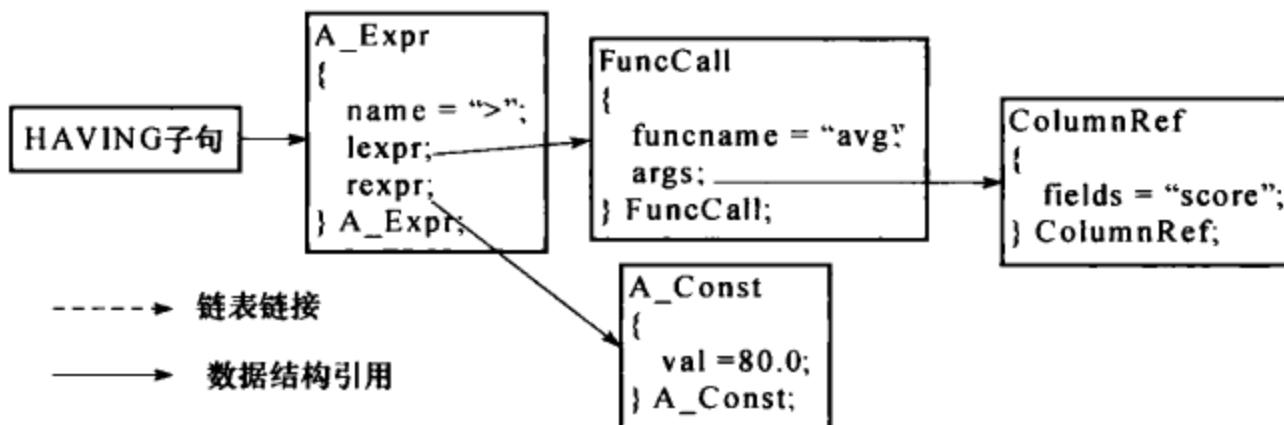


图 5-13 例 5.1 对应的 HAVING 子句的内存组织结构

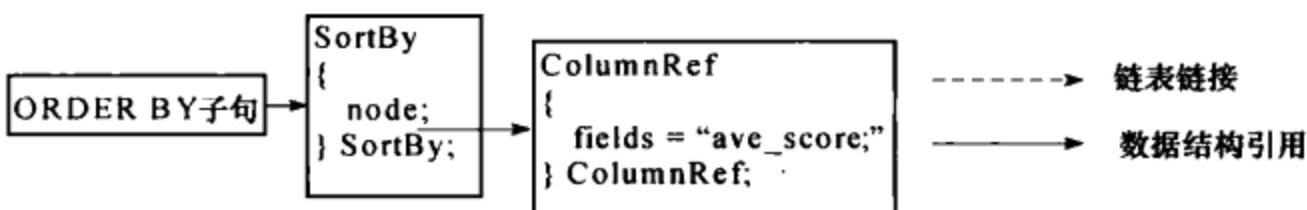


图 5-14 例 5.1 对应的 ORDER BY 子句的内存组织结构

什么情况下 raw_parse-tree_list 中会含有多个分析树？

在某些情况下，用户可能在一个命令字符串中执行多个 SQL 命令，也就是说客户端提交给服务进程的字符串中包含多个 SQL 命令。例如，用户可能输入“CREATE TABLE t (a int); INSERT INTO t VALUES (1); SELECT * FROM t;”这样一个字符串，那么 Postgres 接收到该字符串之后进行词法和语法分析的结果就是三个分析树：CreateStmt、InsertStmt 和 SelectStmt。在返回的 raw_parse-tree_list 中就有三个 ListCell 分别包含上述三个分析树。

5.2.3 语义分析

语义分析阶段会检查命令中是否有不符合语义规定的成分。例如，所使用的表、属性、过程函数等是否存在，聚集函数（如求和函数 SUM、平均函数 AVG 等）是否可以合法使用等。其主要作用在于检查该命令是否可以正确执行。语义分析器会根据分析树中的内容得到更有利于执行的数据，例如，根据表的名字得到其 OID，根据属性名得到其属性号，根据操作符的名字得到其对应的计算函数等。

`exec_simple_query` 在从词法和语法分析模块获取了 `parsetree_list` 之后，会对其中每一棵分析树调用 `pg_analyze_and_rewrite` 进行语义分析和查询重写，而在其中负责语义分析的则是 `analyze.c` 文件中的 `parse_analyze` 函数。`parse_analyze` 会根据分析树生成一个对应的查询树，而查询重写模块则继续对这一查询树进行修改，并且有可能会将这个查询树改写成一个包含多棵查询树的链表。因此，`pg_analyze_and_rewrite` 最终返回给 `exec_simple_query` 的将是一个查询树链表。


```

Node           *p_value_substitute;
bool          p_variableparams;           //为真表示参数的类型还未确定
bool          p_hasAggs;                 //是否有聚集函数
bool          p_hasWindowFuncs;         //是否有窗口函数
bool          p_hasSubLinks;            //是否有子链接
bool          p_is_insert;              //是否为 INSERT 语句
bool          p_is_update;              //是否为 UPDATE 语句
Relation      p_target_relation;        //目标表
RangeTblEntry *p_target_rangetblentry; //目标表在 RangeTable 中对应的项
} ParseState;

```

函数 `parse_analyze` 首先将生成一个 `ParseState` 结构用于记录语义分析的状态，然后通过调用函数 `transformStmt` 来完成语义分析过程。函数 `transformStmt` 会根据不同的查询类型调用相应的函数进行处理。从词法和语法分析中介绍的数据结构可以看到，PostgreSQL 中几乎每一个数据结构的第一个字段都是 `NodeTag` 类型的 `type`。在 PostgreSQL 为了传递参数方便，把很多数据节点的指针都通过指针类型转换成 `Node` 结构的指针，`Node` 结构中只有一个类型为 `NodeTag` 名为 `type` 的字段，通过这样一种方式 PostgreSQL 把大多数要传递的数据结构都包装成了称为“`Node`”（节点）的统一形式，从而通过一套统一的操作函数进行处理。为了能够确定接收到的数据结构到底是哪一种，PostgreSQL 中把 `NodeTag` 设计成了一个枚举类型，每一种数据结构都在其中对应一个值，这些值的名称都由“`T_`”及其数据类型名称组成，例如 `SelectStmt` 数据结构对应的 `NodeTag` 值就是“`T_SelectStmt`”。这样当我们得到一个 `Node` 时，只需要检查其 `type` 字段的值即可确定要处理的是哪种数据结构。

`transformStmt` 函数只有两个参数，一个是 `ParseState`，另一个就是要处理的包装成节点的分析树。通过节点 `type` 字段，`transformStmt` 可以处理七种分析树，相应的处理函数见表 5-9。

这些函数名以“`transform`”开头的函数处理方式大同小异，限于篇幅本书只对其中最为复杂的函数 `transformSelectStmt` 进行分析。`transformSelectStmt` 函数主要是对 SELECT 查询的各个子句分别进行处理，现仍结合例 5.1 来分析其处理过程。

`transform` 函数的主要作用是将分析树转换为查询树（`Query` 结构），查询树是 SQL 语句在 PostgreSQL 中被处理时的内部表现形式，组成语句的每个独立部分（子句）存储在查询树的各个字段中。在 5.3 节中将要讲到的规则也是以查询树的文本形式存储在系统表中。若要将查询树打印到日志文件，可以在配置文件 `postgresql.conf` 中修改参数 `debug_print_parse`、`debug_print_rewritten` 和 `debug_print_plan`，重启数据库之后每次执行命令时都会将查询树以文本的方式打印到系统日志。查询树是以 `Query` 结构存储和组织的，如数据结构 5.6 所示。其中：

表 5-9 分析树的语义分析函数

NodeTag 值	语义分析函数
<code>T_InsertStmt</code>	<code>transformInsertStmt</code>
<code>T_DeleteStmt</code>	<code>transformDeleteStmt</code>
<code>T_UpdateStmt</code>	<code>transformUpdateStmt</code>
<code>T_SelectStmt</code>	<code>transformSelectStmt</code> 或 <code>transformValuesClause</code>
<code>T_DeclareCursorStmt</code>	<code>transformDeclareCursorStmt</code> （游标定义语句）
<code>T_ExplainStmt</code>	<code>transformExplainStmt</code> （ <code>EXPLAIN</code> 语句）
其他	作为 Utility 类型处理，直接在分析树上封装一个 <code>Query</code> 节点返回

数据结构 5.6 表示查询树的 Query 结构

```

typedef struct Query
{
    NodeTag          type;           //节点类型,T_Query
    CmdType          commandType;    //命令类型
    QuerySource      querySource;    //是原始查询还是来自于规则的查询（被规则取代）
    bool             canSetTag;     //查询重写时用到，如果该Query是由原始查询转换而来则此字段为假，如果Query是由查询重写或查询规划时新增加的则此字段为真
    Node             *utilityStmt;   //定义游标或者不可优化的查询语句
    int              resultRelation; //结果关系
    IntoClause       *intoClause;    //处理 SELECT INTO 和 CREATE TABLE AS 的信息
    bool             hasAggs;        //目标属性或 HAVING 子句中是否有聚集函数
    bool             hasWindowFuncs; //目标属性中是否有窗口函数
    bool             hasSubLinks;   //是否有子查询
    bool             hasDistinctOn; //DISTINCT 子句
    bool             hasRecursive;  //公共表表达式中是否允许递归
    List             *cteList;        //WITH 子句，用于公共表表达式
    List             *rtable;         //范围表
    FromExpr         *jointree;      //连接树，描述 FROM 和 WHERE 子句出现的连接
    List             *targetList;    //目标属性
    List             *returningList; //RETURNING 子句
    List             *groupClause;   //GROUP 子句
    Node             *havingQual;   //HAVING 子句
    List             *windowClause; //窗口函数子句
    List             *distinctClause; //DISTINCT 子句
    List             *sortClause;    //ORDER 子句
    Node             *limitOffset;   //OFFSET 子句
    Node             *limitCount;    //LIMIT 子句
    List             *rowMarks;      //行标记链表，用于 FOR UPDATE/FOR SHARE 子句
    Node             *setOperations; //集合操作 (UNION、INTERSECT、EXCEPT)
} Query;

```

1) `commandType`: 查询树对应命令的类型，它是一个枚举类型，说明由哪类命令生成该查询树，包括以下几类：CMD_SELECT、CMD_INSERT、CMD_UPDATE、CMD_DELETE 和 CMD/utility。如果命令类型为 CMD.utility，则查询优化器不会对该查询树进行优化。

2) `rtable`: 范围表，是查询中使用的表的列表。在 SELECT 语句中，在 FROM 子句中出现的表或视图都称为范围表。范围表用编号而不用名字引用，以避免出现重名问题。在生成路径过程中，通过连接操作所创建的新的关系中会记录所引用的范围表的编号。

3) `resultRelation`: 结果关系，是涉及数据修改的范围表（实际使用的是范围表的编号）。这个字段只适用于 INSERT/UPDATE/DELETE 命令，表示这些命令中需要更改的表或视图。SELECT 命令通常没有结果关系，但 SELECT INTO 作为特例处理（因为 SELECT INTO 语句会将查询的结果插

入到另外一个表中), 参考 Query 中的 intoClause 成员。

4) targetList: 目标属性, 用于存放查询结果属性的表达式。目标属性里的每个元素都包含一个表达式。它可以为常量值、某个范围表的一个属性、参数或者由函数调用、常量、变量、操作符等构成的表达式树。分以下四种情况:

①SELECT 语句: 目标属性是位于 SELECT 和 FROM 之间的表达式。

②DELETE 语句: 不需要目标属性, 因为 DELETE 语句不返回任何元组。

③INSERT 语句: 目标属性描述插入到结果关系的元组的属性。这些属性包括表名后指定的要插入值的属性或 INSERT ... SELECT 语句中 SELECT 子句里的表达式。查询重写过程的第一步就是为查询中缺省字段增加目标属性, 剩下的字段 (既无指定值也无缺省值) 将会赋予常量 NULL。

④UPDATE: 目标属性描述被更新的属性, 即 SET 子句中的属性。

5) joontree: 连接树, 查询的连接树显示了 FROM 子句中表的连接情况。对于类似于 SELECT ... FROM a, b, c 的简单查询, 连接树只是 FROM 中表的简单列表, 因为允许以任意顺序连接这些表。如果使用 JOIN 表达式 (尤其是外连接), 必须按照该连接指定的顺序进行连接。此时, 连接树显示了 JOIN 表达式的结构。JOIN 子句上的约束条件 (ON 或 USING 表达式) 作为附加在连接树节点上的条件表达式处理。通常, 还会把顶层 WHERE 表达式作为附加在顶层连接树节点上的条件表达式来处理, 因此, 连接树实际上表示了 SELECT 语句中的 FROM 和 WHERE 子句。

1. transformSelectStmt 函数

transformSelectStmt 函数的作用是根据一个 SelectStmt 结构生成一个查询树, 其主要流程如下:

- 1) 创建一个新的 Query 节点, 设置其 commandType 字段为 CMD_SELECT。
- 2) 调用 transformWithClause 函数处理 WITH 子句。
- 3) 调用 transformFromClause 函数处理 FROM 子句。
- 4) 调用 transformTargetList 函数处理目标属性。
- 5) 调用 transformWhereClause 函数处理 WHERE 子句和 HAVING 子句。
- 6) 调用 transformSortClause 函数处理 ORDER BY 子句。
- 7) 调用 transformGroupClause 函数处理 GROUP BY 子句。
- 8) 调用 transformDistinctClause 函数或 transformDistinctOnClause 函数处理 DISTINCT。
- 9) 调用 transformLimitClause 函数处理 LIMIT 和 OFFSET。
- 10) 处理 INTO 子句。
- 11) 设置 Query 节点的其他标志 (如 hasAggs、hasSubLinks 等字段)。
- 12) 返回 Query 节点。

由此可见, 对于 SelectStmt 的处理过程被分解成对其各个子句部分的处理, 各子句的处理过程相对独立。其他类型的分析树的语法分析过程与此大同小异。这里我们着重对 FROM 子句、目标属性、WHERE 子句的处理进行分析。

2. FROM 子句的语义分析处理

transformFromClause 负责处理 FROM 子句 (SelectStmt 的 fromClause 字段) 的语法分析, 并生成范围表 (Query 节点的 rtable 字段)。该函数以 ParseState 结构和 SelectStmt 结构的 fromClause 字段为参数。fromClause 是一个 List, 其中每一个 ListCell 包含一个 Node, 这些 Node 是由语法分析模块根据 FROM 子句中出现的表、视图、子查询、函数或者连接表达式生成的各类数据结构包装而来。

`transformFromClause` 函数主要调用函数 `transformFromClauseItem` 对 `fromClause` 中的每一个 `Node` 进行处理。处理完每一个 `Node` 之后都会检查该 `Node` 表示的范围表是否已存在（即重名检查，SELECT 语句中同一个表不能定义两次，使用别名除外）。`transformFromClause` 主要用于对结构体各字段进行赋值，从而存储语义分析的相关信息。

函数 `transformFromClauseItem` 的主要工作是根据 `fromClause` 的 `Node` 产生一个或者多个 `RangeTableEntry` 结构（数据结构 5.7，简称 RTE）加入到 `ParseState` 的 `p_rtable` 字段指向的链表中，每一个 `RangeTableEntry` 结构就表示一个范围表，最终生成的查询树的 `rtable` 字段也将指向该链表，所以该函数实际完成了范围表链表的构造。

`transformFromClauseItem` 函数的主要流程如下：

1) 如果要处理的 `Node` 是一个 `RangeVar` 结构（表示一个普通表）：

① 检查该表是不是一个 CTE（Common Table Expression，公共表表达式），如果是则调用 `transformCTEReference` 将它作为一个 RTE 加入到 `ParseState` 的 `p_rtable` 中，当然这个 RTE 的类型会置为 `RTE_CTE`。

② 如果不是一个 CTE，则调用 `transformTableEntry` 将表作为一个 `RTE_RELATION` 类型的 RTE 加入到 `p_rtable` 中。

③ 利用参数 `relnamespace` 将生成的 RTE 传出，并创建一个引用该 RTE 的 `RangeTblRef` 结构返回。这一步骤对子查询、JOIN 表达式等类型的 `Node` 同样需要。

数据结构 5.7 RangeTblEntry

```
typedef struct RangeTblEntry
{
    NodeTag          type;
    RTEKind          rtekind;           //RTE 的类型, RTE_RELATION 表示普通表, RTE_SUBQUERY 表示子
                                         //查询, RTE_JOIN 表示连接, RTE_SPECIAL 表示 NEW 或者 OLD 两个
                                         //虚表, RTE_FUNCTION 表示函数, RTE_VALUES 表示 VALUES 表达式,
                                         //RTE_CTE 表示公共表表达式
    Oid              relid;            //表的 OID, 不是 RTE_RELATION 类型时 relid 为 0
    Query            *subquery;         //子查询的查询树, RTE_SUBQUERY 类型有效
    JoinType         jointype;         //连接类型, 后续字段对 RTE_JOIN 类型有效
    List             *joinaliasvars;    //连接结果中属性的别名
    Node             *funcexpr;         //函数调用的表达式树, 后续字段对 RTE_FUNCTION 类型有效
    List             *funccoltypes;      //函数返回的记录中属性类型的 OID 列表
    List             *funccoltypmods;    //函数返回的记录中属性类型的 typmods 列表
    List             *values_lists;     //VALUES 表达式列表, RTE_VALUES 类型有效
    char             *ctename;          //公共表表达式名称, 后续字段对 RTE_CTE 类型有效
    Index            ctelevelsup;       //公共表表达式的级别
    bool             self_reference;    //是否是递归地自我引用
    List             *ctecoltypes;       //CTE 返回的记录中属性类型的 OID 列表
    List             *ctecoltypmods;     //CTE 返回的记录中属性类型的 typmods 列表
```

```

Alias          *alias;           // 用户定义的别名, 后续字段对于所有 RTE 类型都有效
Alias          *eref;            // 扩展的引用名称
bool           inh;              // 是否需要继承
bool           inFromCl;        // 是否出现在 FROM 子句中
AclMode        requiredPerms;   // 该 RTE 所需的访问权限, 见第 8 章
Oid            checkAsUser;    // 如果是有效值, 则以该 OID 对应的用户身份进行权限检查
Bitmapset      *selectedCols;   // 需要 SELECT 权限的属性集合
Bitmapset      *modifiedCols;   // 需要 INSERT/UPDATE 权限的属性集合
} RangeTblEntry;

```

2) 如果要处理的 Node 是一个 RangeSubselect 结构 (表示一个子查询), 则调用 transformRangeSubselect 将子查询作为一个 RTE_SUBQUERY 类型的 RTE 加入到 p_rtable 中, 当然这里还会递归调用 transformSelectStmt 生成子查询的查询树用于填充该 RTE 的 subquery 字段。

3) 如果要处理的 Node 是一个 RangeFunction 结构 (表示一个函数), 则调用 transformRangeFunction 生成 RTE_FUNCTION 类型的 RTE 加入到 p_rtable 中, 除填充该 RTE 中与函数相关的字段之外, 还会根据该函数是聚集函数或者窗口函数来设置 ParseState 中的相应标志字段。

4) 如果要处理的 Node 是一个 JoinExpr 结构 (表示一个连接表达式, 见数据结构 5.8), 由于 JoinExpr 有左右子树且有不同连接类型, 因此这种情况下的处理要复杂一些。

①对 JoinExpr 的左右子树分别递归调用 transformFromClauseItem 进行语义分析。最简单的情况下, 左右子树分别由一个表构成, 那么这里的递归调用会把这两个表加入到 p_rtable 中。

②进行局部的重名检查, 由于 JoinExpr 结构中可能也涉及多个表, 因此这里需要进行重名检查。注意, 这里的局部重名检查和在 transformFromClause 函数中的重名检查一样, 只有两个 RTE 的表名称及其别名都相同才被认为是重名。检查完成后将左右子树对应的 RTE 加入到参数 rel_namespace 指向的链表中以便传给 transformFromClause 函数进行重名检查。

③如果 JoinExpr 中的 isNatrual 为真 (自然连接), 由于自然连接是以两个表的公共属性为连接属性, 因此自然连接的情况下不需要通过 USING 关键字来指定连接属性, 语义分析模块需要为自然连接找到连接属性 (即公共属性) 并放在 JoinExpr 的 using 字段中 (这种情况下 using 字段原本是为空的)。

④如果 using 字段不为空, 则需要处理由 USING 给出的连接条件。由于 PostgreSQL 的语法规规定, 在连接中 USING 和 ON 是互斥的, 因此 using 字段不为空时, 只需要考虑形如 “USING (a, b, …)” 的情况, 也就是说 USING 后面的属性都是左右两个 RTE 同时具有的。这里将从 using 字段中逐一取出属性名, 然后检查它是否在左右两个 RTE 的属性中存在, 只要在某一方不存在则会报错并终止整个 SQL 命令的执行。在检查过程中将分别为左右 RTE 生成参加连接的属性的列表, 完成这一检查之后, 根据这些属性列表调用函数 transformJoinUsingClause 生成多个 “=” 表达式, 然后将它们作为一个 List 放入 JoinExpr 的 quals 字段中。

⑤如果 quals 字段不为空, 则需要处理由 ON 给出的连接条件。将调用 transformJoinOnClause 来处理 quals 中的内容, 实际上该函数将会调用 transformWhereClause 来完成它的工作。

⑥最后调用 addRangeTableEntryForJoin 从完成的 JoinExpr 生成一个表达该连接的 RTE 加入 p_rtable 中并返回 JoinExpr。

数据结构 5.8 JoinExpr

```

typedef struct JoinExpr
{
    NodeTag      type;
    JoinType     jointype;           //连接类型
    bool         isNatural;         //是否自然连接
    Node         *larg;             //左子树
    Node         *rarg;             //右子树
    List         *using;            //USING 子句
    Node         *quals;             //连接条件
    Alias        *alias;             //别名
    int          rtindex;           //连接表达式在范围表中的索引号,如没有则为0
}RangeTblEntry;

```

在每次调用 `transformFromClauseItem` 处理完一个节点之后, `transformFromClause` 还要调用 `checkNameSpaceConflicts` 对刚处理完的节点中产生的 RTE 进行重名检查(全局重名检查), 其原理就是把 `transformFromClauseItem` 的参数 `relnamespace` 传回的 RTE 逐个和已经通过检查的 RTE 进行名字和别名的比对, 然后将名字和别名没有冲突的 RTE 存放在 `ParseState` 的 `p_relnamespace` 字段所指向的列表中。当然, `checkNameSpaceConflicts` 还会把通过该检查的 RTE 加入到 `p_relnamespace` 中去。最后 `transformFromClauseItem` 的返回值会被以 `Node` 的形式加入到 `ParseState` 的 `p_joinlist` 中。

与词法、语法分析部分的介绍一样, 图 5-16 也给出了例 5.1 中查询树的范围表在内存中的数据组织结构。图中 Subquery 是 FROM 子句中的子查询, 对于该子查询会再创建一个 Query 结构体来组织, 在此不再详述。

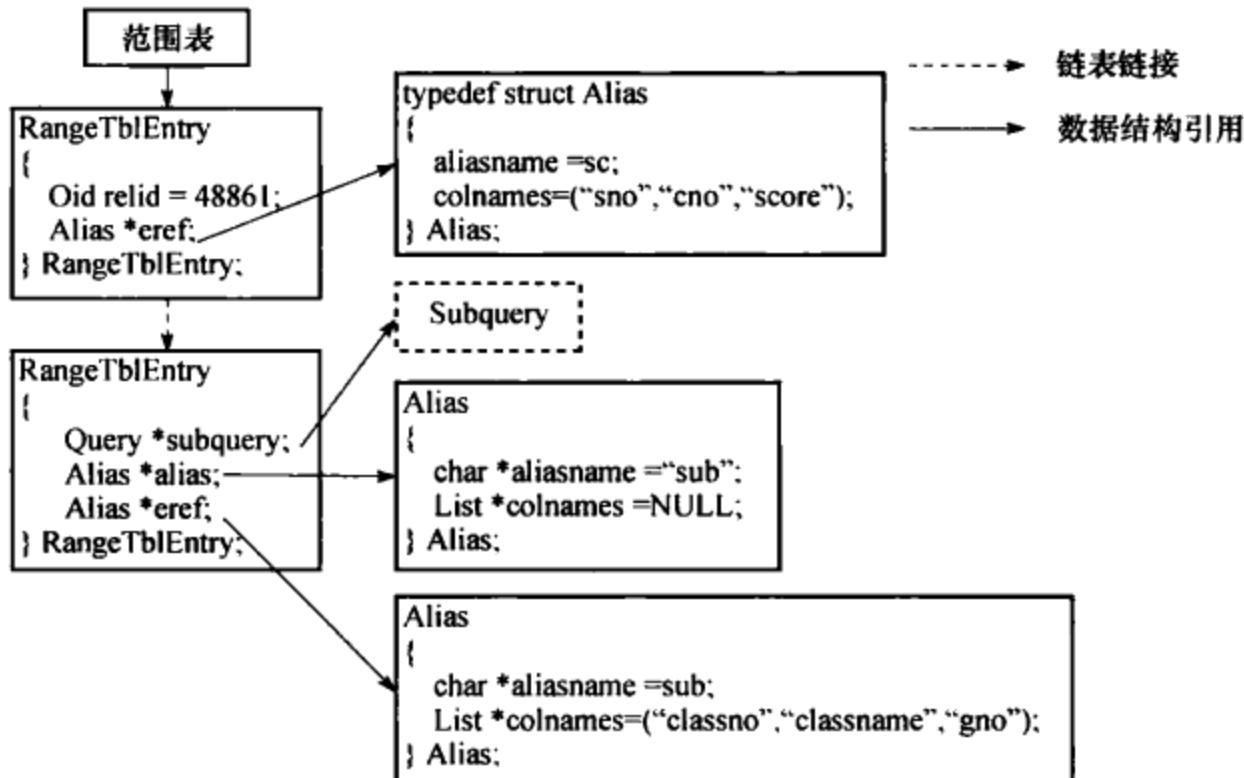


图 5-16 例 5.1 对应查询树的范围表数据组织结构

3. 目标属性的语义分析处理

处理目标属性的人口函数是 transformTargetList 函数，它通过调用函数 transformTargetEntry 来处理分析树的目标属性中的每一项。对于目标属性中的每个项，都会用函数 makeTargetEntry 创建结构体 TargetEntry（数据结构 5.9）来存储和组织。因此，该函数的作用就是将一个 ResTarget 结构体的链表转换成一个 TargetEntry 结构体的链表，每一个 TargetEntry 表示查询树中的一个目标属性。

数据结构 5.9 TargetEntry

```

typedef struct TargetEntry
{
    Expr           xpr;          //表达式头部,严格来说 TargetEntry 并不是一个表达式,因为它不
                                //能被表达式计算函数处理,但 PostgreSQL 仍将它视为一种表达式,
                                //因为在很多地方把目标属性作为一个表达式树进行处理会很方便

    Expr           *expr;        //目标属性中需要计算的表达式

    AttrNumber     resno;        //属性编号,在 SELECT 的目标属性中 resno 对应了属性出现的位置

    char           *resname;     //属性名

    Index          ressortgroupref; //被 ORDER BY 和 GROUP BY 子句引用时使用,0 值表示未被 ORDER BY
                                //和 GROUP BY 子句引用,正值表示被引用,但正值的具体值没有特别
                                //的意义

    Oid            resorigtbl;   //属性所属的表(源表)的 OID

    AttrNumber     resorigcol;   //属性在源表中的属性号

    bool           resjunk;      //表示该属性是否是一个 junk 属性,如果为真表示这个属性是一个
                                //工作属性,在输出结果时应该去除。比如一个仅用于排序的属性其
                                //resjunk 就为真

}TargetEntry;

```

transformTargetList 的参数有两个：ParseState 和分析树的 targetList 字段指向的链表（targetlist 参数）。transformTargetList 的工作流程如下：

- 1) 准备一个空的链表 p_target。
- 2) 对于 targetlist 参数中的每一个节点。
 - ①取出节点中的 ResTarget 结构。
 - ②如果 ResTarget 的 val 字段中是一个“*”号，则调用 ExpandColumnRefStar 或 ExpandIndirectionStar 将其展开，即把星号变成其指代的一系列属性，然后为每一个属性创建 TargetEntry 并加入到 p_target 中。
 - ③如果 ResTarget 的 val 字段中不是一个“*”号，都调用 transformTargetEntry 函数为其生成 TargetEntry 并加入到 p_target 中。transformTargetEntry 函数会调用 transformExpr 函数为 val 字段中的数据类型生成表达式结构。例如，如果是 ColumnRef 结构（表示目标属性就是一个属性），则会为其生成一个 Var 结构（数据结构 5.10）来存储该属性及其源表的相关信息。

数据结构 5.10 Var

```

typedef struct Var
{
    Expr           xpr;
    Index          varno;           //变量(属性)所在的表在范围表中的编号
    AttrNumber     varattno;        //属性编号
    Oid            vartype;         //该属性数据类型在 pg_type 中的 OID
    int32          vartypmod;       //该属性对应的 pg_attribute 元组的 typmod 值
    Index          varlevelsup;     //用于子查询中引用外层表的变量,表示子查询的层数
    Index          varnoold;        //varno 的原始值,用于调试信息
    AttrNumber     varoattno;       //varattno 的原始值,用于调试
    int            location;        //符号出现的位置
}Var;

```

在处理完所有 ResTarget 结构之后, transformTargetList 将返回 p_target, 而 transformSelectStmt 将直接把 transformTargetList 的返回值作为目标属性赋值给查询树的 targetList 字段。例 5.1 中查询树的目标属性在内存中的数据组织结构如图 5-17 所示。

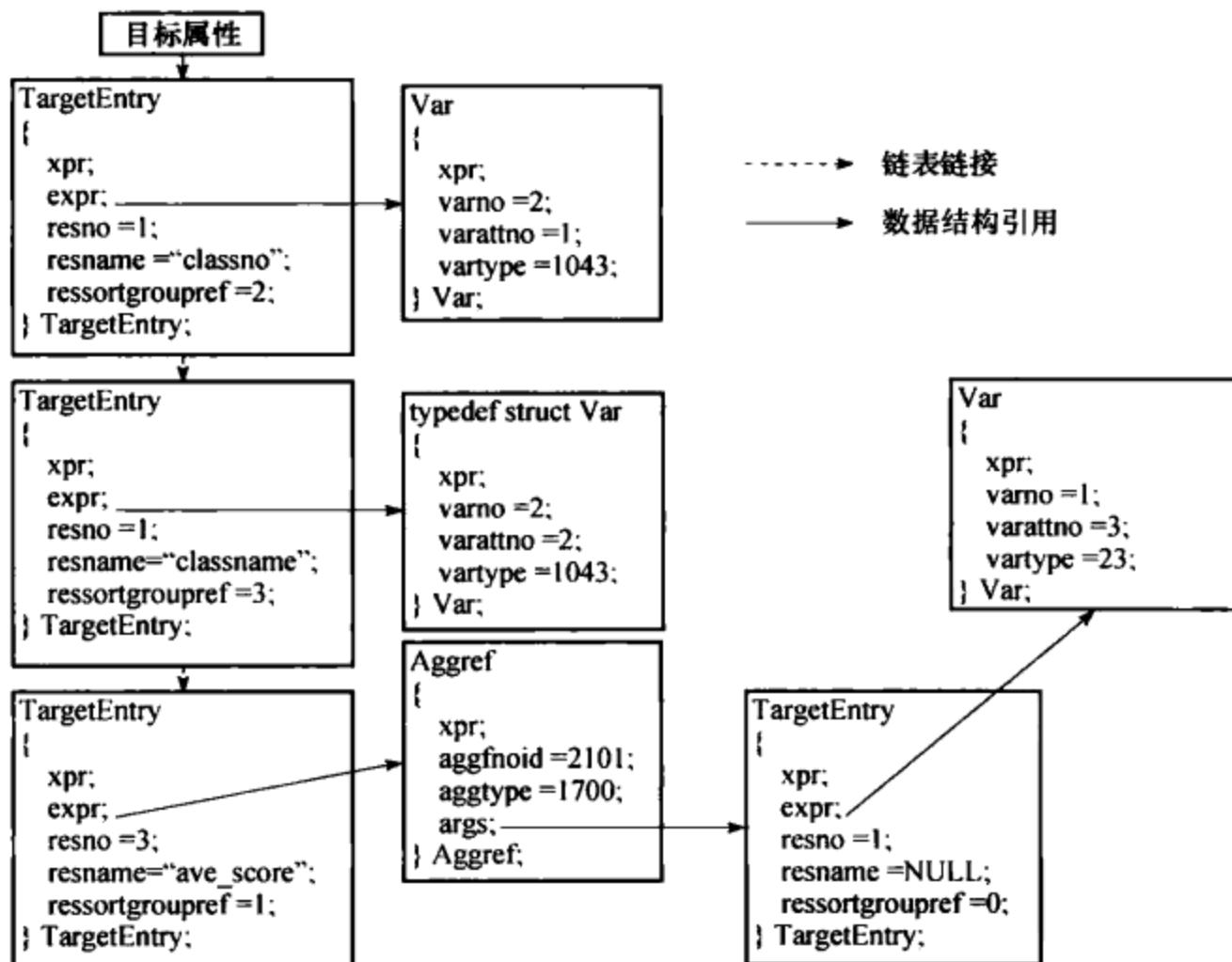


图 5-17 例 5.1 对应查询树的范围表数据组织结构

4. WHERE 子句的语义分析处理

处理 WHERE 子句的人口函数为 `transformWhereClause`。该函数调用函数 `transformExpr` 将分析树

的 whereClause 字段表示的 WHERE 子句转换为一棵表达式树，然后将 ParseState 的 p_joinlist 字段所指向的链表以及从 WHERE 子句得到的表达式树包装成一个 FromExpr 结构存入查询树的 joontree，用以表示连接树。例 5.1 的查询树中连接树的数据组织结构如图 5-18 所示。

至此，我们已经介绍了 SELECT 语句中三个核心子句 SELECT、FROM、WHERE 的语义分析过程，这些知识已经足以让我们了解一个最简单形式的 SELECT 语句会被转换成一个何种结构的查询树。对于 SELECT 语句中的其他部分，例如 DISTINCT、WITH、GROUP BY 等，其语义分析的方式和 SELECT 等基本相同，故此不再详述。

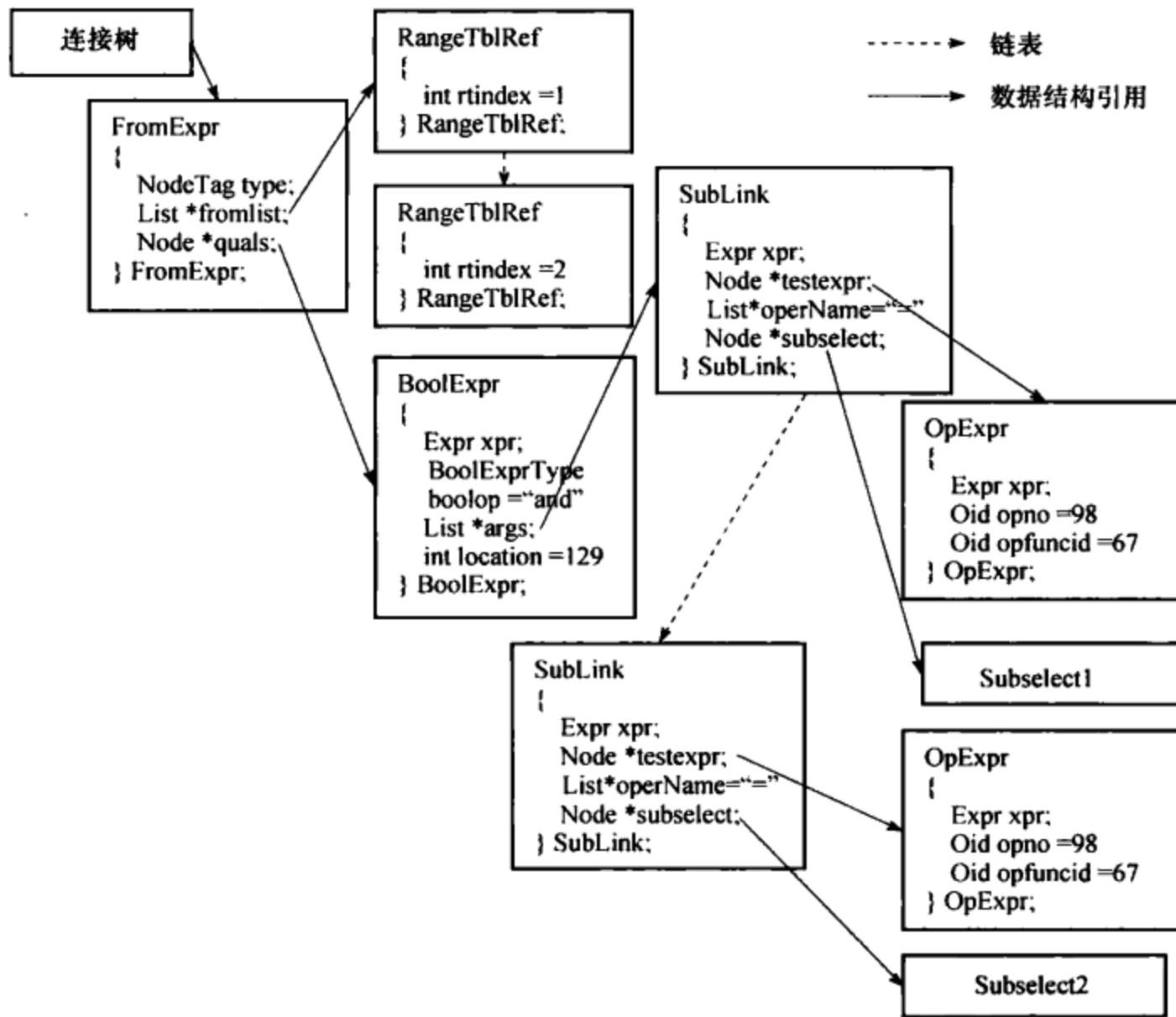


图 5-18 例 5.1 对应查询树的连接树数据组织结构

5.3 查询重写

在完成语义分析步骤得到查询树之后，会立刻对该查询树进行查询重写处理。查询重写的人口函数是 pg_rewrite_query，它在 pg_analyze 之后被 pg_analyze_and_rewrite 调用，且 pg_rewrite_query 的参数就是 pg_analyze 的返回值——查询树。查询重写模块使用规则系统判断来进行查询树的重写，如果查询树中某个目标被定义了转换规则，则该转换规则会被用来重写查询树。查询重写的源代码在 src\Backend\rewrite 文件夹中。

5.3.1 规则系统

查询重写的核心就是规则系统，而规则系统则由一系列的规则组成。系统表 pg_rewrite 存储重写规则，具体内容如表 5-10 所示。pg_rewrite 中的每一个元组代表一条规则。

表 5-10 系统表 pg_rewrite

名字	类型	引用	描述
rulename	name		规则名称
ev_class	oid	pg_class. oid	使用该规则的表名称
ev_attr	int2		规则适用的属性
ev_type	char		规则适用的命令类型：1 = SELECT, 2 = UPDATE, 3 = INSERT, 4 = DELETE
is_instead	bool		如果是 INSTEAD 规则，则为真
ev_qual	text		规则的条件表达式 (WHERE 子句)
ev_action	text		规则动作的查询树 (DO 子句)

从系统表 pg_rewrite 的内容及其含义我们大致可以知道某条规则的工作原理：对于每一条规则（一个 pg_rewrite 元组），该元组的 ev_class 属性表示该规则适用的表名，如果在该表的指定属性（由 ev_attr 属性记录）上执行特定的命令（由 ev_type 属性记录）且满足了规则的条件（由 ev_qual 属性记录）时，用规则的动作（由 ev_action 属性记录）替换原始命令的动作或者将规则的动作附加在原始命令之前（或者之后）。

根据系统表 pg_rewrite 的不同属性，规则可以按两种方式分类：

- 按照规则适用的命令类型分类，可以分成 SELECT、UPDATE、INSERT 和 DELETE 四种。
- 按照规则执行动作的方式分类，可以分成 INSTEAD（替代）规则和 ALSO 规则。

1. SELECT/INSERT/UPDATE/DELETE 规则

SELECT/INSERT/UPDATE/DELETE 四种规则通过其 pg_rewrite 元组的 ev_type 属性值来区分。四种规则有不同的特点：

1) SELECT 规则中只能有一个动作，而且只能是不带条件的 INSTEAD 规则。SELECT 规则的执行效果类似于视图。使用方法可以参考“视图和规则系统”部分的举例。

2) INSERT/UPDATE/DELETE 规则具有以下特性：

- 可以没有动作。
- 可以有多个动作。
- 可以是 INSTEAD 型规则（替代规则）或 ALSO 型规则（缺省）。
- 可以使用伪关系 NEW 和 OLD。
- 可以使用规则条件。
- 不是修改查询树，而是创建零个或多个新查询树并且可能把原始的查询树丢弃。

2. INSTEAD 规则和 ALSO 规则

INSTEAD 规则和 ALSO 规则通过规则的 pg_rewrite 元组的 is_instead 属性值来区分，为真表示 INSTEAD 规则，为假则表示 ALSO 规则。

INSTEAD 规则执行动作的方式很简单，就是用规则中定义的动作替代原始的查询树中的对规

则所在表的引用。

而 ALSO 型规则中原始查询和规则动作都会被执行，但执行顺序根据不同的命令也有所不同：

- INSERT 规则：原始查询在规则动作执行之前完成，这样可以保证规则动作能引用插入的行。
- UPDATE 和 DELETE 规则：原始查询在规则动作之后完成，这样能保证规则动作可以引用将要更新或删除的元组；否则，那些要访问旧版本元组的规则动作就无法完成。

3. 视图和规则系统

PostgreSQL 中的视图通过规则系统实现。在创建视图时，系统会自动按照其定义生成相应的规则，当查询涉及该视图时，查询重写模块都会用对应的规则对该查询进行重写，将对视图的查询改写为对基本表的查询。在生成视图的规则时，规则动作是视图创建命令中 SELECT 语句的拷贝，并且该规则是无条件的 INSTEAD 规则。

下面我们通过一个实际的例子来展示如何通过规则系统来实现视图。在此之前，先介绍 PostgreSQL 中定义规则的命令，其语法为：

```
CREATE[ OR REPLACE] RULE name AS ON event
TO table[ WHERE condition]
DO[ ALSO |INSTEAD ] {NOTHING |command |(command;command... )}
```

其中：

- name：规则名。
- event：SELECT、INSERT、UPDATE、DELETE 命令之一。
- table：规则作用的表或视图（可以有模式修饰）。
- condition：规则条件，满足条件则执行规则动作。只能引用 OLD 和 NEW^①不能引用其他的表，不能有聚集函数。
- INSTEAD：表示用 DO 子句中的命令替换 table。
- ALSO：DO 子句中的命令和 event 所对应的命令都要执行，是缺省值。
- command：规则动作，有效的语句是 SELECT、INSERT、UPDATE、DELETE 和 NOTIFY。

使用规则系统实现视图的具体步骤如下：

(1) 创建需要用到的函数和表

```
-- min 函数用于返回两个整数值中较小值
CREATE FUNCTION min(integer,integer)RETURNS integer AS $$
    SELECT CASE WHEN $1 < $2 THEN $1 ELSE $2 END
$$ LANGUAGE SQL STRICT;

-- 鞋子数据表
CREATE TABLE shoe_data
{
```

① 在 condition 和 command 里，特殊的表名字 NEW 和 OLD 用于指向被引用表里的数值。对于 ON INSERT 规则，NEW 指向要被插入的新行；对于 ON UPDATE 规则，NEW 指向被更新之后的新行，OLD 指向现存的要被更新的行；对于 ON DELETE 规则，OLD 指向要被删除的行。

```

shoename text,      -- 主键
sh_avail integer,   -- 鞋的可用对数
slcolor text,        -- 首选的鞋带颜色
slminlen real,       -- 鞋带最短长度
slmaxlen real,       -- 鞋带的最长长度
slunit  text         -- 长度单位
);

-- 度量单位换算表
CREATE TABLE unit
(
    un_name text,      -- 主键
    un_fact real,      -- 转换成厘米的系数
);

```

(2) 创建需要用到的视图

```

-- 鞋带的视图(将鞋带的长度单位换算成厘米)
CREATE VIEW shoelaceAS
SELECT      s.sl_name,
            s.sl_avail,
            s.sl_color,
            s.sl_len,
            s.sl_unit,
            s.sl_len * u.un_fact AS sl_len_cm
FROM        shoelace_data s,unit u
WHERE       s.sl_unit = u.un_name;

-- 鞋子的视图(将允许搭配的最小、最大鞋带长度的单位换算成厘米)
CREATE VIEW shoeAS
SELECT      sh.shoename,
            sh.sh_avail,
            sh.slcolor,
            sh.slminlen, sh.slminlen * un.un_fact AS slminlen_cm,
            sh.slmaxlen, sh.slmaxlen * un.un_fact AS slmaxlen_cm,
            sh.slunit
FROM        shoe_data sh,unit un
WHERE       sh.slunit = un.un_name;

-- 鞋子和鞋带可配套的数据表(颜色相同、鞋带长度符合鞋子的要求)
CREATE VIEW shoe_ready AS
SELECT      rsh.shoename,
            rsh.sh_avail,
            rsl.sl_name,
            rsl.sl_avail,
            min(rsh.sh_avail,rsl.sl_avail) AS total_avail
FROM        shoe rsh,shoelace rsl

```

```

WHERE      rsl.sl_color = rsh.slcolor
AND        rsl.sl_len_cm >= rsh.slminlen_cm
AND        rsl.sl_len_cm <= rsh.slmaxlen_cm;

```

第一个 CREATE VIEW 命令将创建 shoelace 视图，系统会自动为该视图生成相应的规则，并在 pg_rewrite 表里将该规则作为一个元组插入。如果某查询的范围表里引用了视图 shoelace，系统会自动从 pg_rewrite 中读取 shoelace 上的规则对该查询进行重写，即用规则动作替换查询中对视图 shoelace 的引用。

(3) 系统为视图 shoelace 自动生成的规则的动作部分

```

SELECT s.sl_name,s.sl_avail,
       s.sl_color,s.sl_len,s.sl_unit,
       s.sl_len * u.un_fact AS sl_len_cm
FROM shoelace *OLD*,shoelace *NEW*,
     shoelace_data s,unit u
WHERE s.sl_unit = u.un_name;

```

(4) 对视图 shoelace 上的查询进行查询重写

1) 原始查询语句：

```
SELECT * FROM shoelace;
```

2) 被分析器解释之后的原始查询树（逻辑上的表示）：

```

SELECT shoelace.sl_name,shoelace.sl_avail,
       shoelace.sl_color,shoelace.sl_len,
       shoelace.sl_unit,shoelace.sl_len_cm
FROM   shoelace shoelace;

```

3) 用规则重写之后的查询树（逻辑上的表示）：

```

SELECT shoelace.sl_name,shoelace.sl_avail,
       shoelace.sl_color,shoelace.sl_len,
       shoelace.sl_unit,shoelace.sl_len_cm
FROM(SELECT s.sl_name,s.sl_avail,
       s.sl_color,s.sl_len,s.sl_unit,
       s.sl_len * u.un_fact AS sl_len_cm
    FROM shoelace_data s,unit u
   WHERE s.sl_unit = u.un_name)shoelace;

```

从上面的例子可以看到，PostgreSQL 中的视图其实就是用规则来实现的。创建一个视图时，会按照视图的定义自动创建一个规则。在对视图进行查询时，则用相应的规则将对视图的查询改写成对基表的查询。

4. ALSO 型规则的查询重写

下面我们还是以上一节的例子说明 ALSO 型规则的查询重写过程。假设需要跟踪 shoelace 表中的 sh_avail 字段的变化轨迹，我们创建一个日志表（shoelace_log）和一条规则（log_shoelace），这条规则的作用是每次在用 UPDATE 命令更新 shoelace 表时都往日志表里写一条记录。例子中需要的

日志表和规则定义如下：

```
-- 日志表
CREATE TABLE shoelace_log
(
    sl_name text,          -- 鞋带
    sh_avail integer,      -- 新的可用数量
    log_who text,          -- 插入者
    log_when timestamp     -- 时间
);

-- 有规则条件 ALSO 型的 UPDATE 规则
CREATE RULE log_shoelace AS ON UPDATE
TO shoelace WHERE NEW.sh_avail <> OLD.sh_avail
DO INSERT INTO shoelace_log
VALUES (NEW.sl_name, NEW.sh_avail, current_user, current_timestamp);
```

假设当前在 shoelace 表中 sl_name 为 “sl7” 的元组 sh_avail 属性的值为 10。对该表执行一条更新语句：

```
UPDATE shoelace
SET sh_avail = 6
WHERE sl_name = 'sl7';
```

显然，这条语句满足规则的条件，因为 NEW. sh_avail 是 6，而 OLD. sh_avail 是 10，两者不相等。因此该规则动作和原始查询动作都将被执行。下面是对该 UPDATE 语句的查询重写过程：

1) 重写第一步：原始查询的范围表集成到规则动作查询里；

```
INSERT INTO shoelace_log
VALUES      (*NEW*.sl_name, *NEW*.sh_avail, current_user, current_timestamp)
FROM        shoelace *NEW*, shoelace *OLD*,
            shoelace shoelace;
```

2) 重写第二步：添加规则条件（将结果集限制为 sh_avail 改变的元组）；

```
INSERT INTO shoelace_log
VALUES      (*NEW*.sl_name, *NEW*.sh_avail, current_user, current_timestamp)
FROM        shoelace *NEW*, shoelace *OLD*,
            shoelace shoelace
WHERE       *NEW*.sh_avail <> *OLD*.sh_avail;
```

3) 重写第三步：添加原始查询的查询条件；

```
INSERT INTO shoelace_log
VALUES      (*NEW*.sl_name, *NEW*.sh_avail, current_user, current_timestamp)
FROM        shoelace *NEW*, shoelace *OLD*,
            shoelace shoelace
WHERE       *NEW*.sh_avail <> *OLD*.sh_avail
AND        shoelace.sl_name = 'sl7';
```

4) 重写第四步：把将 NEW 引用替换为来自原始查询树的目标属性或来自结果关系的相匹配的变量引用；

```
INSERT INTO shoelace_log
VALUES      (shoelace.sl_name,6,current_user,current_timestamp)
FROM        shoelace *OLD*,shoelace shoelace
WHERE       6 <> *OLD*.sh_avail
AND         shoelace.sl_name = 's17';
```

5) 重写第五步：用结果关系引用把 OLD 引用替换掉

```
INSERT INTO shoelace_log
VALUES      (shoelace.sl_name,6,current_user,current_timestamp)
FROM        shoelace shoelace
WHERE       6 <> shoelace.sh_avail
AND         shoelace.sl_name = 's17';
```

6) 重写第六步：由于 ALSO 型规则要输出原始查询，因此，从规则系统输出的是两个查询的列表

```
INSERT INTO shoelace_log
VALUES      (shoelace_data.sl_name,6,current_user,current_timestamp)
FROM        shoelace
WHERE       6 <> shoelace.sh_avail
AND         shoelace.sl_name = 's17';

UPDATE      shoelace
SET         sh_avail = 6
WHERE       sl_name = 's17';
```

由于是 ALSO 型的 UPDATE 规则，所以原始查询必须在规则动作之后执行；否则，规则动作将找不到满足规则条件的元组。而对于 INSTEAD 型规则就没有第六步，因为只执行规则动作，不执行原始查询。

5. 规则系统与触发器的区别

从以上的例子可以看出，用规则（特别是 ALSO 规则）重写查询和触发器的功效非常相似，都可以在某种命令和条件下被激活，并且可以执行原始查询之外的动作。但是，两者从本质上还是有区别的。触发器对涉及的每个元组都要执行一次，而规则是对整个查询树进行修改或生成额外的查询。因此，如果在一个语句中涉及多个元组，一个规则通常比触发器的效率高。同时，触发器从概念上要比规则简单，触发器实现的功能可以用规则实现。但是目前某些约束不能用规则实现，特别是外键。

5.3.2 查询重写的处理操作

查询重写部分的处理操作主要包括定义规则、删除规则以及利用规则进行查询重写。下面分别对这些操作进行介绍。

1. 定义重写规则

在使用规则系统之前首先需要定义规则，规则的定义通过 CREATE RULE 命令来完成。而定义重写规则的操作主要由函数 DefineRule 实现，“CREATE RULE”命令被词法和语法分析模块处理之后，相关信息被存储在一个 RuleStmt 结构（数据结构 5.11）中，最后查询执行模块会把该结构交给 DefineRule 来完成规则的创建。DefineRule 的流程如下：

1) 首先调用 transformRuleStmt 对 RuleStmt 进行处理：

① 把要定义规则的表分别以别名“*OLD*”和“*NEW*”封装成 RTE 加入到 ParseState 的 p_rtable（范围表）中。

② 根据规则适用的命令不同将 OLD 或者 NEW 对应的 RTE 加入到 ParseState 的 p_joinlist 中：对于 SELECT 和 DELETE 命令只加入 OLD，对于 INSERT 命令只加入 NEW，对于 UPDATE 命令要将 OLD 和 NEW 都加入。

③ 调用 transformWhereClause 将 RuleStmt 中的 WHERE 子句转换成一个表达式树，并由参数 whereClause 将表达式树的指针传出。

④ 做一系列检查，包括 WHERE 子句中不能引用其他表、不能使用聚集和窗口函数。

⑤ 对于没有动作的规则（RuleStmt 的 actions 字段为 NIL），创建一个命令类型为 CMD_NOTHING 且没有连接树的 Query，并将该 Query 包装成一个 List，通过参数 actions 将 List 的指针传出。

⑥ 将 RuleStmt 的 actions 字段中的每一个节点通过 transformStmt 函数转换成一个 Query 结构，每一个 Query 结构的范围表中都要增加 OLD 和 NEW，最后将这些 Query 放在一个链表中并由参数 actions 传出。

2) 获取要定义规则的表的 OID。

3) 调用函数 DefineQueryRewrite 将已处理好的规则作为一个元组插入到系统表 pg_rewrite 中，DefineQueryRewrite 会把处理好的 WHERE 子句的表达式树以及规则的动作作为其参数之一。

数据结构 5.11 RuleStmt

```
typedef struct RuleStmt
{
    NodeTag      type;          // 节点类型
    RangeVar     *relation;     // 规则关系
    char         *rulename;     // 规则名
    Node         *whereClause;  // 条件字句
    CmdType      event;        // 动作类型
    bool         instead;       // 该规则是否为 instead 类型
    List         *actions;       // 该规则的动作
    bool         replace;       // 创建语句中是否有 OR REPLACE
}RuleStmt;
```

函数 DefineQueryRewrite 的流程如图 5-19 所示。在首先进行的检查工作中将做如下检查：

- 规则只能定义在表或视图上。
- 规则不能定义在系统表上，除非全局变量 allowSystemTableMods 设置为真。
- 规则动作不允许修改 OLD 或 NEW 两个虚表。

以上检查只要有一个不能通过就会报错。如果通过检查，则会调用函数 InsertRule 将处理好的

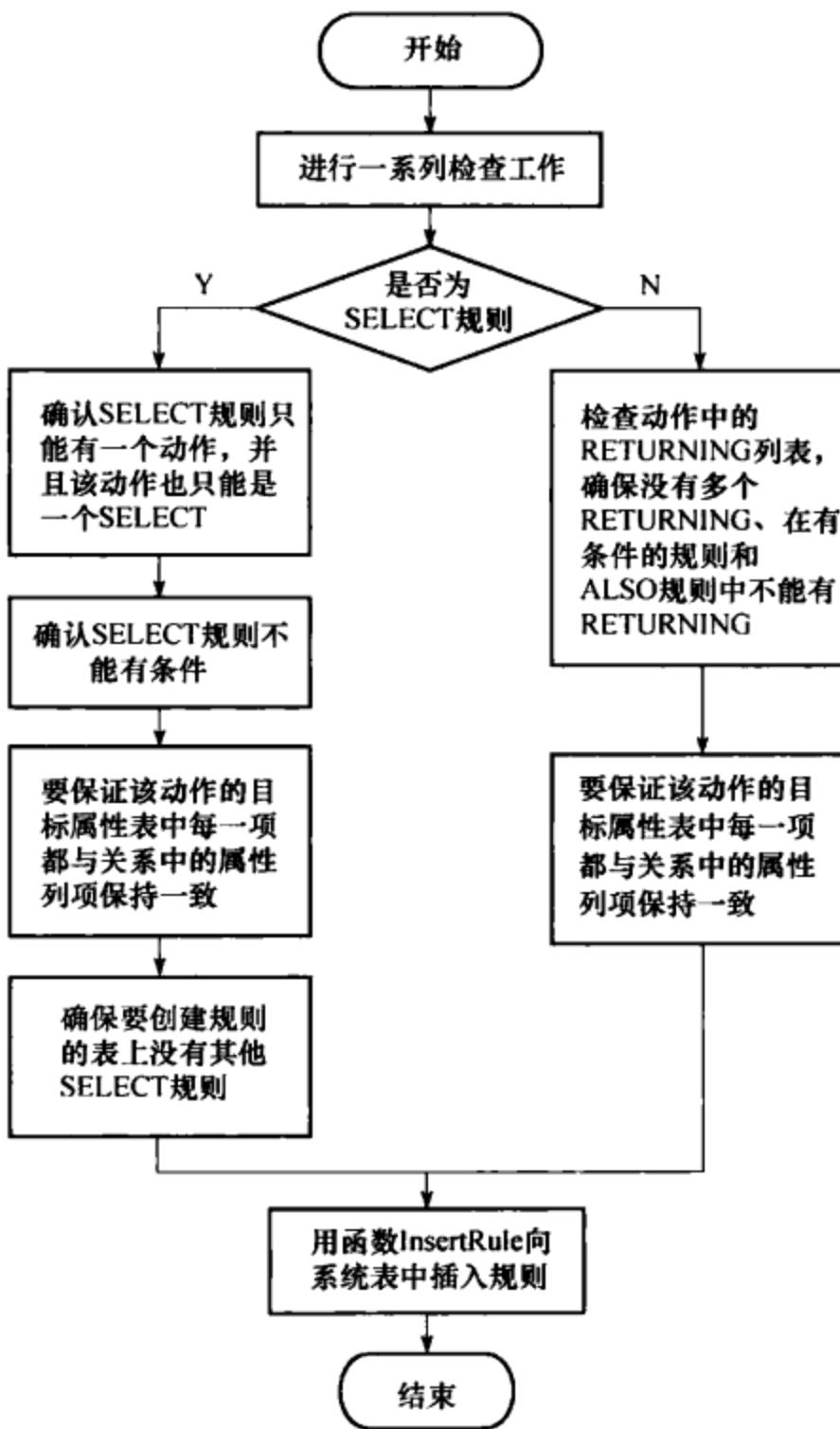


图 5-19 函数 DefineQueryRewrite 处理流程

规则信息组成一个 pg_rewrite 表元组，然后将这个元组插入到系统表 pg_rewrite 中去。将规则成功插入后，该函数返回规则所对应的 OID。规则插入的流程如下：

- 首先检查系统中是否有同名规则，如果有则报错。
- 新建一个 pg_rewrite 元组并对其每一个属性赋值，其中规则条件和规则动作都来自于 transformRuleStmt 传出的 WHERE 子句的表达式树和规则动作列表，并且会将它们转换成文本格式存放在元组的对应属性中。
- 将 pg_rewrite 元组插入该系统表中。
- 最后更新索引并返回规则的 OID。

2. 删除重写规则

在 PostgreSQL 中实现了两种删除规则的方式：第一种是根据规则名删除规则，由函数 Re-

`moveRewriteRule` 实现；第二种是根据规则的 OID 删除规则，由函数 `RemoveRewriteRuleById` 实现。当我们输入“DROP RULE”命令删除规则时，实际的删除工作将由 `RemoveRewriteRule` 完成；而 `RemoveRewriteRuleById` 的作用是在删除其他对象（比如视图）时，用于级联删除与这个对象相关的规则。

(1) 根据名称删除规则

函数 `RemoveRewriteRule` 删除名称由参数 `name` 指定的规则。它首先在 `SysCache` 中查找目标规则对应的 `pg_rewrite` 元组，若果不存在拥有该 `name` 的规则，则返回错误；否则先确认用户是否拥有删除该规则的权限，若用户没有相关权限，则返回错误。通过上述检查后才会将该规则对应的 `pg_rewrite` 元组删除。

`RemoveRewriteRule` 函数的主要流程如图 5-20 所示。

(2) 根据 OID 删除规则

函数 `RemoveRewriteRuleById` 删除给定 OID 所对应的规则。其流程如下：

- 1) 打开系统表 `pg_rewrite`，并对其加 `RowExclusiveLock` 锁，以防止在该表上执行 `INSERT/UPDATE/DELETE` 命令。
- 2) 利用指定规则的 OID 在系统表 `pg_rewrite` 中找到该规则所对应的元组。
- 3) 通过该元组中的 `ev_class` 字段获取该规则所作用的表，并对该表添加 `AccessExclusiveLock` 锁，以防止对其进行修改和删除等操作。
- 4) 删除系统表 `pg_rewrite` 中该规则所对应的元组，并关闭该系统表。
- 5) 更新该规则所作用的表的 `RelationData` 结构，并关闭此表。

3. 对查询树进行重写

`pg_rewrite_query` 中会调用函数 `QueryRewrite` 来完成查询树的重写。`QueryRewrite` 的处理流程如下：

- 1) 用非 `SELECT` 规则将一个查询重写为 0 个或多个查询，此工作通过调用函数 `RewriteQuery` 完成。
- 2) 对上一步得到的每个查询分别用 RIR 规则（无条件 `INSTEAD` 规则，并且只能有一个 `SELECT` 规则动作）重写，通过调用函数 `fireRIRules` 完成。
- 3) 将这些查询树作为查询重写的结果返回。

其中，函数 `RewriteQuery` 适用于 `UPDATE`、`DELETE`、`INSERT` 三种类型。这个函数的流程如图 5-21 所示。

函数 `fireRules` 用于对一系列规则进行循环处理。其中，对于带有条件的 `INSTEAD` 规则，将规则的条件的否定形式加入到原始查询树中，这样，当规则的条件不满足时整个查询可以按照原始查询的方式进行。然后调用 `rewriteRuleAction` 将规则的动作应用到查询树中，最终得到一个查询树的链表返回给 `RewriteQuery`。`RewriteQuery` 得到重写过的查询链表之后，还会对其中的每一个查询树递归地进行重写（可以看成是规则的递归触发），最后对 `RETURNING` 做一些检查之后给 `QueryRewrite` 返回查询树的链表。

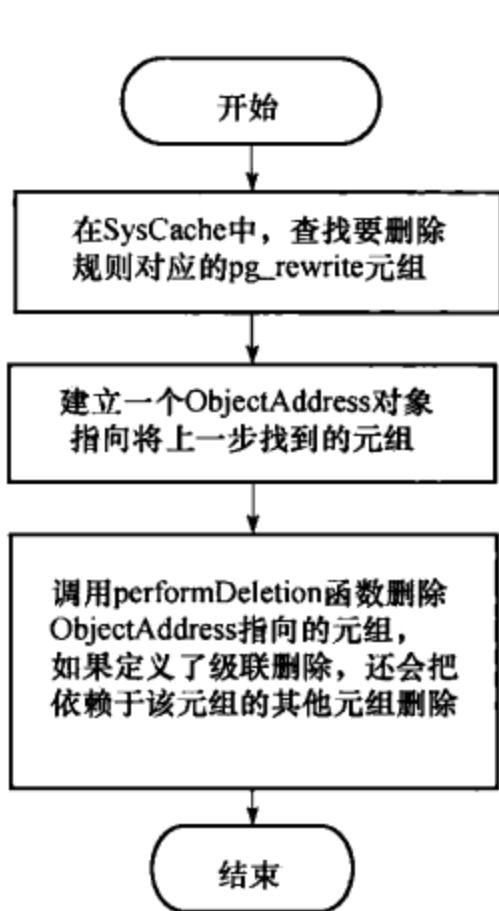


图 5-20 函数 RemoveRewriteRule 的处理流程



图 5-21 函数 RewriteQuery 处理流程

RewriteQuery 处理完所有非 SELECT 规则的重写之后，再由 fireRIRrules 来处理 SELECT 规则的重写。QueryRewrite 会调用 fireRIRrules，对 RewriteQuery 返回的查询树链表中的每一个查询树进行如下处理：

1) 对于该查询树的每一个范围表 (RTE)：

① 如果 RTE 是子查询，则递归调用 fireRIRrules 对其重写。

② 如果 RTE 是一个表，且在查询树中被引用，则找到它符合该查询树命令类型的规则，并调用 ApplyRetrieveRule 应用之。

2) 对于查询树中的每一个公共表表达式 (CTE)，也递归调用 fireRIRrules 对其查询树进行重写。

3) 如果查询树还有子链接 (子查询)，则调用 query_tree_walker 对整个查询树的子链接进行遍历，并用 fireRIROnSubLink 进行重写处理，而该函数实际也是调用 fireRIRrules 来完成重写工作。

4) 返回重写后的查询树。

QueryRewrite 将 fireRIRrules 返回的查询树组织成一个链表，对它们的 canSetTag 字段进行设置之后返回给 pg_rewrite_query。到这里，整个查询重写工作就已经结束了，最终生成的是一个由多棵查询树组成的查询树链表，该链表将会传递给查询规划模块形成执行计划。

5.4 查询规划

在数据库管理系统中，用户的查询请求可以采用不同的方案来执行。尽管不同方案返回给用户

的结果相同，但执行效率却存在差异，查询规划就用于选择一种代价最小的执行方案。因此，查询规划在数据库的查询性能方面起着举足轻重的作用。本节将介绍在 PostgreSQL 中如何生成并选择代价最小的执行方案（也称为执行计划）。

在数据库查询中，最耗时的是表连接，查询优化的核心思想是“尽量先做选择操作，后做连接操作”，因为先做选择操作可以减少后面进行表连接的数据量。PostgreSQL 中提升子链接和提升子查询的目的在于将连接操作尽量推迟，并且在此过程中将子查询的 WHERE 子句与父查询合并。由于 WHERE 子句中各种行约束条件的执行代价不一样，尽量将各个行约束条件合并到同一个 WHERE 子句中，再重新确定其执行顺序，就可以降低选择操作的代价，该问题在“扫描计划”这一节中会介绍。PostgreSQL 将需要进行连接操作的表提升到同一个查询层次后，将考虑如何选择一种代价最小的连接方案，该问题在“生成路径”这一节会讲到，在生成路径时将采用动态规划算法和遗传算法。

5.4.1 总体处理流程

查询规划的最终目的是得到可被执行器执行的最优计划，整个过程可分为预处理、生成路径和生成计划三个阶段。

预处理实际上是对查询树（Query 结构体）的进一步改造，这种改造可通过 SQL 语句体现（可以参考 5.4.8 节）。在此过程中，最重要的是提升子链接和提升子查询。在生成路径阶段，接收到改造后的查询树后，采用动态规划算法或遗传算法，生成最优连接路径和候选的路径链表。在生成计划阶段，用得到的最优路径，首先生成基本计划树（查询语句的 SELECT…FROM…WHERE 部分，在 5.4.4 节和 5.4.5 节介绍），然后添加 GROUP BY、HAVING 和 ORDER BY 等子句所对应的计划节点形成完整计划树（在 5.4.6 节介绍）。

查询分析工作完成之后，其最终产物——查询树链表将被移交给查询规划模块，该模块的入口函数是 pg_plan_queries，它负责将查询树链表变成执行计划链表。pg_plan_queries 只会处理非 Utility 命令，它调用 pg_plan_query 对每一个查询树进行处理，并将生成的 PlannedStmt 结构体（数据结构 5.12）构成一个链表（执行计划链表）返回。pg_plan_query 中负责实际计划生成的是 planner 函数，我们通常认为从 planner 函数开始就进入了执行计划的生成阶段。planner 函数的处理流程和函数调用关系如图 5-22 所示。

planner 函数会调用函数 standard_planner 进入标准的查询规划处理。函数 standard_planner 接收 Query 查询树及外部传递的参数信息，返回 PlannedStmt 结构体，该结构体包含执行器执行该查询所需要的全部信息，包括计划树 Plan、子计划树 SubPlan 和执行所需的参数信息。standard_planner 通过调用函数 subquery_planner 和 set_plan_references 分别完成计划树的生成、优化与清理工作。

函数 subquery_planner 接收 Query 查询树，返回 Plan（计划树），该计划树被包装在 PlannedStmt 中。对含有子查询的情况，可通过递归调用为子查询生成相应的子计划并链接到上层计划。该函数主要调用相关预处理函数，依据消除冗余条件、减少递归层数、简化路径生成等原则对 Query 查询树进行预处理。具体包括以下内容：

- 1) 一些特殊情况的处理，如内联返回值函数、预处理行标记、扩展继承表等。
- 2) 提升子链接，提升子查询。

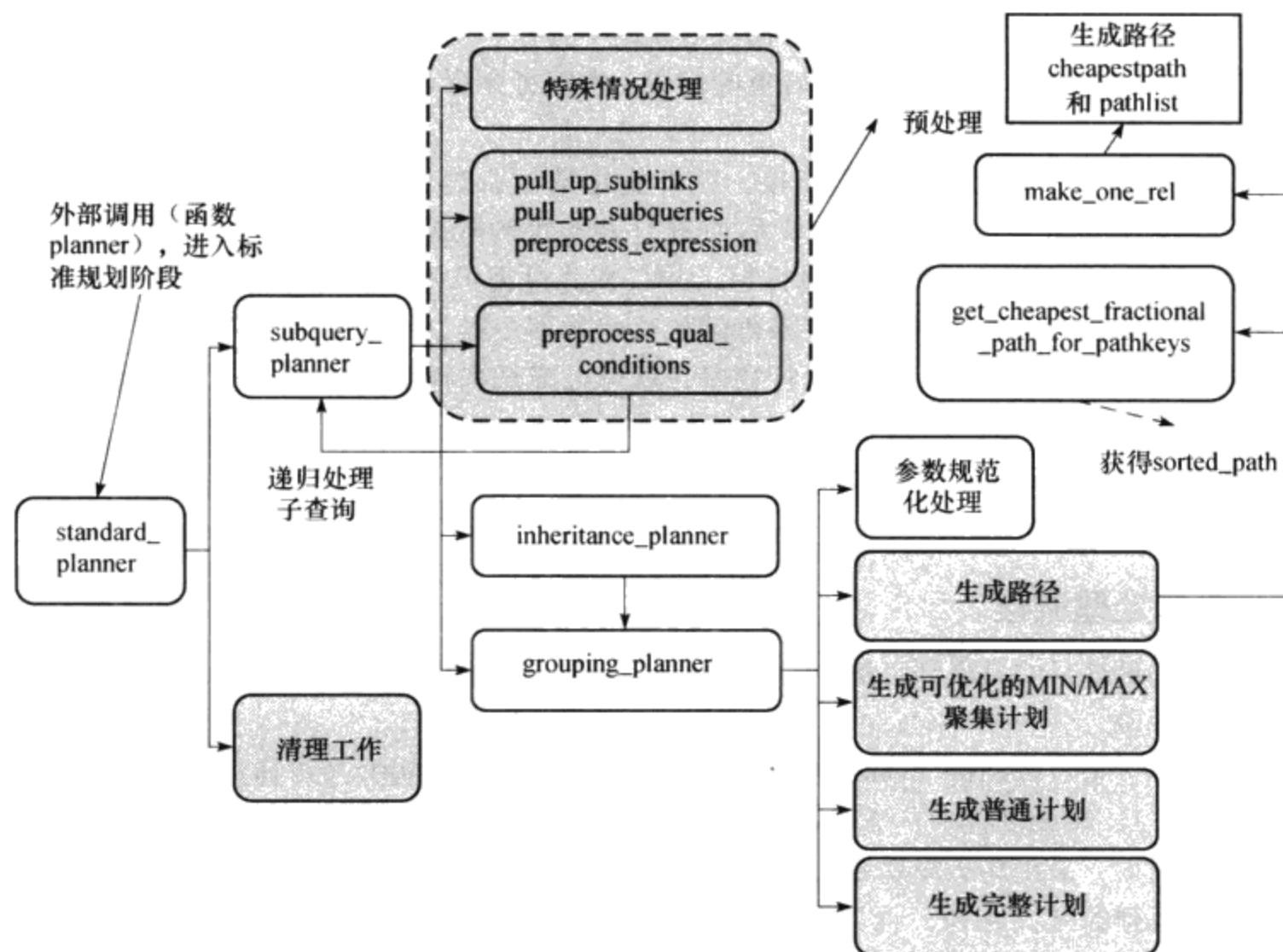


图 5-22 Planner 函数流程及主要函数调用情况

数据结构 5.12 PlannedStmt

```

typedef struct PlannedStmt
{
    NodeTag          type;
    CmdType          commandType;           // 计划所对应的命令类型
    bool             canSetTag;            // 是否需要设置命令结果标志
    bool             transientPlan;        // 当 TransactionXmin 改变时是否要重做计划
    struct Plan      *planTree;            // 计划树
    List             *rtable;              // 范围表
    List             *resultRelations;     // 计划中的结果关系,由结果关系的 RTE 索引构成
    Node             *utilityStmt;         // 定义游标时用来记录游标定义语句的分析树
    IntoClause       *intoClause;          // SELECT INTO 或 CREATE TABLE AS 的目标
    List             *subplans;             // 子计划
    Bitmapset        *rewindPlanIDs;        // 需要回滚的子计划的索引信息
    List             *returningLists;        // 需要返回的目标属性的链表,用于 RETURNING
    List             *rowMarks;             // 用于 SELECT FOR UPDATE 等
    List             *relationOids;          // 该计划所依赖的表的 OID
    List             *invalItems;            // 计划所依赖的其他对象,用 PlanInvalItem 结构表示,
                                         // 其中记录了被依赖的对象所属的 SysCache 的
                                         // ID 以及其对应的系统表元组的 TID
    int              nParamExec;           // 计划所需参数的数目
} PlannedStmt;
  
```

- 3) 预处理表达式。
- 4) 预处理连接约束条件，如果发现连接约束条件中含有子链接/子查询则会递归 subquery_planner 来优先规划子查询。

预处理按照“尽量先做选择操作后做连接操作”的思想改造查询树。在预处理完成后，会依据表是否存在继承关系分别调用函数 inheritance_planner 和 grouping_planner 进行规划处理。当所要处理的表存在继承关系时，函数 inheritance_planner 会将其处理成非继承关系的表，然后调用函数 grouping_planner 来处理。

函数 grouping_planner 在执行过程中不再对查询树做变换。SELECT 查询所包含的信息主要有以下四部分：目标属性、范围表、表连接约束条件以及行显示顺序约束条件。函数 grouping_planner 将这些信息进行规范化，传递给函数 query_planner，生成代价最小路径 cheapest_path 和排序路径 sorted_path（该路径可能为空）。

函数 query_planner 调用 make_one_rel 函数（生成路径主入口函数），即进入查询优化阶段。函数 make_one_rel 依据最优启动代价、最优执行代价、路径中包含的索引是否包含行排序信息等三个条件，保留相应路径到 pathlist 路径链表中，并推荐最优执行代价的路径放入 cheapest_path。之后函数 query_planner 通过调用函数 get_cheapest_fractional_path_for_pathkeys 在 pathlist 中寻找符合排序需求的排序路径 sorted_path。

函数 grouping_planner 得到代价最小路径 cheapest_path 和排序路径 sorted_path 后，依据是否存在 Hash 等条件从中最终确定最优路径 best_path。之后，函数 grouping_planner 调用函数 create_plan 为路径生成基本计划树（只包含目标属性、范围表以及表连接约束信息的连接计划树），最终依据分组和排序等条件（GROUP BY 和 ORDER BY 等）在生成的基本计划树上添加相应的计划节点，生成完整的计划树。

表 5-11 查询优化的各主要函数功能说明

主要函数	说明
planner	优化器的人口函数；输入经过重写器处理的查询树，输出最优的执行计划（如果连接表比较多时，为近似最优）
standard_planner	标准的规划处理；PostgreSQL 中提供了扩展接口，开发人员可自行设计规划器并在 PostgreSQL 中使用之
subquery_planner	优化处理的主体函数，可以递归处理子查询
inheritance_planner	存在继承表的规划处理
grouping_planner	不存在继承表的规划处理
set_plan_references	完成生成执行计划后的清理工作

下面对查询规划处理总体控制的几个函数进行介绍。

(1) standard_planner

standard_planner 的返回值是一个 PlannerStmt 结构，其中包含执行器执行计划时的全部信息，包含了计划树 (Plan)、子计划树和相关参数信息。该函数有三个参数：

- 1) parse：一个 Query 结构，表示需要处理的查询树。
- 2) cursorOption：整型，表示游标选项，在处理与游标操作时用到。游标选项主要有以下几种。
 - CURSOR_OPT_BINARY：游标以二进制的形式而不是文本形式返回结果。
 - CURSOR_OPT_SCROLL：允许游标以非顺序的方式返回结果（例如反向扫描）。

- CURSOR_OPT_NO_SCROLL：不允许以非顺序的方式返回结果。
- CURSOR_OPT_INSENSITIVE：游标建立之后对其所在表上的更新操作不会影响游标的结果。
- CURSOR_OPT_HOLD：把游标一直保持，即使创建它的事务提交之后也仍然可以使用。
- CURSOR_OPT_FAST_PLAN：希望创建一个能快速启动的计划。

3) boundParams：ParamListInfo 结构，记录了所用到的参数信息。

一般情况下，CursorOption、boundParams 参数为零值/空值。standard_planner 的主要流程如图 5-23 所示。

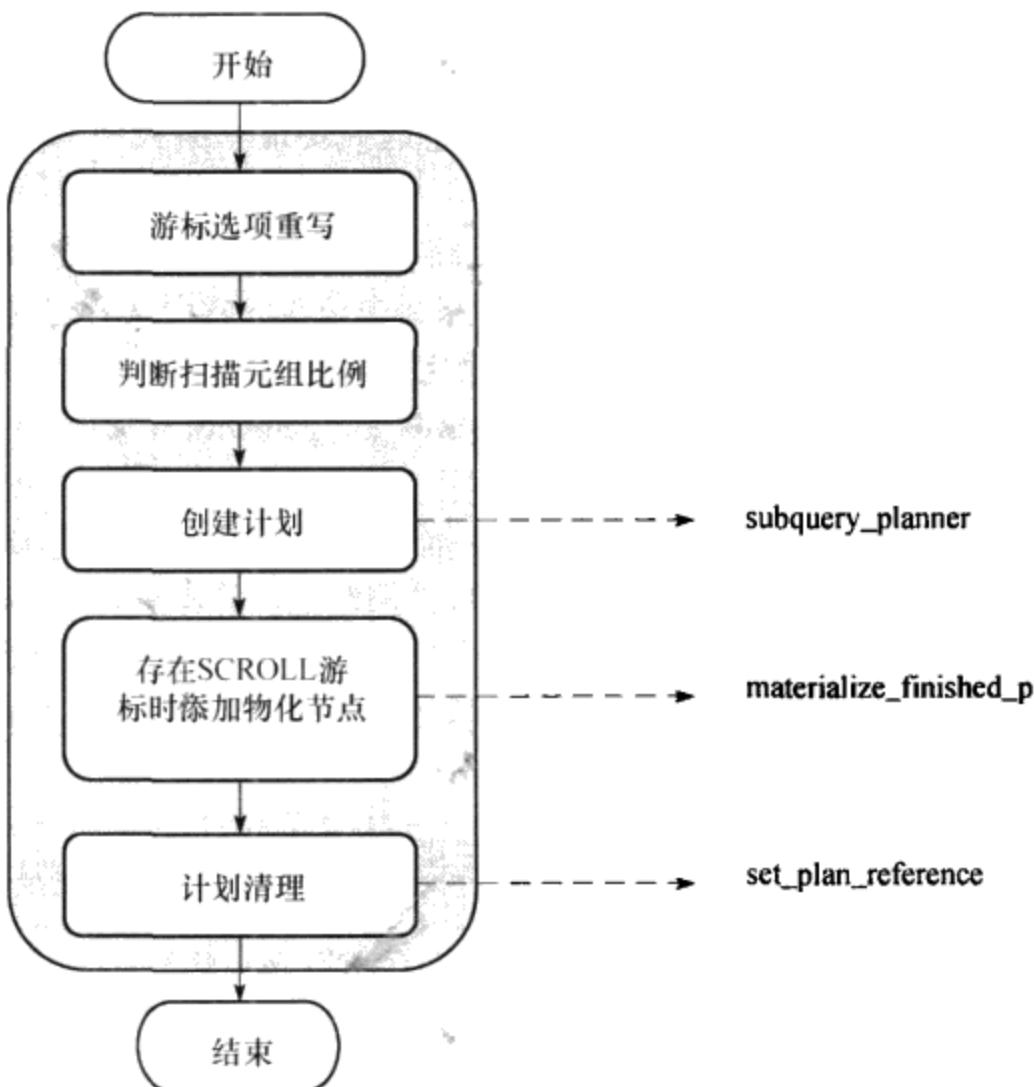


图 5-23 函数 standard_planner 的流程

1) 游标选项可能来自于外界传递的参数，也可能来自游标声明状态，故需要结合查询树 parse 中命令状态变量 utilityStmt 对游标参数重写。

2) 依据游标选项判断需要扫描多少比例的元组。由 `tuple_fraction` 变量记录，如果游标选项设置为 CURSOR_OPT_FAST_PLAN，则将 `tuple_fraction` 设置为 GUC 变量 `cursor_tuple_fraction` 的值（默认为 0.1）。由于该 GUC 变量可以在配置文件（`postgresql.conf`）中设置，因此还需要对 `tuple_fraction` 进行一次检验并调整。扫描比例变量 `tuple_fraction` 的取值范围为 [0 - 1]。默认情况下扫描全部元组，即 `tuple_fraction = 0`；当 `tuple_fraction` 大于等于 1 时（用户在配置文件中设置失误可能导致 `tuple_fraction` 超出其取值范围 [0 - 1]），则将其置为 0；当 `tuple_fraction` 小于 0 时，则将其置为 `1e-10`。如果游标选项设置为其他值，设置 `tuple_fraction` 为 0。

3) 调用函数 `subquery_planner` 生成计划。

4) 如果游标参数是 CURSOR_OPT_SCROLL，且该计划类型不支持反向执行，也就是说该计划

不能支持游标的非顺序返回元组，则调用 `materialize_finished_plan` 函数在原有计划树的上层添加物化节点（物化节点会将可能产生的所有元组全部找到并缓存起来，详见 6.4.3 节）作为新的计划树。这个操作实际上是将原始计划树的结果物化下来，这样游标就可以进行反向返回元组。

5) 调用 `set_plan_references` 函数对计划树做清理操作，该函数只是方便执行计划做相应的变量调整，并不对计划树做本质上的改变。

6) 创建 `Plannedstmt` 类型结点，填入相应变量，并返回。

(2) subquery_planner

`subquery_planner` 负责创建计划，可递归处理子查询。`subquery_planner` 的工作分为两部分：

- 依据消除冗余条件、减少查询层次、简化路径生成的基本思想，调用预处理函数对查询树进行预处理。
- 调用 `inheritance_planner` 或 `grouping_planner` 进入生成计划流程，该过程不对查询树做出实质性的改变。

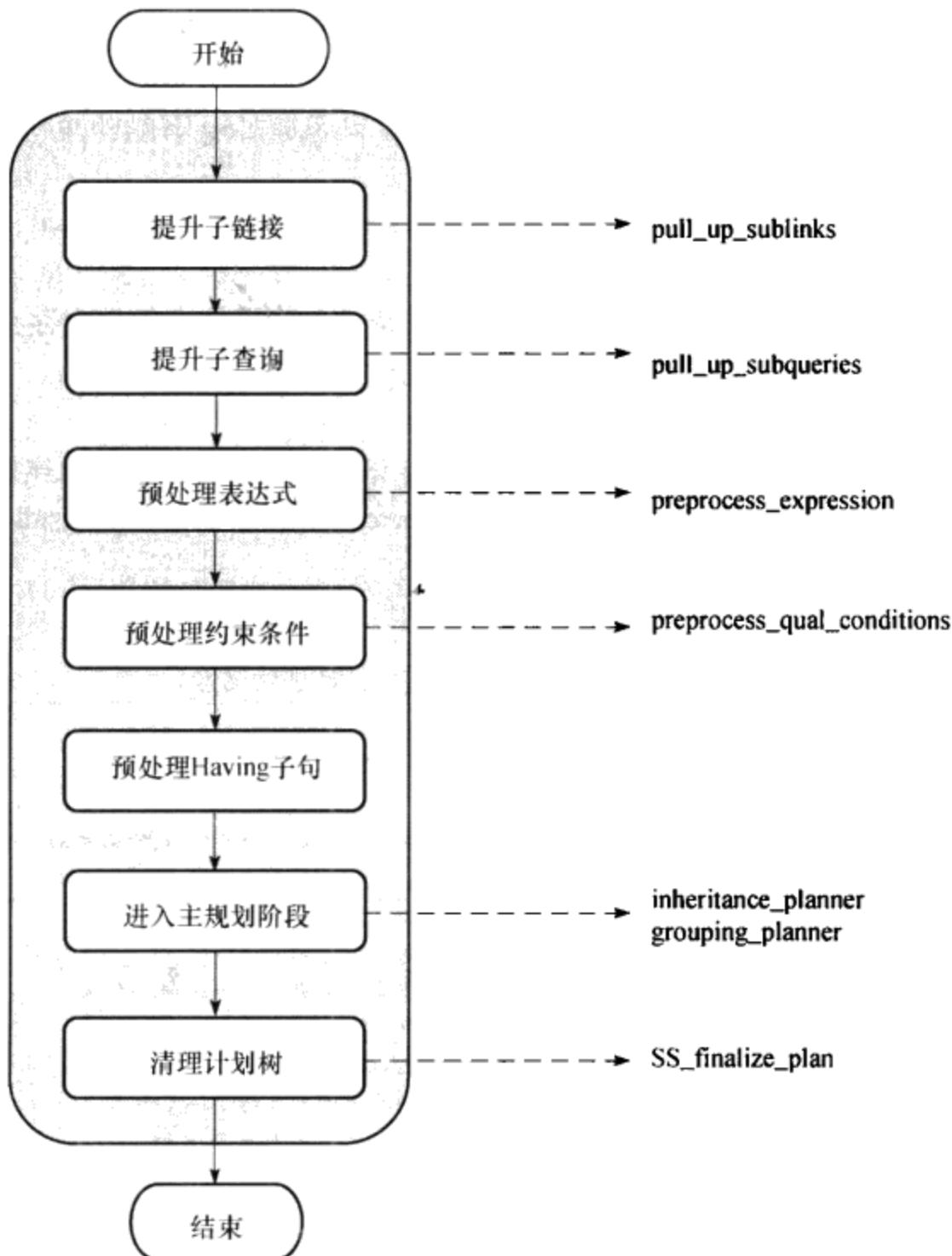


图 5-24 函数 `subquery_planner` 的流程

subquery_planner 的返回值是一个 Plan 结构（计划树，见数据结构 5.13），它有以下几个参数：

- glob：PlannerGlobal 类型指针，记录做计划期间的全局信息，例如子计划、子计划的范围表等，每条查询语句有且仅有一个该变量。
- parse：Query 类型指针，指向要生成计划的查询树。
- parent_root：PlannerInfo 类型指针，指向父查询的规划器相关信息（首次调用为空）。
- subroot：PlannerInfo 类型指针的指针，用于向外部函数传出本层查询的规划器相关信息。
- hasRecursion：bool 类型，如果正在处理的是一个递归的 WITH 查询，则 hasRecursive 为真。
- tuple_fraction：double 类型，表示计划扫描元组的比例。

如图 5-24 所示，函数 subquery_planner 的处理流程如下：

- 1) 初始化信息，为该子查询创建一个 PlannerInfo 结构（数据结构 5.14）。如果 hasRecursive 为真，还要为该子查询新建一个 worktable 变量（见第 6 章中关于 CTEScan 的相关内容），并把其标识符放在 PlannerInfo 的 wt_param_id 字段中。
- 2) 调用 pull_up_sublinks 提升子链接。
- 3) 调用 pull_up_subqueries 提升子查询。
- 4) 根据查询树中的相关信息确定 PlannerInfo 中的相关信息，如 hasJoinRTEs、hasOuterJoins、hasHavingQual 等。
- 5) 调用 preprocess_expression 对查询树的 targetlist、returninglist、havingQual、limitOffset、limitCount 等字段进行处理，同时也对范围表进行预处理。
- 6) 调用 preprocess_qual_conditions 函数预处理约束条件，如果约束条件中存在子链接/子查询，则递归调用 subquery_planner 处理子链接和子查询。
- 7) 将 HAVING 子句提升到 WHERE 条件中去。
- 8) 如果存在 hasOuterJoins，调用 reduce_outer_joins 试图将其简化为简单内连接。
- 9) 进入生成计划树的主要步骤，调用函数 inheritance_planner 和 grouping_planner 生成计划树，用数据结构 Plan 存储和表示。
- 10) 调用 SS_finalize_plan 对计划树进行清理。

数据结构 5.13 Plan (存储执行计划的节点)

```
typedef struct Plan
{
    NodeTag          type;
    Cost             startup_cost;           // 获取任何元组之前的代价，也称为启动代价
    Cost             total_cost;            // 总代价（假设所有的元组都已经获取）
    double           plan_rows;            // 计划期望获取的元组数
    int              plan_width;           // 行的平均字节宽度
    List             *targetlist;           // 该节点中需要计算的目标属性的链表
    List             *qual;                 // 查询条件
    struct Plan     *lefttree;             // 当前计划节点的左子树
    struct Plan     *righttree;            // 当前计划节点的右子树
    List             *initPlan;             // 初始计划节点（非相关子查询）
    Bitmapset        *extParam;             // 子计划用到的所有外部参数，比如相关子计划从父计划得到的参数
    Bitmapset        *allParam;             // 所有参数
} Plan;
```

数据结构 5.14 PlannerInfo

```

typedef struct PlannerInfo
{
    NodeTag          type;
    Query            *parse;                                //需要生成计划的查询树
    PlannerGlobal    *glob;                                //规划器当前运行状态的全局信息
    Index             query_level;                         //当前查询的层次，最外层查询的层次为1
    struct PlannerInfo *parent_root;                      //根节点的生成计划信息
    struct RelOptInfo **simple_rel_array;                 //基本关系的信息，生成路径时用到
    int               simple_rel_array_size;              //simple_rel_array 数组的大小
    RangeTblEntry    **simple_rte_array;                  //范围表
    List              *join_rel_list;                     //进行连接操作的表的信息
    struct HTAB      *join_rel_hash;                    //进行连接操作的表的 Hash 表，当表很多
                                                               //时，可以加快查找
    List              *resultRelations;                  //结果关系
    List              *returningLists;                   //目标属性
    List              *init_plans;                        //查询的初始子计划
    List              *cte_plan_ids;                     //计划的 CTE 子计划
    List              *eq_classes;                        //等值连接的 PathKeyItems 链表
    List              *canon_pathkeys;                   //规范化的 PathKeys 链表
    List              *left_join_clauses;                //左外连接的 RestrictInfo 结构
    List              *right_join_clauses;                //右外连接的 RestrictInfo 结构
    List              *full_join_clauses;                 //全连接的 RestrictInfo 结构
    List              *join_info_list;                   //SpecialJoinInfo 结构
    List              *append_rel_list;                  //AppendRelInfo 结构，追加的表的信息
    List              *placeholder_list;                 //PlaceHolderInfo 结构
    List              *query_pathkeys;                   //函数 query_planner 需要的 pathkeys
    List              *group_pathkeys;                   //GROUP 子句的 pathkeys
    List              *window_pathkeys;                  //窗口函数的 pathkeys
    List              *distinct_pathkeys;                //DISTINCE 子句的 pathkeys
    List              *sort_pathkeys;                     //ORDER BY 子句的 pathkeys
    List              *initial_rels;                     //进行连接操作之前的初始的表
    MemoryContext     planner_ctxt;                    //规划器的内存上下文信息
    double            total_table_pages;                //所有表的页数
    double            tuple_fraction;                  //元组数
    bool              hasJoinRTES;                      //范围表是否为连接类型
    bool              hasHavingQual;                  //是否有 HAVING 子句
    bool              hasPseudoConstantQuals;           //是否有任何一个 RestrictInfo 的
                                                               //pseudoconstant 为真
    bool              hasRecursion;                   //WITH 子句是否允许递归处理，hasRecursion
                                                               //为真时，后续字段有效
    int               wt_param_id;                     //工作表的 PARAM_EXEC ID
    struct Plan      *non_recursive_plan;             //非递归计划
} PlannerInfo;

```

(3) inheritance_planner

在 PostgreSQL 中，计划分为两种类型：继承关系计划（Inheritance—Plan）和一般计划（Grouping—Plan）。前者由函数 `inheritance_planner` 负责，用于处理更新操作（`UPDATE`、`DELETE`、`INSERT`）中的对象存在继承关系的情况。在这种情况下，为了保证一致性，需要为查询请求中涉及的每一个基本关系都生成一个计划。最后，我们得到的是一个计划链。一般计划类型则由 `grouping_planner` 函数负责，用于生成不涉及继承关系的计划。

函数 `inheritance_planner` 的工作其实就是多个一般计划的生成，其过程如下：

- 1) 根据继承关系，找到操作中涉及的所有基本关系。
- 2) 为每个涉及的基本关系生成一个查询树。
- 3) 对于每个生成的子关系查询树，循环调用函数 `grouping_planner`，为它们生成一个计划。
- 4) 将生成的多个计划链接成为一个计划链返回。

(4) grouping_planner

`Grouping-plan` 的生成由函数 `grouping_planner` 实现，其主要流程如图 5-25 所示。`grouping_planner` 要考虑查询中是否需要进行集合操作，该操作可通过查询树中的 `setOperations` 成员变量来判断，如果不为空则表示需要进行集合操作。对于集合操作，遍历 `setOperations`，为其中的每个子查询生成计划。对于非集合操作，计划生成过程如下：

1) 对于含有 `GROUP BY` 子句的查询，调整其 `GROUP BY` 属性的顺序以匹配 `ORDER BY` 子句中的属性顺序，这样可以使用一次排序操作来同时实现排序和分组。

2) 调用 `preprocess_targetlist` 预处理 `INSERT`、`UPDATE`、`DELETE` 及 `FOR UPDATE` 情况下的目标属性。

3) 计算并确定代表排序需求的路径关键字，主要有 `groupClause`、`WindowClause`、`distinctClause` 和 `sortClause`，其优先级依次降低。

4) 调用 `query_planner` 为一个基本查询创建路径，包含 `cheapest_path` 和 `sorted_path` 两条路径。前者为代价最低执行路径，后者为对排序最优执行路径（可能为空）。

5) 确定最优路径 `best_path`，存在 Hash 分组或不存在 `sorted_path` 时采用 `cheapest_path` 为最优路径。

6) 生成可优化的 `MIN/MAX` 聚集计划。

7) 如果上一步没有生成聚集计划，调用 `create_plan` 生成普通计划。

8) 根据生成的计划，结合用户条件进行包装：

①若有 `GROUP BY` 子句，则在已得到的计划上再封装 `GROUP BY` 信息，生成一个 `Group` 型计划。

②若有聚集操作则调用函数 `make_agg` 生成 `Agg` 类型计划。

③若有 `ORDER BY` 子句，则增加专门的排序操作，并生成 `Sort` 计划。这样处理的前提是已得到的计划不是按照该子句信息进行排序，如果现有计划已经能满足排序要求，则这里不需要额外的生成 `Sort` 计划。

④若有 `DISTINCT` 子句，则包装生成 `Unique` 计划。

⑤若有 `LIMIT` 子句，则在计划中增加 `LIMIT` 信息。

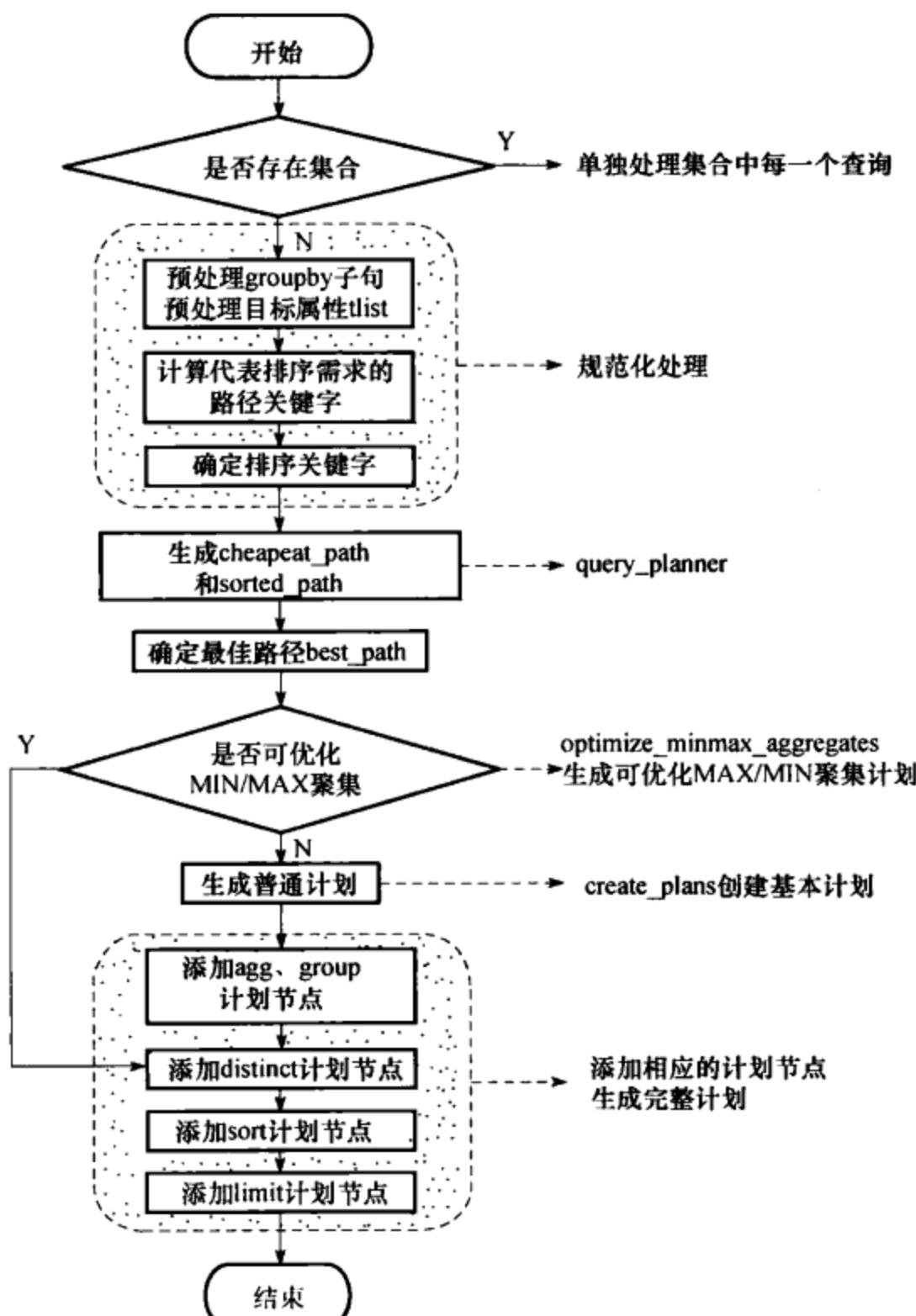


图 5-25 grouping_planner 的处理流程

5.4.2 预处理

前面说过，在实际进行计划生成之前将对查询树做一些预处理，预处理的主要工作是提升子链接和子查询以及预处理表达式和 HAVING 子句等。预处理部分主要是对查询树 Query 中的范围表 rtable 和连接树 jointree 等进行处理。

1. 提升子链接/子查询

一个“SELECT…FROM…WHERE”语句称为一个查询块，将一个查询块嵌套到另一个查询块的 FROM 子句、WHERE 子句或 HAVING 子句中的查询称为嵌套查询，其中被嵌入其他查询块中的查询块称为嵌套子查询。SQL 语句允许多层嵌套子查询，嵌套查询的一般处理方法是由里向外处

理，即每个子查询在其父查询处理之前求解，子查询的结果用于建立父查询的查找条件。在 PostgreSQL 中子链接用来表示出现在表达式中的子查询与普通子查询的联系和区别：子查询是一条完整的查询语句，而子链接是一条表达式，但该表达式内部可以包含查询语句。从直观上来说，子查询是出现在 FROM 子句中的，而子链接则出现在 WHERE 子句或 HAVING 子句中。

按相关性，嵌套查询可分为相关子查询和非相关子查询。相关子查询是指该子查询的执行依赖于外层父查询的某些属性值。它需要接受父查询的参数，因此会设置相关标记，当参数改变的时候，需要重置参数后重新执行一遍子查询得到新的结果。非相关子查询中子查询完全可以独立。

按关键字，嵌套查询可以分为以下几类：

- 1) EXISTS：声明了 EXISTS 的子查询。
- 2) ALL：声明了 ALL 或 NOT IN 的子查询。
- 3) ANY：声明 ANY 或 IN 的子查询。
- 4) EXPR：子查询返回一个参数给外层父查询。
- 5) MULTIEXPR：子查询返回多个参数给外层父查询，例如语句“SELECT *FROM B WHERE (b1,3,’aa’) > SELECT * from A；”中的子查询将向父查询返回多个属性值；
- 6) ARRAY：子查询是将某些值构成数组的表达式，例如：“SELECT ARRAY[1,2,3 + 4]；”。

一个出现在 FROM 子句中的子查询如果使用了诸如聚集函数、分组、DISTINCT 属性，它将被独立规划成一个子计划，在规划父查询的过程中会被当作一个黑盒子。但是，如果子查询仅仅只是一个简单的扫描或是连接，把子查询当作是一个黑盒子可能会产生一个较差的规划。

下面我们结合一个例子来解释上面的问题。例如有一个原始的 SQL 语句如下：

```
SELECT D.dname
  FROM dept D
 WHERE D.deptno IN
      (SELECT E.deptno FROM emp E WHERE E.sal = 10000);
```

对于该语句中的子查询，如果它被独立地规划，那么对于表 dept 中的每一个元组 deptno 值，都要搜索整个表 emp 以确定该元组是否满足 WHERE 条件，显然这样做的代价是非常大的。如果我们把子查询提升并合并到父查询中，其效果将会大不一样。其子查询存在于 WHERE 子句中，在 PostgreSQL 实现层面为一个子连接。

首先，对语句提升子链接：

```
SELECT D.dname
  FROM dept D, (SELECT E.deptno FROM emp E WHERE E.sal = 10000) As Sub
 WHERE D.deptno = Sub.deptno
```

然后，提升子查询：

```
SELECT D.dname
  FROM dept D,emp E
 WHERE D.deptno = E.deptno AND E.sal = 10000;
```

可以看到，提升子查询之后，子查询被合并到父查询中，可以进行统一优化。在这个例子中，

执行时只需要先做一下过滤 ($E.sal = 10000$)，然后把结果和 dept 做一下连接就可以了，查询的效率大大提高。

因此，在查询规划的预处理阶段，优化器会查找简单的子查询，把它们放到整个父查询的连接树中，这个过程称为提升子查询。

如果把子查询合并到父查询中导致了一个 FROM 子句作为另一个 FROM 子句的一部分（在它们之间没有显式的连接标识），那么可以把这两个 FROM 子句列表合并到一起。然而，这可能导致规划时间迅速增长，因为动态规划搜索的运行时间复杂度是以所涉及的 FROM 子句中表的个数为指数增长的。因此，如果合并导致表个数过多，就不会进行 FROM 子句的合并。

subquery_planner 调用 pull_up_sublinks 处理 WHERE 子句和 JOIN/ON 子句中的 ANY 和 EXISTS 类型的子链接。

pull_up_sublinks 调用 pull_up_sublinks_jointree_recurse 递归处理 jointree，函数 pull_up_sublinks_jointree_recurse 处理流程如下：

- 1) 对于 RangeTable 类型，直接返回。
- 2) 对于 FromExpr 类型，递归调用 pull_up_sublinks_jointree_recurse 处理每个节点，然后调用 pull_up_sublinks_qual_recurse 处理约束条件。
- 3) 对于 JoinExpr 类型，递归调用处理左右子树每个节点，然后调用 pull_up_sublinks_qual_recurse 处理约束条件。

提升子查询通过递归调用函数 pull_up_subqueries 来实现，当子查询仅仅只是一个简单的扫描或连接时，规划器就会使用提升子查询操作把子查询或者子查询的一部分合并到父查询中，以便进行统一的优化。提升子查询时分以下三种情况处理：

- 1) 范围表中存在子查询。如果是简单的子查询，则调用函数 pull_up_simple_subquery 直接提升；如果是简单的 UNION ALL 子查询，调用函数 pull_up_simple_union_all 直接提升；其他情况不处理。
- 2) FROM 表达式中的子查询。对于 FROM 列表中的每个节点都调用函数 pull_up_subqueries 递归处理。
- 3) 连接表达式中的子查询。根据不同的连接类型，递归调用函数 pull_up_subqueries 进行处理。

2. 预处理表达式

表达式的预处理工作由函数 preprocess_expression 完成，处理的对象有：目标属性、HAVING 子句、OFFSET 和 LIMIT 子句、连接树。preprocess_expression 采用递归扫描的方式对 PlannerInfo 结构中的表达式进行处理，主要的工作包括：

- 1) 用基本关系变量取代连接别名变量，此工作由函数 flatten_join_alias_vars 完成。
- 2) 进行常量表达式的简化，此工作由函数 eval_const_expressions 完成。主要是将常量表达式的值直接计算出来替换掉原有的表达式，例如 “ $2 + 2 < > 4$ ” 会被 False 替换，“ $x \text{ AND } \text{False} < > \text{False}$ ” 会被 False 替换。
- 3) 对表达式进行规范化，该工作由函数 canonicalize_qual 完成。这个步骤不会对目标属性进行处理，其目标是将表达式转化成为最佳析取范式或合取范式。
- 4) 将子链接 (SubLinks) 转化成为子计划 (SubPlans)，对每一个子链接都调用函数 make_subplan 完成转换。

函数 make_subplan 的处理过程如下：

1) 复制 SubLink 中的 Query。

2) 对于一个 EXISTS 的子计划，调用 `simplify_EXISTS_query` 对 Query 副本进行处理。由于 EXISTS 类型的子查询只需要根据是否有满足条件的元组返回真值或假值，因此可以简化掉与此无关的语法部分，例如可以丢弃掉目标属性、GROUP BY 子句等。

3) 调用 `subquery_planner` 为子链接生成计划，并通过设置其 `tuple_fraction` 参数来告诉低层规划器要取多少元组。如果是 EXISTS 类型的子查询，只需要取第一个元组即可完成查询，设置 `tuple_fraction` 为 1.0；如果是 ALL 或 ANY 类型的子查询，从概率上来说，大约取到 50% 的元组就可以得到结果，因此设置 `tuple_fraction` 为 0.5；对于其他情况，设置 `tuple_fraction` 为 0。

4) 调用 `build_subplan` 将上一步生成的计划转换成 SubPlan 或 InitPlan 的形式，并将当前查询层次的参数列表传递给它的子计划。

5) 如果生成的是 SubPlan 且是一个简单的 EXISTS 类型的子查询，调用 `convert_EXISTS_to_ANY` 尝试将它转换为一个 ANY 类型的子查询，并为其创建计划。

3. 预处理 HAVING 子句

对于 HAVING 子句，除了进行前面所提到的预处理外，我们还需要处理其中的每个条件。如果不含有聚集，则将 HAVING 子句提升到 WHERE 条件中去，否则将它放到 Query 的 HavingQual 字段中。

例如，语句 “`SELECT a FROM A WHERE b < 10 HAVING a > 1;`” 中，HAVING 子句中没有聚集，则应将 “`a > 1`” 提升到 WHERE 条件中形成语句 “`SELECT a FROM A WHERE b < 10 AND a > 1;`”。

其具体步骤为：

1) 初始化空 HAVING 子句链表 newHaving。

2) 扫描 HAVING 子句链表 (Query 的 HavingQual 字段) 中每一条 HAVING 子句。

①如果 HAVING 子句存在包含聚集、易失函数 (volatile function)、子计划任一种情况，则添加至 newHaving，即保留为 HAVING 子句。

②如果不包含 GROUP 子句，则添加至 WHERE 子句中。

③不属于以上两种情况，则将其添加至 WHERE 子句，同时保留至 newHaving。

3) 将处理后的子句链表 newHaving 替代原 HAVING 子句链表。

5.4.3 生成路径

对于 SQL 中的计划命令的处理，无非就是获取一个（或者一系列）元组，然后将这个元组返回给用户或者以其为基础进行插入、更新、删除操作。因此对于一个执行计划来说，最重要的部分就是告诉查询执行模块如何取到要操作的元组。执行计划要操作的元组可以来自于一个基本表或者由一系列基本表连接而成的“连接表”，当然一个基本表也可以看成是由它自身构成的连接表。这些基本表连接成连接表的关系可以从逻辑上表示成一个二叉树结构（连接树），由于表之间不同的连接方式和顺序，同一组基本表形成连接表的连接树会有多个，每一棵连接树在 PostgreSQL 中都称为一条路径。因此，路径在查询的规划和执行过程中表示了从一组基本表生成最终连接表的方式。而查询规划的工作就是从一系列等效的路径中选取效率最高的路径，并形成执行计划。

生成路径的工作是由函数 `query_planner` 来完成的，该函数的处理流程如图 5-26 所示。

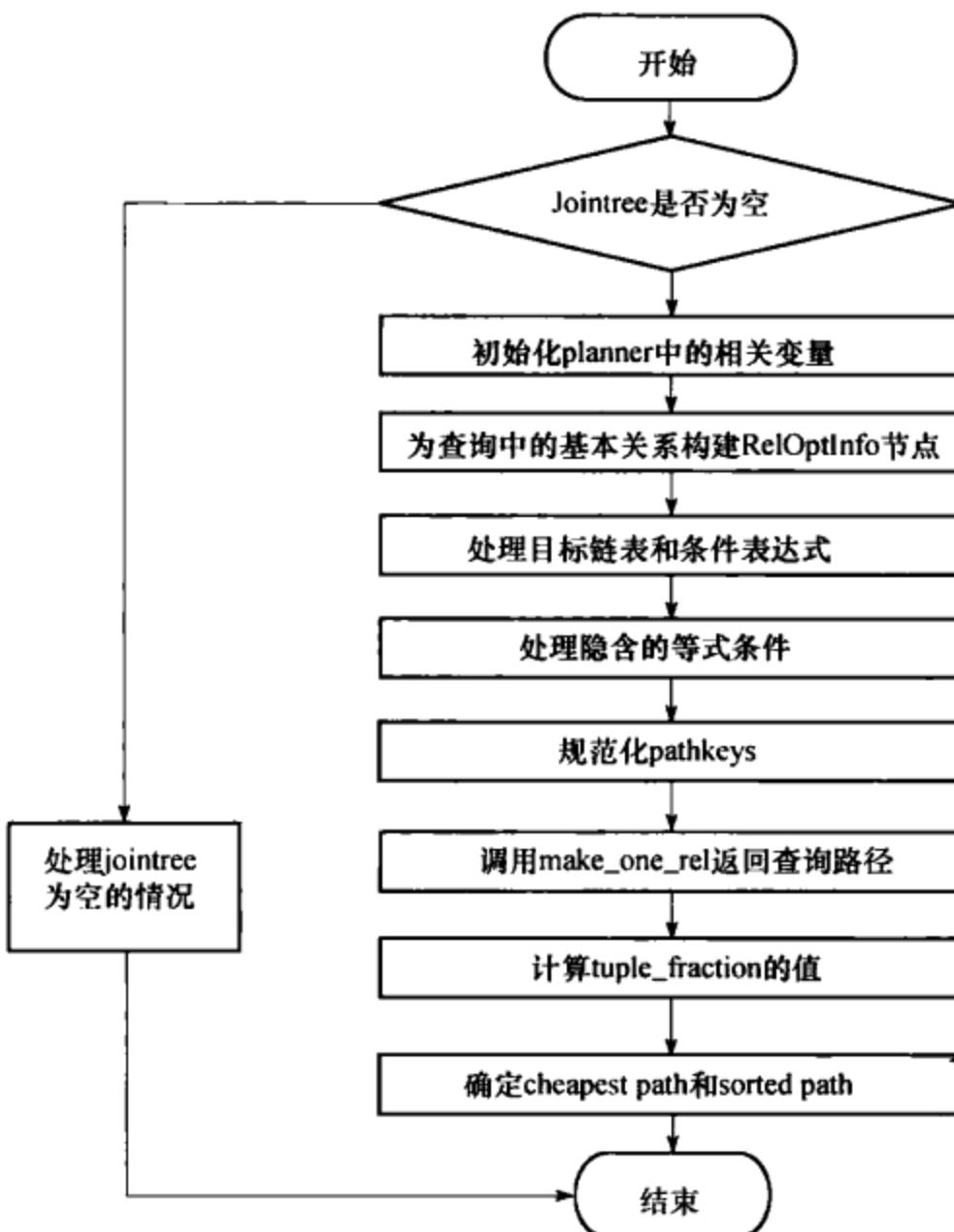


图 5-26 函数 query_planner 处理流程

RelOptInfo 结构（数据结构 5.15）是贯穿整个路径生成过程的一个数据结构，生成路径的最终结果始终存放在其中，生成和选择路径所需的许多数据也存放在其中。路径生成和选择涉及的所有操作几乎都是针对这个结构进行的，因此搞清楚这个结构对于理解整个路径生成过程非常重要。

数据结构 5.15 RelOptInfo

```

typedef struct RelOptInfo
{
    NodeTag      type;           //节点类型
    RelOptKind   reloptkind;     //关系的类型(基本关系类型、连接关系类型、其他成员关系类型)
    Relids       relids;         //关系所涉及的基本表(在范围表中的编号)
    double       rows;           //元组数(估计)
    int          width;          //元组平均字节宽度(估计)
    void*        relcache;       //物化信息
    List        *reltargetlist;  //关系的目标属性
}

```

```

List    *pathlist;           //多个基本关系生成该关系的路径(多条可行路径)
struct Path *cheapest_startup_path; //最小的启动代价的路径
struct Path *cheapest_total_path;   //最小的总代价的路径
struct Path *cheapest_unique_path; //最小的唯一代价的路径
                                //基本关系的信息,对于连接关系(多个基本关系通过连接操作生成的关系)无效

Index      relid;          //当前关系在范围表中的编号
RTEKind    rtekind;        //基本关系的类型(普通关系、子查询等)
AttrNumber min_attr;       //关系中最小的字段编号
AttrNumber max_attr;       //关系中最大的字段编号
Relids     *attr_needed;   //连接关系中那些字段是需要的
int32     *attr_widths;    //每个字段的宽度的估计
List      *indexlist;      //索引关系的信息
BlockNumber pages;         //关系占用的页面数
double    tuples;          //关系的元组数
struct Plan *subplan;      //子查询的计划
List      *subrtable;      //子查询的范围表
                           //各种扫描和连接时使用
List      *baserestrictinfo; //基本关系上的约束信息
QualCost baserestrictcost; //对以上约束信息的代价估计
List      *joininfo;        //连接子句中调用该关系时的约束信息
bool     has_eclass_joins; //是否存在等值连接
                           //对关系进行索引扫描时的缓存信息
Relids    index_outer_relids; //可进行索引扫描的关系
List      *index_inner_paths; //存放最佳索引路径
} RelOptInfo;

```

该结构涉及 baserel（基本关系）和 joinrel（连接关系）的概念。baserel 是一个普通表、子查询或者是范围表中出现的函数。joinrel 是两个或者两个以上的 baserel 在一定约束条件下的简单合并。注意，对任何一组 baserel 仅有一个 joinrel，即对于任何给定 baserel 的集合，只有一个 RelOptInfo 结构。比如，连接 {a, b, c} 是用同一个 RelOptInfo 来表示的，无论 a、b、c 这三个 baserel 连接的顺序如何。因为对于 joinrel 来说，baserel 之间是没有顺序的。至于这三个 baserel 是以什么样的路径（连接顺序、方式）连接成 {a, b, c} 形式的连接表，则记录在 RelOptInfo 的 pathlist 字段中。

在 RelOptInfo 中，pathlist 记录了生成该 RelOptInfo 在某方面较优的路径，其中每一个节点都是一个 Path 结构（数据结构 5.16）的指针，Path 结构也称为路径。路径描述的是扫描表的不同方法（由 Path 结构的 pathtype 字段表示，比如 T_SeqScan 表示顺序扫描、T_IndexScan 表示索引扫描），以及元组排序的不同结果（由 pathkeys 字段表示）。事实上，Path 结构只是路径的一个超类或者是一个头部信息，pathlist 中的 Path 结构指针可以根据 Path 结构中的 pathtype 字段表示的路径类型转换成具体的路径节点，比如 T_IndexScan 类型的具体路径节点对应着 IndexPath 数据结构，而 IndexPath 结构的第一个字段就是 Path 结构，这种处理方式类似于 5.2.3 节中提到的把数据结构包装成 Node 的手段。在 pathlist 中，每一个 Path 节点对应的具体路径节点中则存放了构成该条路径的具体信息，

包括连接的方式、顺序，以及参加连接的 baserel 的访问方式等。

pathkeys 为 Pathkey 类型的链表，它描述一个路径包含的排序信息。顺序扫描的路径的 pathkey 是 NIL，表明没有已知的排序。如果有索引扫描的话，索引扫描路径的 pathkey 描述了所选择的索引的排序键。

数据结构 5.16 Path

```
typedef struct Path
{
    NodeTag      type;          //节点类型
    NodeTag      pathtype;      //标记扫描/连接方法
    RelOptInfo*  *parent;        //此路径所构建的关系
    Cost         startup_cost;   //执行此路径的启动代价估计
    Cost         total_cost;    //执行此路径的总代价估计
    List         *pathkeys;      //此路径的输出的排序信息,
                                //list 元素为 PathKeyItem
} Path;
```

路径是用一个树结构表示的。叶子节点是对基本关系的扫描路径，内部节点是连接路径的节点。例如，图 5-27 中的树结构是图中 SQL 语句生成的路径之一。

```
SELECT *
FROM Student,Course,SC
WHERE Student.sno = SC.sno AND Course.cno = SC.cno AND sname = '张三';
```

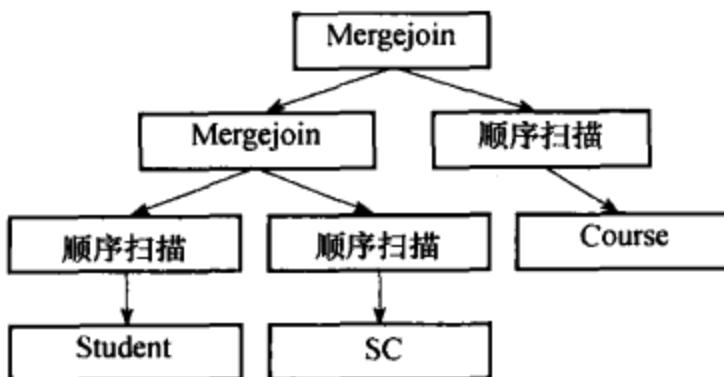


图 5-27 路径实例

在 Path 生成过程中，在某一方面代价较优的路径（启动代价、总代价等）都存放在 RelOptInfo 结构的 pathlist 链表中，如图 5-28 所示。

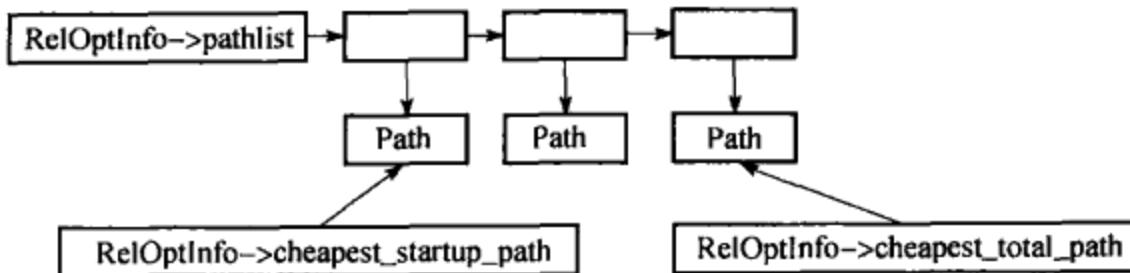


图 5-28 RelOptInfo 中的路径

1. 路径生成算法

路径代表了对一个表或者多个表中数据的访问方式。由于单个表的访问方式（顺序访问、索引访问、TID 访问）、两个表间的连接方式（嵌套循环连接、归并连接、Hash 连接）以及多个表间的连接顺序（左连接、右连接和布希连接）都有多种，因此访问一个表或多个表的路径也会有多种，每个路径都可能是上述访问方式、连接方式和连接顺序的一种组合。查询执行模块只需要执行效率最高的路径。因此在准备计划时，查询规划器需要考虑所有的路径，并从中挑选出最优的路径来生成执行计划，这个生成并挑选最优路径的工作由路径生成算法完成。PostgreSQL 中有两种路径生成算法：动态规划算法和遗传算法。

(1) 动态规划算法

在 PostgreSQL 中，通常情况下是使用动态规划来获得最优的路径的，下面是用动态规划算法实现路径生成的步骤：

1) 初始：初始状态下，为每一个待连接的 baserel 生成基本关系访问路径，选出最优路径。这些关系称为第 1 层中间关系。把所有的由 n 个关系连接生成的中间关系称为第 n 层关系。

2) 归纳：已知第 1 ~ n - 1 层的路径，用下列方法生成第 n 层的关系（n > 1）：

- 将第 n - 1 层的关系分别与每个第 1 层的关系连接，并在各自最优路径的基础上，生成使用不同连接方法的路径。
 - 若 n > 3，将第 n - 2 层的关系分别与每个第 2 层的关系连接，n - 3 层的关系分别与每个第 3 层的关系连接，并在各自最优路径的基础上，生成使用不同连接方法的路径，并依此类推。
- 3) 生成第 n 层关系后，选出每个第 n 层关系的最优路径作为结果。

例如，对 A、B、C 三个 baserel 构成的连接，用动态规划算法生成路径的过程如图 5-29 所示。

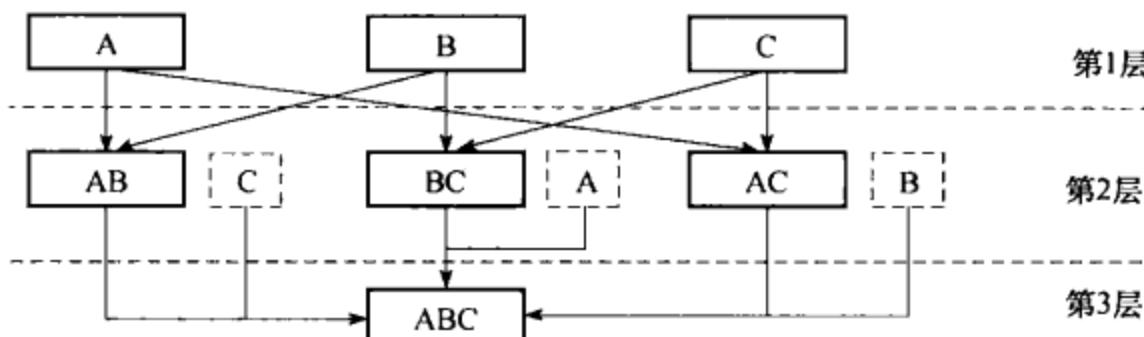


图 5-29 生成 {A, B, C} 的连接路径

在初始步骤要确定 baserel 的访问方式，在归纳步骤则要确定这些 baserel 和 joinrel 的连接顺序和连接方式。

对 baserel 的访问方式有顺序访问（Sequential Access）、索引访问（Index Access）和 TID 直接访问三种。而连接顺序有左连接、右连接和布希连接三种。

- 左连接（left-handed）是指外关系是一个 joinrel，而内关系总是一个 baserel。
- 右连接（right-handed）是指外关系总是一个 baserel，而内关系则是一个 joinrel。
- 布希连接（bushy）则是指内外关系都能自行连接，也就是说内外关系都可以是 joinrel。

例如，当构建 {1, 2, 3, 4} 时：

- 左连接的方式是先连接 {1} 和 {2}，然后连接 {1, 2} 和 {3}，最后连接 {1, 2, 3} 和 {4}。

- 右连接的方式是先连接 {2} 和 {1}，然后连接 {3} 和 {1, 2}，最后连接 {4} 和 {1, 2, 3}。
- 布希连接的方式是先分别连接 {1} 和 {2}、{3} 和 {4}，然后连接 {1, 2} 和 {3, 4}。

连接方式有嵌套循环连接 (Nest-loop join)、归并连接 (Merge Join) 和 Hash 连接 (Hash Join) 三种。

- 嵌套循环连接

这种连接方式对左边关系的每个元组都要扫描右边关系的所有元组进行连接操作。这个策略最容易实现，但是可能会很耗费时间。如果右边的关系可以用索引扫描，那么这个策略可能就是个好策略，我们可以用来自左边关系的当前元组的连接属性值为键值对右边关系做索引扫描。图 5-30 是一个嵌套循环连接的例子。

- 归并连接

归并连接在连接开始之前，会将每个关系都按连接属性进行排序。然后对两个关系做并行扫描，匹配的元组就组合起来形成连接元组。这种组合效率更高，因为每个关系只扫描一次。其中，排序步骤可以通过指定一个排序操作来完成，或者是使用连接属性上的索引按照恰当的顺序扫描关系。图 5-31 是一个归并连接的例子。

- Hash 连接

在 Hash 连接时，首先扫描右边的关系，并把扫描到的元组用连接属性值作为 Hash 键装载进一个 Hash 表，然后扫描左边的关系，并将找到的每个元组的连接属性值用作 Hash 键字来定位 Hash 表里匹配的元组。图 5-32 是一个 Hash 连接的例子。

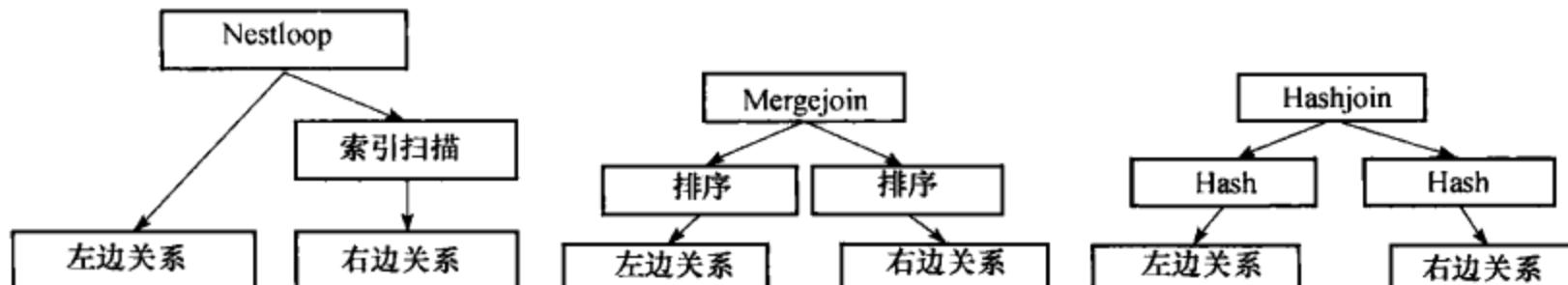


图 5-30 嵌套循环连接

图 5-31 归并连接

图 5-32 Hash 连接

连接的顺序不同会导致中间关系的大小不同，间接导致执行时的磁盘 I/O 时间和 CPU 时间不同；连接的实现方法不同，直接导致磁盘 I/O 时间和 CPU 时间不同。选择路径时，我们需要同时考虑这两个因素，为不同的顺序和不同的连接方法形成的组合分别生成路径。而计算某一路径的代价时，将根据底层关系的大小和连接的实现方法，计算出总的代价，从而来评判路径的优劣。代价评估将在 5.5 节具体介绍。

因此，在路径生成过程中，每生成一个中间关系，我们要估算出它的大小、路径及其代价，供上层路径计算代价。但在具体实现中，我们不能精确判断出中间关系的路径中哪一个是最优的，所以要尽量保留可能的最优路径，用于更高层的路径生成。保留时主要考虑三个指标：路径的启动代价、路径总的执行代价和路径的输出排序键。只要一个路径在一个指标上最优就保留它，因此在每个 RelOptInfo 中都有 cheapest_startup_path、cheapest_total_path 和 cheapest_unique_path 三个字段，分别指向启动代价最优路径、总代价最优路径和最优路径。

(2) 遗传算法

动态规划算法是生成路径的默认算法，但在有些情况下，检查一个查询所有可能的路径会花去很多的时间和内存空间，特别是所要执行的查询涉及大量的关系的时候。在这种情况下，为了在合理的时间里判断一个合理的执行计划，PostgreSQL 将使用遗传算法生成路径。是否启用遗传算法由两个因素决定：在系统配置中是否允许使用遗传算法以及需要连接的基本关系是否超过使用遗传算法的阈值。前者由全局变量 `enable_geqo` 控制，其值来自于配置文件中的 `geqo` 配置项，默认情况下是允许；后者由全局变量 `geqo_threshold` 控制，其值来自于配置文件中的 `geqo_threshold` 配置选项，默认值是 12，即参加连接的基本关系数大于或等于 12 时就会采用遗传算法来生成路径。

遗传算法（GA）是一种启发式的优化方法（heuristic optimization method），它通过既定的随机搜索进行操作。优化问题的可能解的集合被认为是个体（individuals）组成的人群（population）。一个个体对它的环境的适应程度由它的适应值（fitness）表示。一个个体在搜索空间里的参照物是用染色体（chromosomes）表示的，实际上这是一套字符串。一个基因（gene）是染色体的一个片段，基因是被优化的单个参数的编码。对一个基因的典型的编码可以是二进制（binary）或整数（integer）。通过仿真进化过程的重组（recombination）、变异（mutation）和选择（selection）找到新一代的搜索点，它们的平均适应值要比它们的祖先好。

在 PostgreSQL 中，遗传算法将路径作为个体，将个体以某种方式编码（个体体现了连接顺序），然后通过重组得到后代，考虑连接代价来计算后代的适应值，再选择合适的后代进行下一次迭代。当到达一定的迭代次数之后，遗传算法终止。选择遗传算法可以减少生成路径的时间，但是遗传算法并不一定能找到“最好”的规划，它只能找到相对较优的路径。PostgreSQL 中的遗传算法将在 5.6 节中详细介绍。

2. 路径生成总体流程

生成路径的工作由函数 `make_one_rel` 实现，该函数只有一个类型为 `PlannerInfo` 的参数 `root`。函数 `make_one_rel` 的工作步骤可分为两个阶段：

- 1) 调用函数 `set_base_rel_pathlists` 生成基本关系的访问路径（即为每个基本关系生成一个 `RelOptInfo` 结构并生成路径），放在其 `RelOptInfo` 的 `pathlist` 字段中，对应动态规划算法中的第 1 层。
- 2) 调用函数 `make_rel_from_joinlist` 生成最终路径。返回一个代表连接所有基本关系的 `RelOptInfo` 结构，它的 `pathlist` 字段就是所有的最终路径组成的链表，对应第 2 层到第 n 层。

3. 生成基本关系访问路径

函数 `set_base_rel_pathlist` 负责为基本关系生成访问路径。它将扫描范围表中的所有基本关系，并考虑顺序扫描和所有可用的索引扫描，最后把每个“值得保留”的路径加入到基本关系的 `RelOptInfo` 结构的 `pathlist` 中。要强调的是，为了减小搜索空间，对每个基本关系或者连接关系，都只在其 `RelOptInfo` 中保留“值得保留”的路径。所谓“值得保留”的路径应满足以下条件之一：

- 总代价最优。
- 启动代价最优。
- 有最好的排序结果。

如果有两个路径的总代价相同且排序效果不同，则两者都会被保留。

如果一个基本关系是子关系（RTE 的 `inh` 字段为真），则调用 `set_append_rel_pathlist` 为其建立

访问路径，否则根据其 RTE 的 rtekind 字段分以下四种情况进行处理：

- RTE_SUBQUERY：表示当前的基本关系是一个子查询，调用函数 set_subquery_pathlist 为子查询建立访问路径。
- RTE_FUNCTION：表示当前的基本关系是一个函数，调用函数 set_function_pathlist 为之建立访问路径。
- RTE_VALUES：表示当前的基本关系是一个 VALUES 列表，调用函数 set_values_pathlist 为 VALUES 项建立访问路径。对于 INSERT 型语句，如果有多个 VALUES 项，系统会为之生成一个单独的 RTE^①。
- RTE_CTE：表示当前的基本关系是一个公共表表达式，如果它是递归的（self_reference 为真），调用函数 set_worktable_pathlist 为其建立访问路径；否则调用函数 set_cte_pathlist 为其生成路径。
- RTE_RELATION：表示当前的基本关系是一个普通表，调用 set_plain_rel_pathlist 为其生成路径。

上面提到的各个函数的作用如下。

(1) set_append_rel_pathlist 函数

该函数为子关系建立访问路径。由于子关系的元组的获得必须要依赖于其父关系的元组，因此要将其父关系的访问路径作为一种 AppendPath（附加路径）加入到子关系的路径节点下层。该函数的流程如下：

1) 寻找子关系的父关系：在此之前的处理步骤中已经将所有需要附加的父关系放在了 Planner-Info 结构的 append_rel_list 中，并且每一个父关系的数据结构中都有一个字段指向了查询中与之相关的子关系，只需对 append_rel_list 做一个简单的循环查找就可以找到父关系。

2) 为父关系对应的基本关系生成 RelOptInfo，并在其中生成父关系的路径。

3) 调用 create_append_path 函数将父关系的总代价最优路径包装成一个 AppendPath 并将之作为子关系的路径加入到其 RelOptInfo 中，加入操作由函数 add_path 完成。

4) 为子关系选择最优路径。

(2) set_subquery_pathlist 函数

该函数为子查询的 RelOptInfo 建立访问路径。如果子查询被附加了约束条件（baserestrictinfo 字段中指定），则首先将它们下推到到子查询本身的 WHERE 或者 HAVING 语句中，采用这种方法可以得到更好的子计划。该函数的流程如下：

1) 检查 baserestrictinfo 中的每一个约束条件，如果通过 subquery_is_pushdown_safe 检查该约束条件是可以下推的，则下推该条件到子查询中。

2) 调用函数 subquery_planner 生成此子查询的子计划，并放入子查询的 RelOptInfo 结构的 subplan 字段中。

3) 调用 create_subqueryscan_path 生成类型为 T_SubqueryScan 的路径，其执行代价与 RelOptInfo 结构 subplan 字段中记录的子计划相同。

4) 将上一步生成的路径加到子查询的 pathlist 中，即此子查询只有唯一的路径。

^① RTE，即 Range Table Entry，指范围表中一个条目。

(3) set_function_pathlist 函数

该函数为函数类型的 RTE 建立访问路径。它先调用函数 set_function_size_estimates 用评估的输出行、宽度来标记该 RTE 的 RelOptInfo 结构，然后调用函数 add_path 将 create_functionscan_path 函数产生的 T_FunctionScan 路径加入到 pathlist，最后调用函数 set_cheapest 选择代价最小的路径。

(4) set_values_pathlist 函数

该函数为 VALUES 列表生成路径，其过程和 set_function_pathlist 函数类似，只是在这里要调用 create_valuesscan_path 来生成一个 T_ValuesScan 路径。

(5) set_worktable_pathlist 函数和 set_cte_pathlist 函数

当一个 CTE 是递归的，则要单独为它生成一个 T_WorkTableScan 路径以配合 RecursiveUnion 节点的扫描（见 6.5.1 小节和 6.5.2 小节），set_worktable_pathlist 将首先找到 CTE 对应的非递归计划，然后调用 set_cte_size_estimates 根据该计划的信息来设置 CTE 的 RelOptInfo，之后调用 create_worktablescan_path 创建一个 T_WorkTableScan 路径并加入到 pathlist 中。

如果 CTE 是非递归的，则可以把每一个 CTE 表达式作为一个独立存在“虚表”，该函数对其中每一个 CTE 表达式检查是否存在计划之后，设置代价信息并调用 create_ctescan_path 创建一个 T_CteScan 路径并加入到 pathlist 中。

(6) set_plain_rel_pathlist 函数

该函数用于生成普通关系（非子查询，且无继承关系）的路径。其流程如下：

- 1) 函数 set_baserel_size_estimates 设定给定基本关系对应的 RelOptInfo 结构中 rows、width、baserestrictcost 三个字段，即此基本关系将输出的元组数、元组的平均大小和 baserestrictinfo 的代价。

- 2) 函数 check_partial_indexes 检查基本关系上是否有局部索引，且对这个查询有效，将可用的部分索引的 predOK 字段置为真。

- 3) 函数 create_or_index_quals 判断是否可以将 AND 子句中的 OR 子句外提，如果可以，将外提的约束条件加入这个基本关系对应的 RelOptInfo 结构中的 baserestrictinfo 链表中，并重新估算基本关系对应的 RelOptInfo 结构中的 rows、width、baserestrictcost 三个字段，并且重新检查是否有部分索引可用。

- 4) 调用 create_seqscan_path 生成一个顺序扫描路径，并调用 add_path 加入到 pathlist 中。

- 5) 调用 create_index_paths 检查基本关系上是否有索引可用，为可用的索引建立索引扫描路径并尝试用 add_path 加入到 pathlist；同样调用 create_tidscan_paths 尝试生成 TID 扫描路径。注意，这里生成的索引扫描路径和 TID 扫描路径不一定会保留在 pathlist 中，除非它们在代价或者排序上具有优势；生成索引扫描路径和 TID 扫描路径的过程将分别在后面介绍。

- 6) 调用函数 set_cheapest 设置最小启动代价和最小总代价。

3. 生成索引扫描路径

当关系上涉及索引扫描时，要生成相应的索引路径。生成索引路径的函数有 create_index_paths 和 find_usable_index。

(1) create_index_paths 函数

该函数为参数 rel 中指定的 RelOptInfo 结构添加可能的索引路径。一个索引必须满足以下条件之一才会被考虑用于生成索引扫描路径：

- 能匹配一个或者多个约束条件（来自于 baserestrictinfo 字段中的 RestrictInfo 结构）。
- 能匹配查询中的连接条件。

- 能匹配查询的排序要求。
- 能匹配查询的条件 (WHERE 子句中除了连接条件的其他条件)。

函数 `find_usable_indexes` 用来寻找基本关系上满足上述要求的索引，并为找到的索引生成索引扫描路径，将生成的索引扫描路径链接成一个 `List` 返回。然后用 `add_path` 尝试将具有 `amgettuple` 方法的索引扫描路径添加到 `rel` 的 `pathlist` 中。有些索引虽然不支持 `amgettuple` 方法，但是能支持 `amgetitem` 方法，则将它们保留，暂时记录在 `bitindexpaths` 链表中。既不支持 `amgettuple` 方法也不支持 `amgetitem` 方法的索引扫描路径将被抛弃。

接下来会调用函数 `generate_bitmap_or_paths` 为 `baserestrictinfo` 中的 OR 子句生成 `BitmapOrPath` 路径，并将它们加入到 `bitindexpaths` 链表中。最后，先调用函数 `create_bitmap_heap_path` 将 `bitindexpaths` 中所有的位图路径放在一个 `BitmapAndPath` 路径中，再调用 `create_bitmap_heap_path` 函数在 `BitmapAndPath` 路径之上生成一个 `BitmapHeapPath` 路径，然后尝试将之加入到 `pathlist`。

(2) `find_usable_indexes` 函数

该函数的主要功能是对一个基本关系，给定一系列约束条件，查找所有可能有用的索引并为之创建索引扫描路径。它在与索引有关的路径的生成过程中都会使用，其主要流程如下：

1) 对于每一个索引调用 `group_clauses_by_indexkey` 找到所有适合的约束条件，将这些约束条件放在 `restrictclauses` 变量中。

2) 调用 `build_index_pathkeys` 为该属性计算描述其排序信息的 `pathkeys`，然后调用 `truncate_useless_pathkeys` 检查这些 `pathkeys` 对该查询的有效性，只保留有效的 `pathkeys`。函数 `truncate_useless_pathkeys` 会返回一个 `pathkey` 组成的链表，该链表满足下列条件之一：①链表中每个 `pathkey` 都出现在涉及该基本关系的等值连接子句中；②该链表的前 m 个 `pathkey` 与 `PlannerInfo` 中 `query_pathkeys` 字段指向的链表相同。

3) 为被 `truncate_useless_pathkeys` 保留下来的每个索引产生一个向前扫描的索引扫描路径（调用 `create_index_path`）。

4) 若 `index` 是有序的，则以步骤 2 和步骤 3 类似的方法创建向后扫描的索引扫描路径。

5) 将所有生成的索引扫描路径组织成一个链表返回。

(3) `cost_index` 函数

找出可用的索引并生成相应的索引路径之后，就需要评价路径优劣，其标准就是该路径的代价。评估路径的代价这个功能在 `Costsize.c` 中实现。每一类路径都有各自的评估函数，下面仅介绍索引路径的代价评估，其他类型路径的评估读者可自行分析。

函数 `cost_index` 用于评估索引路径代价。其中，在估算 I/O 代价（即取出的页面数）时，使用了函数 `index_pages_fetched`，该函数的实现原理是 Mackert 和 Lohman 提出的估计公式。在该公式中，根据实际应取出的索引基表的元组数，便可估计取出这些元组花费的 I/O 次数。当索引的排序与基表的排序无关时，根据 Mackert 和 Lohman 提出的估计公式，所取的页面数（`PFMax`）可以按如下方式计算：

$$PFMax = \begin{cases} \min\left(\frac{2 \times T \times N \times S}{2 \times T + N \times S}, T\right), & T \leq b \\ \frac{2 \times T \times N \times S}{2 \times T + N \times S}, & T > b \text{ 且 } N \times S \leq \frac{2 \times T \times b}{2 \times T - b} \\ b + \left(N \times S - \frac{2 \times T \times b}{2 \times T - b}\right) \times \frac{T - b}{T}, & T > b \text{ 且 } N \times S > \frac{2 \times T \times b}{2 \times T - b} \end{cases} \quad (\text{公式 5.1})$$

其中，T 表示基表的页面数；N 表示基表的元组数；s 表示选择度（selectivity），即所取元组在元组总数中占的比例；b 表示系统可用的缓冲区（Buffer）数。

函数 cost_index 的主要流程如下：

- 1) 调用索引的 amcostestimate 方法评估索引本身的启动代价、总代价、选择度和相关度。
- 2) 用变量 startup_cost 和 run_cost 分别表示路径的启动代价和运行代价（两者之和表示总代价）。这里先设置 startup_cost 为索引的启动代价，而 run_cost 设置为索引的运行代价（索引总代价减去索引启动代价）。
- 3) 估算基表要取出的元组数，记录在变量 tuples_fetched 中。
- 4) 如果这个索引路径将用于连接（即存在外关系 outer_rel），估算在整个连接中索引路径需要取出的页面数，然后分别估算最大和最小 I/O 代价。估算最大 I/O 时使用 $N = \text{tuples_fetched} * \text{outer_relrow}$ ；估算最小 I/O 时使用 $N = \text{tuples_fetched} * \text{outer_relrow} * \text{索引选择度}$ 。
- 5) 如果是不直接用于连接的索引路径，估算其最大和最小 I/O 代价。
- 6) 用 max_IO_cost 和 min_IO_cost 来表示前面计算的最大和最小 I/O 代价，然后根据索引的相关度计算总的 I/O 代价并加到 run_cost 中。
- 7) 将约束条件的启动代价加入到 startup_cost 中，然后将全局变量 cpu_tuple_cost（表示一个元组需要的平均 CPU 代价）加上约束条件对每个元组产生的代价得到 cpu_per_tuple（表示该索引路径每个元组的实际 CPU 代价），用 cpu_per_tuple 乘以 tuples_fetched 得到该索引路径的总 CPU 代价，并将之加到上一步得到的 run_cost 中形成最终的运行代价。
- 8) 根据最后的 startup_cost 和 run_cost 的值设置索引路径的启动代价和总代价。

4. 生成 TID 扫描路径

在 PostgreSQL 中，TID 用来表示元组的物理地址，一个 TID 从逻辑上可以看成是一个三元组（表 OID，块号，元组的块内偏移量）。由于 TID 可以让执行器快速定位到磁盘上的元组，因此通过 TID 来访问元组非常高效。在 PostgreSQL 的查询中也允许在 WHERE 子句中使用 TID 作为查询条件，对于这种情况，查询规划器将为之生成 TID 扫描路径。TID 扫描路径生成由函数 creat_tidscan_paths 实现，该函数分为两步：

- 1) 从表的 RelOptInfo 结构的 baserestrictinfo 中提取 TID 条件。
- 2) 根据 TID 条件调用 create_tidscan_path 生成 TID 扫描路径，并尝试用 add_path 加入到 pathlist 中。

函数 create_tidscan_path 先创建一个 T_TidScan 类型的路径节点，然后将提取到的 TID 条件放在该路径节点的 tidquals 字段中，其重点工作在于创建完路径节点后估算其代价，这个工作由 cost_tidscan 完成，其流程如下：

- 1) 首先估算要取多少元组。
- 2) 估算 I/O 代价，这里实际估算的是 I/O 代价的最大值，也就是认为要取出的页面数就是要取出的元组数（一个页面中只命中一个元组，这是最坏的情况），然后将 I/O 代价加入到运行代价中。
- 3) 估算 CPU 代价。CPU 代价包括要取出的元组的代价和基本关系上约束信息的代价，并加到运行代价中。
- 4) 为 TID 扫描路径设置启动代价和总代价。

5. 生成最终路径

所有基本关系的路径都创建好之后，就可以把它们连接起来形成最终的连接关系，最终路径的生成由函数 `make_rel_from_joinlist` 实现，其函数原型为：

```
RelOptInfo*make_rel_from_joinlist (PlannerInfo*root, List*joinlist)
```

其中，参数 `joinlist` 链表中每一个节点都代表一个基本关系，该基本关系可能表示一个范围表或者是一个子查询。该函数会将 `joinlist` 中的基本关系连接起来生成最终的连接关系，并且为最终的连接关系建立 `RelOptInfo` 结构，其 `pathlist` 字段就是最终路径组成的链表。该函数的实质就是选择不同的连接方式作为中间节点，选择不同的连接顺序将代表基本关系的叶子节点连接成一棵树，使其代价最小。这是一个递归的生成二叉树的过程，使用的是 PostgreSQL 中标准的递归处理二叉树的方法。`make_rel_from_joinlist` 函数的流程如下：

- 1) 由于 `joinlist` 中的基本关系可能是一个子查询，因此首先需要递归调用 `make_rel_from_joinlist` 为子查询生成路径，然后将子查询的 `RelOptInfo` 和其他基本关系的 `RelOptInfo` 一起作为动态规划算法第一层关系。动态规划算法所需层数与 `joinlist` 链表的长度相同。将所有第一层关系的 `RelOptInfo` 结构放入链表 `initial_rels` 中。
- 2) 如果算法所需层数等于 1，表示 `FROM` 子句中只有一项，则可以将 `initial_rels` 中第一个节点对应的 `RelOptInfo` 作为最终的连接关系路径返回。
- 3) 如果算法所需层数大于 1，那么需要考虑关系的不同连接顺序。这种情况下，连接路径的生成有三种解决方案：一种是当启用了遗传算法并且动态规划算法所需层数超过了遗传算法的触发值时，将利用遗传算法完成路径生成；另一种方案就是使用动态规划算法进行路径生成；第三种方案是使用全局变量 `join_search_hook` 定义的“第三方”函数来进行路径生成。但在 PostgreSQL 8.4.1 版本中，`join_search_hook` 设置为空，该变量是为了以后便于扩展新的生成路径的方法。简而言之，如果待连接的基本关系太多则会使用遗传算法，否则使用动态规划算法。

遗传算法将在 5.6 节单独介绍，这里重点介绍动态规划算法。动态规划算法由函数 `standard_join_search` 实现，其原型为：

```
RelOptInfo*standard_join_search (PlannerInfo*root, int levels_needed, List*initial_rels)
```

该函数从第一层的关系（参数 `initial_rels`）开始，基于前面层次的结果依次生成上层的连接关系，最后返回由所有基本关系连接而成的关系。该函数首先为 `root` 的 `join_rel_level` 字段分配内存，使之成为一个长度为 `levels_needed + 1` 的链表数组。除第一个元素外，数组中的每一个元素都指向一个链表，该链表用于存放相应层次的关系，层号等于该数组元素的下标（C 语言数组下标从 0 开始，因此第一个数组元素并未使用）。另外，此数组中的所有关系将被链接在一个链表中，放入 `root` 的 `join_rel_list` 字段（该链表上建有 Hash 表用于查找）。显然，`initial_rels` 中的关系将被放到数组的第二个链表中，代表第一层关系，即算法的初始状态。然后，调用函数 `join_search_one_level` 依次生成第 2 到第 `levels_needed` 层的关系，每次调用该函数只能生成一层的关系。函数 `join_search_one_level` 的流程如下（假设要生成第 n 层的关系）：

- 1) 先把第 $n - 1$ 层的关系与第 1 层的关系两两组合连接起来。对每个第 $n - 1$ 层的关系：
 - ①如果它没有任何连接子句，则调用函数 `make_rels_by_clauseless_joins` 将它与每个未和它连接的第 1 层关系做笛卡儿积，从而生成深度较低的路径树。

②如果它有连接子句，调用 `make_rels_by_clause_joins` 将它与其连接子句中的第 1 层关系连接起来。其中，`make_rels_by_clause_joins` 中会调用 `make_join_rel` 来生成实际的连接路径。

2) 以同样的方法考虑其他层次之间的连接（第 2 层与第 $n - 2$ 层、第 3 层与第 $n - 3$ 层等），将层次之间通过连接子句指定的关系调用 `make_join_rel` 连接起来。

3) 如果所有可用的连接子句都已处理完，则这里将强制生成一些笛卡儿积。这是为了处理一种特殊情况：关系上有连接子句，但目前还不能使用这些子句。例如，有一个查询：“`SELECT * FROM a, b, c WHERE (a.f1 + b.f2 + c.f3) = 0;`”，连接子句必须在第三层才能被应用，那么在第二层我们只能进行笛卡儿积。此步骤保证了关系生成的顺利进行，在这里我们只考虑第 $n - 1$ 层与第 1 层的乘积，通过对每个第 $n - 1$ 层的关系调用 `make_rels_by_clauseless_joins` 实现其与第 1 层关系的笛卡儿积。

其中：

- 函数 `make_rels_by_clause_joins` 调用 `make_join_rel` 连接关系，连接时要求两个关系中的基本关系不存在交集，且存在连接子句时连接。
- 函数 `make_rels_by_clauseless_joins` 调用 `make_join_rel` 连接关系，连接时要求两个关系中的基本关系不存在交集，且不存在连接子句时连接。
- 函数 `make_join_rel` 把给定的两个关系连接起来，生成连接后的关系的 `RelOptInfo` 结构，若此前该结构已经生成，则只在其 `pathlist` 域中增加利用两个关系的路径生成的新的路径，并返回此 `RelOptInfo` 结构。

函数 `make_join_rel` 的主要功能是生成连接路径，该函数的主体工作是调用函数 `add_paths_to_joinrel` 实现的。为了处理三种不同的连接方式，`add_paths_to_joinrel` 将分别调用 `sort_inner_and_outer`、`match_unsorted_outer`、`hash_inner_and_outer` 这三个函数。`add_paths_to_joinrel` 的流程如下：

1) 从连接约束条件（`restrictlist` 参数）中，选出所有可用于生成 `MergeJoin` 路径的条件（`mergeclause_list`）。

2) 调用 `sort_inner_and_outer` 生成需要先将左右关系排序，然后进行归并连接的路径。

3) 调用 `match_unsorted_outer` 生成不需对左关系进行重新排序的路径，包括 `Merge Join` 路径和 `Nest-loop Join` 路径。

4) 调用 `hash_inner_and_outer` 生成需要先将左右关系进行 Hash 后，再形成 Hash 连接的路径。

下面详细介绍其中各个生成连接路径的步骤：

(1) `sort_inner_and_outer` 函数

该函数会为待连接的关系（左右关系）生成按某给定的 `pathkeys` 排序，然后进行归并连接的路径。我们先将 `mergeclause_list` 中的每个 `merge` 子句所生成的 `pathkey`（键值），放入链表 `all_pathkeys`。分别将此链表的每个 `pathkey` 生成按此 `pathkeys` 排序的 `MergeJoin` 路径。最后调用 `create_mergejoin_path` 生成路径节点。此时，若左右的关系已经按期望的方式排序，则将 `outersortkeys` 或 `innersortkeys` 置空。

(2) `match_unsorted_outer` 函数

该函数生成不需对左关系排序的路径，包括嵌套循环连接路径和归并连接路径。该函数对于左关系的 `RelOptInfo` 结构中记录的每个路径按下列步骤进行（假设被处理路径的 `pathkeys` 为 `p`）处理：

1) 生成三种嵌套循环连接路径：分别选用右关系的最小总代价路径、最小启动代价路径以及

最优的作为内关系时的索引扫描路径（如果存在，在右关系的内连接中）进行连接。注意，生成的路径的 pathkeys 与该左关系的路径相同。

- 2) 选用右关系的最小总代价路径，并对其按 p 重排序，然后生成归并连接路径。
- 3) 选出以 p 的前 m ($m = 1, 2, \dots, \text{length}(p)$) 个排序键排序且总代价最小的右关系路径，找出其中总代价最小和启动代价最小的右关系路径，然后和左关系路径生成归并连接路径。
- 4) 以上所有路径生成之后将立刻通过 add_path 函数尝试加入到 pathlist 中，该函数会对路径的代价进行比较，保留较优的路径。

(3) hash_inner_and_outer 函数

该函数生成 hash 连接路径。该函数对每个可以作为 Hash 连接的约束子句分别利用左关系的最小总代价路径和最小启动代价路径与右关系的最小总代价路径连接，生成 Hash 路径。通过 add_path 把这些路径尝试添加到 pathlist。

以上三种连接方式分别有自己的代价评估函数，但是原理相似，有兴趣的读者可以自行分析比较。

至此，一个最终连接关系的路径树就建立起来了，根据代价评估的结果以及是否有 Hash 分组和排序操作，选择最优路径传给计划生成器。

5.4.4 生成可优化的 MIN/MAX 聚集计划

完成路径生成之后，查询规划器就可以开始进行计划的生成。首先，规划器会处理一种比较特殊的查询：查询中含有 MIN/MAX 聚集函数，并且聚集函数使用的属性上建有索引或者属性恰好是 ORDER BY 子句中指定的属性。在这种特殊情况下，可以直接从索引或者已排序好的元组集中取到含有最大值或最小值的元组，从而避免了扫描全表带来的开销。规划器会先检查一个查询是否可以优化到不对全表扫描而直接读取元组，如果可以则生成可优化的 MIN/MAX 聚集计划，否则（需要对全表扫描）生成普通计划。

生成可优化的 MIN/MAX 聚集计划的主函数是 optimize_minmax_aggregates，该函数从一个选定的路径生成一个计划，如果该路径对应的查询满足可优化的 MIN/MAX 聚集计划的条件，该函数会返回一个计划；否则该函数返回空值。如果该函数能够生成一个非空的计划，则后续生成普通计划的步骤就不再进行，将这个计划进行完善后作为最终的计划。该函数的主要步骤如下：

- 1) 通过查询树（Query 结构体）中的字段 hasAggs 判断该查询（或子查询）中是否存在聚集函数，如果存在则进行第 2 步；否则退出。
- 2) 通过查询树中的字段 groupClause 和 hasWindowFuncs 分别判断查询（或子查询）中是否存在分组（GROUP BY）或窗口函数，如果不存在则进行第 3 步；否则退出。因为如果存在分组或窗口函数，则必须对表中所有元组进行扫描，不可能优化到直接获取单个元组。
- 3) 通过检查该查询（或子查询）中范围表的个数来判断是否为单表查询，如果存在 2 个或 2 个以上的范围表，则退出；否则，通过函数 planner_rt_fetch 读取该范围表，并进行第 4 步。
- 4) 检查该范围表是否为普通关系类型（基本表）以及是否存在继承关系，如果是非继承的普通关系则通过函数 find_base_rel 访问该关系，并进行第 5 步；否则退出。
- 5) 通过函数 find_minmax_aggs_walker 查找目标属性和 HAVING 子句中出现的所有聚集函数，

如果都是 MIN/MAX 聚集函数，进行第 6 步；否则退出。函数 `find_minmax_aggs_walker` 递归扫描目标属性或 HAVING 子句中的聚集节点（用数据结构 `Aggref` 表示的节点），检查是否都为 MIN/MAX 型聚集节点。如果是，则将所有 MIN/MAX 型聚集节点保存到变量 `aggs_list` 中，并返回 FALSE；如果在查找过程中发现了非 MIN/MAX 型的聚集节点，则返回 TRUE，说明存在其他类型的聚集函数，此种情况下，直接退出，不会进行优化处理。

6) 对于变量 `aggs_list` 中保存的每个 MIN/MAX 聚集函数，通过函数 `build_minmax_path` 查找可用哪个索引对其进行优化，并计算其优化后的代价（利用索引访问一条元组的代价，不包括对目标属性的代价评估）。如果都可被优化，进行第 7 步；否则退出。函数 `build_minmax_path` 对于给定的 MIN/MAX 聚集节点 `info`，在关系 `rel` 中试图找到一个可以使之优化的索引，并创建最优的索引路径。其处理过程如下：

① 判断聚集与索引属性是否匹配，并且确定了索引扫描方向。主要是通过函数 `match_agg_to_index_col` 来判断是否匹配，如果聚集的排序操作符与索引中前向扫描的操作符类（operator class）匹配则返回 `ForwardScanDirection`，进行前向扫描；未匹配，则判断聚集的排序操作符与索引中后向扫描的操作符类是否匹配，如果匹配则返回 `BackwardScanDirection`；如果不匹配则返回 `NoMovementScanDirection`，表明不进行扫描。

② 提取约束信息（即按照怎样的约束条件通过索引获取元组）。

③ 创建索引访问路径。通过调用函数 `create_index_path` 来实现，具体实现可参考前面创建索引路径的章节。

④ 进行代价评估，并选取所有索引扫描路径中代价最小的路径。

7) 通过函数 `cost_agg` 评估不优化 MIN/MAX 聚集函数时的代价。并将优化前后的代价进行比较，如果优化后的代价比优化前的小，则进行第 8 步；否则，退出。

8) 通过函数 `make_agg_subplan` 为每个 MIN/MAX 聚集函数生成优化后的子计划。在此过程中，会将子计划转换成初始计划 `initPlan` 保存到结构体 `MinMaxAggInfo`（该结构体保存了可优化的 MIN/MAX 聚集节点）的成员变量 `param` 中，作为子计划的输出参数。

9) 通过函数 `replace_aggs_with_params_mutator` 将目标列和 HAVING 子句中的 MIN/MAX 聚集函数调用替换成相应子计划的输出参数（即初始计划）。

10) 通过函数 `mutate_eclss_expressions` 将等值类（即 equivalence class，包括目标列和等值表达式）中的 MIN/MAX 聚集函数替换成相应子计划的输出参数（即初始计划）。

11) 通过函数 `make_result` 生成最终的输出计划节点。

12) 通过函数 `cost_qual_eval` 对目标属性进行代价评估，并更新计划节点中的代价评估值。

5.4.5 生成普通计划

如果 `optimize_minmax_aggregates` 返回的是空值，则需要继续生成普通计划。生成普通计划的人口函数是 `create_plan`，该函数为最优路径创建计划，依据其路径节点类型的不同，分别调用不同的函数生成相应的计划。

普通计划分为以下几类：

- 扫描计划：主要有顺序扫描、索引扫描等计划类型。
- 连接计划：主要有嵌套循环连接、Hash 连接、归并连接等计划类型。

- 其他计划：如 Append 计划、Result 计划、物化计划等。

其中，顺序扫描和嵌套循环连接计划属于简单的计划类型，我们通过分析这两种简单计划的创建过程来说明在生成计划过程中做了哪些主要工作。

1. 顺序扫描计划

函数 `create_seqscan_plan` 用于生成顺序扫描计划（SeqScan），该函数返回一个类型为 SeqScan 的结构（数据结构 5.17），其参数如下：

- `root`: 类型为 `PlannerInfo` 指针，指向当前规划器状态信息。
- `best_path`: 类型为 `Path` 指针，指向最优路径。
- `tlist`: 类型为 `List` 指针，指向目标属性链表，链表中每一个节点都是一个 `Var` 结构。
- `scan_clauses`: 为 `List` 指针，指向扫描约束信息。

数据结构 5.17 SeqScan

```
typedef struct Scan
{
    Plan      plan;           // 该顺序扫描节点对应的计划信息
    Index     scanrelid;      // 要扫描的表在范围表中的索引号
} Scan;
typedef Scan SeqScan;
```

其中，`tlist` 和 `scan_clauses` 均从 `best_path` 的 `parent` 字段中得到，将其提取出来方便统一处理。如图 5-33 所示，生成顺序扫描计划主要完成以下工作：

1) 对扫描约束信息按执行器最佳执行顺序排序。此处没有考虑选择度，仅仅依据代价估计排序。例如，表 A 有 a, b, c 三个属性，其扫描约束条件为 $a = b \text{ AND } c = \text{Const}$ 。依据执行代价，在做相等比较之前，条件 $a = b$ 需要读两次数据，而条件 $c = \text{Const}$ 只需读一次数据，那么就应先执行后者。

2) 对约束信息的排序是完整的 `RestrictInfo` 类型。当对约束信息排序完成后，需要获得实际的约束子句。

3) 创建顺序扫描计划，将目标属性、扫描表、约束子句赋值给 `SeqScan` 计划节点，并将其左右子树置为空（扫描计划节点均以叶子形式存在）。

4) 复制代价估计信息，该部分是对 `Path` 中代价估计的一个拷贝。

对于基本表的扫描还有另外一种重要的形式——索引扫描。其生成计划的过程与上类似，即从 `indexPath` 获得相应索引信息、扫描约束信息并排序，最终赋值给 `IndexPlan` 计划节点，将左右子树置空。

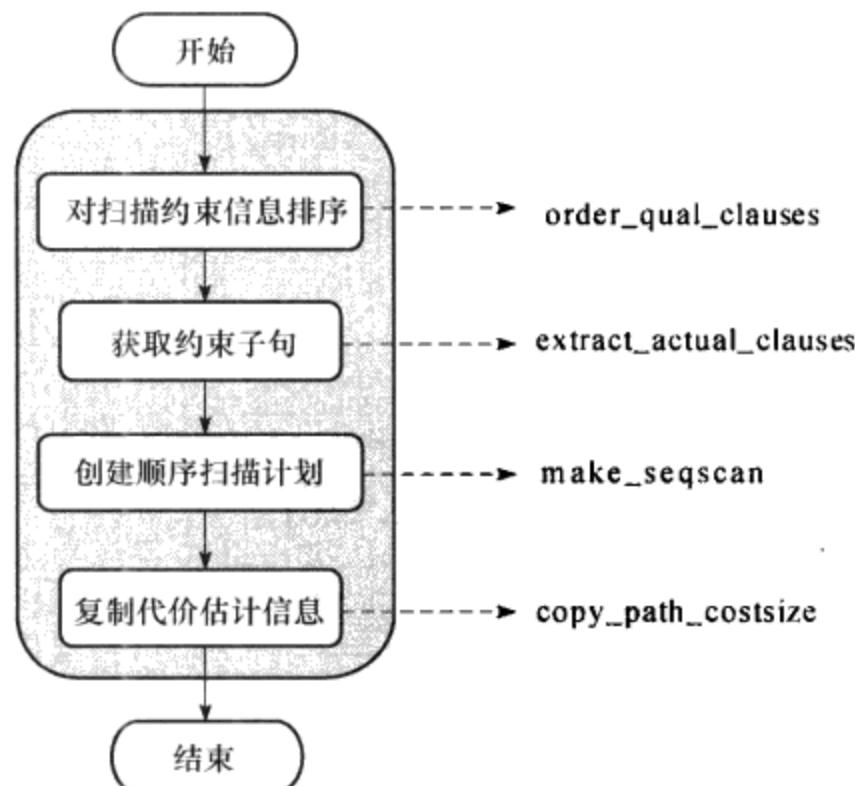


图 5-33 顺序扫描计划的生成过程

2. 嵌套循环连接计划

嵌套循环连接是三种连接方式中最简单的，嵌套循环连接计划由函数 `create_nestloop_plan` 生成，该函数除了具有和 `create_seqscan_plan` 相同的两个参数 `root` 和 `best_path` 之外，还有两个特有的参数 `outer_plan` 和 `inner_plan`，分别表示外关系和内关系对应的计划（已经由外部调用者创建好）。该函数将创建一个嵌套循环连接计划节点（数据结构 5.18），并连接其子计划树（`outer_plan` 和 `inner_plan`）到该计划节点，最终形成嵌套循环连接计划并返回。

如图 5-34 所示，创建嵌套循环连接计划的函数流程如下：

- 1) 从 `best_path` 中获取目标属性及连接约束信息。
- 2) 删除冗余的连接信息。如果内关系路径为嵌套循环索引扫描，会将连接约束信息作为索引条件，这种情况下不需要再重复检查。
- 3) 对连接约束信息按执行器最佳执行顺序排序。同样，该部分只考虑其代价估计信息，如对于表 A(a, b, c) 和表 B(a, b)，其连接约束条件有 $A.b + A.c = B.b$ AND $A.a > B.a$ ，则执行器会先执行后者。

数据结构 5.18 NestLoop

```
typedef struct Join
{
    Plan      plan;
    JoinType jointype;
    List     *joinqual; //连接约束信息
} Join;

typedef struct NestLoop
{
    Join      join;
} NestLoop;
```

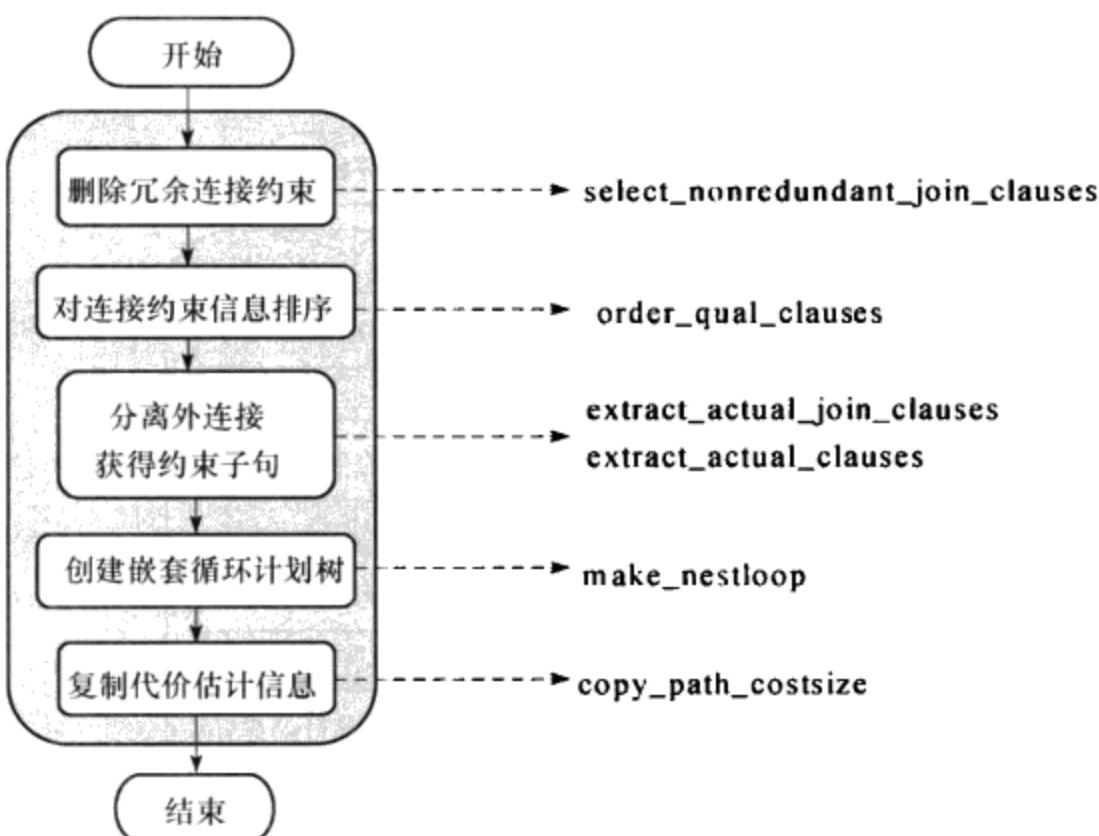


图 5-34 创建嵌套循环连接计划的流程

- 4) 存在外连接则分离外连接获得约束子句。
- 5) 创建嵌套循环连接计划，将目标属性、连接子句、外连接子句赋值给 `NestLoop` 计划节点，并将其外关系计划树和内关系计划树作为其左右子树。
- 6) 复制代价估计信息，该部分是对 `Path` 中代价估计的一个拷贝。

这里需要特别解释一下外连接子句的处理。我们知道，外连接的结果元组某一侧可以为空值。PostgreSQL 并没有将外连接作为连接子句去处理，而是将其赋值给 NestLoop 节点包含的 plan 中的 qual 字段（在顺序扫描的情况下，我们可以看到该变量存放的是顺序扫描的约束条件），采取一种类似于基本表的扫描方式。

除嵌套循环连接外，还有归并连接和 Hash 连接两种连接方式。但其生成计划的过程类似，主要过程都包含提取有效连接约束信息、对连接约束信息排序、分离外连接、创建相应类型的计划树、复制代价估计等步骤。但对于归并连接，还需要考虑内外关系的排序情况；对于 Hash 连接，要从连接约束子句中提取出 Hash 连接子句，将其赋值给计划节点的不同字段。

5.4.6 生成完整计划

上节介绍了依据最优路径 best_path 生成计划树的过程。前面讲过，生成路径仅仅考虑基本查询语句信息，并没有保留诸如 GROUP BY、ORDER BY 等信息。grouping_planner 函数调用 create_plan 生成基本计划树后，则会依据查询树相关约束信息在前面生成的普通计划之上添加相应的计划节点生成完整计划。

1. 创建聚集计划节点

函数 make_agg 用于生成聚集计划节点，函数返回一个类型为 Agg 的结构（数据结构 5.19），其参数除了当前规划器状态信息、目标列、聚集约束条件外，还有以下几个特有参数：

- aggstrategy：类型为 AggStrategy，指该聚集计划采用的聚集策略。聚集策略有以下三种：AGG_PLAIN（普通聚集）、AGG_SORTED（排序聚集）、AGG_HASHED（Hash 聚集）。
- numGroupCols：类型为 int，需要聚集的属性数目。
- grpColIdx：类型为 AttrNumber 指针类型，代表聚集属性的索引。
- grpOperator：类型为 OID 指针类型，代表用于分组时的排序操作符。
- numGroup：类型为 long，代表预计的分组数目。
- numAggs：类型为 int，代表聚集数目。
- lefttree：为 Plan 的指针类型（可以是实际的计划类型强制转换成 Plan 类型），代表添加聚集计划节点前的计划树，该计划树将作为 Agg 的左子树。

数据结构 5.19 Agg

```
typedef struct Agg
{
    Plan      plan;
    AggStrategy aggstrategy; //聚集策略
    int       numCols;      //分组属性列数
    AttrNumber *grpColIdx; //分组列的索引
    Oid        *grpOperators; //分组操作符
    long       numGroups;   //预计分组数目
} Agg;
```

如图 5-35 所示，函数 make_agg 执行流程如下：

- 1) 初始化计划节点。
- 2) 调用 copy_plan_costsize 函数获得左子树代价。
- 3) 调用 cost_agg 函数计算聚集代价。
- 4) 调用 cost_qual_eval 函数计算 HAVING 子句代价。
- 5) 最后复制总代价信息、目标属性，连接原计划树至 Agg 节点的左子树指针。

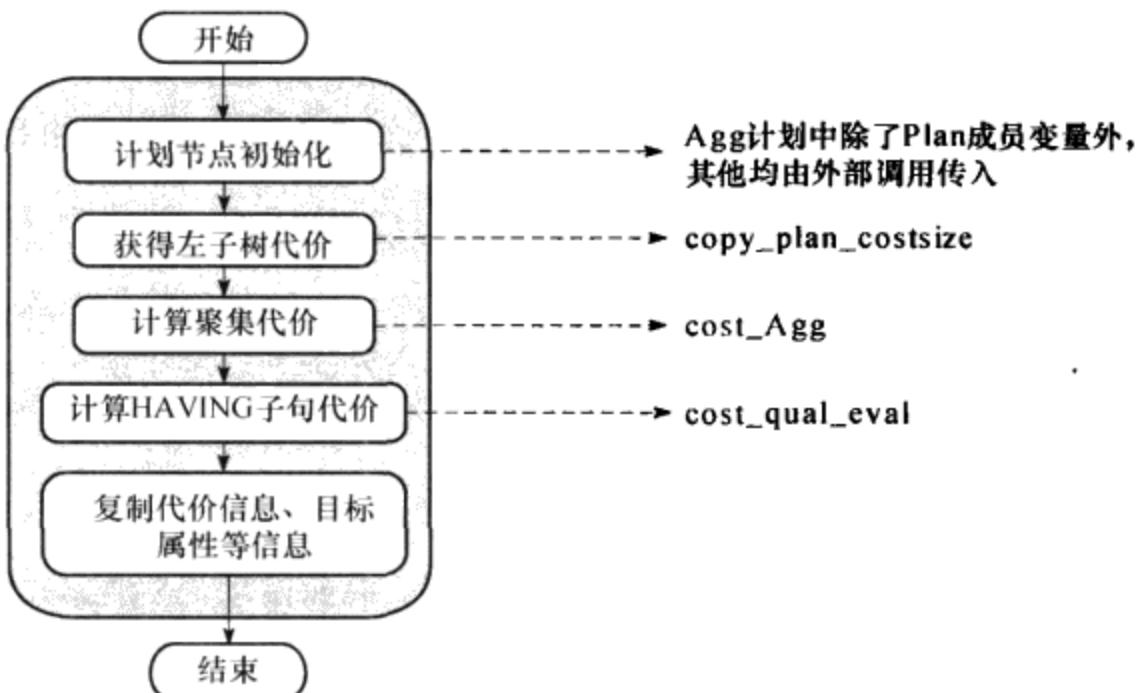


图 5-35 函数 make_agg 的执行流程

通过数据结构可以看到，Agg 计划节点除了包含 Plan 类型成员变量外，还添加了和聚集分组相关的一些成员变量，包括分组属性数目、分组属性的索引、分组的操作符等，这些扩展成员变量由外部调用传入参数直接赋值；获得左子树的代价（仅保留大小估计，时间代价会被覆盖），调用 cost_agg 计算其聚集代价，对于存在 HAVING 子句的情况调用 cost_qual_eval 估计其代价；最后将代价、目标属性、HAVING 子句赋给 plan 相应的成员变量（具体成员变量解释见顺序扫描计划），并将原有计划树作为其左子树添加进来形成完整的计划树。

2. 创建排序计划节点

函数 make_sort 用于生成排序计划节点，函数返回一个类型为 Sort 的结构（数据结构 5.20），参数主要有：

- numCols：类型为 int，代表排序的属性数目。
- sortColIdx：类型为 AttrNumber 指针类型，代表排序的属性索引。
- sortOperators：类型为 Oid 指针类型，代表排序的操作符。

如图 5-36 所示，函数 make_sort 的执行流程如下：

- 1) 调用 copy_plan_costsize 函数获得左子树的代价信息。
- 2) 调用 cost_Sort 函数计算排序代价。
- 3) 最后复制相应信息。

数据结构 5.20 Sort

```

typedef struct Sort
{
    Plan      plan;
    int       numCols;           //排序属性列数
    AttrNumber* sortColIdx;     //排序属性列索引
    Oid       *sortOperators;   //排序操作符
    bool      *nullsFirst;
} Sort;

```

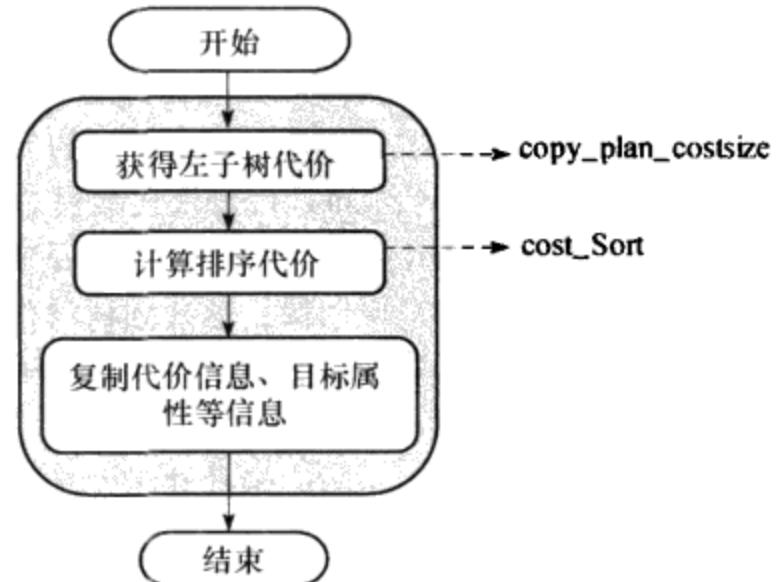


图 5-36 函数 make_sort 的执行流程

通过数据结构可以看出，Sort 计划节点除了包含 Plan 类型成员变量外，添加了和排序相关的一些成员变量，包括排序属性索引、排序操作符等，这些扩展成员变量由外部调用传入参数直接赋值。这里需要额外做的工作是调用 cost_Sort 估计排序代价，并将原有计划树连接至左子树，形成完整计划树。

由以上计划节点生成过程可以得知，后期完整计划树的封装主要是计算其相应代价信息，并将相关信息填入计划节点。对于不同语句，前期还会生成相应的物理执行策略。例如，聚集的时候会依据数据的统计信息决定采取普通聚集、排序聚集或 Hash 聚集，不同的物理执行策略就会产生不同的代价。

5.4.7 整理计划树

生成的完整计划经过计划树整理之后就可以交给查询执行器去执行了，负责整理工作的主函数是 set_plan_references。整理计划树是查询规划器处理工作的最后一步，主要是为了方便执行器的执行，对计划树一些表达上的细节做最后的调整。例如，将上层的 Var 结构变为对子计划的输出结果的引用、获取操作符的 OID 等。同时，这一步也会删除那些没有任何用处的子查询扫描计划节点。在进行整理工作时会使用不同的函数来进行不同的整理动作，其主要函数及其功能如表 5-12 所示。

表 5-12 计划树整理函数功能介绍

主要函数	功能介绍
fix_expr_references	通过调用 fix_expr_references_walker，完成表达式（目标属性表或条件表达式）的清理工作
set_subqueryscan_references	在 SubqueryScan 上进行 set_plan_references 的操作，即试图去除掉 SubqueryScan 节点
fix_opfuncids	通过调用 fix_opfuncids_walker，在一个表达式的树中对操作符调用的表达式节点（OpExpr node）的值进行补全
trivial_subqueryscan	检测一个 SubqueryScan 是否可以从计划树中删除
adjust_plan_varnos	通过 rtoffset 调整 varnos 和其他所涉及的表的索引在计划树中的偏移量
adjust_expr_varnos	通过 rtoffset 来调整变量中的 varnos 在一个表达式中的偏移量
fix_expr_references	做表达式（如目标属性表或约束表达式）最后的清理工作
set_join_references	通过设置 varnos 为内连接或外连接和设置 attno 的值为内连接或外连接的结果的元组的数目来修改 join 结点的目标属性表和约束表达式
set_inner_join_references	处理出现在内部索引连接的表达式的 join 节点
set_uppernode_references	根据左子树子计划的返回元组来更新上层的目标变量列表和表达式
build_tlist_index	为一个子节点的目标属性表建立一个索引的结构
search_indexed_tlist_for_var	在一个索引的列表中寻找一个变量 var，假如找到的话，返回它的一份拷贝，假如没有找到，则返回 null
search_indexed_tlist_for_non_var	在一个索引的列表中寻找一个非变量的节点，假如存在返回一个 Var 指向这个 item，假如不存在，则返回 null
join_references	通过修改句子中的 varno/varattno 的值为外连接或内连接的目标属性表来建立一个 join 表达式或目标属性表的集合
replace_vars_with_subplan_refs	这个程序修改表达式树，使得所有的变量节点都与子计划的目标节点相关联。它主要用来处理上层计划节点中的非连接表达式或目标属性表

5.4.8 实例分析

在 5.2 节的例 5.1 中介绍了如何对 SQL 语句进行查询分析最终生成查询树，现在其基础上进一步介绍如何对查询树进行规划处理。

1. 查询分析之后的查询树

图 5-37 中给出了例 5.1 对应的原始查询树在数据结构上的组织形式。其中 jointree 中的节点可能有以下三种类型：

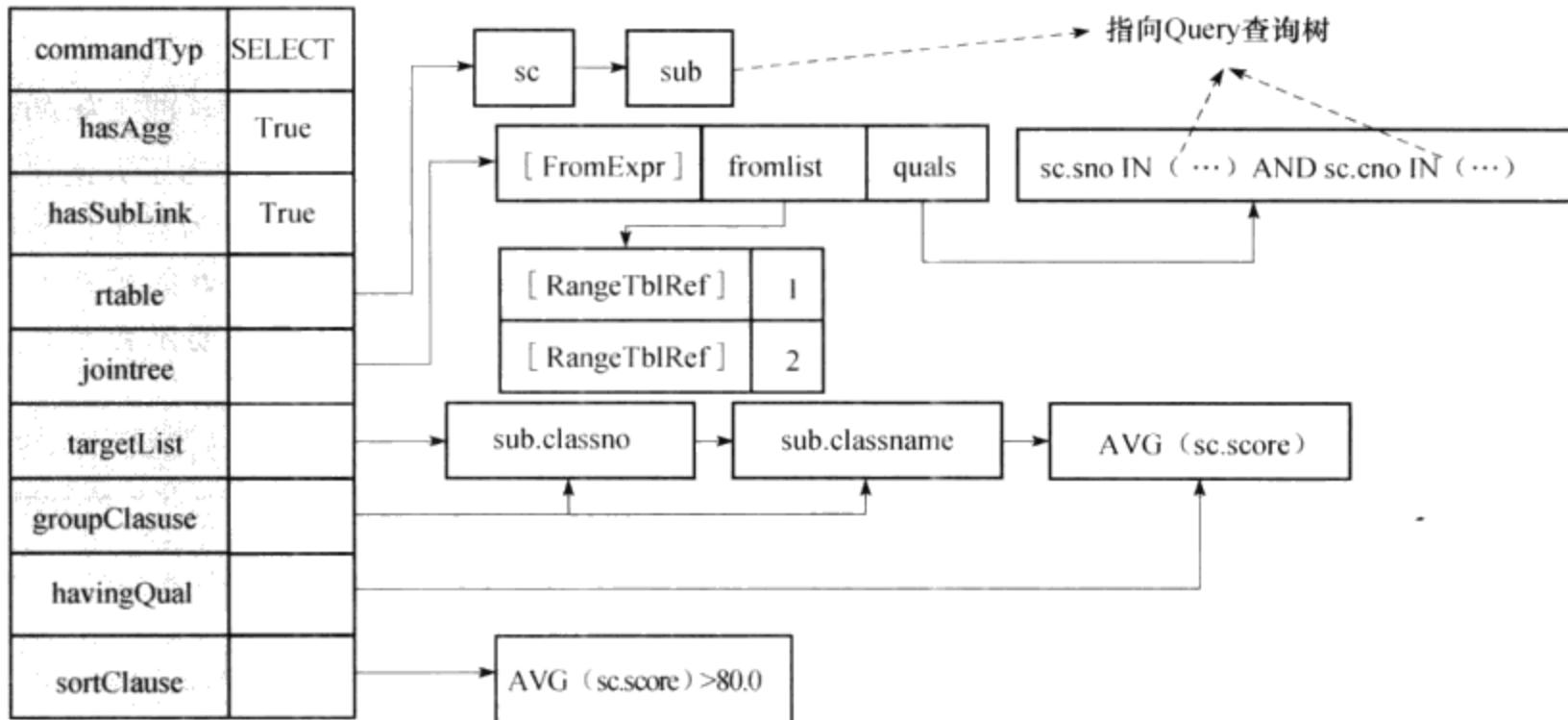


图 5-37 例 5.1 对应的查询树

1) 范围表索引 (RangeTblRef)：实际就是 1、2、3 等整数值，代表引用查询树中 rtable 链表中相应索引号的范围表。

2) FROM 表达式 (FromExpr)：FROM 表达式数据结构包含两个主要字段：

- fromlist：连接子树链表，成员类型可为范围表索引、FROM 表达式、JOIN 表达式。
- quals：约束条件，作用于 FROM 表达式的约束条件，来自于 WHERE 子句。

3) JOIN 表达式 (JoinExpr)：JOIN 表达式数据结构包含三个主要字段：

- larg：左连接子树。
- rarg：右连接子树。
- quals：连接约束条件。

对于 jointree，其顶层节点不能是 JOIN 表达式类型，否则查询分析器会在 jointree 顶层强制增加 FROM 表达式节点，并将其 quals 设置为空。

对于例 5.1 中的连接树，其顶层为 FROM 表达式类型，范围表索引 1、2 分别代表 sc 和子查询 sub，其连接条件为 “sc.sno IN (SELECT sno FROM student WHERE student.classno = sub.classno) AND sc.cno IN (SELECT course.cno FROM course WHERE course.cname = '高等数学')”。

2. 提升子链接

由原始查询树可知，WHERE 子句中存在两个子链接，可以将其提升成为范围表（如果子链接不能被提升则将作为子查询处理）。提升子链接主要有以下步骤：

- 1) 扫描连接树，发现存在非相关子链接。
- 2) 在范围表中添加相应子查询 ANY_subquery - 索引号 3。
- 3) 构建 JoinExpr 节点，将原始 jointree 作为其左子树；将范围表索引 3 作为其右子树；将原有的“sc.cno IN (……)”转换为“sc.cno = ANY_Subquery.cno”，作为新建 JoinExpr 的连接约束条件。
- 4) 新建空的 FromExpr，将连接树作为其 fromlist，约束条件为空。
- 5) 删除原有子链接。

提升子链接后，例 5.1 对应的查询树转变成图 5-38 所示的样子，所对应的查询语句如下所示：

```
SELECT sub.classno, sub.classname, AVG(sc.score) AS avg_score
FROM sc, (SELECT* FROM class WHERE class.gno = '2005') AS sub,
      (SELECT course.cno FROM course WHERE course cname = '高等数学') AS ANY_subquery
WHERE sc.sno IN (SELECT student.sno FROM student WHERE student.classno = sub.classno)
      AND sc.cno = ANY_subquery.cno
GROUP BY classno, classname
HAVING AVG(score) > 80.0
ORDER BY avg_score;
```

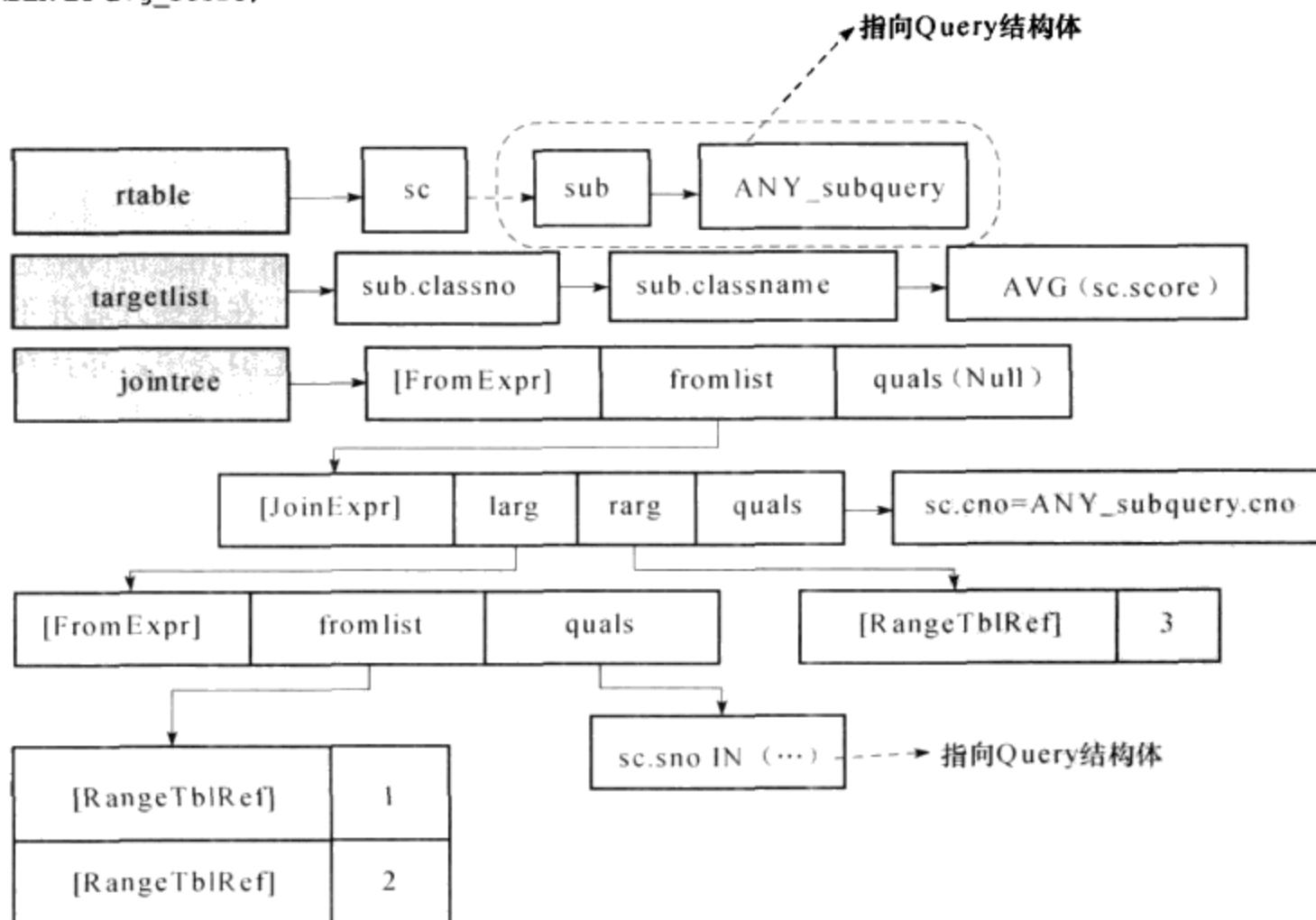


图 5-38 例 5.1 提升子链接后的查询树

3. 提升子查询

从提升子链接后的查询树和 SQL 语句可以看出，FROM 子句的 3 个范围表有 2 个是子查询（非相关）的形式，因此可以将其提升到父查询。

对图 5-38 中范围表索引为 2 的子查询进行提升需进行以下步骤：

- 1) 扫描 jointree 发现存在子查询 2-sub 和 3-ANY_subquery。
- 2) 在范围表中添加子查询中所包含的物理表，即 class，其索引为 4。
- 3) 创建 FROM 节点，依据原有子查询 2 填充 FromExpr——将其 fromlist 设为表 4 的索引，约束条件即子查询的约束条件 “class.gno = ‘2005’”。
- 4) 删除原来对范围表索引 2 的引用，将新建的 FromExpr 节点连接到链表上。
- 5) 调整所有引用了范围表索引为 2 (sub) 的地方，将其改为 4。

对于图 5-38 中范围表索引为 3 的子查询，由于其没有受影响的目标属性，按前三个步骤处理即可。所对应的查询语句如下所示：

```
SELECT class.classno, class.classname, AVG(sc.score) AS avg_score
FROM sc, (SELECT* FROM class WHERE class.gno = '2005') AS sub,
      (SELECT course.cno FROM course WHERE course cname = '高等数学') AS ANY_subquery, class, course
WHERE sc.sno IN (SELECT student.sno FROM student WHERE student.classno = class.classno)
      AND sc.cno = course.cno
      AND course cname = '高等数学'
      AND class.gno = '2005'
GROUP BY classno, classname
HAVING AVG(score) > 80.0
ORDER BY avg_score;
```

需要指出的是：PostgreSQL 在查询规划过程中，考虑到后续步骤仍可能引用之前的信息，范围表中的子查询 sub（范围表编号为 2）和 ANY_subquery（范围表编号为 3）在被提升后并不会被删除。从子查询中提升出来的基本表 class 和 course 会直接添加到父查询的范围表中，并且在父查询的范围表中编号为 4 和 5。提升子查询之后对应的查询语句如下所示：

```
SELECT classno, classname, AVG(sc.score) AS avg_score
FROM sc, class, course
WHERE sc.sno IN (SELECT student.sno FROM student WHERE student.classno = class.classno)
      AND sc.cno = course.cno
      AND course cname = '高等数学'
      AND class.gno = '2005'
GROUP BY classno, classname
HAVING AVG(score) > 80.0
ORDER BY avg_score;
```

提升子查询后例 5.1 对应的查询树如图 5-39 所示。

4. 生成路径

查询树改造完毕之后将先生成路径，然后再进行计划的生成。路径生成按照先基本关系后连接

关系的顺序进行。

(1) 基本关系的路径

首先处理的是相关子查询，它只针对基本关系 student，不存在关系的连接。所以对此 student 的扫描形成了此相关子查询的最终路径链。由于在 student 上并没有索引，因此 student 的路径只有一种情况：顺序扫描。记 student 的顺序扫描路径为 R1。

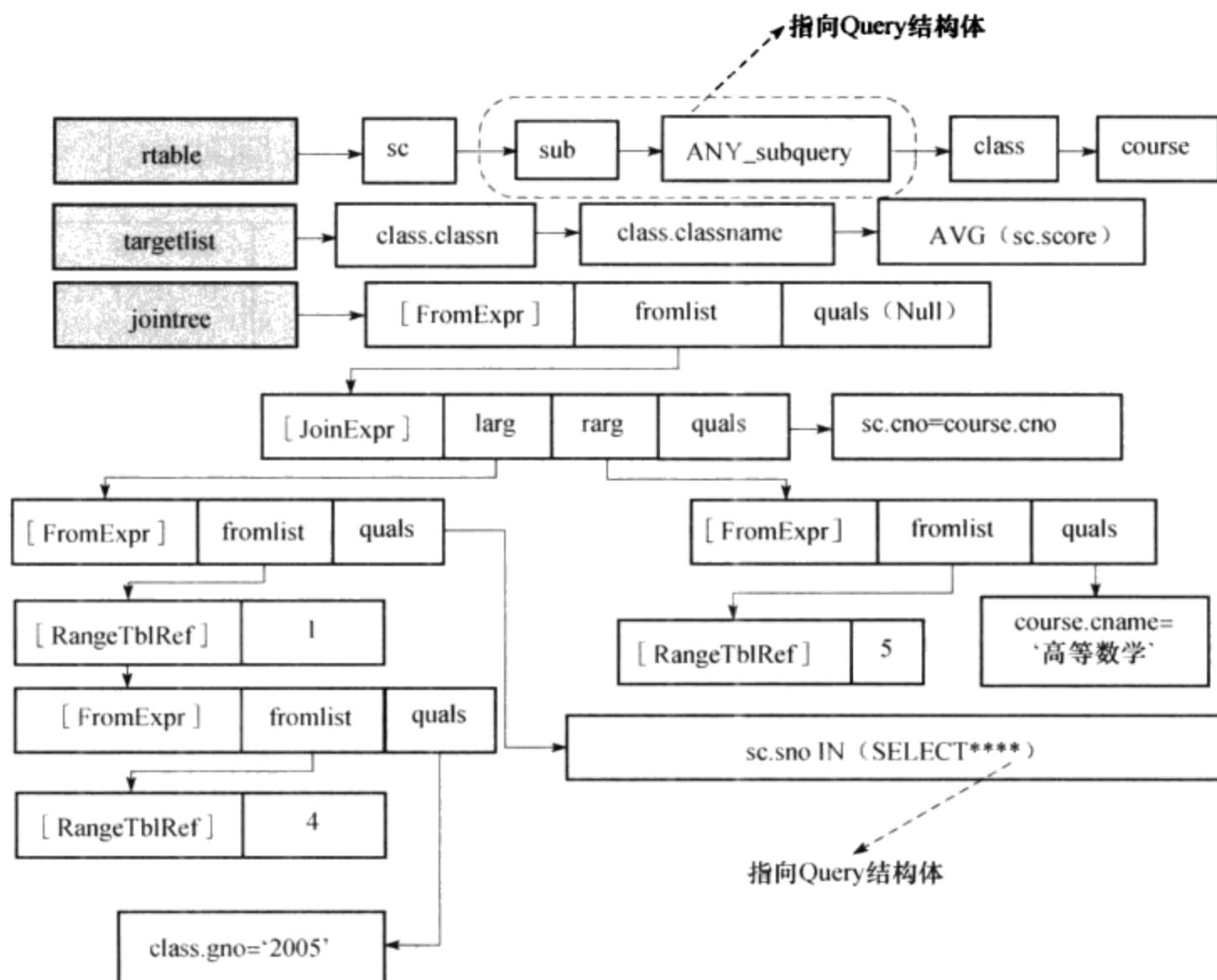


图 5-39 例 5.1 提升子查询后的查询树

相关子查询路径生成之后将为父查询生成基本关系的路径，父查询包含三个基本关系，分别为 sc、class 和 course，在范围表中对应的索引号为 1、4、5。开始得到的 R1 生成计划后得到的子计划 subplan 将作为 sc 及 class 的连接约束信息放在这两个基本关系的 joininfo 字段内。例 5.1 中基本关系的路径如图 5-40 所示。

(2) 连接关系的路径

基本关系的路径生成之后，就需要利用动态规划算法来生成由 sc、class、course 连接而成的连接关系的路径，生成的过程如下：

第一层即基本关系的扫描。依据最佳启动代价、最佳执行代价、键值信息三个筛选依据，结果如图 5-41 所示。其中表 sc 和表 class 的 pathlist 中仅有一条路径，而表 course 存在两条路径，其中索引扫描在键值上符合外部的排序需求，被保留下。

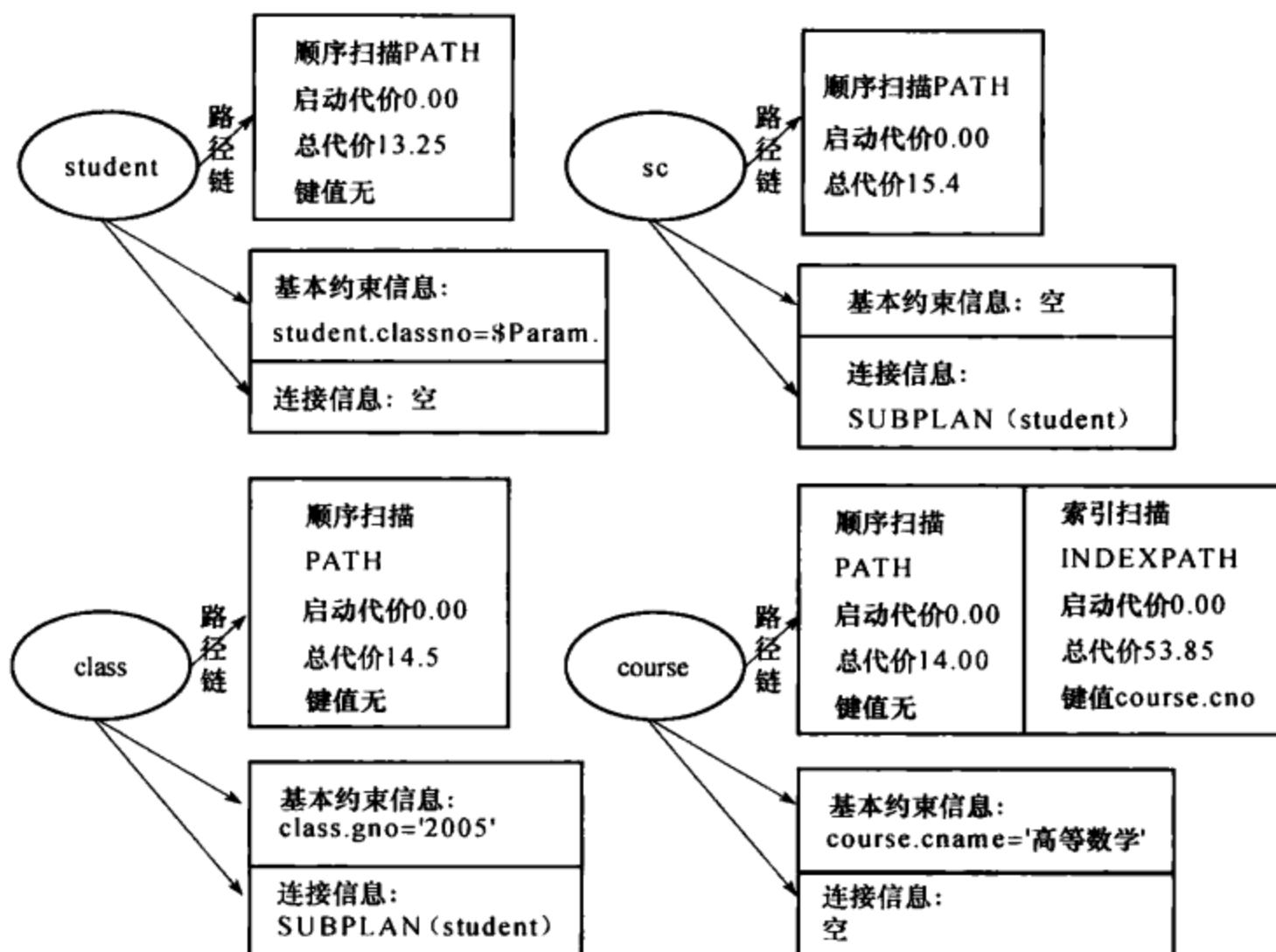


图 5-40 例 5.1 中基本关系的路径

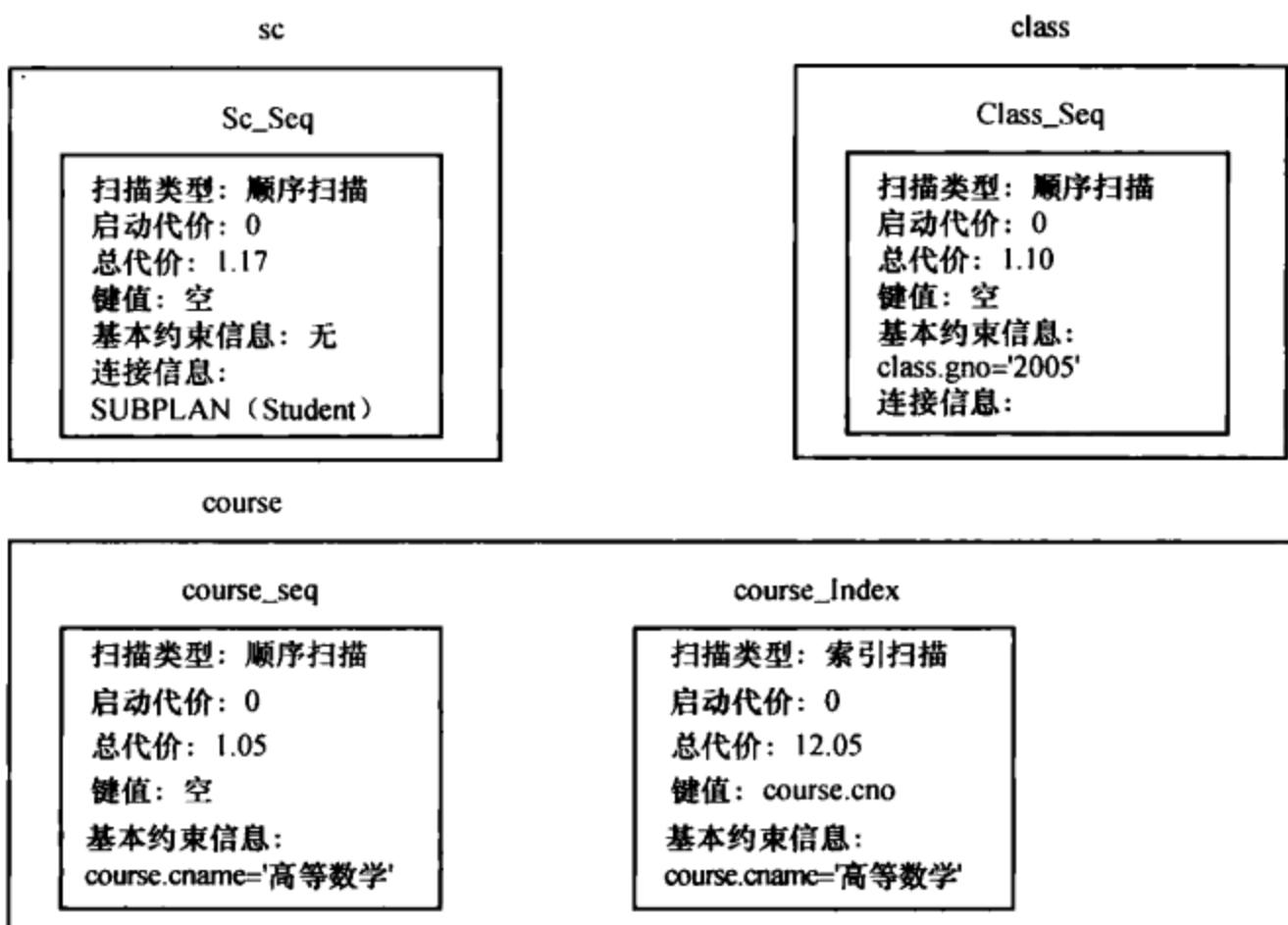


图 5-41 动态规划第一层的结果

第二层进行基本关系表连接。第二层的结果如图 5-42 所示。采取动态规划算法生成路径，保留具有“优势”的路径。表 sc 和表 class 的连接存在两条路径，连接类型均为嵌套循环连接，不同的是，在 sc_class_nest2 中，对 class 表顺序扫描后添加物化计划节点。

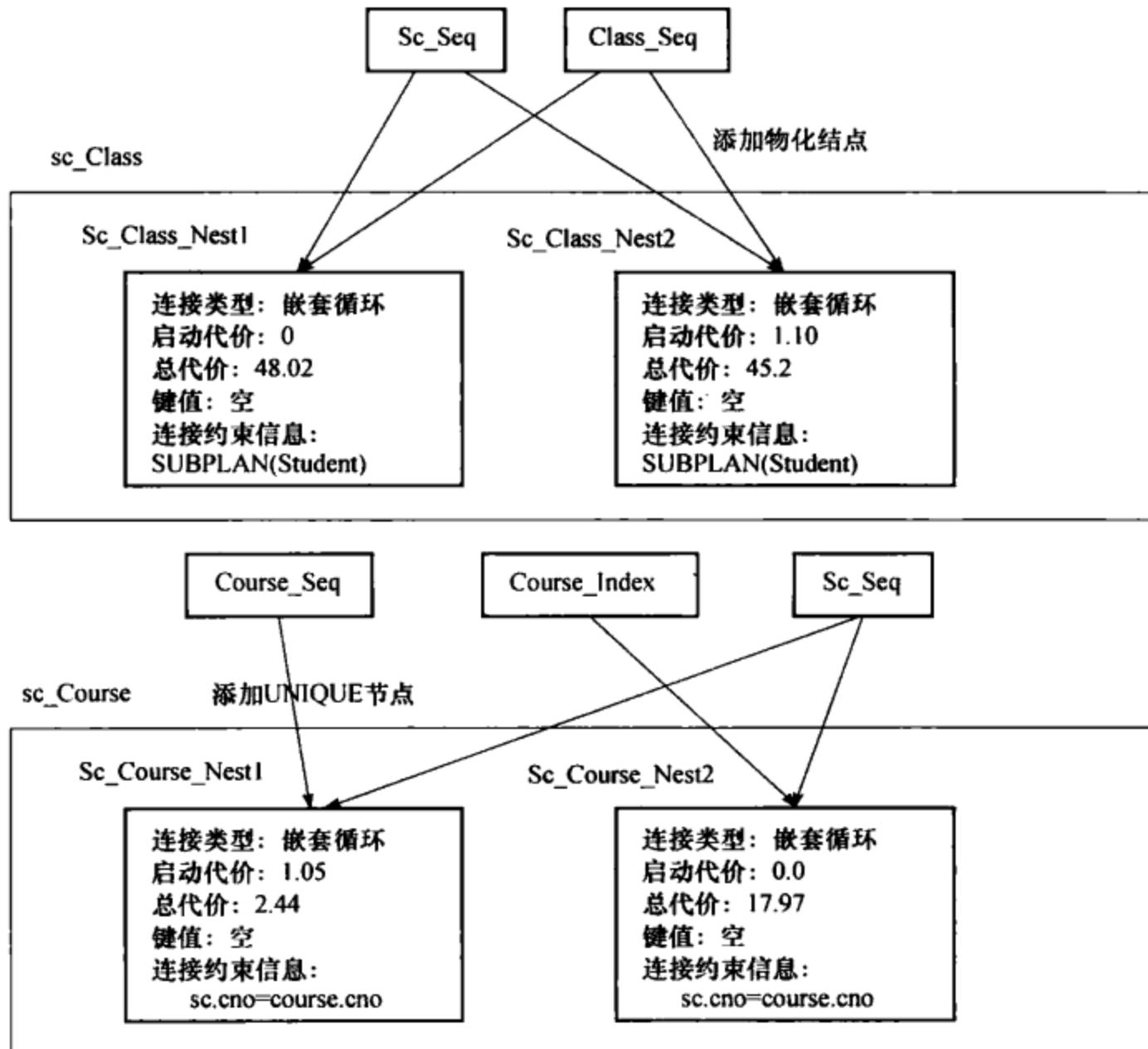


图 5-42 动态规划第二层的结果

对于 sc 和 course 的连接同样存在两条路径，其中，sc_course_nest1 节点的生成过程是：先为 course 的顺序扫描节点添加 UNIQUE 计划节点，然后与 sc 顺序扫描节点进行嵌套循环连接得到 sc_course_nest1 节点。sc_course_nest2 为 course 索引扫描与 sc 顺序扫描进行嵌套循环连接。

表 class 与 course 是完全不相关的，进行连接的话只能求笛卡儿积（耗时，且没有必要）。

第三层进一步对基本关系或连接关系连接。结果如图 5-43 所示，在 sc_class_course 连接的 path-list 中存在三条路径。最终的最佳执行路径为 sc_class_nest1。

5. 生成计划

由上步生成的最优路径，按照后序遍历路径树，提取相应的约束信息（并按最佳执行顺序排序）、基本关系或连接关系的 ID、目标属性等信息，创建相应的计划节点，生成计划树。由最优路径只能生成普通计划树，后期再进行包装成完整计划树。生成的计划树如图 5-44 所示。

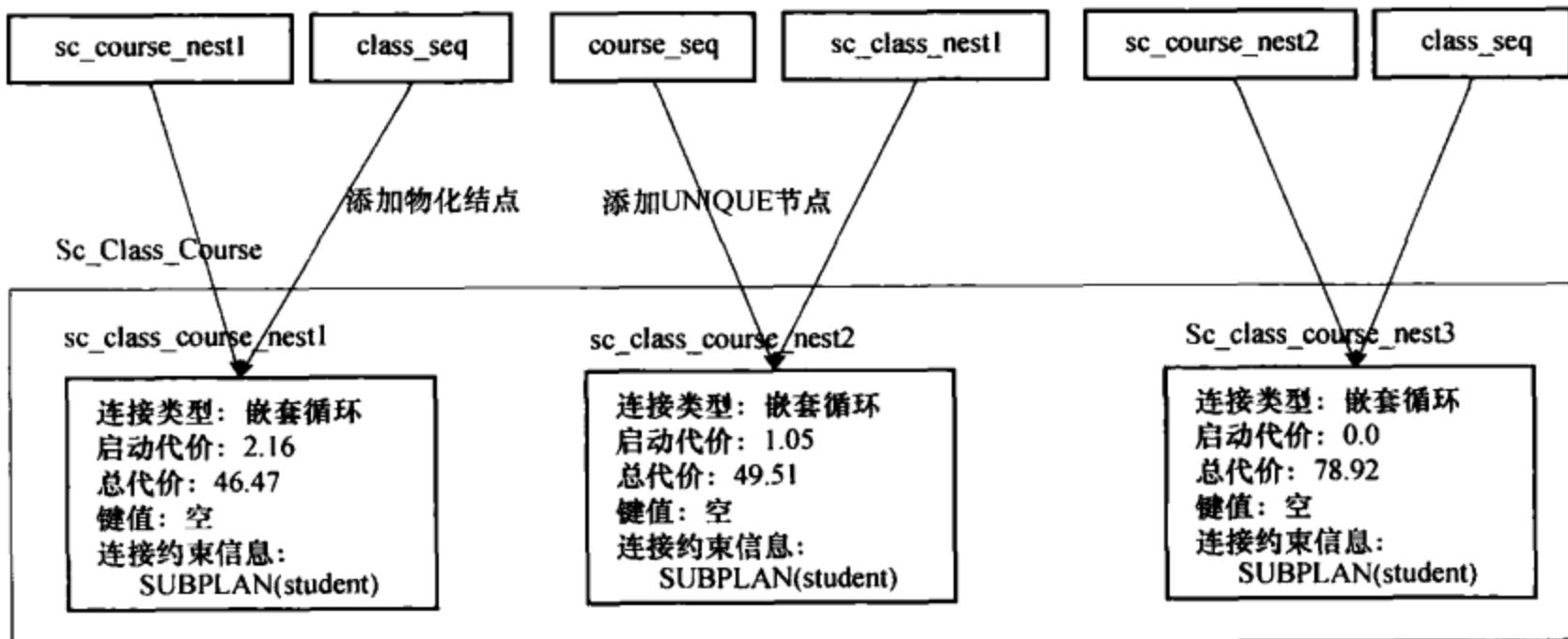


图 5-43 动态规划第三层结果 (最终结果)

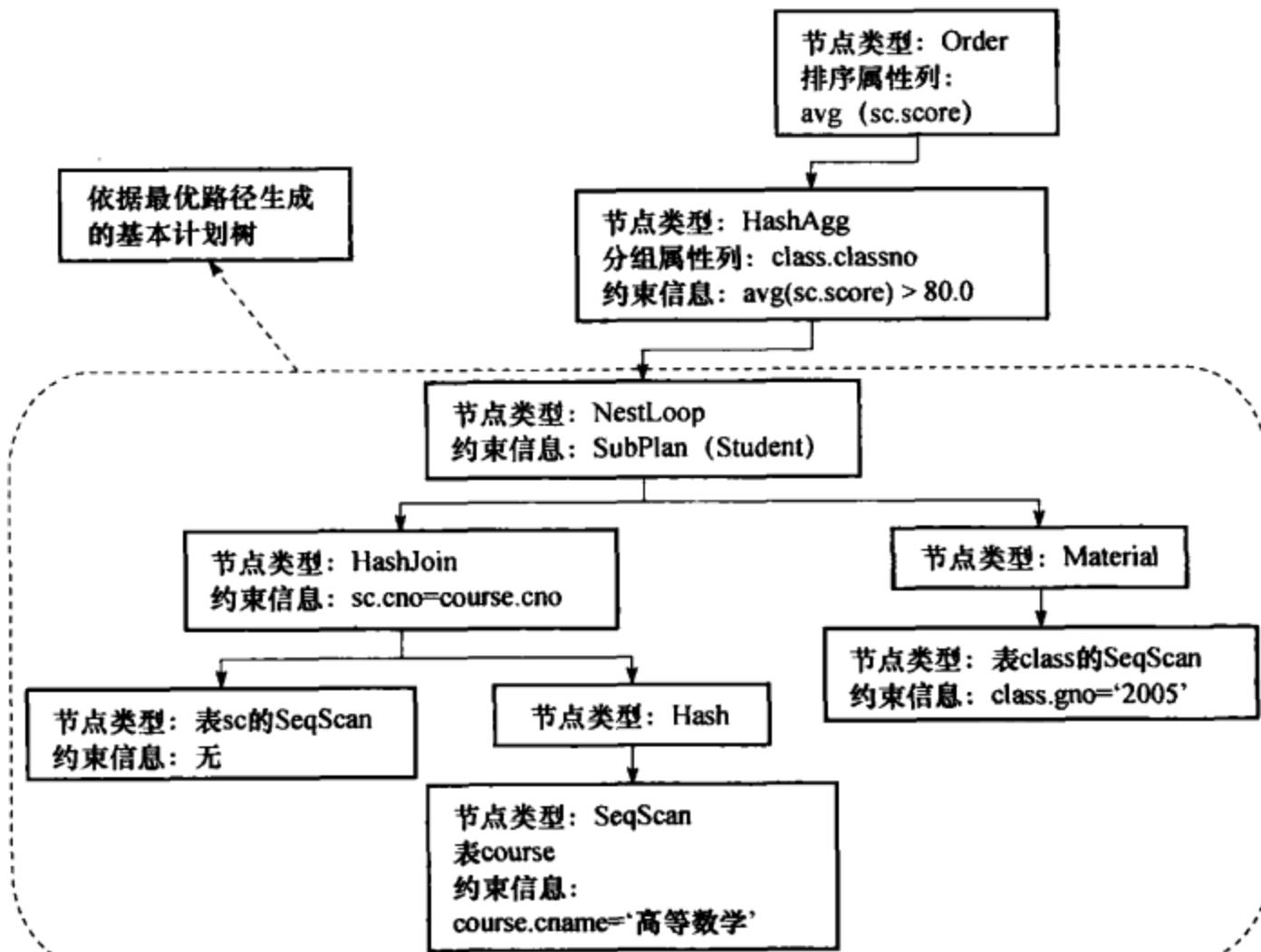


图 5-44 例 5.1 对应的计划树

5.5 代价估计

在 PostgreSQL 的查询规划过程中，查询请求的不同执行方案是通过建立不同的路径（Path）来

表达的。在生成了许多符合条件的路径之后，要从中选择出代价最小的路径，把它转化为一个计划，传递给执行器执行。因此，规划器的核心工作就是建立多条路径，然后从中找出最优的那一条。同一个查询请求有不同路径主要是因为：表的不同访问方式、表之间不同的连接方式、表之间不同的连接顺序等因素造成的。而评价路径优劣的依据是用系统表 pg_statistic 中的系统统计信息估计出的不同路径的代价（Cost）。

某个路径的代价要考虑 CPU 代价和磁盘存取代价两方面。磁盘代价以从磁盘顺序存取一个页面的代价为单位，所有其他形式的代价计算都是相对磁盘存取代价来计算的。用于估算代价的参数有：

- seq_page_cost：顺序存取页的代价，值为 1.0。
- random_page_cost：非顺序存取页的代价，值为 4.0。
- cpu_tuple_cost：典型的 CPU 处理一个元组的代价，值为 0.01。
- cpu_index_tuple_cost：典型的 CPU 处理一个索引元组的代价，值为 0.005。
- cpu_operator_cost：CPU 处理一个典型的 WHERE 操作的代价，值为 0.0025。
- effective_cache_size：用来量度 PostgreSQL 和 OS 缓存的磁盘页的数量，值为 16384。

一个路径的代价由三部分组成：启动代价（Startup Cost）、总代价（Total Cost）、执行结果的排序方式（PathKeys）。计算启动和总代价所用的参数如下：

1) 基本参数：

- 表元组数 ntuples。
- 表磁盘块数 nblocks。
- 符合选择条件的元组数 nrows。

2) 统计信息：查询规划器需要估计一个查询检索的元组的数目，这样才能选择正确的查询规划，统计信息的主要内容就是每个表和索引中的元组总数，以及每个表和索引占据的磁盘块数。这个信息保存在系统表 pg_class 的 reltuples 和 relpages 属性中^①。

3) 直方图信息（存在系统表中）：各个属性值出现次数的统计信息。

路径的代价估算的基本步骤都是相似的：首先根据统计信息和查询条件，估算出这次查询要进行的 I/O 次数以及要取出的元组个数，然后根据元组个数（分为表元组和索引元组）计算出 CPU 代价，最后综合考虑 CPU 代价和 I/O 次数（磁盘代价）即可得到最后的代价。不同的路径类型有不同的代价估算方式，在 5.4.3 节中我们已经介绍过索引扫描路径的代价估算函数，其他类型路径的代价估算函数读者可以参考索引扫描路径进行分析，在本节中我们将给出影响代价估算的重要因素及代价估算的一般性公式。

5.5.1 代价估算公式

代价估算从 I/O 次数和 CPU 开销两个方面考虑，估算公式为 $P + W * T$ 。其中，P 表示在执行时所要访问的页面数，反映了磁盘 I/O 次数；T 表示在执行时所要访问的元组数，反映了 CPU 开销；W 表示磁盘 I/O 代价和 CPU 开销间的权重因子。通过对估算公式的分析可以发现，计算代价时只需考虑访问的页面数和元组数。

^① <http://www.pgsqldb.org/mwiki/index.php/规划器使用的统计信息>

5.5.2 选择度

选择度 (selectivity) 用来定量地描述前述代价估算公式中的权重因子 W。选择度的计算要综合考虑以下各类参数：约束条件中的操作符、约束常量、索引中的元组数、某字段的最大值和最小值等。表 5-13 中列出了对于选择度计算需要考虑的因素。

表 5-13 选择度

约束条件	选择度
r. field = value	1/字段 r. field 上所定义的索引关系的元组数
r. field > value	(字段 r. field 的最大值 - value)/(字段 r. field 的最大值 - 字段 r. field 的最小值)
r. field < value	(value - 字段 r. field 的最小值)/(字段 r. field 的最大值 - 字段 r. field 的最小值)

5.5.3 单个表的扫描代价

计算单个表的扫描代价时，需要将表的大小与各个约束条件所对应的选择度相乘，如表 5-14 所示。其中，NumPages 表示表的页面数；NumTuples 表示表的元组数；ITuples 表示索引表的元组数；F 表示多个约束条件组合后的选择度。

表 5-14 扫描代价计算方法

扫描方式	P	T
顺序扫描	NumPages	NumTuples
一级索引扫描	NumPages * F	NumTuples * F
二级索引扫描	NumPages * F + ITuples * F	ITuples * F + NumTuples * F

5.5.4 两个表的连接代价

两个表的连接代价计算方法如表 5-15 所示。其中：

- Couter 表示扫描外连接表的代价。
- Cinner 表示扫描内连接表的代价。
- Csortouter 表示将外连接表进行排序的代价（使用临时存储空间）。
- Csortinner：表示内连接表进行排序的代价（使用临时存储空间）。
- Ccreatehash 表示对内连接表进行 Hash 的代价（使用临时存储空间）。
- Chash 表示单独的 Hash 的代价。
- Nouter 表示外连接表的大小。

表 5-15 两个表的连接代价计算方法

连接方式	代价计算公式
嵌套循环连接	$C_{outer} + N_{outer} * C_{inner}$
归并连接	$C_{outer} + C_{sortouter} + C_{inner} + C_{sortinner}$
Hash 连接	$C_{outer} + C_{createhash} + N_{outer} * C_{hash}$

5.6 PostgreSQL 中的遗传算法

遗传算法 (GA) 是一种启发式的优化法，它通过既定的随机搜索进行操作。在 PostgreSQL 中，遗传算法主要用在连接路径的生成操作中，其流程如图 5-45 所示。

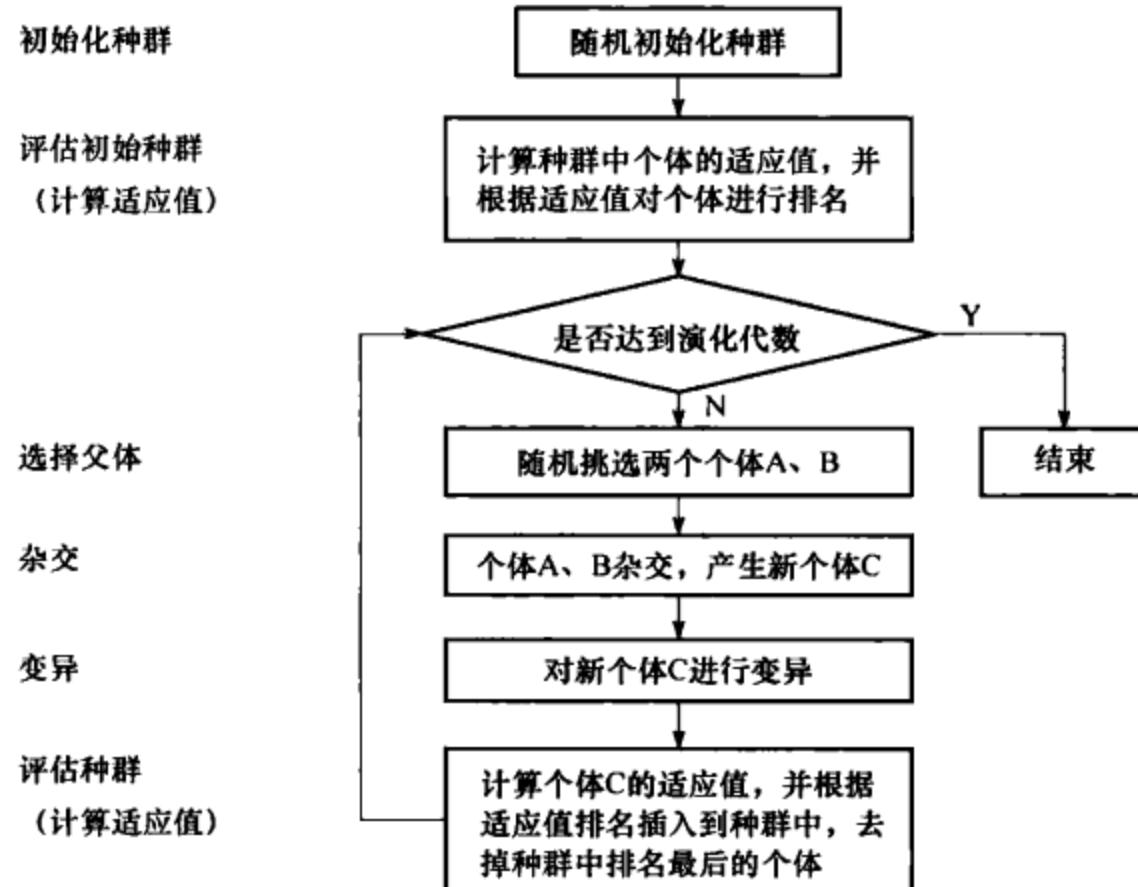


图 5-45 PostgreSQL 中遗传算法处理流程

遗传算法让 PostgreSQL 查询规划器可以通过非穷举搜索有效地支持很多表的连接路径的生成。下面将对其中的个体编码方式、适应值、父体选择策略等进行介绍。

5.6.1 个体编码方式及种群初始化

PostgreSQL 中用遗传算法解决表连接问题的方式类似于解决 TSP 问题。可能的连接路径被当作整数串进行编码。每个串代表查询中可能的一种连接顺序。例如，图 5-46 中的查询树是用整数串“4132”编码的，也就是说，首先联接表“4”和“1”，得到的结果表再和表“3”连接，最后再和表“2”连接。这里的 1、2、3、4 都是 PostgreSQL 里的关系标识 (relids)。

遗传算法执行的第一步是随机初始化种群，种群大小的计算方法参考 5.6.6 节。假设种群大小为 n，首先随机初始化 n 个排列数，即 n 个类似于上面提到的数字串“4132”，每一个排列数就是一个个体。然后基于 n 个排列数生成基本表的连接路径，在此过程中会进行代价评估，将最后的代价作为适应值来衡量该个体的优劣。每一个个体都用 Chromosome 结构 (数据结构 5.21) 表示，其中记录了该个体的排列数和代价。

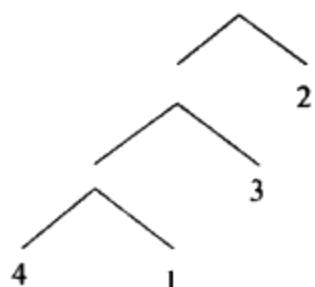


图 5-46 表连接

数据结构 5.21 Chromosome

```

typedef int Gene;
typedef struct Chromosome
{
    Gene      *string;          //染色体的数据值
    Cost      worth;           //对该染色体的代价评估值
}Chromosome;

```

5.6.2 适应值

个体的适应值等于该个体中 N 个表的连接代价（参考表间的连接方式）。适应值计算由函数 geqo_eval 实现。包括以下步骤：

- 1) 检查个体的有效性。
- 2) 确定个体的连接次序以及连接方式。
- 3) 计算个体的适应值。

计算个体的适应值时，首先要检查个体是否有效。也就是说，对一个给定的个体，能不能把这个个体中的表按照某种次序连接起来，因为有些表之间是不能连接的。如果一个个体按任何次序都不能连接，那么这个个体是无效的。对于有效的个体，还要确定连接次序和连接方式。在连接次序和连接方式确定之后才能计算个体的适应值。

5.6.3 父体选择策略

父体选择策略采用基于排名的选择策略，选择概率函数如公式 5-2 所示。其中，max 是个体总数，bias 的默认值为 2.0，geo_rand 是 0.0 ~ 1.0 之间的随机数，f(geo_rand) 表示当前个体在种群中的编号（该编号是根据当前个体的适应值在种群中的排名来确定）。该函数表明，排名越靠前的个体被选择的概率越大。不过在 PostgreSQL 8.4.1 的实现中，父体选择时没有考虑 A 小于 0 的情况。

$$f(\text{geo_rand}) = \begin{cases} \max * \frac{\text{bias} - A}{2 * (\text{bias} - 1)} & (A < 0) \\ \max * \frac{\text{bias} - \sqrt{A}}{2 * (\text{bias} - 1)} & (A \geq 0) \end{cases} \quad (\text{公式 5-2})$$

其中, $A = \text{bias} * \text{bias} - 4 * (\text{bias} - 1) * \text{geo_rand}$

5.6.4 杂交算子

PostgreSQL 中的遗传算法提供了边重组杂交、部分匹配杂交、循环杂交、基于位置的杂交和顺序杂交等多种杂交算子，用于从父辈种群中产生新的后代个体，默认使用的是边重组杂交算法。

1. 边重组杂交 ERX (edge recombination crossover)

边重组杂交算子由函数 gimme_edge_table 和 gimme_tour 实现。函数 gimme_edge_table 用来计算边关系，如步骤 1 和步骤 2 所示；函数 gimme_tour 由边关系得到后代，如步骤 3 所示。

边重组杂交的过程如下：

- 1) 两个父体中的基因构成循环队列，如图 5-47 所示。

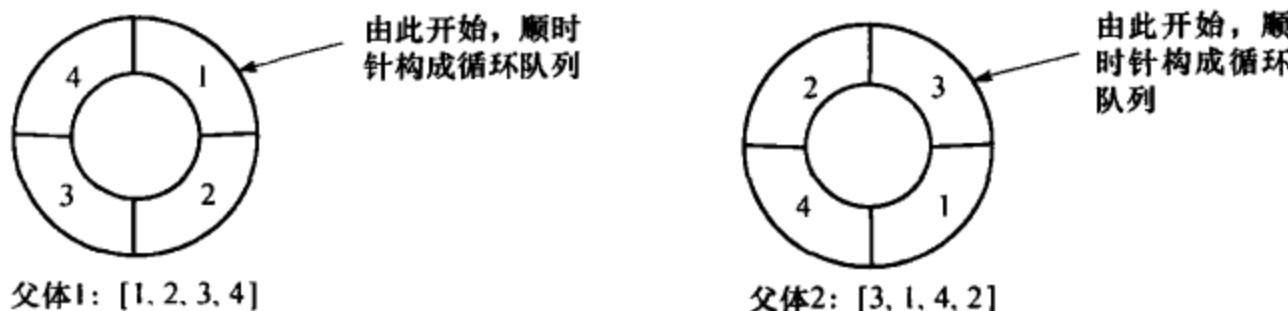


图 5-47 边杂交重组—循环队列

2) 确定父体中的边关系。在步骤1的循环队列中，任意一个基因和相邻的基因构成“边关系”。如果某“边关系”在父体1和父体2中同时存在，则称为“共享边”，用负值表示。例如，在图5-48中描述了图5-47中各个基因之间的边关系。

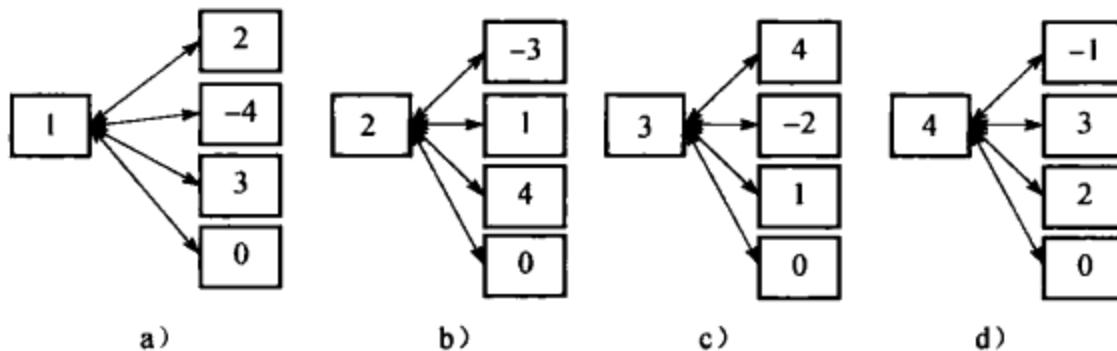


图 5-48 基因的边关系

- a) 从基因1起始的边关系；b) 从基因2起始的边关系；
- c) 从基因3起始的边关系；d) 从基因4起始的边关系

每一个基因的边关系存储在数据结构 Edge (数据结构5.22) 中，所有基因的 Edge 结构按照编号顺序排列组成一个数组 edge_table (边表)。Edge 中的 edge_list 字段是一个表示边关系的数组，长度为 4 (每个基因在两个父体中最多只有 4 个边关系)。如果该基因 X 与某个基因 Y 有边关系，则在基因 X 的 edge_list 中找一个没有使用的元素填上基因 Y 的编号，如果 X 和 Y 的边关系在两个父体中都存在，则在 X 的 edge_list 中填 Y 的编号时要加上负号。edge_list 中没有使用的元素都填 0 值。Edge 中的 unused_edges 表示未使用边的数目，所谓“未使用”是指还没有被用来杂交的边，一旦一个边在杂交过程中被使用过，就会被从 Edge 中清除。例如图 5-48a 中从基因 1 起始的边关系在数据结构层面的表示为：

```
edge_table[1].edge_list[1] = 2;
edge_table[1].edge_list[2] = -4;
edge_table[1].edge_list[3] = 3;
edge_table[1].edge_list[4] = 0;
edge_table[1].total_edges = 3;
edge_table[1].unused_edges = 3;
```

这说明 1 与 2、4、3 都存在边关系，并且 1 与 4 构成了“共享边”。

3) 由边关系得到后代。边重组杂交的基本思想是：随机地选择一个基因作为起始点，顺着它的边关系找到下一个基因（优先考虑共享边），再顺着找到的基因的边关系找到第三个基因，直到找到的基因能够构成一个个体为止，最后将找到的基因按找到的顺序组成一个个体即可。具体的杂交流程如图 5-49 所示。按照以上方法，图 5-48 中的两个父体可产生的一个后代为 [1, 4, 3, 2]。

数据结构 5.22 Edge

```
typedef struct Edge
{
    Gene     edge_list[4];           // 边关系列表
    int      total_edges;          // 总的边数目
    int      unused_edges;         // 未使用边的数目
} Edge;
```

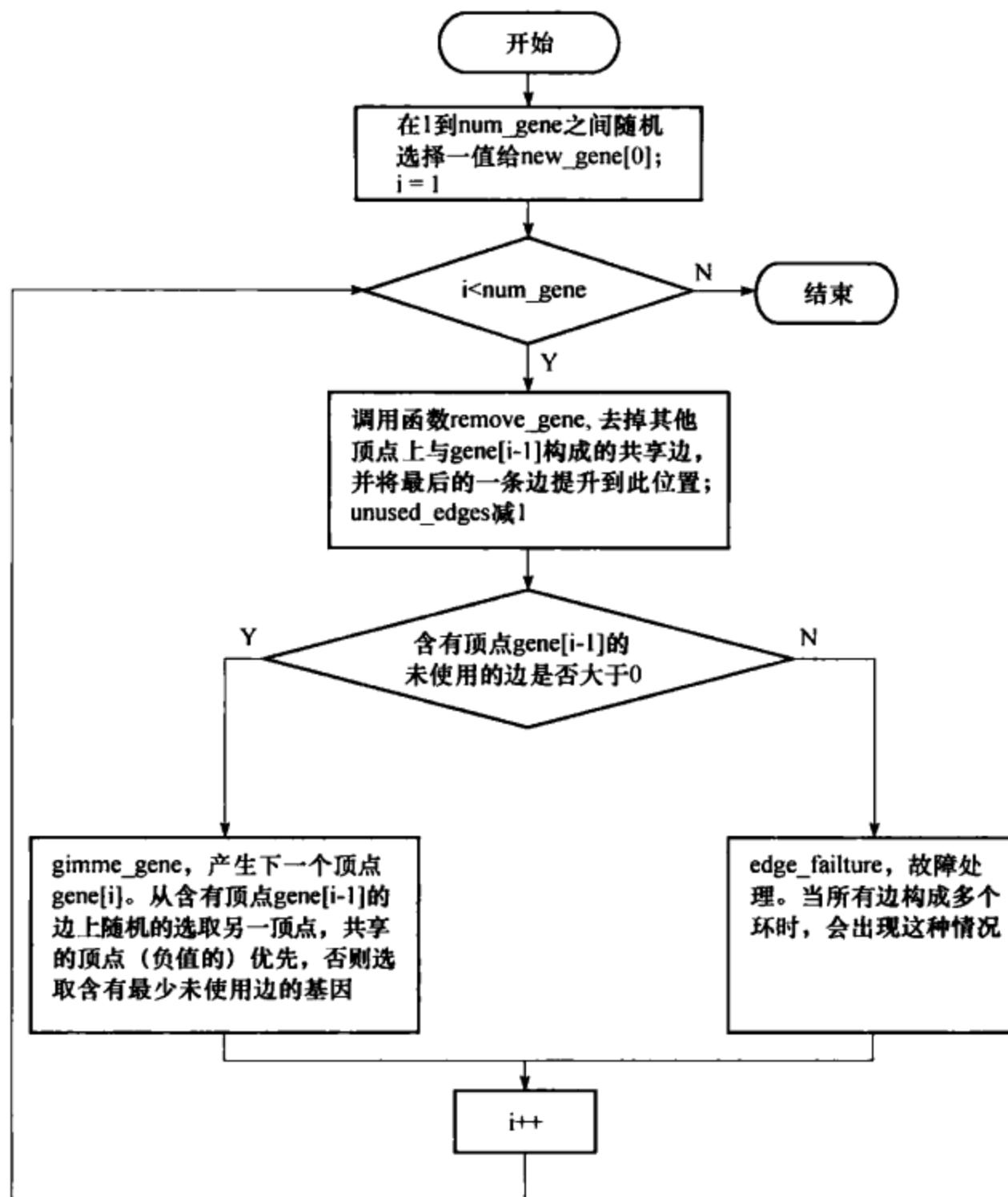


图 5-49 边重组杂交流程

2. 部分匹配杂交

部分匹配杂交 (Partially Matched Crossover, PMX) 算子由函数 pmx 实现，该函数的流程如下：

1) 在字串上均匀随机地选择两点，由这两点确定的子串称为映射段，定义两个整型变量 left 和 right ($left < right$) 表示选取的映射段的起始边界。

- 2) 用父体 2 的映射段替换父体 1 的映射段产生原始后代。
- 3) 确定两映射段之间的映射关系。
- 4) 根据映射关系将后代合法化。

例如，有两个父体，对它们用 pmx 进行杂交的过程如下：

- 1) 随机选择两个点 (下划线部分就是映射段)：

父体 1： [1, 2, 3, 4, 5, 6, 7, 8, 9]

父体2: [5, 4, 6, 9, 2, 1, 7, 8, 3]

2) 交换双亲中的子串:

后代: [1, 2, 6, 9, 2, 1, 7, 8, 9]

3) 确定映射关系。根据交换部分的上下对应关系所生成映射关系: “ $1 \longleftrightarrow 6 \longleftrightarrow 3, 2 \longleftrightarrow 5, 9 \longleftrightarrow 4$ ”。

4) 用映射关系将产生的新个体变成合法的个体(个体中的数字不能重复出现)。例如, 数字2出现了两次, 产生了冲突, 由于2和5是映射的, 因此将非交换部分的2替换成5, 其他的冲突做类似处理。得到最终的后代:

后代: [3, 5, 6, 9, 2, 1, 7, 8, 4]

3. 循环杂交

循环杂交算子由函数 cx 实现, 该方法从一个双亲中取一些基因, 而其他的基因则取自另外一个双亲。该方法首先随机确定一个初始位置, 作为当前位置, 把父体1当前位置下的基因(即编号)赋值给当前位置下的子代, 并标记此基因已被使用。然后推进当前位置到父体2当前位置(未修改前的位置)下的基因在父体1中的位置, 同样把父体1当前位置下的基因赋值给当前位置下的子代。同理推进当前位置, 生成子代的基因片段, 直到循环到初始位置下父体1的基因与父体2中当前位置下的基因一样时循环结束。如果循环结束后仍有部分基因没有被使用, 则将父体2中的这些基因按在父体2中存在的位置赋值到子代中同样位置。例如有两个父体, 对它们用 cx 进行杂交的过程如下:

1) 随机选择两个父体:

父体1: [1, 2, 3, 4, 5, 6, 7, 8, 9]

父体2: [4, 1, 2, 8, 7, 6, 9, 3, 5]

2) 循环查找:

①首先随机选择一个位置, 假设为1。将父体1中1号位置的基因(1)放入后代的1号位置, 此时后代为(x表示该位置上还没有放置基因):

[1, x, x, x, x, x, x, x, x]

②以上一步获得的基因编号(1)为位置在父体2中取出相应的基因(4), 并以其编号为位置放入后代中, 即将基因4放入后代的4号位置, 此时后代为:

[1, x, x, 4, x, x, x, x, x]

③在父体1中取出4号位置的基因4放入后代, 由于基因4已经在后代中了, 因此实际并不对后代进行修改, 此时后代仍为:

[1, x, x, 4, x, x, x, x, x]

④在父体2中取出4号位置的基因8放入后代的8号位置, 此时后代为:

[1, x, x, 4, x, x, x, 8, x]

⑤在父体1中取出8号位置的基因3放入后代的3号位置, 此时后代为:

[1, x, 3, 4, x, x, x, 8, x]

⑥在父体2中取出3号位置的基因2放入后代的2号位置, 此时后代为:

[1, 2, 3, 4, x, x, x, 8, x]

⑦在父体1中取出2号位置的基因2放入后代的2号位置, 此时后代为:

[1, 2, 3, 4, x, x, x, 8, x]

⑧最后由于父体 2 的 2 号位置的基因为 1，与最初放入后代的基因一致，因此循环结束，得到的后代为：

[1, 2, 3, 4, x, x, x, 8, x]

3) 用父体 2 对应位置的基因填满后代中的空位得到后代：

后代：[1, 2, 3, 4, 7, 6, 9, 8, 5]

4. 基于位置的杂交

基于位置的杂交算子由函数 px 实现，该函数的处理流程如下：

1) 根据基因数目 gene，首先随机选择一个区间 $[gene/3, 2/3 * gene]$ 中的整数 p，从父体 1 中随机选择 p 个基因，设为集合 A。

2) 对于在 A 中的基因，从父体 1 拷贝到后代中对应的位置上。

3) 对于不在 A 中的基因，按它们在父体 2 中的顺序拷贝到后代剩余的位置上。

例如，有两个父体，对它们用 px 进行杂交的过程如下：

1) 随机选择两个父体：

父体 1：[5, 3, 2, 1, 4, 6]

父体 2：[2, 4, 6, 5, 1, 3]

2) 基因数目为 6，随机生成一个 [2, 4] 之间的整数 p，假设为 3，从父体 1 中随机的选择 3 个基因是 2、4、6，将它们拷贝到后代中对应的位置上：

后代为：[x, x, 2, x, 4, 6]

3) 剩余的位将从父体 2 中产生，并按它们在父体 2 中的顺序排列：

最后的后代为：[5, 1, 2, 3, 4, 6]

5. 顺序杂交 (ox1)

顺序杂交算子有两种：OX1 和 OX2。OX1 算子由函数 ox1 实现，可以将它看做是一种带有不同修复程序的 PMX 的变型。函数 ox1 的处理流程如下：

1) 从第一个父体中随机选择一个子串。

2) 将子串复制到一个空串的相应位置，产生一个原始后代。

3) 删去第二个父体中子串已有的基因，得到原始后代需要的其他基因的顺序。

4) 按照这个基因顺序，从左到右将这些基因复制到后代的空缺位置上。

例如，有两个父体，用 ox1 对它们进行杂交的过程如下：

1) 随机选择两个父体。

父体 1：[1, 2, 3, 4, 5, 6, 7, 8, 9]

父体 2：[5, 7, 4, 9, 1, 3, 6, 2, 8]

2) 从父体 1 中选择一个子串，假设为 3, 4, 5, 6。将这个子串复制到原始后代对应的位置中。

后代：[x, x, 3, 4, 5, 6, x, x, x]

3) 删除父体 2 中从父体 1 选择的子串，得到的基因顺序是 7、9、1、2、8，将剩余的基因按从左到右的顺序复制到后代空缺的位置上得到最后的后代。

后代：[7, 9, 3, 4, 5, 6, 1, 2, 8]

6. 顺序杂交 (ox2)

OX2 算子由函数 ox2 实现，其处理流程如下：

1) 根据基因的数目 gene, 随机选择一个 $[gene/3, 2/3 * gene]$ 之间的整数 p, 从父体 1 中随机选择 p 个基因, 设为集合 A。

2) 对于父体 2 中不在 A 中基因, 拷贝到后代对应的位置中。

3) 对在 A 中的基因, 按它们在父体 1 中的顺序依次拷贝到后代剩余的位中。

例如, 有两个父体, 对它们用 ox2 进行杂交的过程如下:

1) 随机选择两个父体:

父体 1: [2, 3, 5, 4, 1, 6]

父体 2: [5, 3, 1, 4, 6, 2]

2) 基因数目为 6, 因此随机生成一个 [2, 4] 之间的整数 p, 假设为 3。从父体 1 中随机选择的 3 个基因是 5、1、6, 那么对于父体 2 中 2、3、4, 它们的位置不变:

后代为: [x, 3, x, 4, x, 2]

3) 对于 5、1、6, 按照他们在父体 1 中的顺序拷贝到剩余的位中。

最后后代为: [5, 3, 1, 4, 6, 2]

5.6.5 变异算子

变异算子由函数 geqo_mutation 实现, 该函数随机地从父体中产生两个变异位, 交换这两个变异位的值, 执行 num_gene (基因数目) 次这样的交换操作。

对一个父体进行变异的过程如下:

1) 选择一个父体:

[1, 2, 3, 4, 5, 6, 7, 8, 9]

2) 随机选择两个变异位 (有下划线的位):

[1, 2, 3, 4, 5, 6, 7, 8, 9]

3) 交换这两个基因位, 得到后代:

[1, 2, 6, 4, 5, 3, 7, 8, 9]

4) 继续执行步骤 2、3 的操作, 做 num_gene-1 次。

5.6.6 终止条件

遗传算法采用设定演化代数的方法, 但演化到一定数量的代数时, 就停止演化。默认的演化代数是种群的大小 (pool_size, 缓冲池的大小)。演化代数的计算涉及下面两个参数:

1) geqo_effort: 整型变量, 是用于限制种群大小的影响因子。取值范围是 [1, 10], 默认值为 5。

2) geqo_pool_size: 整型变量, 表示缓冲池 (用于存储种群中的个体) 大小, 缓冲池的大小和种群大小相同, 其值至少为 2。

种群大小 (缓冲池的大小) pool_size 由函数 gimme_pool_size 确定, 其参数 nr_rel 为查询中表的数量。计算方法如下:

1) 计算上限值 maxsize 和下限值 minsize:

maxsize = 50 * geqo_effort

minsize = 10 * geqo_effort

2) 计算基准大小 $\text{size} = \text{pow}(2.0, \text{nr_rel} + 1.0)$ 。

3) 如果基准大小位于上限值和下限值之间，则取基准大小作为种群大小；如果低于下限值，则取下限值；如果高于上限值，则取上限值。

5.6.7 基于排列生成路径

在遗传算法中由排列生成连接路径是按照如图 5-50 所示的算法实现的。其中变量 `rels` 中保存了按照排列数所对应的各个基本关系（即基本表），变量 `rels_temp` 用来临时保存当前不可连接的关系（即后文提到的代码中的 `clumps` 链）。如果两个表存在如下几种情况：内连接、左外连接、右外连接、全连接以及 IN/EXISTS 子链接（子链接也是转换成表的连接来处理），则认为它们是可连接的（joinable）。而不可连接指的是只能用笛卡儿积的连接。

该算法处理过程为：依次取 `rels` 中的所有基本表，在 `rels_temp` 中依次寻找可与其连接的表，如果存在可连接的表，则把这两个表进行连接生成新的表，并从头开始继续在 `rels_temp` 中寻找可连接的表，并将其与新生成的表连接，一直到找不到可连接的表为止，最后将最终生成的新表插入到 `rels_temp` 结尾；继续在 `rels` 中取下一个表，重复以上过程。

5.6.8 实例分析

在 5.4.8 节中我们已经针对例 5.1 的查询树介绍了利用动态规划算法进行路径生成的方法，这里我们将介绍一下用遗传算法来对例 5.1 的查询树生成路径的过程。在例 5.1 中，连接路径的生成涉及了表 `sc`、`class` 及 `course`，其范围表的索引分别为 1、4、5。为了方便起见，这里给它们重新编号为 1 (`sc`)、2 (`class`)、3 (`course`)。

根据 5.6.6 节中的计算方法，可以得到迭代次数（缓冲池大小）为 50。接下来将随机生成 1、2、3 构成的各种排列（每个排列是一个个体）来填满缓冲池。注意，由于遗传算法启动的第一条件是连接的表个数超过 12 个，而我们这个例子中只有 3 个表，因此它们的全部排列也只有 6 种，是不可能填满缓冲池的。但为了更容易让读者理解遗传算法的工作过程，我们假设在这个例子上能够使用遗传算法来生成路径。

(1) 第一步：计算连接代价最小的排列

首先进行第一次迭代过程如下：

1) 利用公式 5.2 计算出要选择的父体在缓冲池中的索引。得到的第一次迭代的两个排列为：

由排列生成连接路径的算法 `creat_path(rels[rels_num])`

```

1 rel=NULL
2 rel_temp=NULL
3 rels_temp[rels_num]→NULL;
4 n→rels_num;
5 FOR i→1 TO rels_num DO
6   rel→rels[i];
7   IF rels_temp is null THEN
8     rels_temp[1]→rel
9   ELSE
10    FOR j→1 TO n DO
11      rel_temp→rels_temp[j]
12    IF rel and rel_temp are joinable THEN
13      rel→rel join rel_temp;
14        delete rel_temp from rels_temp;
15      j→1; n→n-1;
16    CONTINUE;
17  ENDIF
18 ENDFOR
19 insert rel into the end of rels_temp;
20 ENDIF
21 ENDFOR
22 RETURN rel

```

图 5-50 遗传算法中由排列生成连接路径的算法

[1, 3, 2] 和 [3, 2, 1]。

2) 为选出的父体准备边表，得到边表如图 5-51 所示。

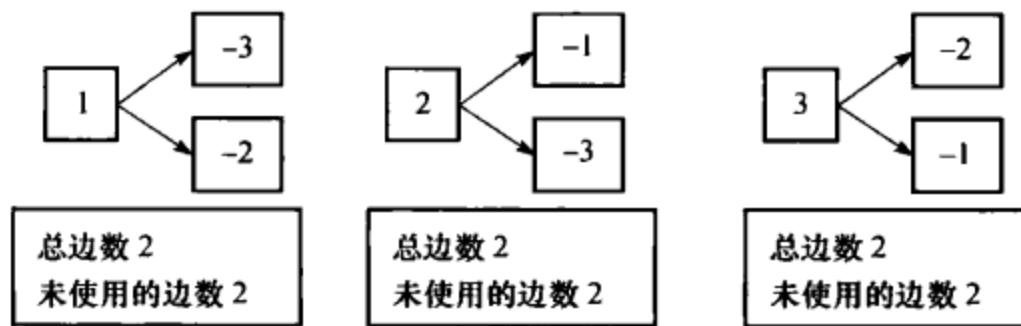


图 5-51 第一次迭代时的初始边表

3) 使用边重组进行杂交。随机选择一个初始起点，假设为 2。删除其他边表中以 2 为终点的边并将未使用边减 1。此时边表如图 5-52 所示。

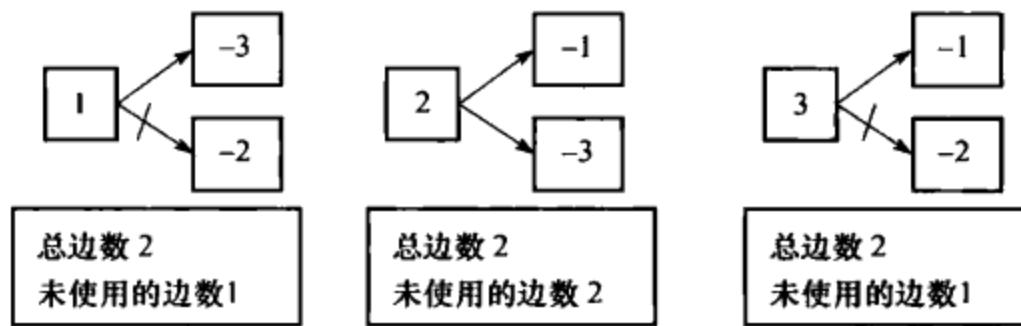


图 5-52 边重组杂交找到第一个基因后的边表

4) 在以 2 为起点的边表中选择一条边，即找出下一节点。规则如下：

- 如果有负值的终点，则优先选择。因为这种边是在两个父体中共享的。说明此两个节点对应的表间的联系比较紧密。
- 当无负值的终点时，找出终点对应的边表的未使用边数最少的终点。如果有多个则随机选择一个。未使用边最少，在某种程度上说明此点对应的表不是与其他很多表都存在某种联系，只与某几个表联系紧密，所以要首先选择出来。

所以，第二个节点是 1。置 2 为起点的边表未使用边为 -1，并删除在其他边表中以 1 为终点的边。此时的边表如图 5-53 所示。

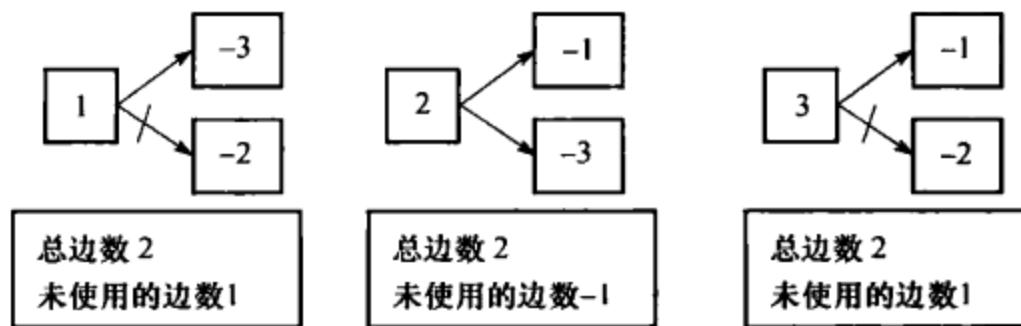


图 5-53 边重组杂交找到第二个基因后的边表

同理，第三个节点是 3。此时边重组完成，边表如图 5-54 所示。生成的新的排列（即后代）为 [2, 1, 3]。

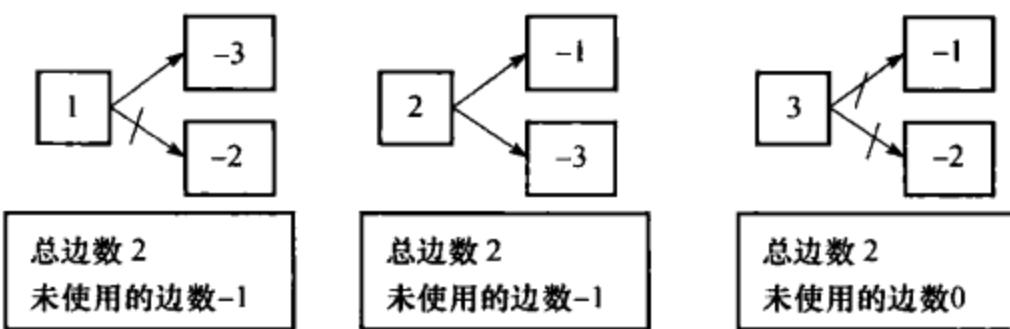


图 5-54 边重组杂交找到第三个基因后的边表

5) 计算后代的代价 (这将在下面给出如何计算代价), 并插入到缓冲池中用以替换缓冲池中代价最高的个体。

到此, 一次迭代完成。持续此过程 50 次。然后从缓冲池中取出第一个个体 (即代价最小的一个), 假设为 [1, 3, 2], 接下来就可以根据该个体中关系的排列来生成连接路径。

(2) 第二步: 基于代价最小的排列生成路径

下面给出排列 [1, 3, 2] 生成路径的过程。首先选择 1 节点对应的表, 即 sc 表。将其关系信息插入到一个 clumps 链中。然后选择 3 节点对应的表, 即 course 表, 与 clumps 中的 sc 比较, 看是否可构成连接。这需要到 PlannerInfo 的 join_info_list 链中寻找满足的特殊连接信息。如图 5-55 所示, 有一个特殊的连接信息是 IN 子查询的, 对应着 “sc.cno IN (SELECT course.cno FROM course WHERE course cname = '高等数学')” 这一条件。

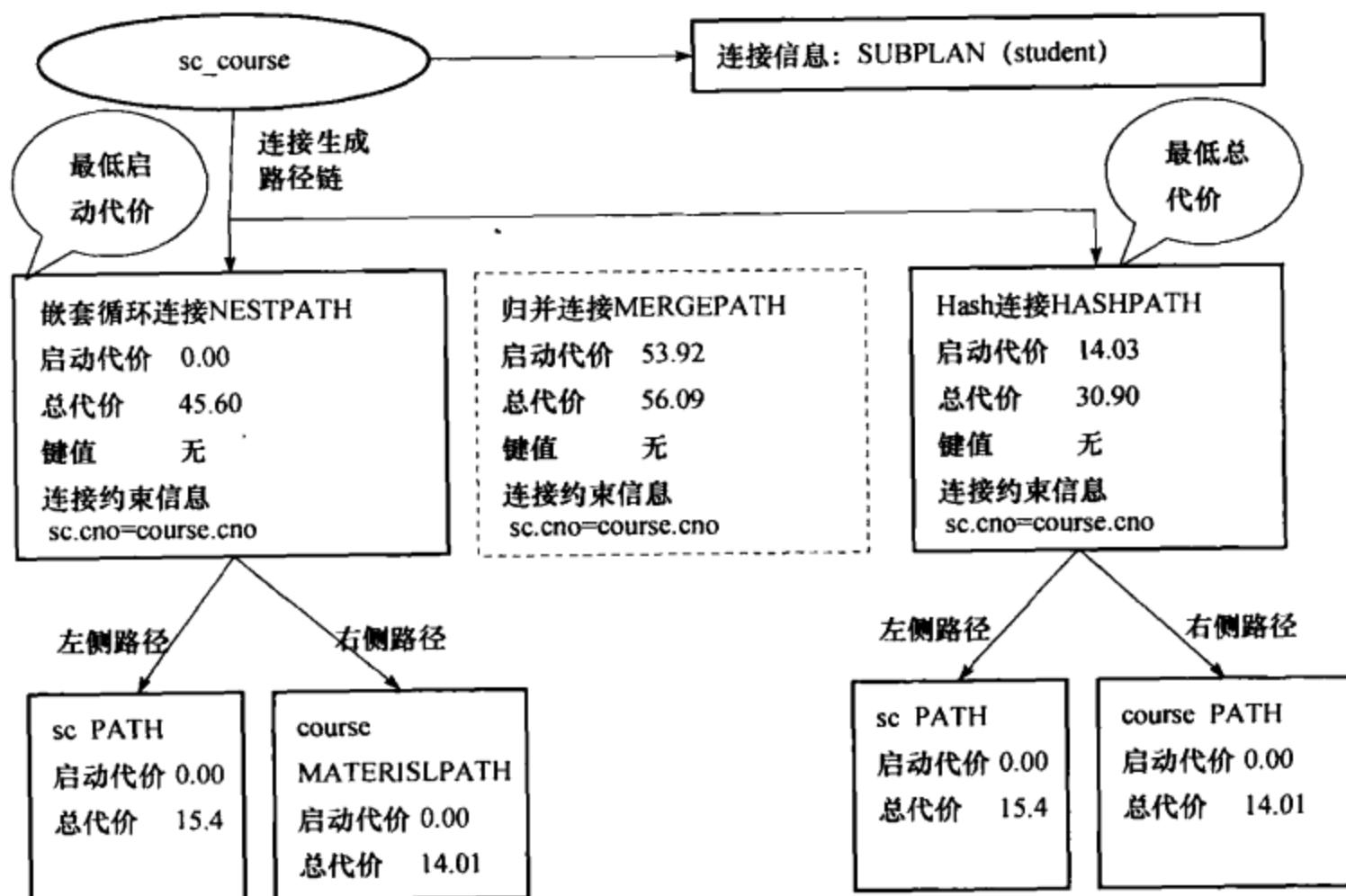
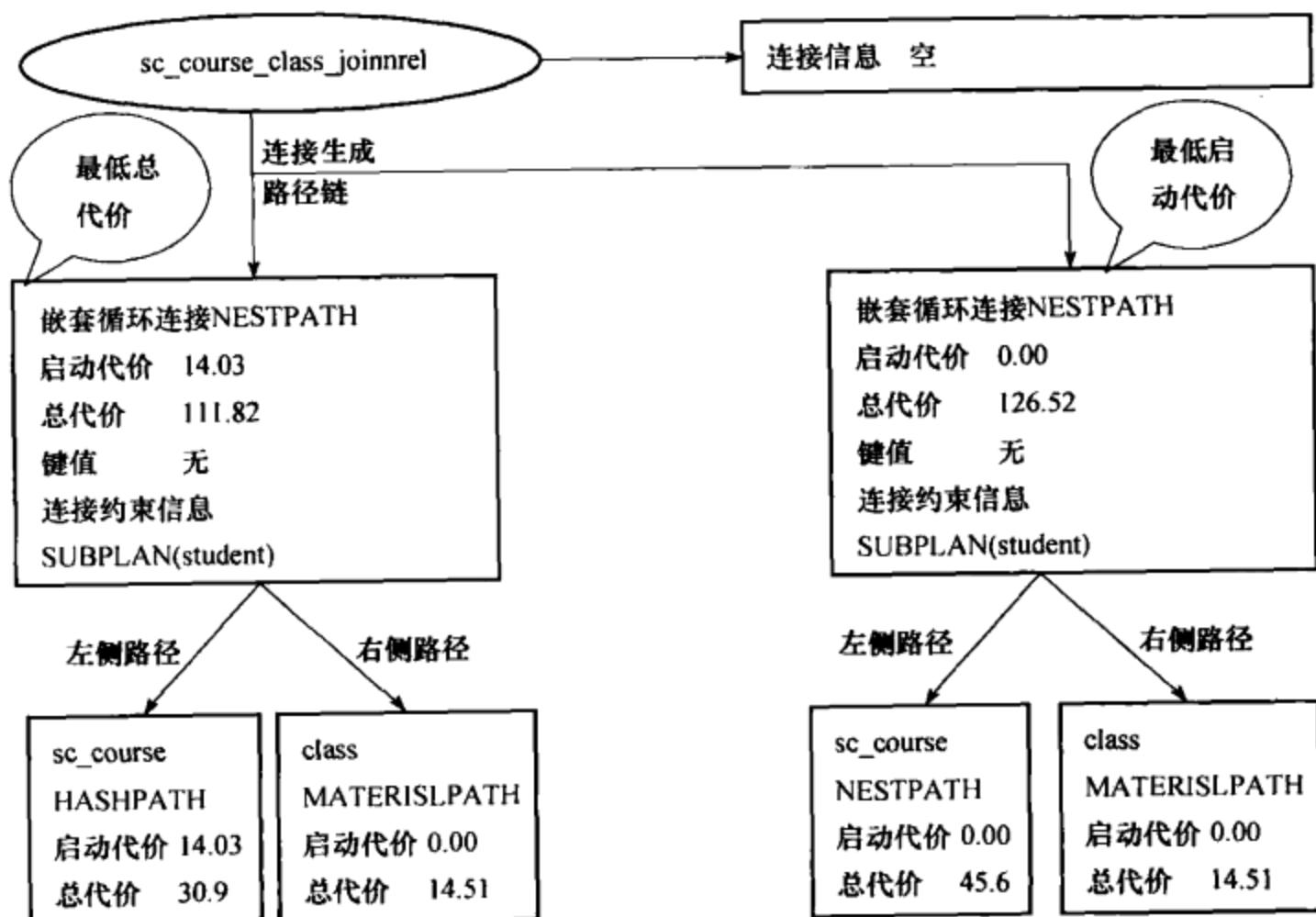
```
{SPECIALJOININFO
:min_lefthand (b1) //左侧是sc表
:min_righthand (b5) //右侧是course表
:syn_lefthand (b14) //整个查询语句中处于左边的有sc及class
:syn_righthand (b5) //整个查询语句中处于右边的有course
:jointype 4          //IN (select) 连接类型
```

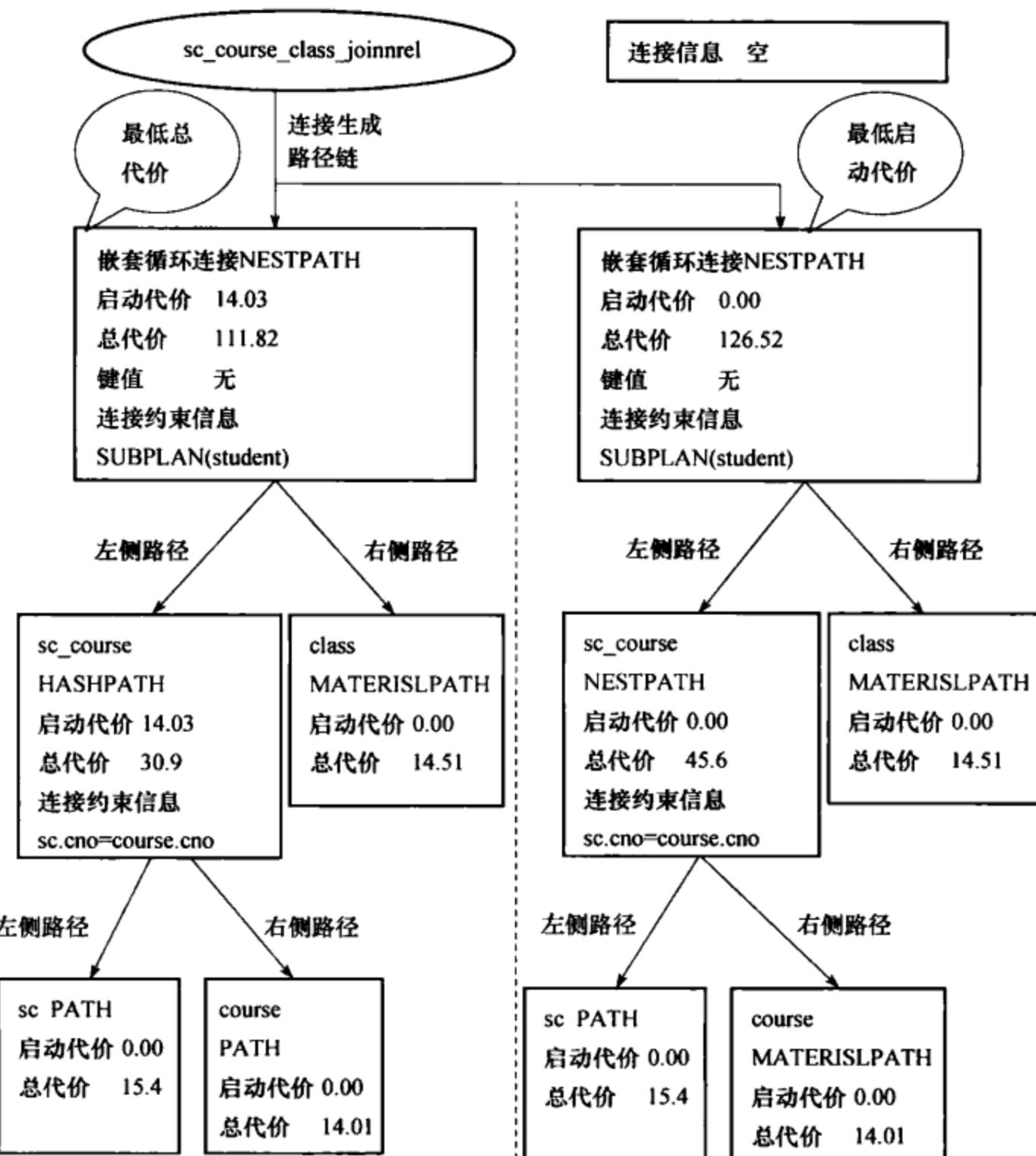
图 5-55 特殊连接信息实例

因此, sc 及 course 的连接是合法的。图 5-56 中计算了表 sc 和 course 的三种连接方式的代价, 其中嵌套循环连接和 Hash 连接分别在启动代价和总代价上占有优势, 而归并连接的这两个代价都比较大, 因此, 将嵌套循环连接和 Hash 连接都保存作为备选路径, 而归并连接被淘汰掉。

将此连接表 (sc_course) 插入到 clumps 链中, 选择排列数 2 所对应的基本表 (即 class 表), 检查是否可以与 clumps 中每个表连接。与特殊连接信息比较后, 发现只可能与新生成的表 sc_course 进行内连接。如图 5-57 所示, 左边和右边分别是表 sc 和表 course 采用嵌套循环连接和 Hash 连接时, 表 sc_course 和表 class 的连接路径。表 sc_course 和表 class 只采用嵌套循环的连接方式。

将表 sc_course 的生成路径添加到图 5-57 中, 最后得到的总的路径如图 5-58 所示。可以发现, 最后生成两条备选路径中, 第一条有最低总代价, 第二条有最低启动代价, 它们有各自的优势。在最后, 会将 pathkey (依据 pathkey 确定元组的排序输出顺序) 考虑进来, 确定最终路径。

图 5-56 表 `sc` 和 `course` 的连接路径图 5-57 表 `class` 和连接表 `sc_course` 的连接路径

图 5-58 表 `sc`、`course`、`class` 的连接路径

5.7 小结

数据库管理系统中，对性能影响最大的是查询处理器。查询处理器由查询编译器和执行器两部分组成，而查询编译器又包括查询分析器、查询预处理器和查询优化器。规划器/优化器的任务是创建一个优化了的执行计划。它首先生成完成查询所有可能的路径。这样创建的所有路径都产生相同的查询结果，而优化器的任务就是计算每个路径的开销并且找出开销最小的那条路径。

从本章介绍的几种优化算法的比较中可以知道，PostgreSQL 采用的优化算法是以穷尽搜索算法为基础的，并在此基础上进行了改进，同时还采用了遗传算法作为辅助，系统能够有选择地采用不同的优化算法，从更大程度上达到查询优化的目的。从实际应用来看，这种查询优化算法在 Post-

greSQL 上能够很好地工作，也能取得很高的优化效率。

习题

习题 5.1 基于 5.3.1 节中提供的函数、表和视图，在 Linux 环境下用 GDB 调试执行以下查询命令：

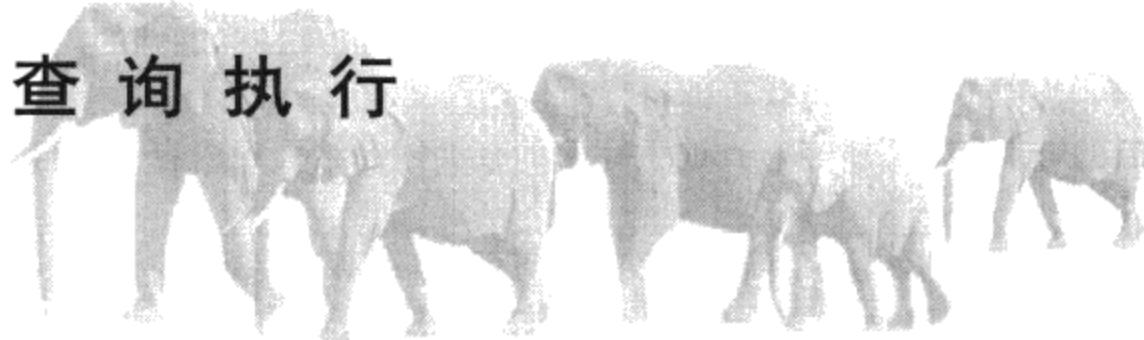
```
SELECT* FROM shoelace;
```

并完成以下问题：

- 1) 经过查询分析后，生成的查询树中（Query 结构），范围表、连接树以及目标属性的内容，并解释各项的含义和用途。（提示：为便于查看，可将查询树打印到日志文件。）
- 2) 基于问题 1 中查询分析后的查询树，分析经过查询重写后的查询树的结构和内容，并与查询重写之前的进行比较，有何不同？
- 3) 执行该查询命令使用了哪种路径生成算法？并详细说明其代价估计的计算过程。

第6章

查询执行



查询编译器将用户提交的 SQL 查询语句转变成执行计划之后，由查询执行器继续执行查询的处理过程。在查询执行阶段，将根据执行计划进行数据提取、处理、存储等一系列活动，以完成整个查询执行过程。查询执行过程更像一个结构良好的裸机，执行计划为输入，执行相应的功能。因此，本章重点介绍查询执行器的框架结构、执行方式，结合实例为读者说明数据库执行计划的相关步骤，帮助读者进一步理解 PostgreSQL 的查询执行过程。

查询执行器的框架结构如图 6-1 所示。同查询编译器一样，查询执行器也是被函数 `exec_simple_query` 调用，只是调用的顺序上查询编译器在前，查询执行器在后。从总体上看，查询执行器实际就是按照执行计划的安排，有机地调用存储、索引、并发等模块，按照各种执行计划中各种计划节点的实现算法来完成数据的读取或者修改的过程。

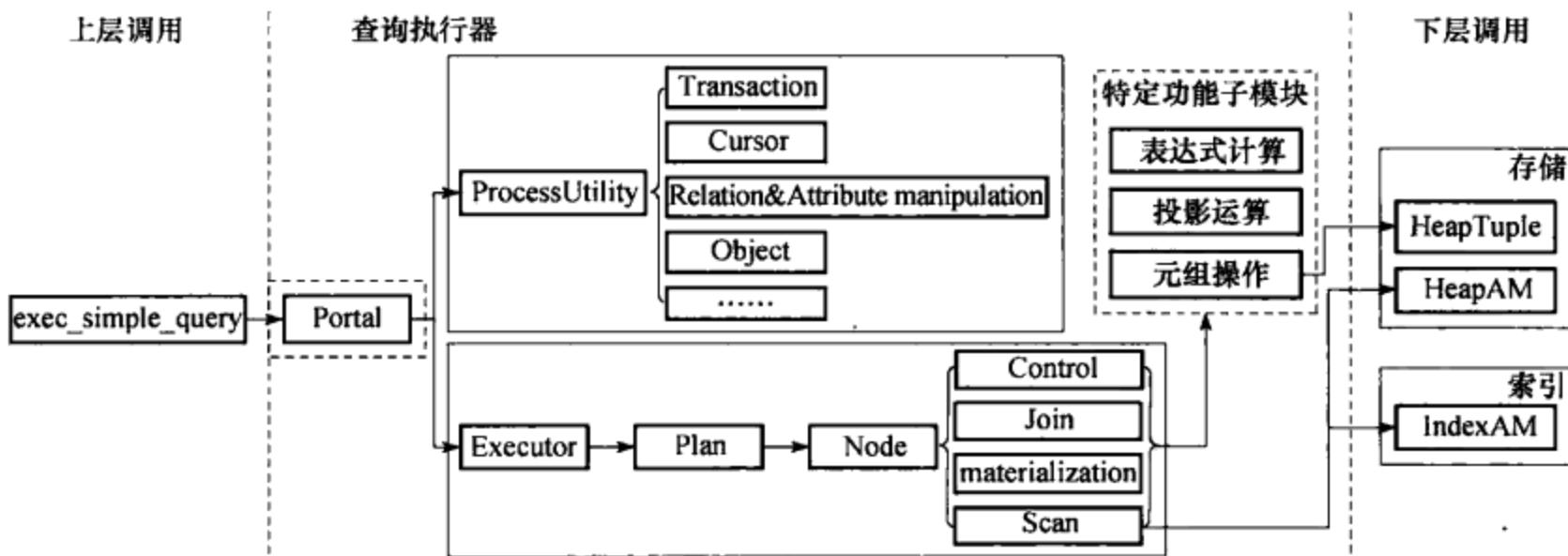


图 6-1 查询执行器的框架结构

如图 6-1 所示，查询执行器有四个主要的子模块：Portal[⊖]、ProcessUtility、Executor 和特定功能

⊖ PortalData 是此部分的核心数据结构，而 PostgreSQL 实现采用了层次模型，因此，此部分被称为 Portal 部分或 Portal 层。PostgreSQL 中很多模块以核心数据结构命名。

子模块部分。由于查询执行器将查询分为两大类别，分别由子模块 ProcessUtility 和 Executor 负责执行，因此查询执行器会首先在 Portal 模块根据输入执行计划选择相应的处理模块（Portal 模块也称为策略选择模块）。选择执行策略后，会将执行控制流程交给相应的处理部件（即 ProcessUtility 或 Executor），两者的处理方式迥异，执行过程和相关数据结构都有很大的不同。Executor 输入包含了一个查询计划树（Plan Tree），用于实现针对于数据表中元组的增删查改等操作。而 ProcessUtility 处理其他各种情况，这些情况间差别很大（如游标、表的模式创建、事务相关操作等），所以在 ProcessUtility 中为每种情况实现了处理流程。当然，在两种执行模块中都少不了各种辅助的子系统，例如执行过程中会涉及表达式计算、投影运算以及元组操作等，这些功能相对独立，并且在整个查询执行过程中会被重复调用，因此将其单独划分为一个模块（特定功能子模块）。

本章将采用由外到内的顺序介绍各个模块。首先介绍策略选择，然后讲解两种执行方式的原理与实现，最后介绍内部各种独立子功能。每个部分都将会是以基础知识和数据结构为基础，然后讲述实现的原理和过程，读者可根据自身需求进行选择阅读。

6.1 查询执行策略

执行流程进入查询执行阶段后，会为每种执行计划选择相应的处理过程，执行相应的处理，最后根据要求返回结果。

图 6-2 显示了 PostgreSQL 计划生成与处理的基本过程。SQL 语句会被查询编译器转换成为两种基本类型的数据结构（执行计划树和非计划树操作），并在执行过程中为其选择合适的执行部件（Executor 或 ProcessUtility）进行处理。DML 语句（SELECT、INSERT、UPDATE、DELETE）会被查询编译器转换成执行计划树，因为 DML 语句的执行过程十分相近，并且都可以使用统一的数据结构加以表示和处理，所以被作为一种特殊的方式，单独进行优化和处理。而其他类型的语句被归为非计划树操作，使用另一处理流程进行处理。然而，有些复杂的 SQL 语句不能简单地归类为一种基本处理类型，在 PostgreSQL 实现中，会将其解析成两种基本类型数据结构的序列（实际是按顺序构成一个链表），并将其顺序执行来完成用户的请求操作。

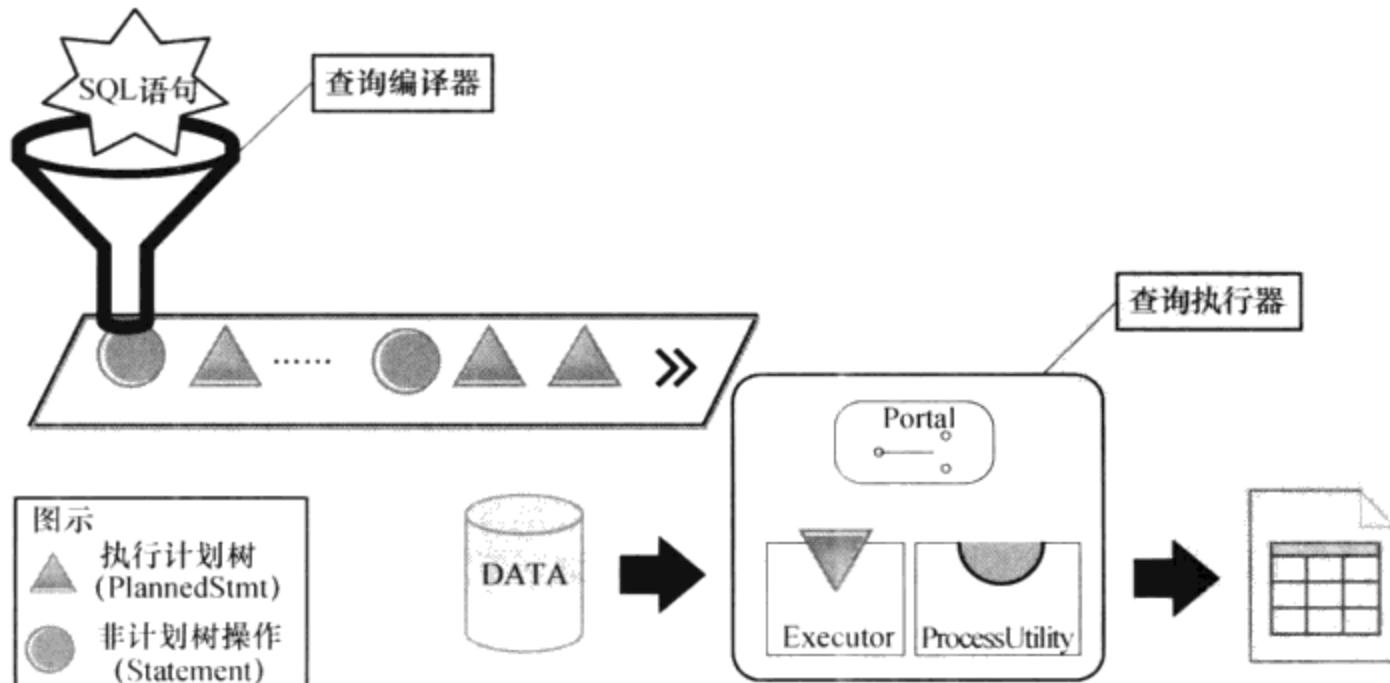


图 6-2 查询语句处理过程

从查询编译器输出执行计划，到执行计划被具体的执行部件处理这一过程，被称作执行策略的选择过程，负责完成执行策略选择的模块称为执行策略选择器。该部分完成了对于查询编译器输出数据的解析，选择预先设定好的执行流程。本节将重点介绍几种执行策略及执行策略的选择过程。

6.1.1 可优化语句和数据定义语句

从上面的介绍可以看出，在 PostgreSQL 数据库中，用户输入的 SQL 语句被分成两种类型，并分别被两种不同的执行部件处理。这两种语句分别被称为：可优化语句（Optimizable statement）和数据定义语句（DDL statement）。

可优化语句主要包括 DML 语句。这类语句的特点是均需要查询相关满足条件的元组，然后将这些元组返回给用户（有可能是经过计算的元组）或者在这些元组上进行某些操作之后写回到磁盘。因此在经过查询编译器处理后，会为其生成一个或多个执行计划树（Plan Tree）[⊖]，用于查询满足相关条件的元组并做相应处理。由于在执行计划树的生成过程中会根据查询优化理论进行重写和优化，以加快查询速度，因此，这类语句被称作可优化语句。

数据定义语句则包括数据表创建等操作，这类语句包含查询数据元组以外的各种操作。语句之间功能相对独立，所以也被称作功能性操作。

针对这两种类型的语句，PostgreSQL 实现了两个处理模块：执行器（Executor）和功能处理器（Utility Processor）。

1) 执行器处理可优化语句，可优化语句包含一个或多个经过重写和优化过的查询计划树，执行器会严格根据计划树进行处理。执行器函数名为 ProcessQuery，主要代码被放置在 src/backend/executor/ 目录中。

2) 功能处理器用于处理数据定义语句，可根据不同类别的功能调用相应的处理函数。功能处理器函数名为 ProcessUtility，各种数据定义语句的具体实现被放置在 src/backend/commands/ 目录中。

6.1.2 四种执行策略

在 PostgreSQL 的实现中，简单的 SQL 语句会被查询编译器转化为一个执行计划树或者一个非计划树操作，而比较复杂的 SQL 语句则可能会被转化成多个执行计划树或者非计划树操作的序列。如图 6-3 所示，左侧的 SQL 语句将被查询编译器转换为右边的三个操作。我们把由查询编译器输出的每一个执行计划树或者一个非计划树操作所代表的处理动作称为一个原子操作。由于一个语句可能被转换为一组原子操作，而这些操作不能简单地使用同一个处理部件进行处理，此时执行策略选择器需要根据原子操作的不同调用相应的处理部件。

此外，有些 SQL 语句虽然仅被转换为一种原子操作，但是其执行中由于各种原因需要能够缓存语句执行的结果，等到 SQL 语句执行完成之后，再返回被缓存的结果：

1) 对于可优化语句而言，当执行修改元组的操作时，希望能够返回被修改的元组。由于原子操作的处理过程不能够被可能有问题的输出过程终止（这是有风险的，一旦输出通道传输出现问题，执行过程将被影响），因此不能边执行边返回结果。此时需要有一个缓存结构来临时存放执行结果，等执行完成后进行返回。

[⊖] PostgreSQL 数据库中用于表示查询计划的一种树状结构，将在 6.3 节中详细介绍。

2) 数据定义语句一般是不包含返回结果的，但是 EXPLAIN、SHOW 等 SQL 语句需要返回结果，此时也需要一个缓存结构暂存执行结果，然后由 Portal 为其调用输出过程。

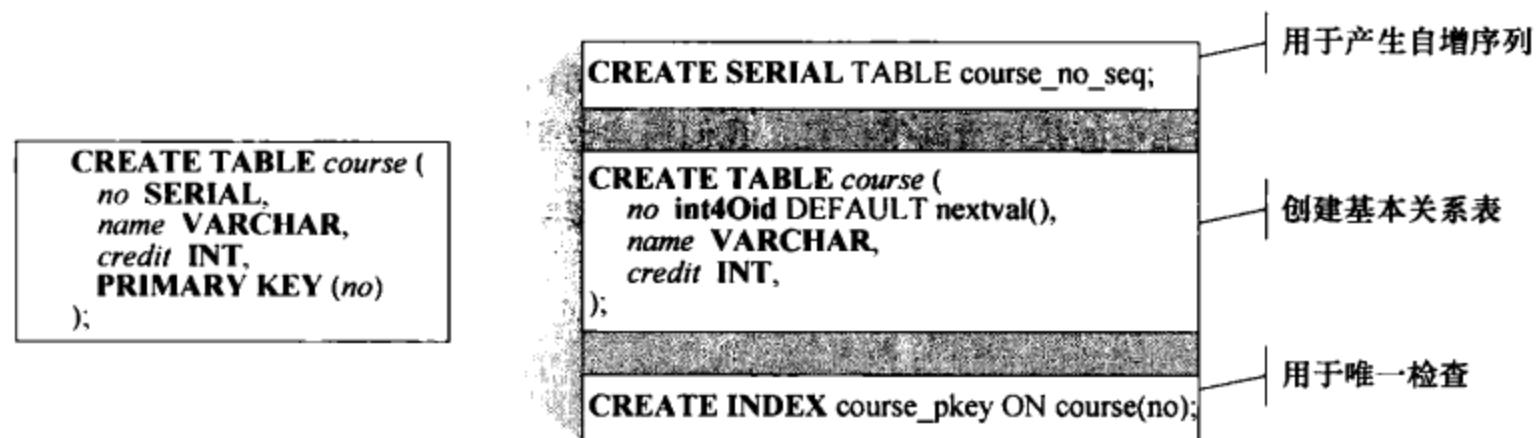


图 6-3 由多个原子操作构成的 SQL 语句例子

针对以上情况，PostgreSQL 实现了不同的执行流程，共分为四类，我们称之为执行策略：

1) PORTAL_ONE_SELECT：如果用户提交的 SQL 语句中仅包含一个 SELECT 类型查询，则查询执行器会使用执行器来处理该 SQL 语句。换句话说，这种策略用于处理仅有一个可优化原子操作的情况。

2) PORTAL_ONE_RETURNING：如果用户提交的 SQL 语句中包含一个带有 RETURNING 子句的 INSERT/UPDATE/DELETE 语句，查询执行器会选择这种策略。因为处理此类语句应先完成所有操作（对元组的修改操作），然后返回结果（例如操作是否成功、被修改的元组的数量等），以减少出错的风险。查询执行器在执行时，将首先处理所有满足条件元组，并将执行过程的结果缓存，然后将结果返回。

3) PORTAL_UTIL_SELECT：如果用户提交的 SQL 语句是一个功能类型语句（数据定义语句），但是其返回结果类似 SELECT 语句（例如 EXPLAIN 和 SHOW），查询执行器将选择这种策略。以此种策略执行时，同样首先执行语句获取完整结果，并将结果缓存起来，然后将结果返回给用户。

4) PORTAL_MULTI_QUERY：用于处理除以上三种情况之外的情况。从其名称中的“MULTI”就能够看出，这个策略更具有普遍性，能够处理一个或多个原子操作，并根据操作的类型选择合适的处理部件。

PORTAL_ONE_RETURNING 和 PORTAL_UTIL_SELECT 两种策略的共同点在于，它们所对应的 SQL 语句的执行过程本来是没有返回结果的，但是由于特殊的需要，会向用户生成类似 SELECT 查询结果的输出，因此这两种策略下的执行需要缓存结果。例如，用 UPDATE 语句更新元组，根据 UPDATE 的语义只需要找到需要更新的元组，然后对其进行修改并写回磁盘即可，并不存在任何返回值。但为了向用户报告执行的状态，PostgreSQL 中的 UPDATE 语句在执行完后会向用户返回此次更新操作修改的元组的数目。而两者的不同在于使用的处理部件不同，前者将使用执行器，后者将使用功能处理器。

简单来说，我们可以认为：

- 1) PORTAL_ONE_SELECT 是用来处理 SELECT 语句的，该策略会调用执行器来执行。
- 2) PORTAL_ONE_RETURNING 面向 UPDATE/INSERT/DELETE 等需要进行元组操作且需要缓存结果的语句，该策略也是调用执行器执行。
- 3) PORTAL_UTIL_SELECT 面向单一数据定义语句，该策略调用功能处理器来执行。
- 4) PORTAL_MULTI_QUERY 是前面三种策略的混合类型，它可以处理多个原子操作。

6.1.3 策略选择的实现

执行策略选择器的工作是根据查询编译器给出的查询计划树链表来为当前查询选择四种执行策略中的一种。在这个过程中，执行策略选择器将会使用数据结构 PortalData（数据结构 6.1）来存储查询计划树链表以及最后选中的执行策略等信息，我们通常也把这个数据结构称为“Portal”。

数据结构 6.1 PortalData（主要字段）

```
typedef struct PortalData
{
    const char *name;           // Portal 的名称
    const char *sourceText;     // 原始 SQL 语句
    List *stmts;               // 查询编译器输出的查询计划树链表
    PortalStrategy strategy;   // 为当前查询选择的执行策略
    PortalStatus status;       // Portal 的状态
    QueryDesc *queryDesc;      // 查询描述符，存储执行查询所需的所有信息
    TupleDesc tupDesc;         // tupDesc 描述可能的返回元组的结构
    .....
}PortalData;
```

查询执行器执行一个 SQL 语句时都会以一个 Portal 作为输入数据，Portal 中存放了与执行该 SQL 语句相关的所有信息（包括查询树、计划树、执行状态等），Portal 及其主要字段之间的结构如图 6-4 所示。其中，stmts 字段是由查询编译器输出的原子操作的链表，在图 6-4 的示例中仅给出了两种可能的原子操作 PlannedStmt 和 Query，两者都能包含查询计划树，用于保存含有查询的操作。当然，有些含有查询计划树的原子操作不一定是 SELECT 语句，例如游标的声明（utilityStmt 字段不为空），以及 SELECT INTO 类型的语句（intoClause 字段不为空）。对于 UPDATE、INSERT、DELETE 类型，含有 RETURNING 子句时 returningList 字段不为空。

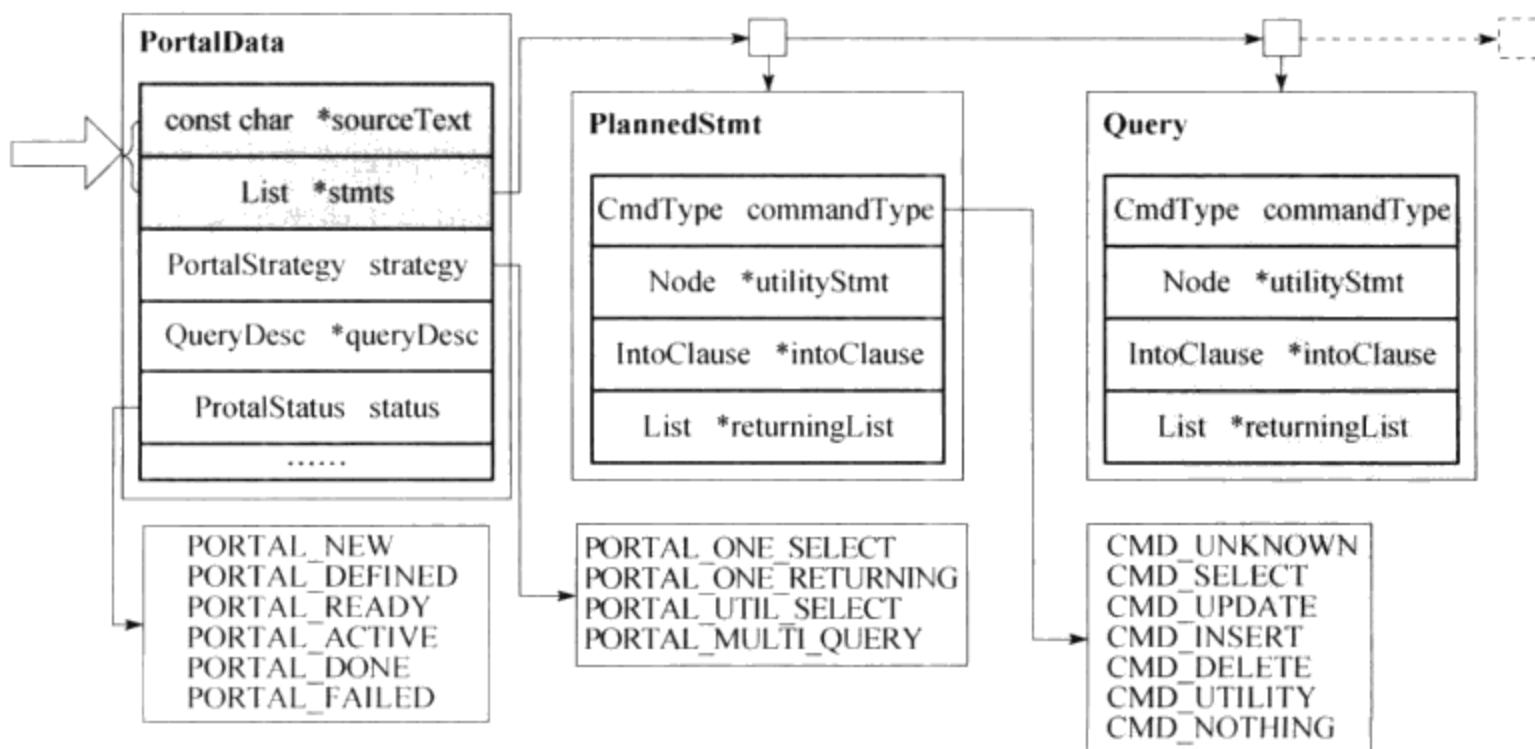


图 6-4 Portal 及其主要字段的结构

PostgreSQL 主要根据原子操作的命令类型以及 stmts 中原子操作的个数来为 Portal 选择合适的执行策略。

由查询编译器输出的每一个查询计划树中都包含有一个类型为 CmdType 的字段，用于标识该原子操作对应的命令类型。命令类型分为六类，使用枚举类型定义：

- CMD_UNKNOWN 表示未定义。
- CMD_SELECT 类型表示 SELECT 查询类型。
- CMD_UPDATE 类型表示更新操作。
- CMD_INSERT 类型表示插入操作。
- CMD_DELETE 类型表示删除操作。
- CMD/utility 类型表示功能性操作（数据定义语句）。
- CMD_NOTHING 类型用于由查询编译器新生成的操作，即如果一个语句通过编译器的处理之后需要额外生成一个附加的操作，则该操作的命令类型就被设置为 CMD_NOTHING。

选择 PORTAL_ONE_SELECT 策略应满足以下条件：

- stmts 链表中只有一个 PlannedStmt 类型或是 Query 类型的节点。
- 节点是 CMD_SELECT 类型操作。
- 节点的 utilityStmt 字段和 intoClause 字段为空。

选择 PORTAL_UTIL_SELECT 策略应满足以下条件：

- stmts 链表仅有的一个 Query 类型的节点。
- 节点是 CMD/utility 类型操作。
- 节点的 utilityStmt 字段保存的是 FETCH 语句（类型为 T_FetchStmt）、EXECUTE 语句（类型为 T_ExecuteStmt）、EXPLAIN 语句（类型为 T_ExplainStmt）或是 SHOW 语句（类型为 T_VariableShowStmt）之一。

而 PORTAL_ONE_RETURNING 策略适用于 stmts 链表中只有一个包含 RETURNING 子句（returningList 不为空）的原子操作。其他的各种情况都将以 PORTAL_MULTI_QUERY 策略进行处理。

执行策略选择器的主函数名为 ChoosePortalStrategy，其输入为 PortalData 的 stmts 链表，输出的是预先定义的执行策略枚举类型 PortalStrategy。该函数的执行流程如图 6-5 所示。

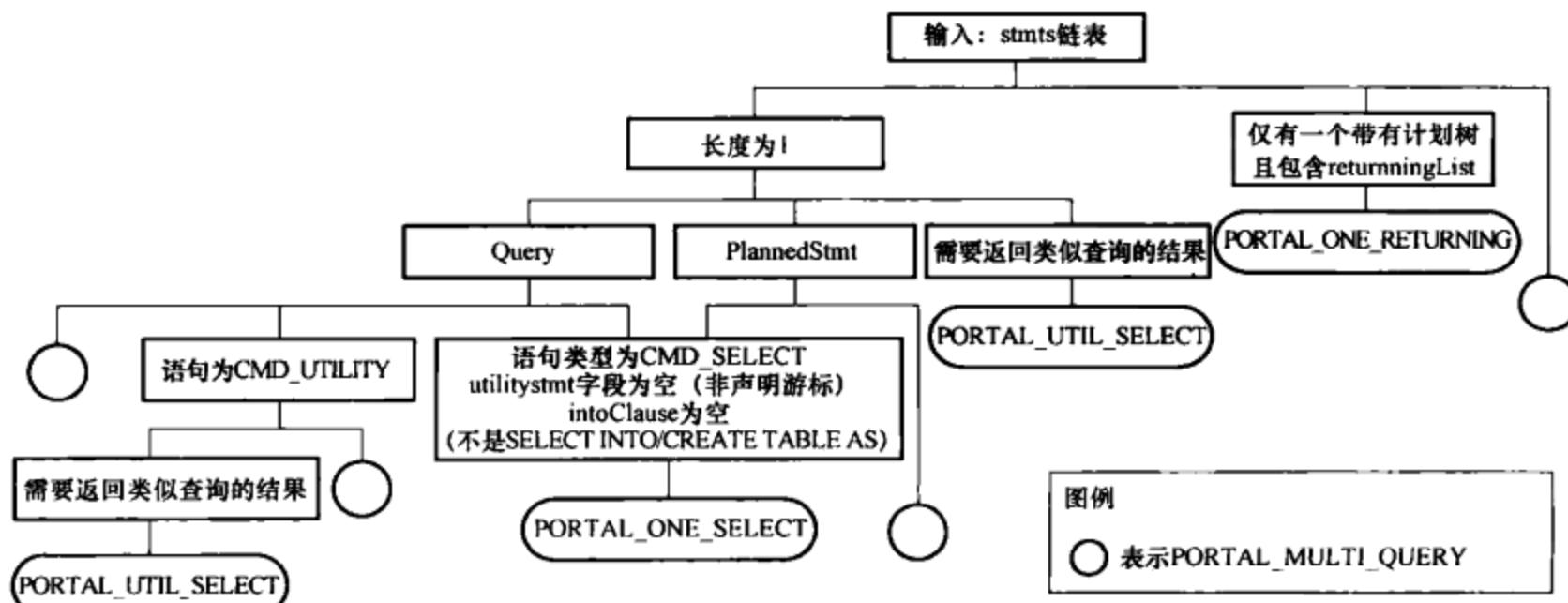


图 6-5 执行策略选择流程

6.1.4 Portal 执行的过程

Portal 是查询执行器执行一个 SQL 语句的“门户”，所有 SQL 语句的执行都从一个选择好执行策略的 Portal 开始。所有 Portal 的执行过程都必须依次调用 PortalStart（初始化）、PortalRun（执行）、PortalDrop（清理）三个过程，PostgreSQL 为 Portal 提供的几种执行策略实现了单独的执行流程，每种策略的 Portal 在执行时会经过不同的处理过程。Portal 的创建、初始化、执行及清理过程都在 `exec_simple_query` 函数中进行，其过程如下：

1) 调用函数 CreatePortal 创建一个干净的 Portal，其中内存上下文、资源跟踪器、清理函数等都已经设置好，但 sourceText、stmts 等字段并没有设置。

2) 调用函数 `PortalDefineQuery` 为刚创建的 Portal 设置 `sourceText`、`stmts` 等字段，这些字段的值都来自于查询编译器输出的结果，其中还会将 Portal 的状态设置为 `PORTAL_DEFINED` 表示 Portal 已被定义。

3) 调用函数 PortalStart 对定义好的 Portal 进行初始化，初始化工作主要如下：

- ①调用 ChoosePortalStrategy 为 Portal 选择策略。
 - ②如果选择的是 PORTAL_ONE_SELECT 策略，调用 CreateQueryDesc 为 Portal 创建查询描述符。
 - ③如果选择的是 PORTAL_ONE_RETURNING 或者 PORTAL_UTIL_SELECT 策略，为 Portal 创建元组的描述符。
 - ④将 Portal 的状态设置为 PORTAL_READY，表示 Portal 已经初始化好，准备开始执行。

4) 调用函数 PortalRun 执行 Portal，该函数将按照 Portal 中选择的策略调用相应的执行部件来执行 Portal。

5) 调用函数 PortalDrop 清理 Portal, 主要是对 Portal 运行中所占用的资源进行释放, 特别是用于缓存结果的资源。

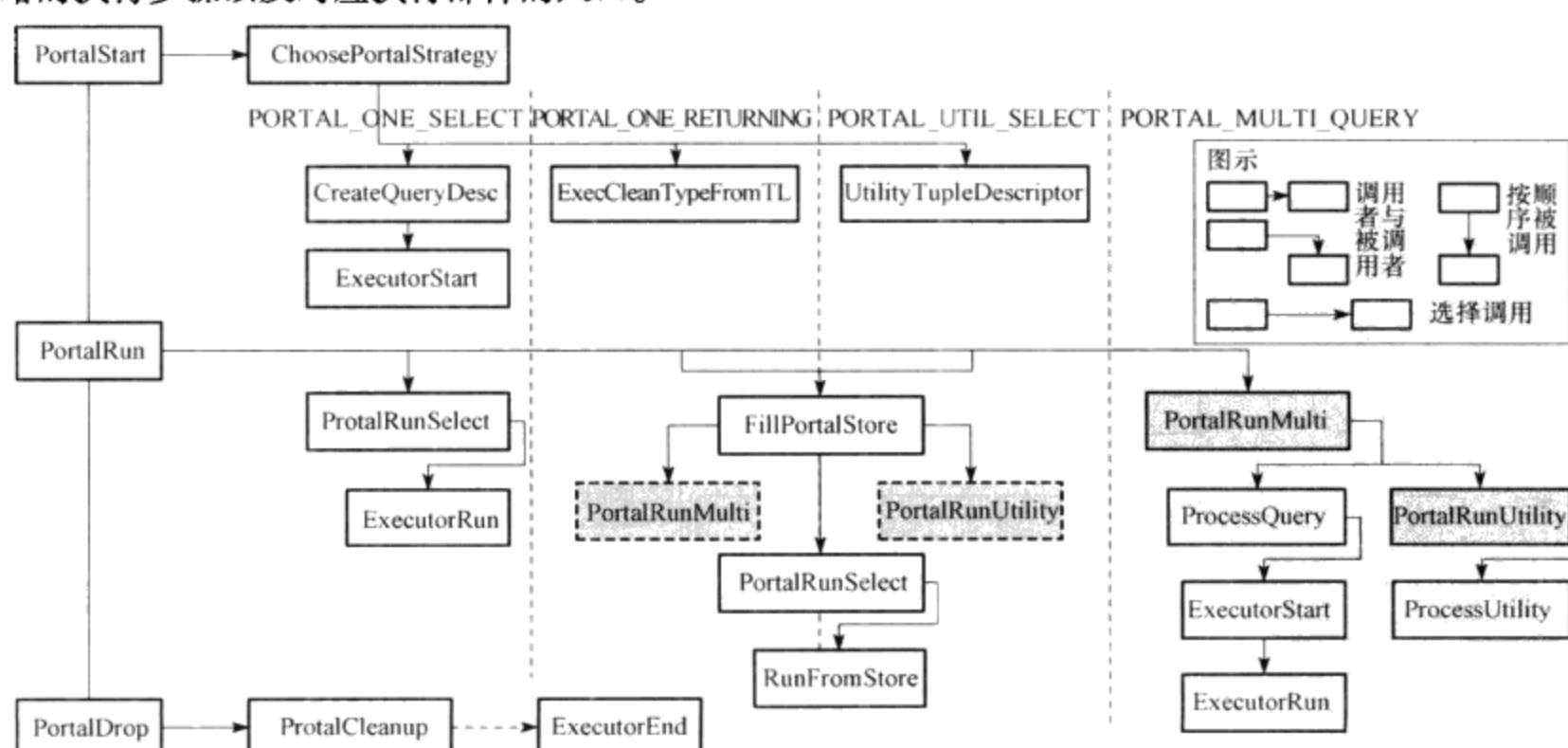


图 6-6 Portal 执行时的函数调用关系

对于 PORTAL_ONE_SELECT 策略的 Portal，其中包含一个简单 SELECT 类型的查询计划树，在 PortalStart 中将调用 ExecutorStart 进行 Executor（执行器）初始化，然后在 PortalRun 中调用 ExecutorRun 开始执行器的执行过程。

PORTAL_ONE_RETURNING 和 PORTAL_UTIL_SELECT 策略需要在执行后将结果缓存，然后将缓存的结果按要求进行返回。因此，在 PortalStart 中仅会初始化返回元组的结构描述信息。接着 PortalRun 会调用 FillPortalStore 执行查询计划得到所有的结果元组并填充到缓存中，然后调用 RunFromStore 从缓存中获取元组并返回。从图 6-6 中可以看到，FillPortalStore 中对于查询计划的执行会根据策略不同而调用不同的处理部件，PORTAL_ONE_RETURNING 策略会使用 PortalRunMulti 进行处理，而 PORTAL_UTIL_SELECT 使用 PortalRunUtility 处理。Portal_MULTI_QUERY 策略在执行过程中，PortalRun 会使用 PortalRunMulti 进行处理。

6.2 数据定义语句执行

数据定义语言是一类用于定义数据模式、函数等的功能性语句。不同于元组增删查改的操作，其处理方式是为每一种类型的描述语句调用相应的处理函数。

数据定义语句的处理过程比较简单，其执行流程最终会进入到 ProcessUtility 处理器，然后执行语句对应的不同处理过程。由于数据定义语句的种类很多，因此整个处理过程中的数据结构和方式种类繁冗、复杂，但流程相对简单、固定。由于处理过程和数据结构的类型众多，我们将会重点介绍表创建过程以及相关数据结构，展示 PostgreSQL 中针对数据描述语句的具体处理过程。

6.2.1 数据定义语句执行流程

由于 ProcessUtility 需要处理所有类型的数据定义语句，因此其输入数据结构的类型也是各种各样，每种类型的数据结构表示不同的操作类型。ProcessUtility 将通过判断数据结构中 NodeTag 字段的值来区分各种不同节点，并引导执行流程进入相应的处理函数。图 6-7 展示了 ProcessUtility 的总体流程。

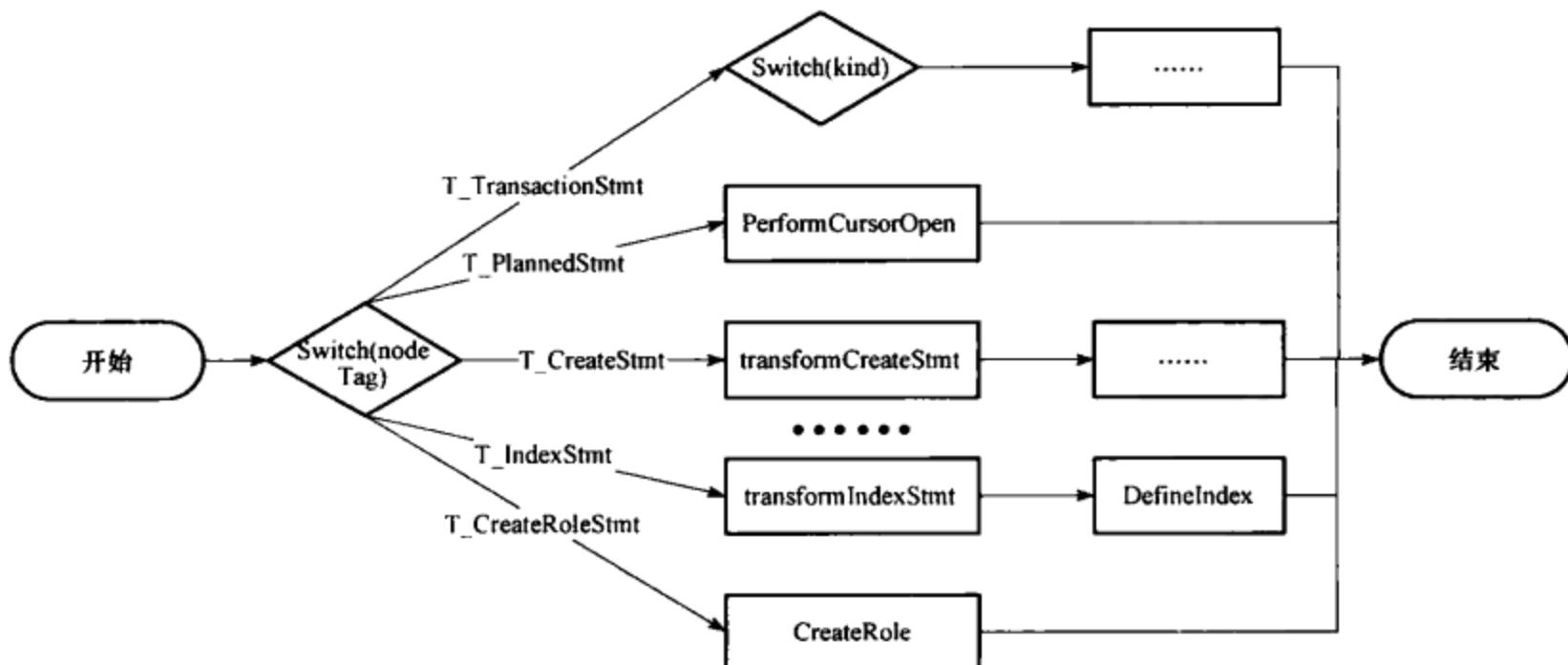


图 6-7 数据定义语句处理过程 ProcessUtility

针对各种不同的查询树，查询编译器在执行处理前会做一些额外的处理对查询树进行分析、处理与转换。例如，创建表的语句会用函数 `transformCreateStmt` 进行查询树的处理。这些处理过程可能会在当前操作之前和之后增加一些新的操作（例如在创建表的操作之前增加创建 `serial` 序列表操作、之后增加创建触发器用于外键约束操作等），也可能会执行对数据结构的处理操作（例如将 `CreateStmt` 节点 `tableElts` 字段中 `CONST_CHECK` 类型的 `Constraint` 节点转存到 `CreateStmt` 的 `constraints` 链表中等）。由于这些处理过程会产生一些新的操作，因此最终会生成一个由多个操作构成的链表。因此，执行过程需要依次扫描该链表，为每一个原子操作调用相应的处理函数。

相同类别的语句处理过程涉及内容相近，实现思想和主要过程相似。例如，事务类处理主要是对于当前事务的状态的判断和转换；游标类处理的主要思想是首次将执行一个查询计划树（`Plan tree`），将结果缓存在 `Portal` 指向的特殊结构中，然后按照要求获取元组数据；表、属性管理类主要涉及权限管理、修改相应系统表以及关系表的存储类别操作等。由于数据定义语句的种类多达上百种，我们下面将以一个创建表的例子来介绍数据定义语句的执行流程。

6.2.2 执行实例

例 6.1 创建一个名为 `course` 的数据表，此表有三个属性：编号（`no`，自增属性）、姓名（`name`）非空、学分（`credit`）非负。其中，包含了一个约束定义，主键被定义为编号（`no`）。对应的 SQL 语句如下：

```
CREATE TABLE course (
    no SERIAL,
    name VARCHAR,
    credit INT,
    CONSTRAINT con1 CHECK(credit >= 0 AND name <> ''),
    PRIMARY KEY(no)
);
```

系统首先会对查询语句进行词法和语法分析，将查询语句构造为查询树的链表。然后，针对链表中的每一个查询树进行如下的处理过程（例 6.1 仅有一个 `T_CreateStmt` 类型的查询树）：

- 1) 分析和重写查询树。
- 2) 生成查询计划。
- 3) 创建及初始化 `Portal`。
- 4) 调用 `Portal` 执行过程。
- 5) 调用 `Portal` 清理过程。

图 6-8 给出了上述查询语句执行时的主要函数调用流程。

在例 6.1 的情况下，查询编译器会生成一个仅包含一个 `T_CreateStmt` 类型节点的查询树链表，因此对应的 `Portal` 的 `stmts` 字段中也只包含一个 `T_CreateStmt` 类型节点。`ChoosePortalStrategy` 函数根据 `stmts` 字段值选择策略时会选择 `PORTAL_MULTI_QUERY` 策略。在接下来的 `PortalRun` 函数中将会调用 `PortalRunMulti` 来执行 `PORTAL_MULTI_QUERY` 策略，将会把处理流程引导到 `ProcessUtility` 中。`ProcessUtility` 将首先调用函数 `transformCreateStmt` 对 `T_CreateStmt` 节点进行转换处理，流程如图 6-9 所示。该过程会做如下转换：

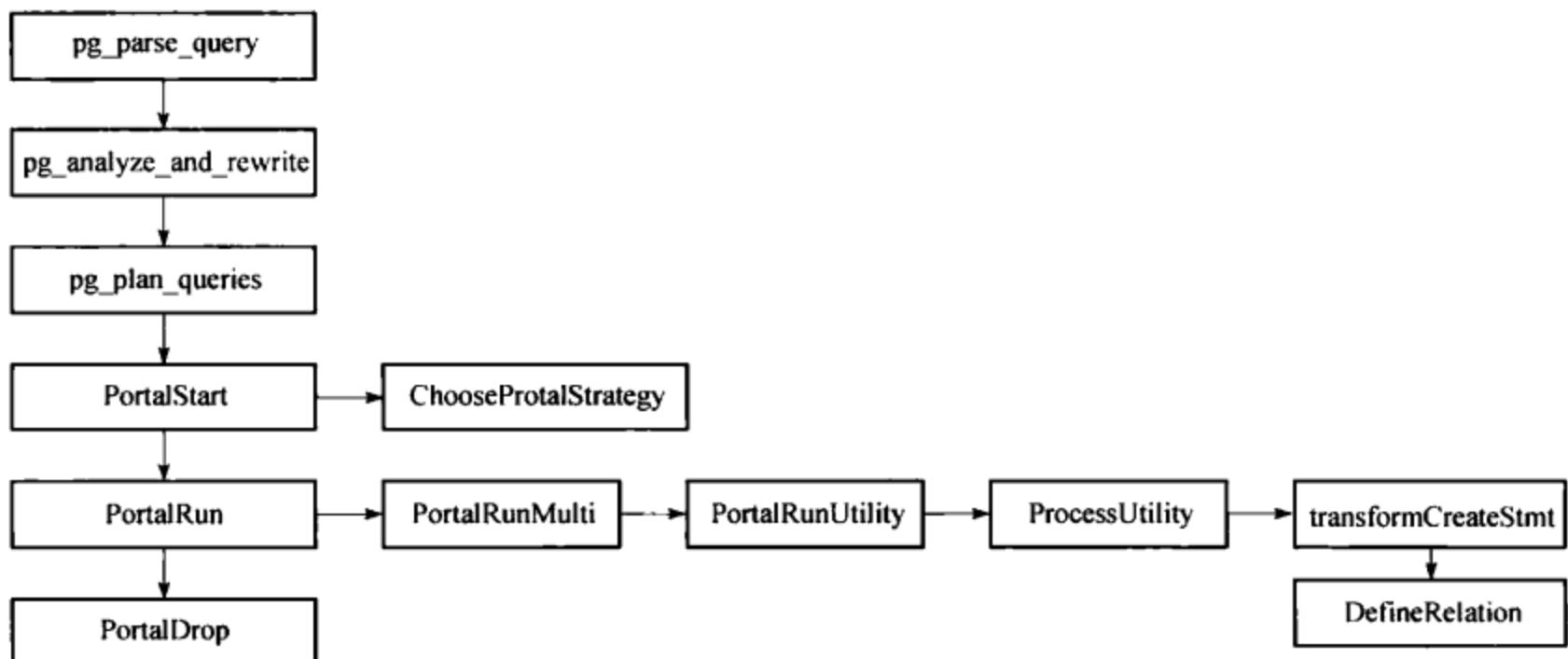


图 6-8 表创建流程

- 将主键约束改为创建唯一索引 (T_IndexStmt 节点)。
- 将自增类型转换为 int4Oid，并附加创建专用的 SERIAL 表（用于记录自增字段，将形成一个 T_CreateSeqStmt 节点）操作。
- 增加 CONSTR_DEFAULT 类型约束作为默认值（被定义为调用函数 nextval）。

创建 SERIAL 表 (T_CreateSeqStmt 节点) 的操作会被放在 stmts 链表中 T_CreateStmt 节点之前的位置，创建唯一约束索引 (T_IndexStmt 节点) 的操作被放置在 T_CreateStmt 节点之后。最后还会将单独定义或与属性同时定义的 CONSTR_CHECK 类型约束全部转移到 T_CreateStmt 节点的 constraints 字段所指向的链表中。

最后，transformCreateStmt 将原有的 T_CreateStmt 操作转换为一个操作序列：依次为 T_CreateSeqStmt（创建序列表）、T_CreateStmt（创建数据表）、T_IndexStmt（创建唯一约束索引）。

```

CREATE SERIAL TABLE course_no_seq; -- 用于产生自增序列
CREATE TABLE course (
    noint4Oid DEFAULT nextval(),
    nameVARCHAR,
    creditINT,
    CONSTRAINT con1 CHECK (credit >= 0 AND name <> ""),
);
CREATE INDEX course_pkey; -- 用于唯一检查

```

之后 ProcessUtility 将逐个对序列中的操作进行处理。对

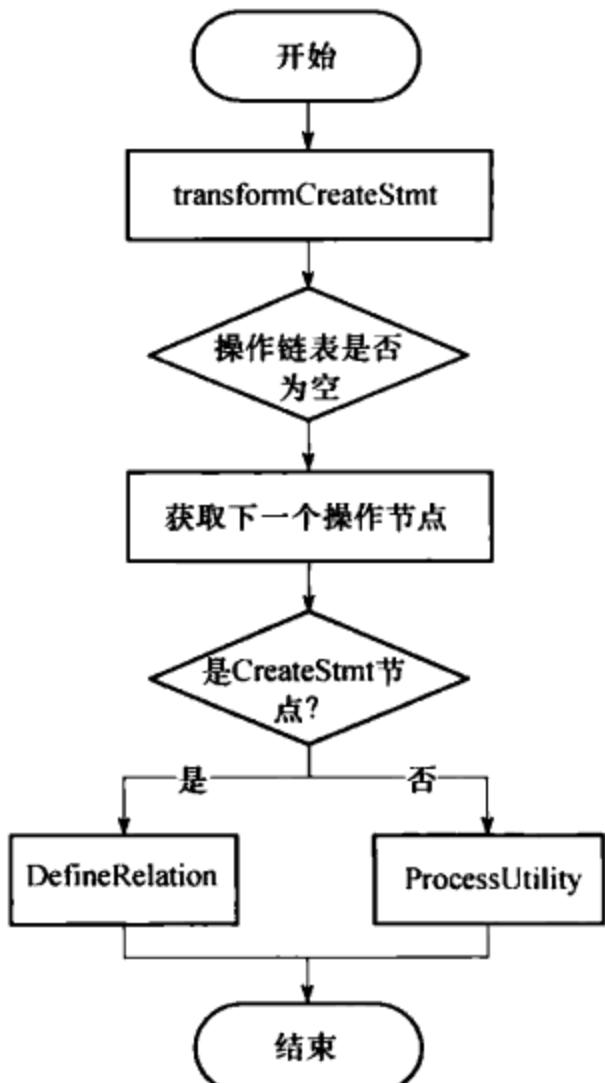


图 6-9 ProcessUtility 对 T_CreateStmt 节点的处理过程

`T_CreateStmt` 操作将会调用 `DefineRelation` 进行数据表的创建，而其他节点则会通过递归调用 `ProcessUtility` 进入相应的处理过程。图 6-10 展示了 `T_CreateStmt` 操作的处理过程。

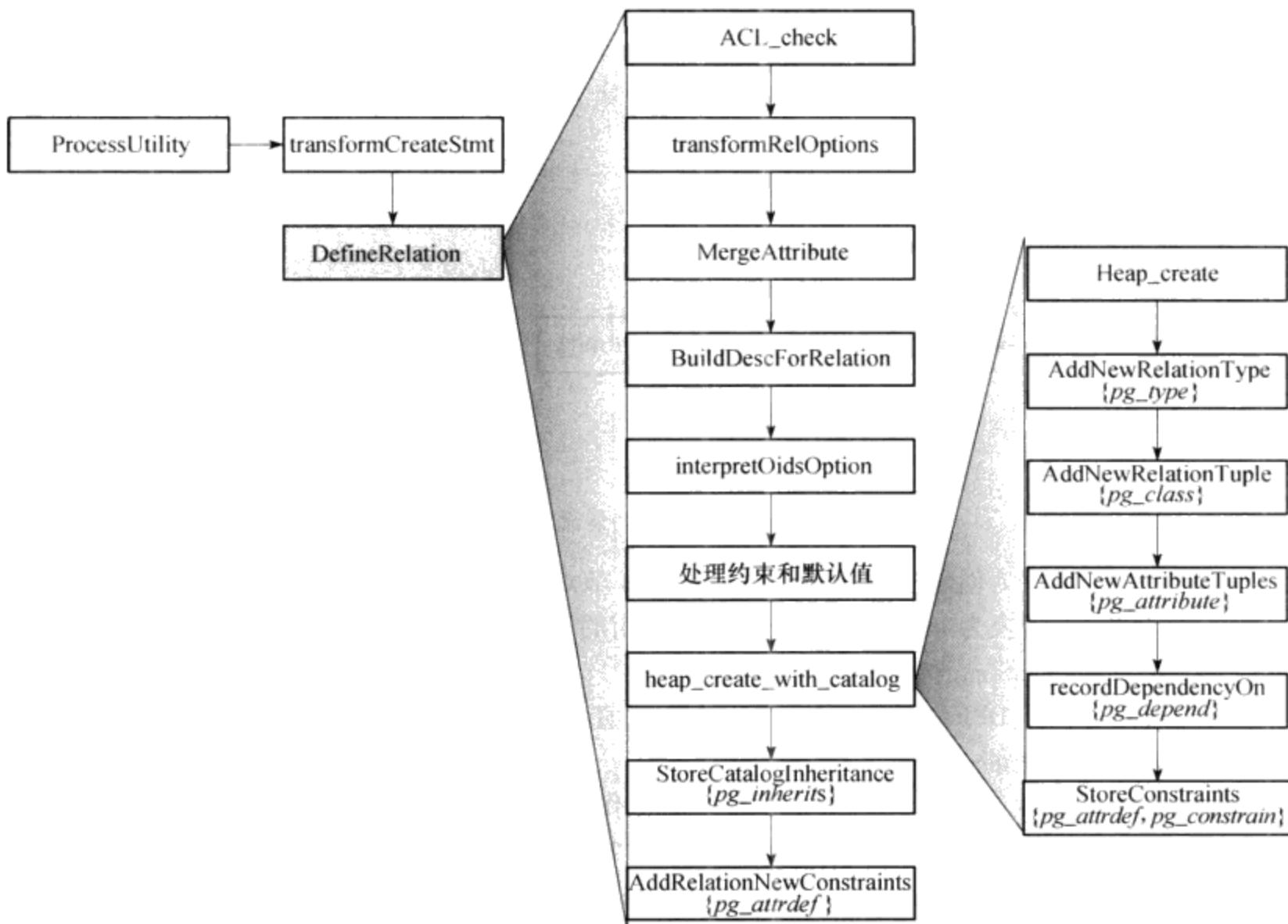


图 6-10 表创建函数 `DefineRelation` 执行流程

创建表的过程由函数 `DefineRelation` 完成，其流程如下：

- 1) 进行权限检查，确定当前用户是否有权限创建表。
- 2) 对表创建语句中的 WITH 子句进行解析（`transformRelOptions`）。
- 3) 调用 `heap_reloptions` 对参数进行合法性验证。
- 4) 使用 `MergeAttributes`，将继承的属性合并到表属性定义中。
- 5) 调用 `BuildDescForRelation` 利用合并后的属性定义链表创建 `tupleDesc` 结构（这个结构用于描述元组各属性结构等信息）。
- 6) 决定是否使用系统属性 OID（`interpretOidsOption`）。
- 7) 对属性定义链表中的每一个属性进行处理，查看是否有默认值、表达式或约束检查。
- 8) 使用 `heap_create_with_catalog` 创建表的物理文件并在相应的系统表中注册。
- 9) 用 `StoreCatalogInheritance` 存储表的继承关系。
- 10) 处理表中新增的约束与默认值（`AddRleationNewConstraints`）。

表创建函数的主要功能是由 `heap_create_with_catalog` 完成的，之前的各种操作主要是构造 `heap_`

`create_with_catalog` 所需要的参数。例如，WITH 子句处理主要完成其中存储相关参数的处理，以便存入 `pg_class` 系统表的 `reloptions` 字段中；`BuildDescForRelation` 主要处理表定义中属性名、类型、非空约束以便构造 `pg_attribute` 系统表相关内容。

`heap_create_with_catalog` 函数首先会根据要创建表的属性描述信息、表的名称、命名空间等使用 `heap_create` 创建一个 `RelationData` 结构并放入 `RelCache`，并根据这些信息通过调用 `RelationCreateStorage` 函数创建物理文件。然后调用 `AddNewRelationType`，向 `pg_type` 中增加一条关于该表的记录。`AddNewRelationTuple` 则会将表的相关信息插入 `pg_class` 系统表中，而 `AddNewAttributeTuples` 将表的每个属性记录到 `pg_attribute` 系统表中。最后还需要通过调用 `StoreConstraints` 将约束和默认值分别存储到 `pg_constraint` 和 `pg_attrdef` 中。

6.2.3 主要的功能处理器函数

从创建表的例子可以看到，功能处理器（`ProcessUtility`）本身只作为入口选择函数，它会根据输入的节点类型调用相应的处理过程。除了创建表的处理过程之外，表 6-1 中还列出了几种常见的输入节点类型，并给出了其对应处理函数以及其功能简介。

表 6-1 常用数据定义语句处理函数

节点类型	核心处理函数	功能
<code>T_TransactionStmt</code>	<code>BeginTransactionBlock</code>	标记事务开始
	<code>EndTransactionBlock</code>	结束事务
	<code>DefineSavepoint</code>	定义保存点
	<code>RollbackToSavepoint</code>	回滚到保存点
<code>T_PlannedStmt</code>	<code>PerformCursorOpen</code>	打开游标操作，初始化查询
<code>T_ClosePortalStmt</code>	<code>PerformPortalClose</code>	关闭游标，释放游标占用资源
<code>T_FetchStmt</code>	<code>PerformPortalFetch</code>	执行 FETCH/MOVE 操作
<code>T_CreateStmt</code>	<code>DefineRelation</code>	创建关系表
<code>T_CreateTableSpaceStmt</code>	<code>CreateTableSpace</code>	创建 tablespace，记录在 <code>pg_tablespace</code> 系统表中
<code>T_DropStmt</code>	<code>RemoveRelations</code>	删除关系表
	<code>RemoveTypes</code>	删除自定义类型
<code>T_CommentStmt</code>	<code>CommentObject</code>	记录注释信息到 <code>pg_description</code>
<code>T_CopyStmt</code>	<code>DoCopy</code>	完成 COPY 命令操作
<code>T_AlterTableStmt</code>	<code>AlterTable</code>	实现 ALTER TABLE 命令
<code>T_IndexStmt</code>	<code>DefineIndex</code>	创建索引
<code>T_CreateSeqStmt</code>	<code>DefineSequence</code>	创建一个用于自增属性的关系表
<code>T_ExplainStmt</code>	<code>ExplainQuery</code>	执行 EXPLAIN 命令

6.3 可优化语句执行

可优化语句的共同特点是它们被查询编译器处理后都会生成查询计划树，这一类语句由执行器（`Executor`）处理。该模块对外提供了三个接口：`ExecutorStart`、`ExecutorRun` 和 `ExecutorEnd`，其输入是包含查询计划树的数据结构 `QueryDesc`，输出则是相关执行信息或结果数据。如果希望执行某个

计划树，仅需构造包含此计划树的 QueryDesc，并依次调用 ExecutorStart、ExecutorRun、ExecutorEnd 三个过程即能完成相应的处理过程。从图 6-6 可以看到，执行器的三个接口函数都是在 Portal 的相关函数中调用的，分别负责执行器的初始化、执行和清理工作，Portal 在处理时也使用了同样的方式，这样可以把资源分配回收工作与执行过程独立开，能够简化执行过程，更是一种很好的资源管理方式。

执行器对于查询计划树的处理，最终被转换为针对计划树上每一个节点的处理。每种节点表示一种物理代数（Physical Algebra）操作，PostgreSQL 对其进行初始化、处理、清理的过程。节点的处理被设计为需求驱动模式，父节点使用孩子节点提供的数据作为输入，并向其上层节点返回处理结果。实际执行时，从根节点开始处理，每个节点的执行过程会根据需求自动调用孩子节点的执行过程来获取输入数据（一般为元组），从而层层递归执行，实现整个计划树的遍历执行过程。初始化和清理也采用相同的设计模式，这种设计模式使得节点处理的代码结构简洁统一、语义明确，且实现方式简单有效。

接下来将对执行器部分的各种原理、实现做进一步的介绍，以帮助读者了解执行器的内部实现和具体执行过程。

6.3.1 物理代数与处理模型

数据库的查询逻辑使用逻辑代数（例如：关系代数）来表示。例如，PostgreSQL 中使用 SQL 语句，然而执行时需要使用物理代数（Physical Algebra），例如，PostgreSQL 中的查询计划。对于一个可优化语句，PostgreSQL 在执行前会给出与之等价的用物理代数表示的查询计划，然后按照查询计划进行执行。

《数据库系统实现》一书中提到：“物理查询计划由操作符构造，每一个操作符实现计划中的一步。物理操作符常常是一个关系代数操作符的特定实现。但是，我们也要用物理操作符来完成另一些与关系代数操作符无关的任务。”此处提到操作符构造的物理计划就是物理代数表达的查询计划，操作符在 PostgreSQL 系统中使用一种节点定义，而无关的任务在该书中提到了表的扫描、排序，PostgreSQL 系统也实现了这些操作的节点，但是在 PostgreSQL 中“逻辑操作符与物理操作符并不是简单的一一对应关系”，它们之间的转换过程实际就是查询编译器所做的工作，本章主要介绍各种物理操作符的实现和执行。

PostgreSQL 中的物理操作符被定义为有 0 ~ 2 个输入和一个输出，这是为了在实现中能够对应二叉树结构：所有的物理操作符被组织为一个二叉树，每一个物理操作符对应于树中的一个节点，下层节点的输出作为上层节点输入。数据（元组）从底层节点向上层节点流动，直至根节点，而根节点的输出即为整个查询的结果。如图 6-11 所示，Join 类型操作符有两个输入，Sort 操作符仅有一个输入，Scan 操作符则没有输入。元组被 Scan 操作符取出后经过连接操作，然后排序后获得结果。

在 PostgreSQL 的实现中，上层函数通过 ExecInitNode、ExecProcNode、ExecEndNode 三个接口函数来统一对节点进

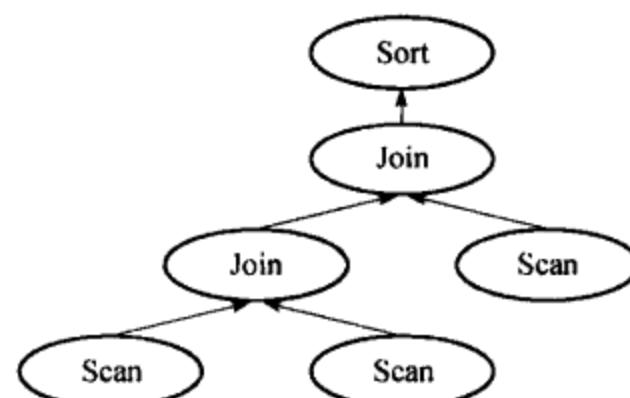


图 6-11 查询计划实例

行初始化、执行和清理，这三个函数会根据所处理节点的实际类型调用相应的初始化、执行、清理函数，例如，若 ExecInitNode 处理的是 SeqScan 节点，则调用 SeqScan 节点的初始化函数 ExecInitSeqScan 来实际完成该节点的初始化工作。这三个函数也是递归执行的：对根节点的初始化会递归地对下层的节点进行初始化；在根节点调用 ExecProcNode 获取结果时，也会递归地对下层节点进行执行以获取上层节点所需的输入数据；执行完成后只需对根节点进行清理，下层节点也会被递归地清理。

由此可见，查询计划树上的节点就构成了物理元组到执行结果的管道，因此查询计划树的执行过程可以看成是拉动元组穿过管道的过程。PostgreSQL 采用了一次一元组的执行模式，每个节点被执行一次仅向上层节点返回一条元组[○]。因此，对于整个查询计划树的执行也是一次一元组的模式。这种模式有很多的优点：

- 减少了返回元组的延迟。
- 对于某些操作（例如游标、LIMIT 子句等）不需要一次性获取所有的元组，节省了开销。
- 减少了实现过程中缓存结果带来的代码复杂性和执行过程中临时存储的开销。

6.3.2 物理操作符的数据结构

从前面的介绍我们已经看到查询计划树是由各种物理操作符（也简称为计划节点）构成，那么在 PostgreSQL 中是如何存储和表示各类节点的呢？图 6-12 给出了 Hash 连接（HashJoin）类型节点的数据结构表示。

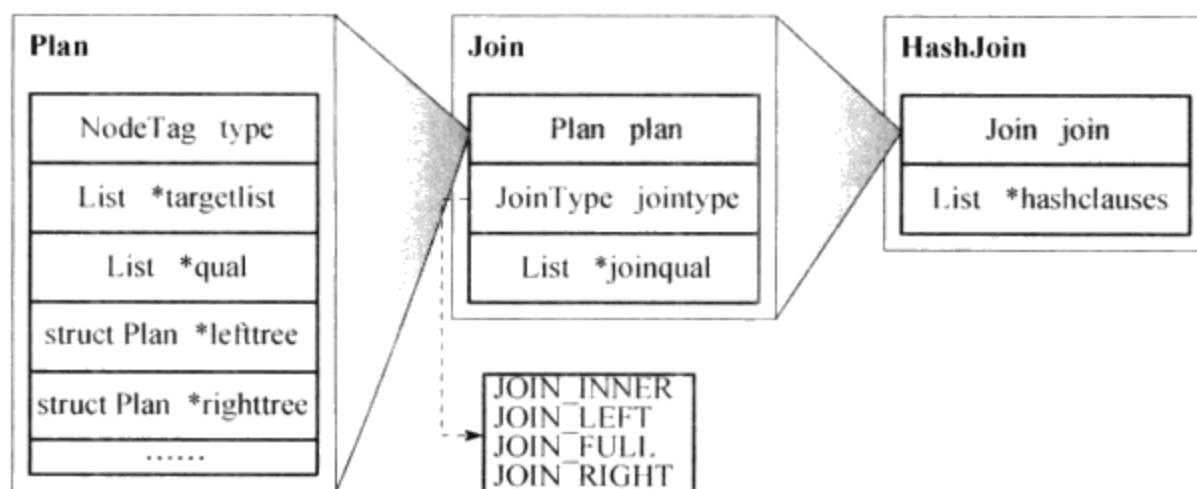


图 6-12 Hash 连接（HashJoin）节点结构定义

所有物理操作符节点的数据结构都以一个 Plan 类型的字段开头，这有点像类的继承：把 Plan 看成一个父类，其他物理操作符节点都是它的直接或者间接子类。如图 6-12 所示，Join 节点是 Plan 的子类，从 Plan 中继承了左右子树指针（lefttree, righttree）、节点类型（type）、选择表达式（qual）、投影链表（targetlist）等公共字段，并有自己的扩展字段连接类型（jointype）和连接条件（joinqual）；HashJoin 节点则是 Join 节点的子类，有自己的扩展字段 hashclauses。

PostgreSQL 系统中将所有的计划节点按功能分为四类：控制节点（control node）、扫描节点

[○] 有一类特殊的节点，称为物化节点（Materialization Node），必须首先计算出所有结果，并将其缓存，实现上采取执行时每次从缓存中获取一条元组。例如排序节点就属于物化节点。

(scan node)、连接节点 (join node) 和物化节点 (materialization node)，并分别为扫描、连接节点类型定义了公共父类 Scan、Join。Hash 连接属于连接节点，因此 Hash 连接继承于 Join 节点。连接节点类型的公共父类定义了连接的类型以及连接的条件。作为 Hash 连接节点，需要使用 Hash 函数，所以 HashJoin 节点扩展定义了 hashclauses 字段来存储相关信息，其中包括需要做 Hash 的属性以及使用的 Hash 函数等。

Plan 的众多子类节点通过 lefttree 和 righttree 字段构成了整个查询计划树，其根节点指针被保存在 PlannedStmt 类型的数据结构中，其中包含了语句的类型 (commandType)、查询计划树根节点 (planTree)、查询涉及的范围表 (rtable)、结果关系表 (resultRelation)[⊖]。PlannedStmt 结构则被放在 QueryDesc 中，QueryDesc 结构的基本定义如图 6-13 所示。

作为执行器的输入，QueryDesc 中包含查询计划树 (plannedstmt 字段)、功能语句相关执行计划 (utilitystmt 字段)、执行器全局状态 (estate 字段) 以及计划节点执行状态 (planstate 字段) 等。从图 6-13 可以看出，执行器全局状态 estate 中保存了查询涉及的范围表 (es_range_table)、Estate 所在的内存上下文 (es_query_ctx)，也是执行过程中一直保持的内存上下文)、用于在节点间传递元组的全局元组表 (es_tupleTable) 和每获取一个元组就会回收的内存上下文 (es_per_tuple_exprContext)。

执行器初始化时，ExecutorStart 会根据查询计划树构造执行器全局状态 (estate) 以及计划节点执行状态 (planstate)。在查询计划树的执行过程中，执行器将使用 planstate 来记录计划节点执行状态和数据，并使用全局状态记录中的 es_tupleTable 字段在节点间传递结果元组。执行器的清理函数 ExecutorEnd 将回收执行器全局状态和计划节点执行状态。

图 6-14 给出了 PostgreSQL 中用于计划节点执行状态记录的数据结构与计划节点之间的对应关系。与图 6-12 类似，PostgreSQL 为每种计划节点定义了一种状态节点。所有的状态节点均继承于 PlanState 节点，其中包含辅助计划节点指针 (Plan)、执行器全局状态结构指针 (state)、投影运算相关信息 (targetlist)、选择运算相关条件 (qual)，以及左右子状态节点指针 (lefttree、righttree)。状态节点之间通过 lefttree 和 righttree 指针组织成和查询计划树结构类似的状态节点树，同时，每个状态节点都保存了指向其对应的计划节点的指针 (PlanState 类型中的 Plan 字段)。图 6-14 中展示了

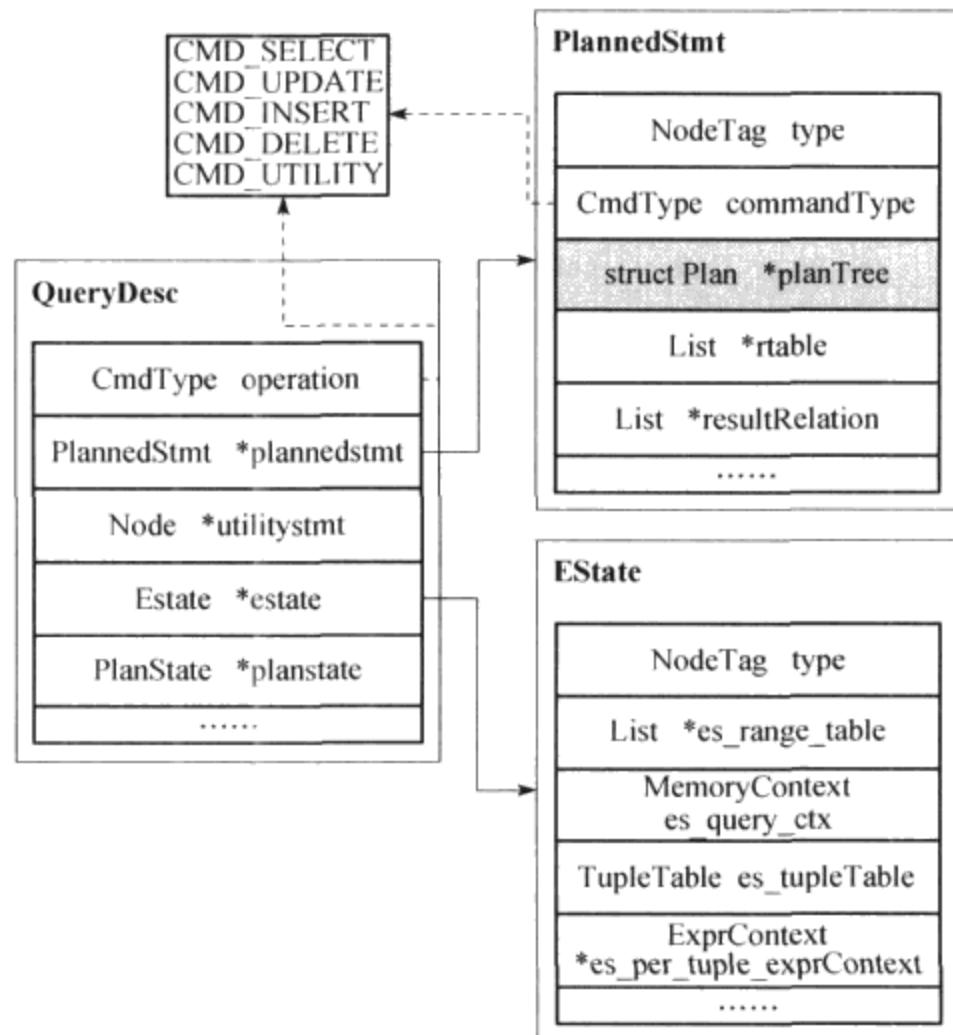


图 6-13 QueryDesc 数据结构组织

[⊖] 链表结构，保存了用于 Insert、Update、Delete 操作的目标关系表在 PlannedStmt 结构中链表 rtable 中的索引位置。

连接节点状态的公共父类 `JoinState`，它继承于 `PlanState`，扩展了连接类型（`jointype`）和连接条件（`joinqual`）属性。而 `HashJoinState` 继承于 `JoinState` 并扩展了更多的属性，包括 Hash 函数相关内容（`hashclauses`、`hj_HashTable`、`hj_OuterHashKeys`、`hj_InnerHashKeys`、`hj_HashOperators`）、左子节点返回元组指针（`hj_OuterTupleSlot`）、右子节点返回元组指针（`hj_HashTupleSlot`）等。

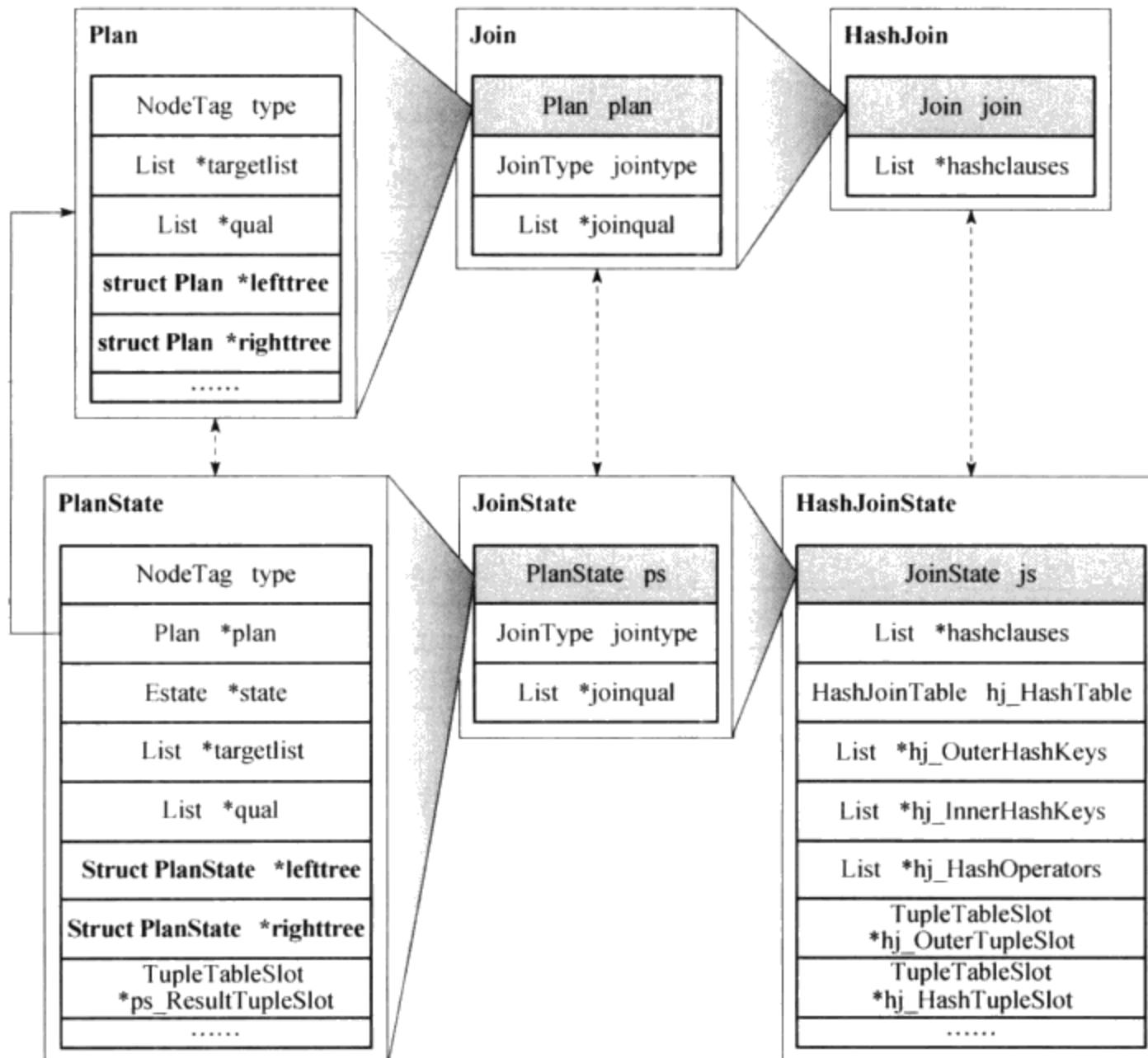


图 6-14 计划节点与节点执行状态

至此，执行器执行过程中涉及的主要各种数据结构已经介绍完毕。执行器的输入是 `QueryDesc`，它包含了存储查询计划树根节点指针的 `PlannedStmt` 结构。执行器执行时，首先构造全局状态记录 `Estate` 结构，并为每个计划节点（`Plan`）构造对应的状态节点（`PlanState`），然后在执行中使用相关结构存储执行状态，执行完毕后释放相关的数据结构。

6.3.3 执行器的运行

在 PostgreSQL 中提供了三个接口函数用于调用执行器，分别为 `ExecutorStart`、`ExecutorRun` 和 `ExecutorEnd`。当需要使用执行器来处理查询计划时，仅需依次调用三个函数即可完成执行器的整个执行过程。

执行器运行时的函数调用关系如图 6-15 所示，ExecutorStart 通过调用 standard_ExecutorStart 对执行器进行必要的初始化，主要工作包括构造 EState 结构和查询计划树的初始化（即构造对应的 PlanState 树，由 InitPlan 函数完成）。ExecutorRun 的功能由 standard_ExecutorRun 实现，在执行过程中会调用 ExecutePlan 完成查询计划的执行。ExecutorEnd 由 standard_ExecutorEnd 函数完成，通过调用 ExecEndPlan 处理执行状态树根节点释放已分配的资源，最后释放执行器全局状态 EState 完成整个执行过程。

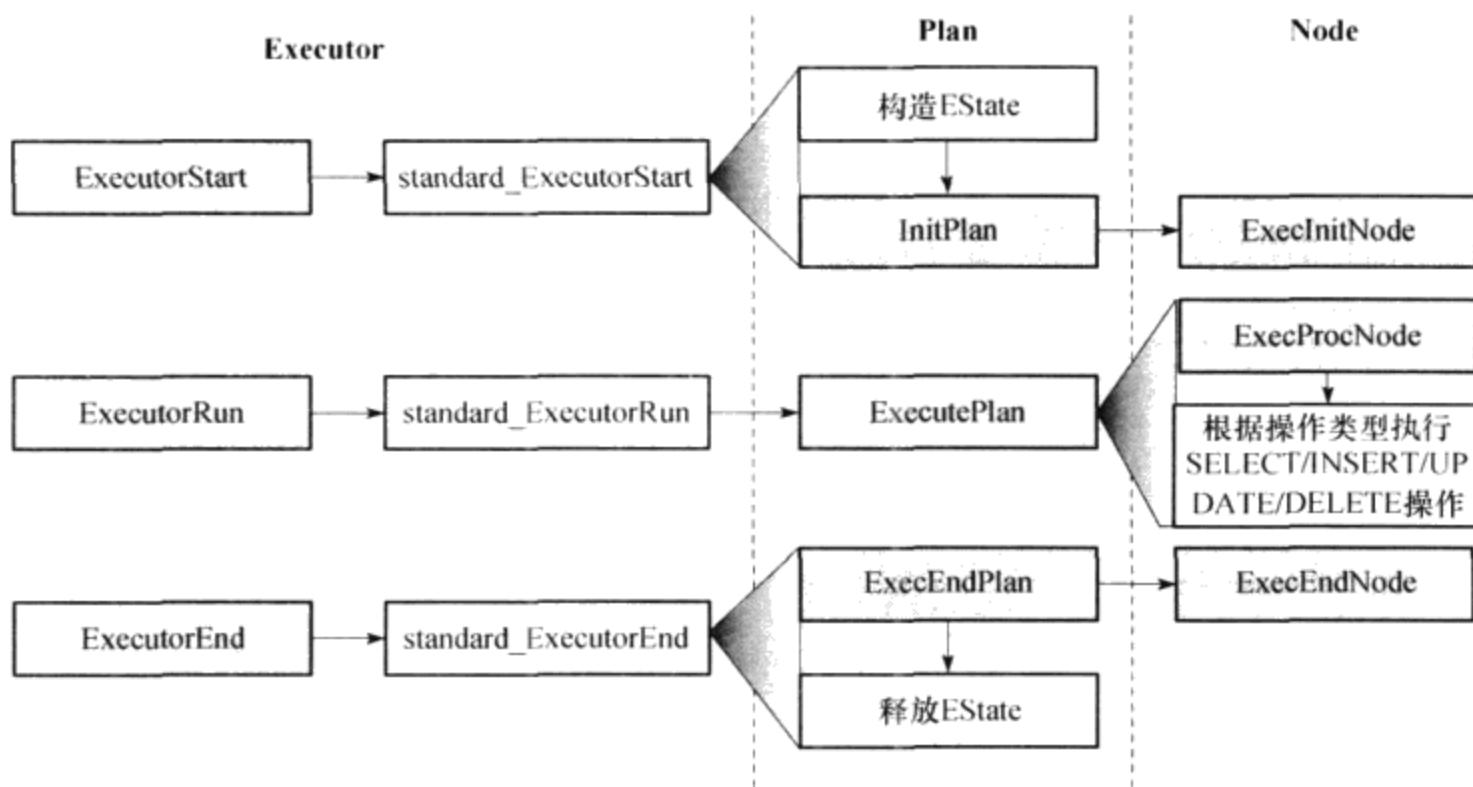


图 6-15 执行器函数调用

(1) 初始化查询计划树

执行器中对查询计划树的初始化都是从其根节点开始，并递归地对其子节点进行初始化。计划节点的初始化过程一般都会经历如图 6-16 所示的几个基本步骤，该过程在完成计划节点的初始化之后会输出与该计划节点对应的 PlanState 结构指针，计划节点的 PlanState 结构也会按照查询计划树的结构组织成计划节点执行状态树。对计划节点初始化的主要工作是根据计划节点中定义的相关信息，构造对应的 PlanState 结构并对相关字段赋值。

计划节点的初始化由函数 ExecInitNode 完成，该函数以要初始化的计划节点为输入，并返回该计划节点所对应的 PlanState 结构指针。在 ExecInitNode 中，通过判断计划节点的类型来调用相应的处理过程，每一种计划节点都有专门的初始化函数，且都以“ExecInit 节点类型”的形式命名。例如，NestLoop 节点的初始化函数为 ExecInitNestLoop。在计划节点的初始化过程中，如果该节点还有下层的子节点，则会递归地调用子节点的初始化函数来对子节点进行初始化。图 6-17 中展示了对 NestLoop 节点进行初始化时对其子节点进行初始化的递归调用过程。如图 6-17 所示，ExecInitNode 函数会根据计划节点的类型 (T_NestLoop) 调用该类型节点的初始化函数 (ExecInitNestLoop)。由于 NestLoop 节点有两个子节点，因此 ExecInitNestLoop 会先调用 ExecInitNode 对其左子节点进行初始化，并将其返回的 PlanState 结构指针存放在为 NestLoop 构造的 NestLoopState 结构的 lefttree 字段中；然后以同样的方式初始化右子节点，将返回的 PlanState 结构指针存放于 NestLoopState 的 righttree 字

段中。同样，如果左右子节点还有下层节点，初始化过程将以完全相同的方式递归下去，直到到达查询计划树的叶子节点。而在初始化过程中构造的 PlanState 子树也会层层返回给上层节点，并被链接在上层节点的 PlanState 结构中，最终构造出完整的 PlanState 树。

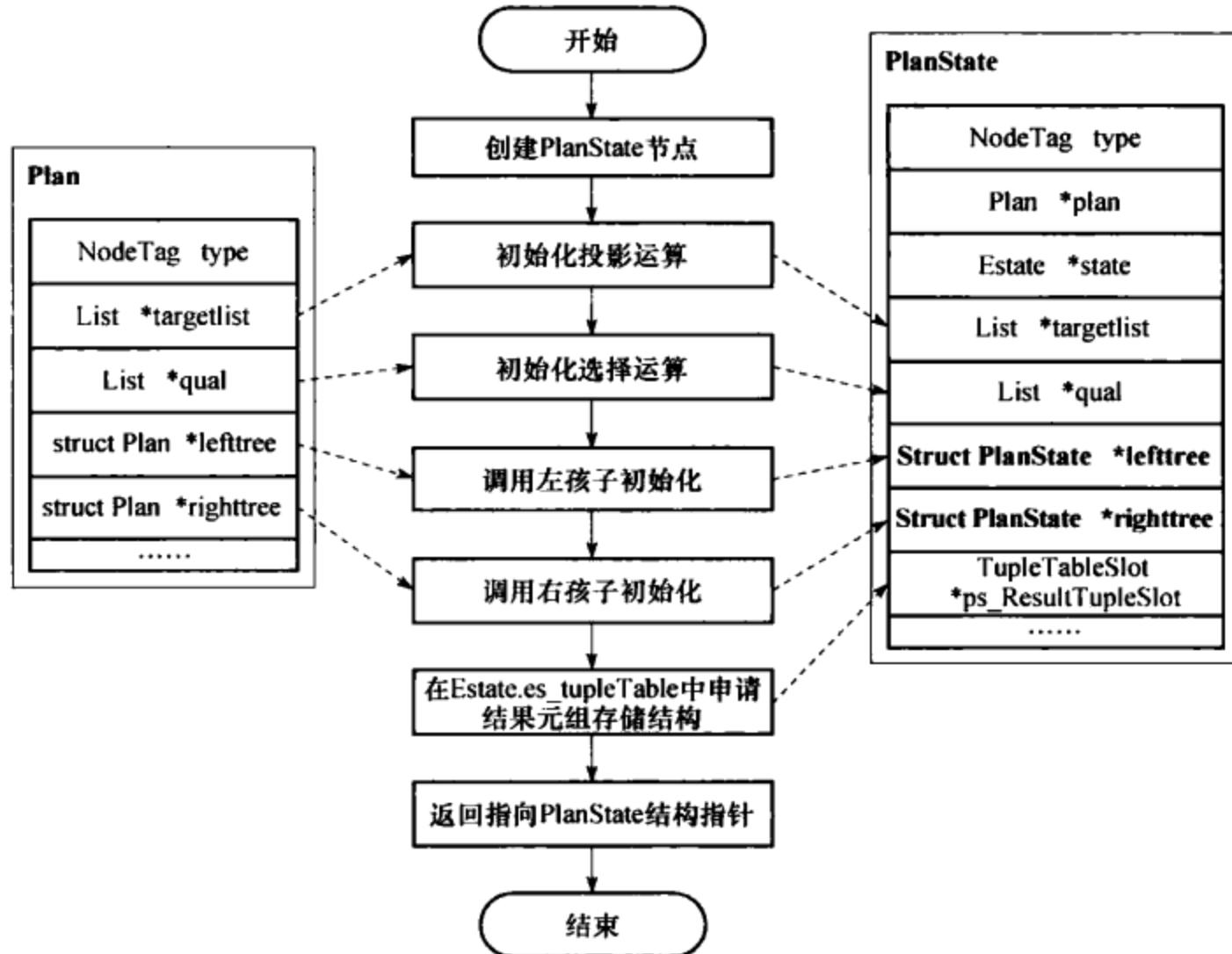


图 6-16 计划节点初始化的一般步骤

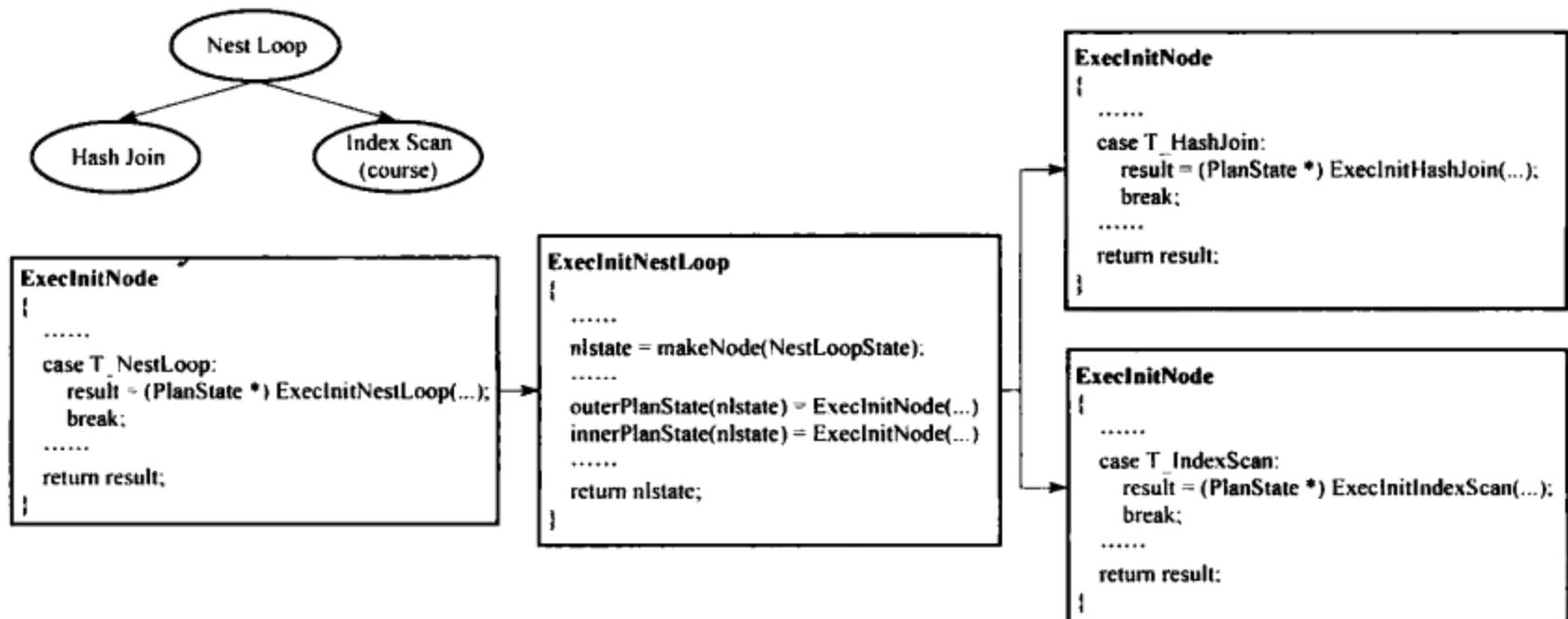


图 6-17 NestLoop 节点初始化的递归调用过程

(2) 查询计划执行

查询计划的实际执行由函数 `ExecutePlan` 完成，该函数的主体部分是一个大的循环，每一次循环都通过 `ExecProcNode` 函数从计划节点状态树中获取一个元组，然后对该元组进行相应的处理（增删查改），然后返回处理的结果。当 `ExecProcNode` 从计划节点状态树中再也取不到有效的元组时结束循环过程。

如图 6-18 所示，`ExecProcNode` 的执行过程也和 `ExecInitNode` 类似：从计划节点状态树的根节点获取数据，上层节点为了能够完成自己的处理将会递归调用 `ExecProcNode` 从下层节点获取输入数据（一般为元组），然后根据输入数据进行上层节点对应的处理，最后进行选择条件的运算和投影运算，并向更上层的节点返回结果元组的指针。同 `ExecInitNode` 一样，`ExecProcNode` 也是一个选择函数，它会根据要处理的节点的类型调用对应的处理函数。例如，对于 `NestLoop` 类型的节点，其处理函数为 `ExecNestLoop`。`ExecNestLoop` 函数同样会对 `NestLoop` 类型的两个子节点调用 `ExecProcNode` 以获取输入数据。如果其子节点还有下层节点，则以同样的方式递归调用 `ExecProcNode` 进行处理，直到到达叶子节点。每一个节点被 `ExecProcNode` 处理之后都会返回一个结果元组，这些结果元组作为上层节点的输入被处理形成上层节点的结果元组，最终根节点将返回结果元组。

每当通过 `ExecProcNode` 从计划节点状态树中获得一个结果元组后，`ExecutePlan` 函数将根据整个语句的操作类型调用相应的函数进行最后的处理。主要处理函数有四个，分别为 `ExecSelect`、`ExecInsert`、`ExecDelete`、`ExecUpdate`。

选择语句的处理较为简单，通过 `ExecSelect` 直接输出查询结果。对于插入语句，则首先需要调用 `ExecConstraints` 对即将插入的元组进行约束检查，如果满足要求，`ExecInsert` 会调用函数 `heap_insert` 将元组存储到存储系统。对于删除和更新，则分别由 `ExecDelete` 和 `ExecUpdate` 调用 `heap_delete` 和 `heap_update` 完成。

(3) 执行器清理

当执行器处理完所有能够获得的元组之后，由执行器清理函数 `ExecutorEnd` 负责善后工作。该函数调用 `ExecEndPlan` 对计划节点执行状态树进行清理。对计划节点执行状态树的清理和执行状

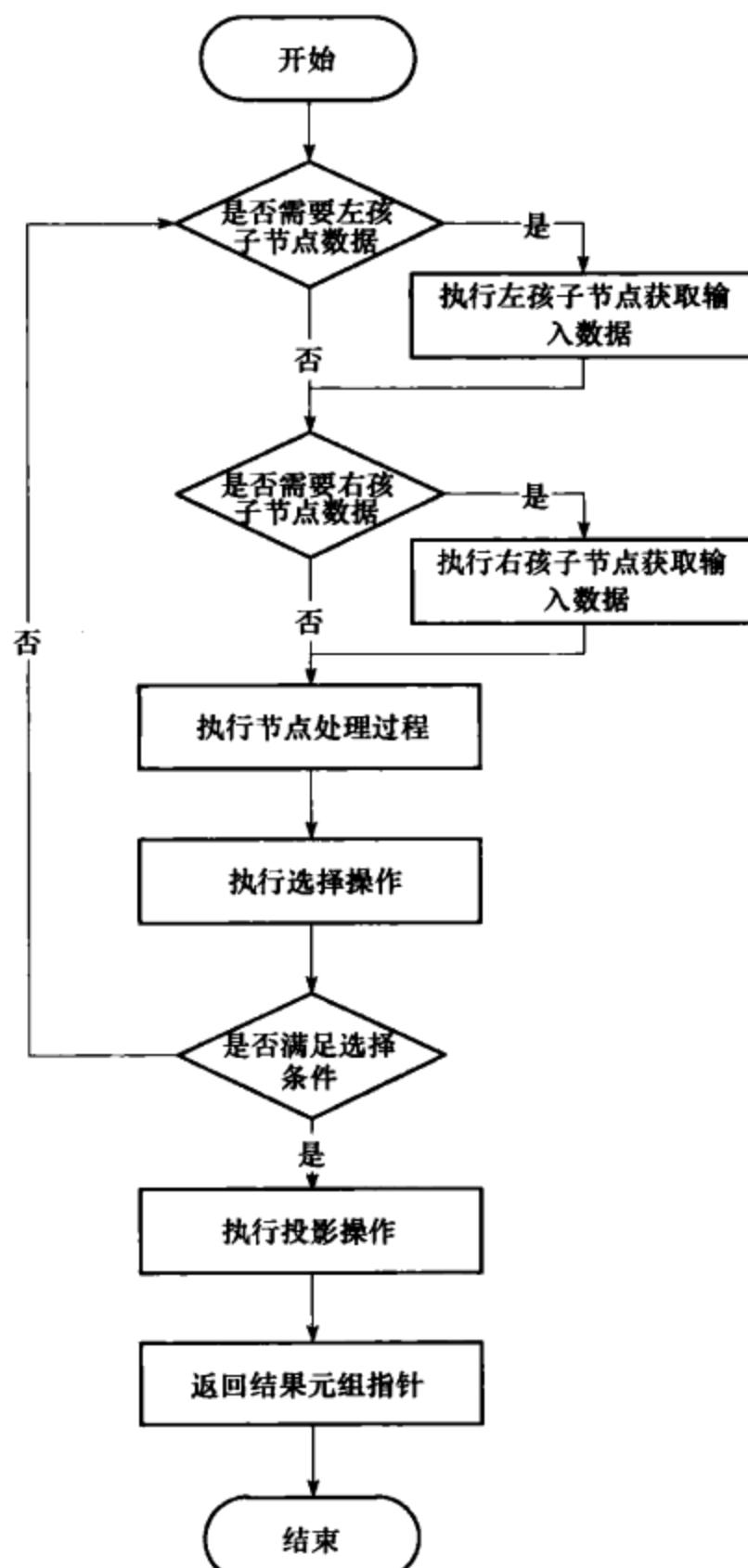


图 6-18 计划节点执行的一般过程

态树的初始化、执行相类似：从根节点开始递归调用 ExecEndNode 对每一个计划节点的执行状态节点进行清理。同样，ExecEndNode 只是一个选择函数，针对不同类型的节点有相应的清理函数。例如，NestLoop 节点的清理函数是 ExecEndNestLoop。如图 6-19 所示，清理过程的任务主要是回收初始化过程中分配的资源、投影和选择结构的内存、结果元组存储空间等，计划节点执行状态树清理完之后，ExecutorEnd 还将调用 FreeExecutorState 清理执行器全局状态。

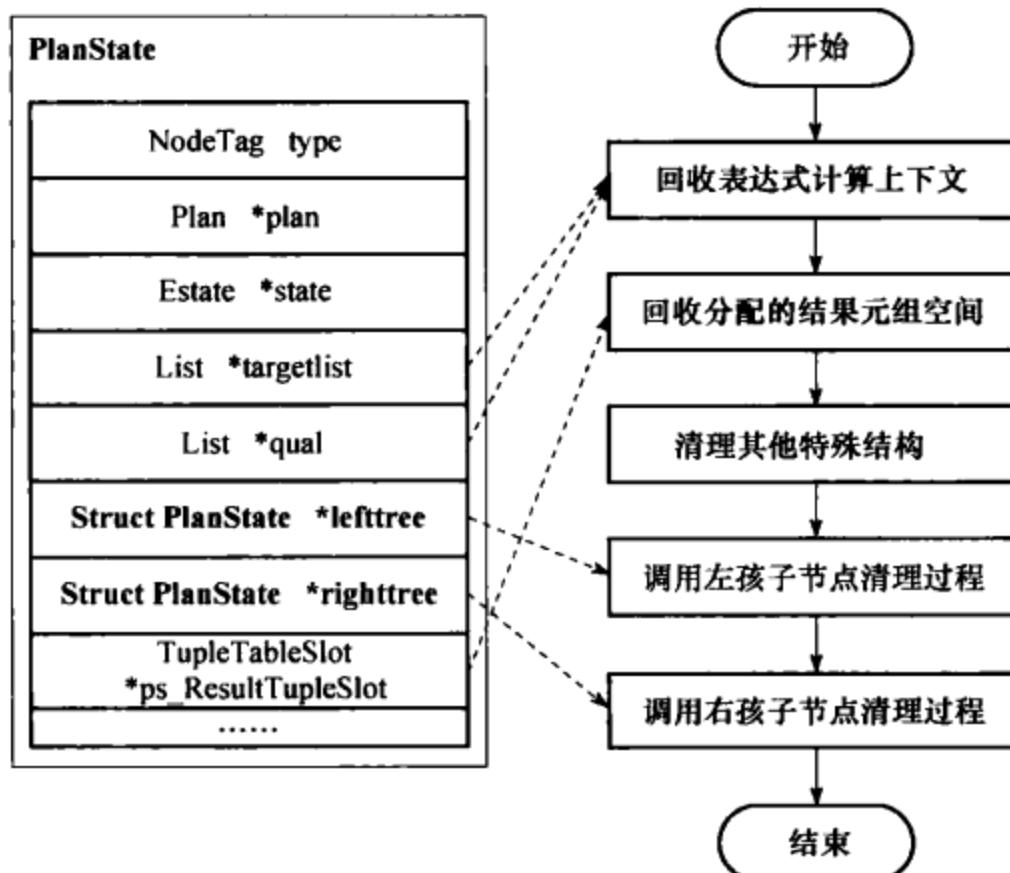


图 6-19 计划节点的清理过程

6.3.4 执行实例

例 6.2 查询 Jennifer 老师所授 Database System 课程学生人数。其中涉及的几个表的模式定义如下：

- course(no, name, credit)：记录课程信息，在课程号上建有索引。
- teacher(no, name, sex, age)：保存教师信息。
- teach_course(tno, cno, stu_num)：记录教师任课信息。

所用到的 SQL 语句为：

```

SELECT t.name, c.name, stu_num
FROM course AS c, teach_course AS tc, teacher AS t
WHERE c.no = tc.cno AND tc.tno = t.no AND c.name = 'Database System' AND t.name = 'Jennifer';

```

PostgreSQL 首先将对该查询语句进行词法和语法分析，将查询语句转换为分析树的链表。然后针对链表中的每一个分析树进行如下的处理（例 6.2 中的分析树链表仅有一个 SelectStmt 节点）：

- 1) 分析和重写分析树。
- 2) 生成执行计划。

- 3) 创建 Portal 数据结构。
- 4) 调用 Portal 初始化过程。
- 5) 调用 Portal 执行过程。
- 6) 调用 Portal 清理过程。

图 6-20 给出执行例 6.2 中 SQL 语句时的函数调用流程。在 PortalStart 过程中，将为该语句选择执行策略，由于仅包含一个 PlannedStmt 节点，且命令类型为 CMD_SELECT，因此 ChoosePortalStrategy 会选择 PORTAL_ONE_SELECT 策略。接下来，PortalStart 会调用 CreateQueryDesc 根据 PlannedStmt 结构构造 QueryDesc 结构体，并将其作为 ExecutorStart 的输入完成执行器初始化工作。PortalRun 会根据选择的执行策略，调用 PortalRunSelect 函数，它调用 ExecutorRun 完成计划的执行。最后，清理函数 PortalDrop 会通过 PortalCleanup 函数来调用 ExecutorEnd 完成执行器内部的清理，最后释放 QueryDesc 结构。

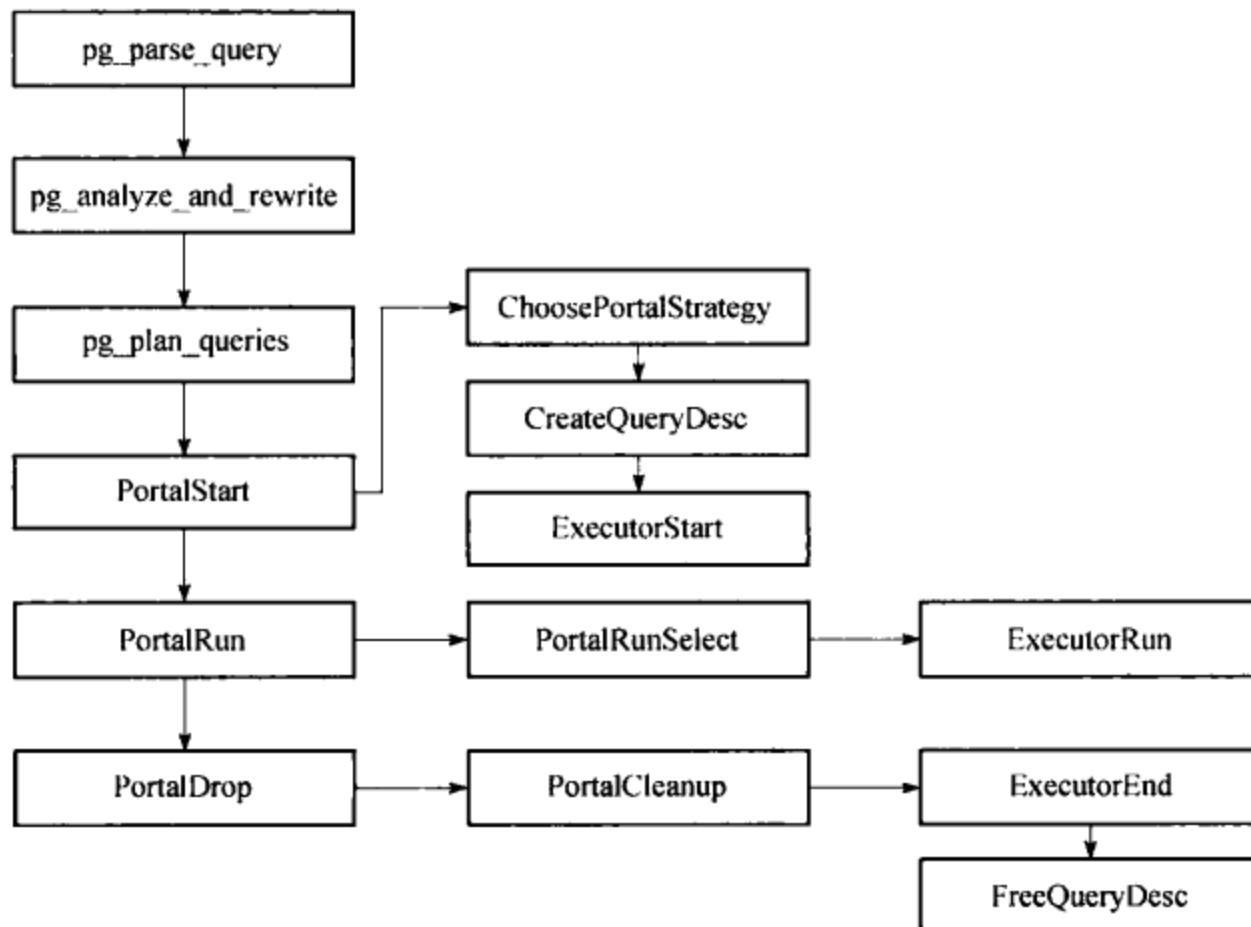


图 6-20 例 6.2 执行时的函数调用

在图 6-21 中展示了例 6.2 中查询语句对应的查询计划树。该查询计划包含三个扫描节点：SeqScan(teach_course)、SeqScan(teacher)、IndexScan(course)，分别从 teach_course、teacher、course 表中扫描元组。扫描节点是一种物理操作符，用于扫描表以获取其中的元组，通常有顺序扫描 (SeqScan) 和索引扫描 (IndexScan) 两种方式。SeqScan 表示从头到尾顺序扫描一个表中的元组，然后返回其中满足条件的元组。而 IndexScan 表示使用索引辅助查找满足条件的元组，但前提是在条件涉及的属性上建有可用的索引。除此之外，图 6-21 中还涉及两个连接节点：HashJoin 和 NestLoop，HashJoin 节点使用 Hash 方法实现连接操作，而 NestLoop 使用嵌套循环连接方法实现连接操作。

当需要从查询计划树的根节点 NestLoop 中取得元组时，实际是从执行状态树的根节点中取元

组, NestLoop 节点需要首先从左子树 (HashJoin 节点) 获取一个元组。而 HashJoin 则会先执行 Hash 节点构造 Hash 表, Hash 表的构造会依次从 SeqScan (teacher) 获取元组直至扫描结束, 至此 Hash 表构造完毕。然后 HashJoin 会从 SeqScan (teach_course) 获得一个元组, 经过 Hash 后与对应桶中的元组进行匹配, 然后将匹配到的元组连接后返回给 NestLoop 节点。从左子树获得一个元组后, NestLoop 还需要从右子树 (IndexScan 节点) 中获取一个元组用来和左子树元组连接, 因此将执行 IndexScan 节点扫描 course 表, 得到相应的元组。IndexScan 会使用左子树元组的 cno 属性值来扫描索引, 并检查元组是否同时满足 name = 'Database System', 满足条件的元组将返回给 NestLoop 节点作为右子树元组。最后, NestLoop 将左右节点提供的元组进行连接形成结果元组返回。这里需要注意的是, 对于每一个得到的左子树元组, 只有在将它和所有能找到的右子树元组连接之后才会去取得下一个左子树元组。也就是说, 下一次从 NestLoop 取元组时, NestLoop 还是会使用上一次得到的左子树元组, 但会尝试从右子树中取一个新的右子树元组与之相连接, 如果不能取得新的右子树元组, NestLoop 才会取下一个左子树元组并重复上述的过程。其他计划节点的执行也遵循类似的方式。

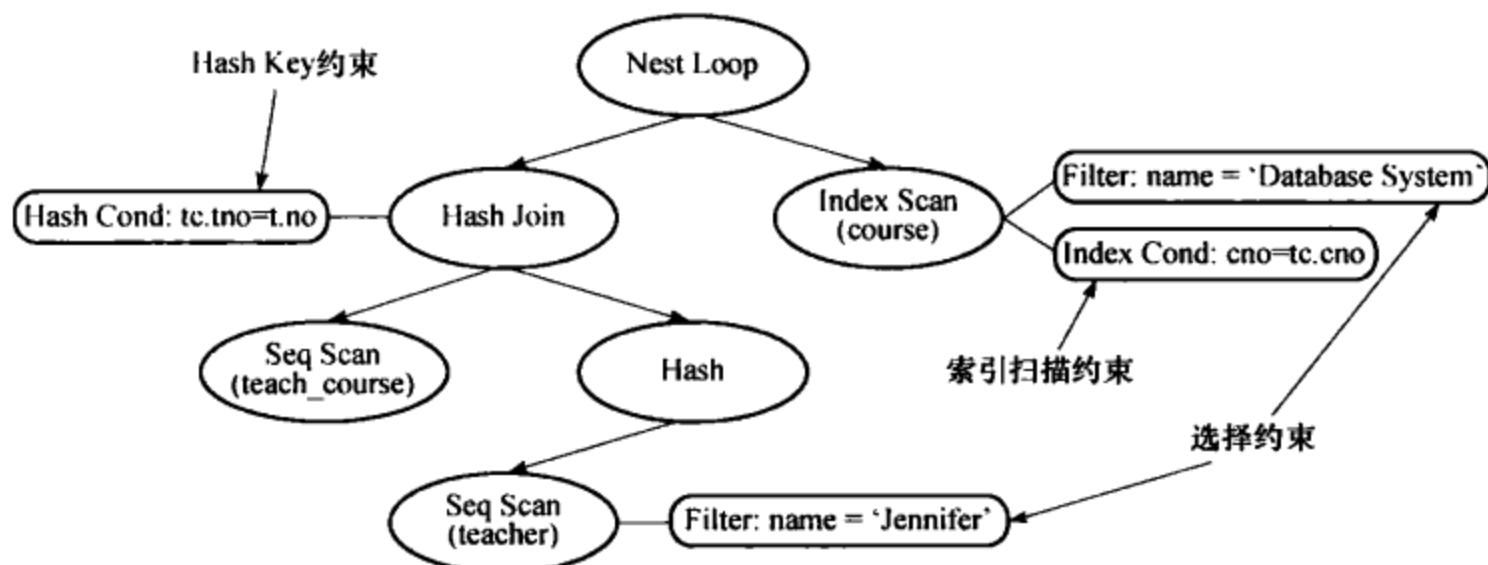


图 6-21 例 6.2 对应的查询计划树

如果从查询计划树的根节点中再也取不到有效的元组, 则说明查询计划的执行过程已经完成, 之后需要清理过程来进行资源回收等工作, 其过程与初始化类似, 这里不再赘述。

6.4 计划节点

从前面介绍的可优化语句处理相关的背景知识、实现思想和执行流程, 不难发现可优化语句执行的核心内容是对于各种计划节点的处理, 由于使用了节点表示、递归调用、统一接口等设计, 计划节点的功能相对独立、代码总体流程相似, 以下将重点介绍执行器中各种计划节点的相关执行过程。

在 PostgreSQL 中, 计划节点分为四类, 分别是控制节点 (Control Node)、扫描节点 (Scan Node)、物化节点 (Materialization Node)、连接节点 (Join Node)。

- **控制节点:** 是一类用于处理特殊情况的节点, 用于实现特殊的执行流程。例如, Result 节点可用来表示 INSERT 语句中 VALUES 子句指定的将要插入的元组。
- **扫描节点:** 顾名思义, 此类节点用于扫描表等对象以从中获取元组。例如, SeqScan 节点用于顺序扫描一个表, 每次扫描一个元组。

- 物化节点：这类节点种类比较复杂，但它们有一个共同特点，即能够缓存执行结果到辅助存储中。物化节点会在第一次被执行时生成其中的所有结果元组，然后将这些结果元组缓存起来，等待其上层节点取用；而非物化节点则是每次被执行时生成一个结果元组并返回给上层节点。例如，Sort 节点能够获取下层节点返回的所有元组并根据指定的属性进行排序，并将排序结果全部缓存起来，每次上层节点从 Sort 节点取元组时就从缓存中按顺序返回下一个元组。
- 连接节点：此类节点对应于关系代数中的连接操作，可以实现多种连接方式（条件连接、左连接、右连接、全连接、自然连接等），每种节点实现一种连接算法。例如，HashJoin 实现了基于 Hash 的连接算法。

6.4.1 控制节点

控制节点用于完成一些特殊的流程执行方式。由于 PostgreSQL 为查询语句生成二叉树状的查询计划，其中大部分节点的执行过程需要两个以内的输入和一个输出。但有一些特殊的功能为了优化的需要，会含有特殊的执行方式和输入需求（例如需要两个以上的输入），这一类的节点被称为控制节点。例如，为了能让 UNION 操作在一个计划节点就执行多个表（大于 2）的合并，Append 节点并未把 UNION 涉及的多个表放在孩子节点中，而是将这些表组成一个链表放在 Append 节点的 appendplans 字段中，处理时依次处理该链表中的节点获取输入。表 6-2 给出了 PostgreSQL 8.4.1 中的控制节点列表。

表 6-2 控制节点

节点类型	文件	描述
Result	nodeResult.c	处理含有仅需一次计算的条件表达式或 INSERT 中仅有一个 VALUES 子句
Append	nodeAppend.c	用于表示和组织需要包含多个（大于 2）子查询执行流程
BitmapAnd	nodeBitmapAnd.c	用于需要对两个或多个位图进行并操作的流程 ^①
BitmapOr	nodeBitmapOr.c	用于需要对两个或多个位图进行或操作的流程
RecursiveUnion	nodeRecursiveUnion.c	用于处理 WITH 子句中递归定义的 UNION 子查询

1. Result 节点

Result 节点有两种用途，第一种是不包含表扫描的情况：

```
SELECT 1 * 2;
INSERT INTO course (name, credit) VALUES ('Database System', 4);
```

当 SELECT 查询中没有 FROM 子句时或者 INSERT 语句只有一个 VALUES 子句时，执行查询计划不需要扫描表，执行器会直接计算 SELECT 的投影属性或者使用 VALUES 子句构造元组。

Result 计划节点的另一种用途是优化包含仅需计算一次的过滤条件：

```
SELECT * FROM course WHERE current_date > '2011-1-1';
```

由于 PostgreSQL 采用了一次一元组的方式执行查询计划树，而 WHERE 子句中的条件表达式结果是常量，只需计算一次即可。因此，Result 节点可针对此种情况进行优化，避免重复计算这类表达式。这时，Result 节点仅有一个子节点（左子节点）。

^① 使用多维索引扫描表时，会使用位图来表示表，并标记扫描到的所有元组，然后对位图进行“与/或”操作，这个功能由 BitmapAnd 或 BitmapOr 节点完成。

为了能够处理以上两种情况，Result 节点被定义成如图 6-22 所示的样子，除了继承 Plan 节点的基本属性外，还扩展定义了 resconstantqual 字段。顾名思义，该字段保存只需计算一次的常量表达式。

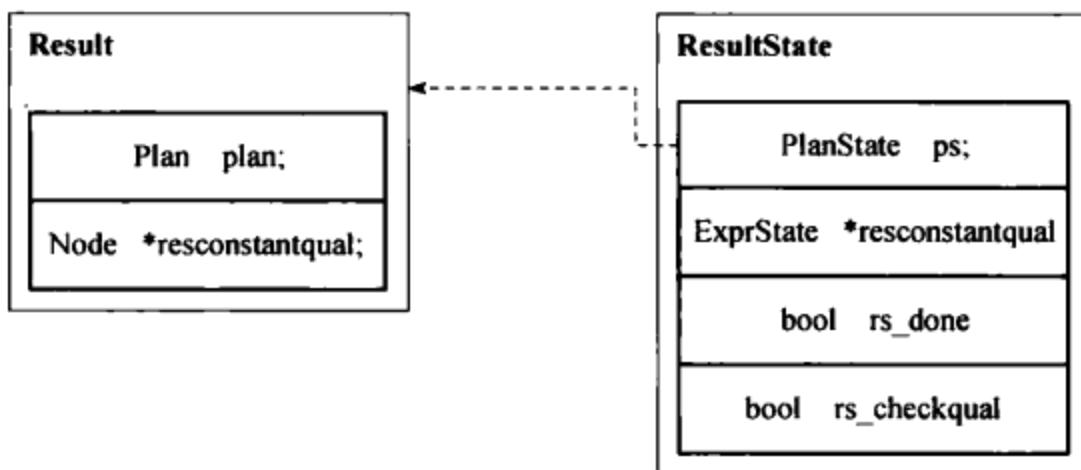


图 6-22 Result 节点相关数据结构

Result 节点初始化过程（ExecInitResult 函数）会初始化 ResultState 结构，如果有常量表达式，则放置在 resconstantqual 字段中，而 rs_checkqual 则用于标记是否需要计算常量表达式，rs_done 表示 Result 节点是否已经处理完所有元组。Result 初始化过程在计划节点的标准初始化过程之外增加了几项检查：

- 1) 保证无右子节点。
- 2) 检查是否有常量表达式，有则将执行状态节点的 rs_checkqual 设置为真，否则设置为假。

Result 节点的执行函数 ExecResult 的流程如图 6-23 所示。在 ExecResult 中，首先判断 rs_checkqual 是否为真，为真表示需要进行常量表达式计算，计算完成后将 rs_checkqual 设置为假，下次执行不会再进行计算。如果表达式计算结果为假，表示没有满足条件的结果，则 Result 节点直接返回 NULL。若表达式计算结果为真，则会检查是否有左子树，有则从其中获取元组；没有左子树表示为 VALUES 子句，则直接设置 rs_done 为真，因为只有一个元组需要输出（这里的 VALUES 子句中只允许有一个元组）。最后对结果元组执行投影操作并返回。

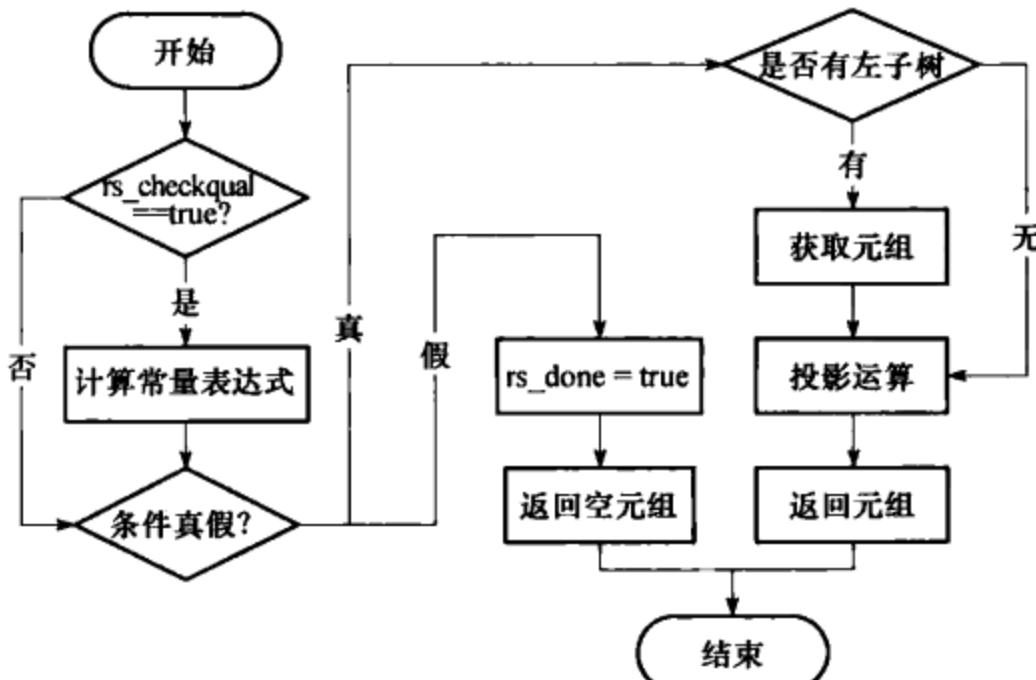


图 6-23 Result 节点的执行流程

对于 Result 节点的清理过程（ExecEndResult 函数），由于不存在右子节点，所以只需调用左子节点的清理过程。

2. Append 节点

Append 节点用于处理包含一个或多个子计划的链表。Append 处理过程会逐个处理这些子计划，当一个子计划返回了所有的结果后，会执行链表中的下一个子计划，直到链表中的所有子计划都被执行完。从这种执行方式可知，Append 节点不使用左右孩子节点。其特点就在于依次执行子计划链表，可用于处理包含多个子查询的 UNION 操作等。

Append 节点的定义如图 6-24 所示，它在 Plan 节点的基础上扩展定义了 appendplans 和 isTarget 字段，其中 appendplans 用于存储子计划链表，isTarget 则用于表示 Append 的子计划链表中的表是否会被修改（即 INSERT/UPDATE/DELETE 操作的对象表）。

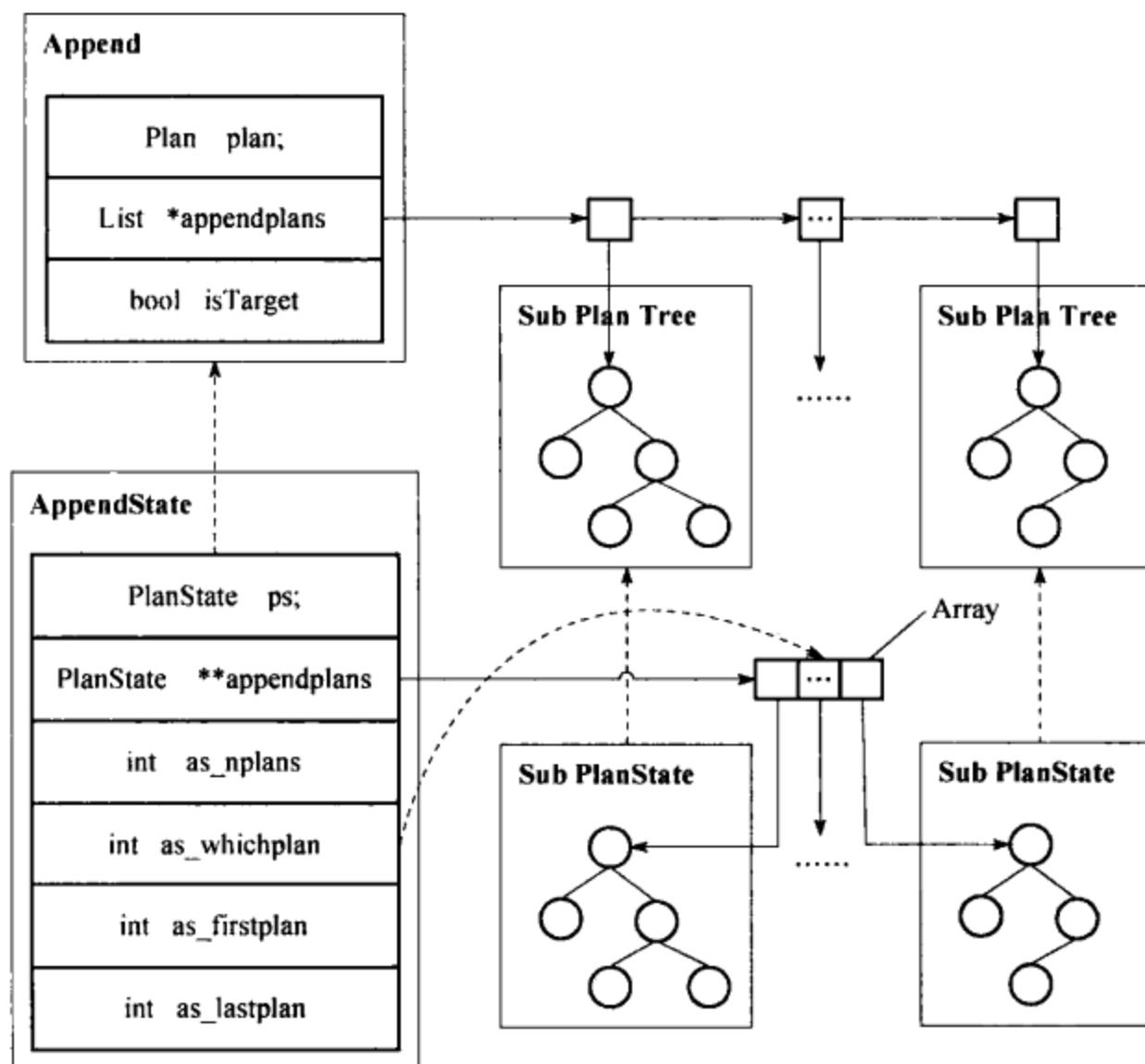


图 6-24 Append 节点相关数据结构

Append 节点的初始化过程（ExecInitAppend 函数）会初始化 AppendState 节点，然后初始化子计划链表中的每一个子计划，并将它们的状态记录节点组织成一个数组，使 AppendState 节点的 appendplans 字段指向该数组。as_nplans 记录 Append 节点中子计划链表的长度，as_whichplan 则用于在执行中指示当前处理的子计划在链表中的偏移量，初始化时一般设置为 as_firstplan 的值。而 as_firstplan 用于标识 Append 节点处理的第一个子计划在链表中的偏移量（一般为 0，表示第一个），as_lastplan 用于

标识 Append 节点处理的最后一个子计划在链表中的偏移量（一般为链表长度减 1，表示最后一个）。

如图 6-25 所示，Append 节点的执行过程由 ExecAppend 函数完成。它会首先从 as_whichplan 标记的子计划开始执行，如果返回的元组非空，则直接返回结果。如果当前子计划返回的是空元组，那么会将 as_whichplan 移动到下一个子计划（因为执行过程可以向前和向后扫描，向前扫描为加 1，反之减 1），并继续执行当前子计划。倘若已经没有下一个子计划可以处理，则直接返回空元组。

由于没有左右孩子节点，Append 节点清理工作（ExecEndAppend 函数）不需要递归调用左右子节点的清理过程，而是扫描整个子计划链表，依次调用子计划的清理函数。

3. BitmapAnd/BitmapOr 节点

BitmapAnd 和 BitmapOr 都是位图（Bitmap）类型节点，用于位图计算。

首先介绍一个位图运算的例子。在使用索引扫描表时，可用位图来标识满足扫描条件的元组，即将满足条件的元组的对应位置 1，这样就可以使用位图来表示整个表中满足扫描条件的元组。当一个表上有多个属性约束且都有索引时，可以利用各个索引分别扫描得到满足对应属性约束的结果位图，然后根据多个属性约束上的逻辑关系对位图进行与或运算得到最终的结果位图。

如图 6-26 所示，假设表有属性 attrA 和 attrB，且分别建立了索引。如果需要进行的查询中有条件“（涉及 attrA 的条件） AND （涉及 attrB 的条件）”，则可以首先利用第一个条件扫描属性 attrA 上的索引并构建位图 Bitmap A，其中每一位对应表中的一个元组，如果元组对应的位为 1 表示该元组满足条件。用同样的方法利用 attrB 上的索引构建位图 Bitmap B。由于两个属性的约束间是“与”关系，因此可以对两个位图进行 AND 操作，从而得到 Bitmap AB，其中值为 1 的位表示对应的元组同时满足 A、B 两属性上约束。

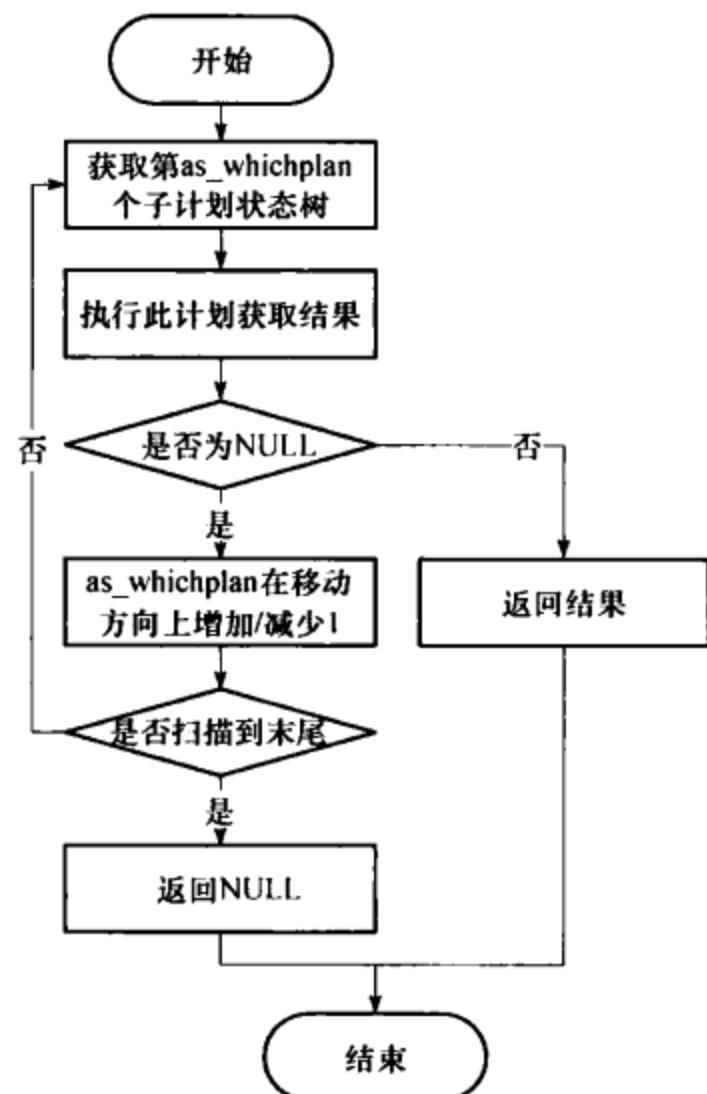


图 6-25 Append 节点执行流程

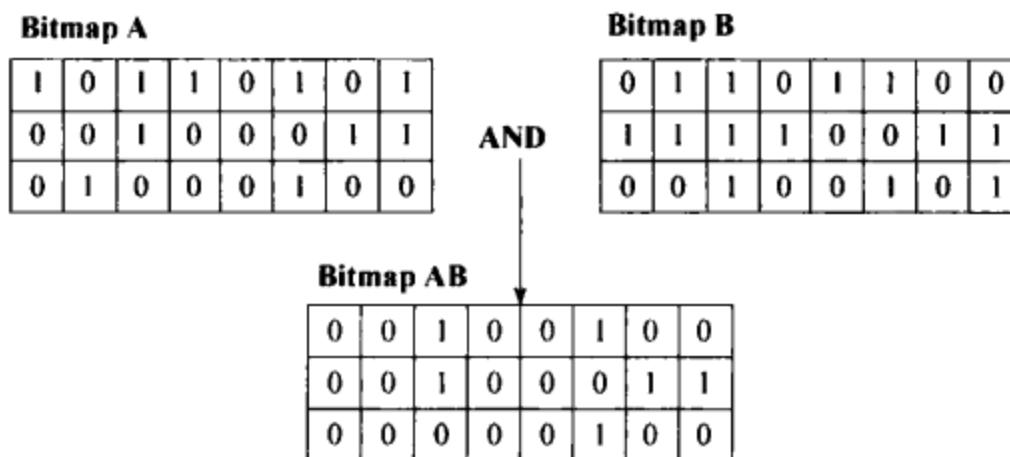


图 6-26 位图操作示例

BitmapAnd 和 BitmapOr 节点实现了两个或多个位图的“与”和“或”运算。这两个节点的数据组织类似于 Append 节点，将产生每一个位图的子计划放在一个链表里，在执行过程中先执行子计划节点获取位图，然后进行“与”（“或”）操作。图 6-27 中给出了 BitmapAnd 和 BitmapOr 节点的数据结构，它们都有一个子计划的链表（bitmapplans 字段）。但和 Append 节点不同之处在于这两个节点的子计划返回的是位图，而不是元组。

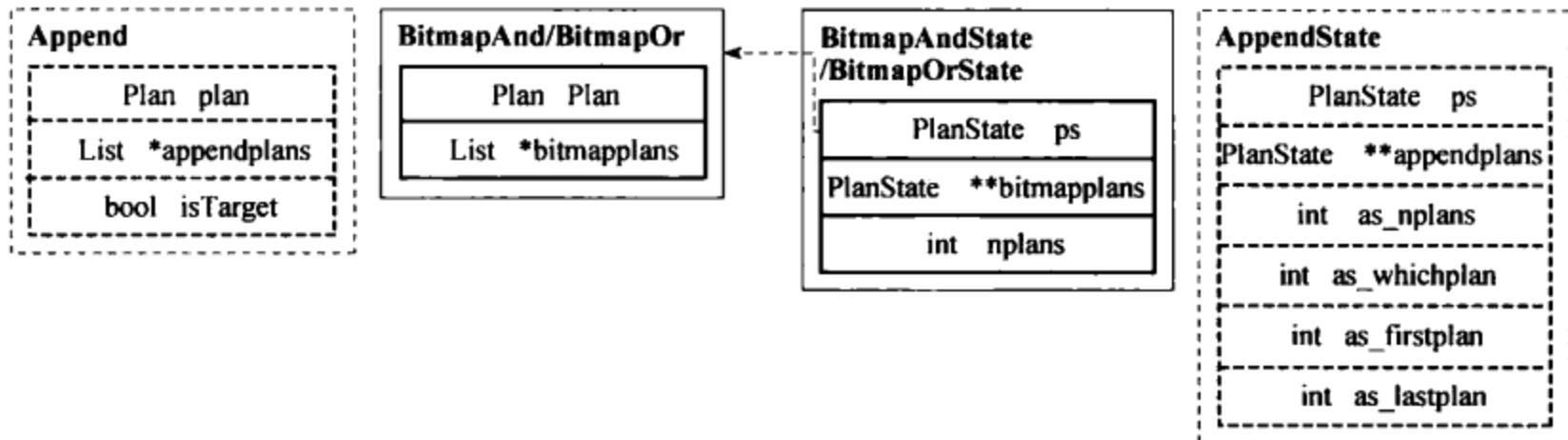


图 6-27 BitmapAnd/BitmapOr 节点相关数据结构

Bitmap 类节点的初始化过程（ExecInitBitmapAnd 和 ExecInitBitmapOr 函数）将初始化 BitmapAndState/BitmapOrState 节点，其过程与 AppendState 节点类似，状态节点定义中也扩展了子计划链表对应的状态节点指针数组 bitmapplans 以及子计划个数 nplans。Bitmap 类节点执行时也是依次获取每个子计划的位图，进行与/或操作，并将结果位图返回。其清理过程也是依次调用每个子计划的清理过程，而不会处理未使用的左右子节点。

4. RecursiveUnion 节点

RecursiveUnion 节点用于处理递归定义的 UNION 语句。下面给出一个例子，使用 SQL 语句定义从 1 到 100 求和，查询语句如下：

```

WITH RECURSIVE t (n) AS (
    VALUES (1)
    UNION ALL
    SELECT n + 1 FROM t WHERE n < 100
)
SELECT sum(n) FROM t;

```

该查询定义了一个临时表 t，并将 VALUES 子句给出的元组作为 t 的初始值，“SELECT n + 1 FROM t WHERE n < 100” 将 t 中属性 n 小于 100 的元组的 n 值加 1 并返回，UNION ALL 操作将 SELECT 语句产生的新元组集合合并到临时表 t 中。以上的 SELECT 和 UNION 操作将递归地执行下去，直到 SELECT 语句没有新元组输出为止。最后执行“SELECT sum (n) FROM t”对 t 中元组的 n 值做求和操作。

上述查询由 RecursiveUnion 节点处理，其过程如图 6-28 所示。有一个初始输入集作为递归过程的初始数据（如上例中“VALUES

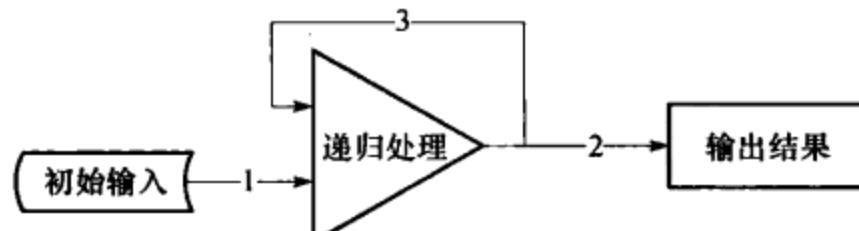


图 6-28 递归合并操作处理过程

(1)”), 然后进行递归部分（上例中的“SELECT n + 1 FROM t WHERE n < 100”）的处理得到输出，并将新得到的输出与初始输入合并后作为下次递归的输入（例如第一次递归处理时新的输出集为 {2}，与初始输入合并后为 {1, 2}）。

RecursiveUnion 节点的数据结构如图 6-29 所示，它在 Plan 的基础上扩展定义了 wtParam 字段用于与 WorkTableScan（负责对临时表的扫描，即递归部分的查询语句的执行，将在 6.4.2 节中介绍）传递参数。执行中 RecursiveUnion 与 WorkTableScan 之间的参数传递是通过 Estate 中 es_param_exec_vals 指向的参数数组完成的，wtParam 表示传递的参数在该数组中的偏移量。其他四个扩展属性（numCols、dupColIdx、dupOperators 和 numGroups）用于 UNION 时不包含关键字 ALL 的情况[⊖]，numCols 存储了用于去重判断的属性个数、dupColIdx 数组记录用于去重判断的属性号、dupOperators 数组则记录了用于判断各属性是否相同的函数的 OID，numGroups 则记录了结果元组数的估计值。

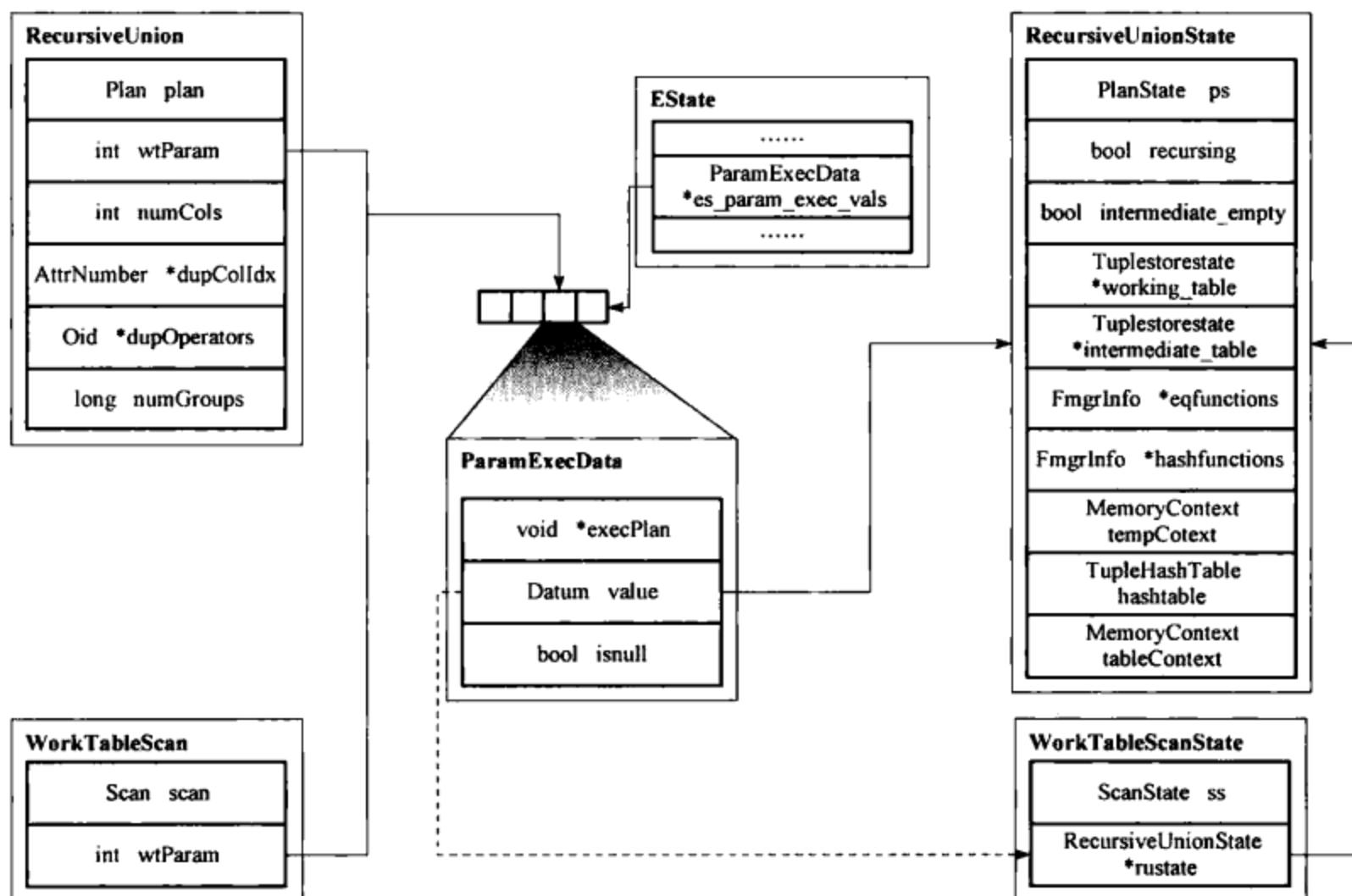


图 6-29 Recursive Union 操作相关数据结构

RecursiveUnion 节点的初始化过程（ExecInitRecursiveUnion 函数）除标准初始化过程外，首先会根据 numCols 是否为 0 来辨别是否需要在合并时进行去重。如果需要去重，则根据 dupOperators 字段初始化 RecursiveUnionState 节点的 eqfunctions 和 hashfunctions 字段，同时申请两个内存上下文用于去重操作，并创建去重操作用的 Hash 表。另外，还将为 RecursiveUnionState 节点的 working_table 和 intermediate_table 字段初始化元组缓存结构。

RecursiveUnion 节点的执行过程（ExecRecursiveUnion 函数）如图 6-30 所示，首先执行左子节点

[⊖] UNION (不含 ALL 关键字) 操作会对结果进行去重，类似于 DISTINCT 关键字。含有 ALL 关键字时，不会进行去重处理。

计划（初始集子计划），将获取的元组直接返回。如果需要去重，则需要把返回的元组加入到 Hash 表（hashtable）中。同时，将初始集存储在 working_table（图中缩写为 WT）指向的元组缓存结构中。当处理完毕所有的左子节点计划后，会执行右子节点计划以获取结果元组。其中，右子节点计划中的 WorkTableScan 会以 working_table 为扫描对象。右子节点计划返回的结果元组将作为 RecursiveUnion 节点的输出，同时缓存在另一个元组存储结构 intermediate_table（图中缩写为 IT）中。每当 working_table 被扫描完毕，RecursiveUnion 节点的执行流程会将 intermediate_table 赋值给 working_table，然后再次执行右子节点计划获取元组，直到无元组输出为止。同样，如果需要去重，右子节点计划的所有输出也会被存入同一个 Hash 表（hashtable），若重复则不会输出。

在 RecursiveUnion 节点的清理过程中，除了对子节点完成递归清理和状态节点的清理外，还会针对初始化中创建的各种数据结构进行清理，例如清理 working_table 和 intermediate_table 两个元组缓存结构，以及去重操作时用到的内存上下文。

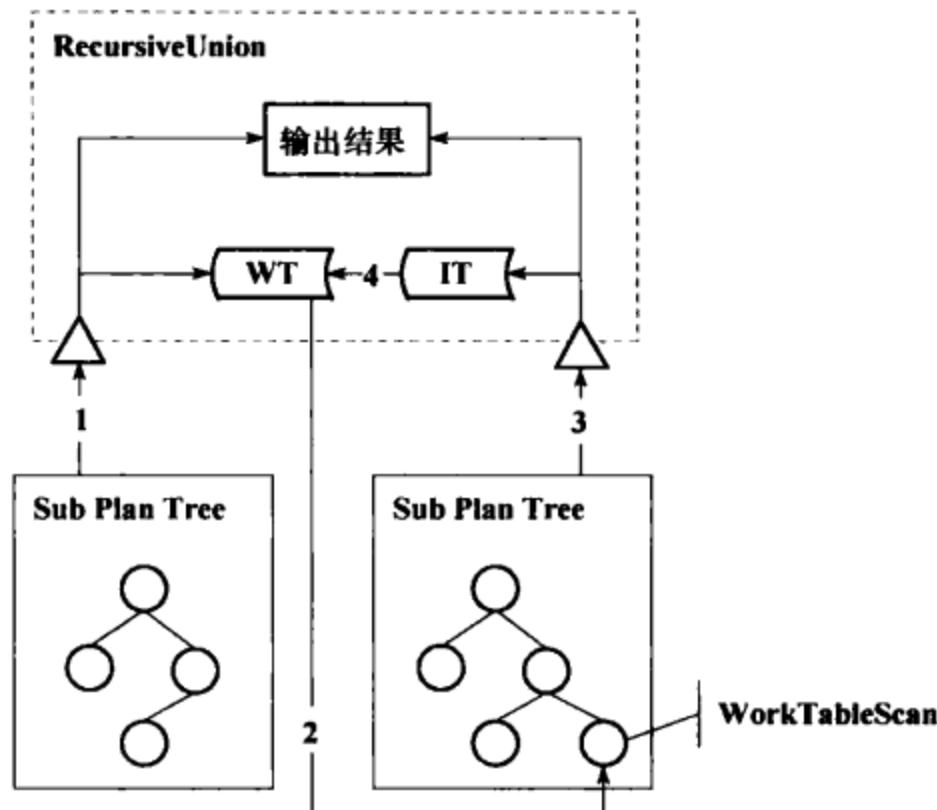


图 6-30 RecursiveUnion 处理过程

6.4.2 扫描节点

扫描节点的作用是扫描表，每次获取一条元组作为上层节点的输入。扫描节点普遍存在于查询计划树的叶子节点，它不仅可以扫描表，还可以扫描函数的结果集、链表结构、子查询结果集等。

出于优化的需要，扫描节点中有一类使用索引进行扫描的节点，前面介绍的 BitmapAnd 和 BitmapOr 节点所需要的位图也是通过索引扫描节点得到的。表 6-3 列出了 PostgreSQL 8.4.1 提供的扫描节点。

表 6-3 扫描节点

节点类型	文件	功能
SeqScan	nodeSeqScan.c	顺序扫描
IndexScan	nodeIndexScan.c	利用索引扫描
BitmapHeapScan	nodeBitmapHeapScan.c	利用 Bitmap 结构获取元组
BitmapIndexScan	nodeBitmapIndexScan.c	利用索引结构获取满足选择条件的 Bitmap
TidScan	nodeTidScan.c	用于扫描一个元组 TID 数组
SubqueryScan	nodeRecursiveUnion.c	扫描一个子查询
FunctionScan	nodeFuncs.c	处理范围表中含有函数的扫描
ValuesScan	nodeValueScan.c	用于扫描 Values 链表
CteScan	nodeCtescan.c	用于扫描 CommonTableExpr ^①
WorkTableScan	nodeWorktableScan.c	扫描 RecursiveUnion 迭代的中间数据 ^②

① 用于存储 WITH 子句的数据结构。

② 见 6.4.3 节，控制节点 RecursiveUnion 处理过程中 WT（ Tuplestorestate 结构），是一个存储迭代中间结果的数据结构，能够使用内存和物理文件多种方式存储。

所有扫描节点都使用 Scan 作为公共父类，Scan 不仅继承了 Plan 的所有属性，还扩展定义了 scanrelid 用于记录被扫描的表在范围表中的序号。扫描节点的执行状态节点都以 ScanState 作为公共父类，ScanState 除了继承 PlanState 的所有属性之外，还扩展定义了 ss_currentScanDesc（记录扫描的位置、关系等信息）和 ss_ScanTupleSlot（记录扫描到的结果）。

扫描节点有各自的执行函数，但是这些执行函数都由公共的执行函数 ExecScan 来实现（参见图 6-31）。ExecScan 需要两个参数，一个是状态节点，另一个是获取扫描元组的函数指针（AccessMtd，由于每一种扫描节点扫描的对象不同，因此函数都不同），ExecScan 迭代地扫描对象，每次执行返回一条结果。ExecScan 会使用 accessMtd 获取元组，存放在 ScanState 的 ss_ScanTupleSlot 中，然后进行过滤条件判断和投影操作，最终返回元组。

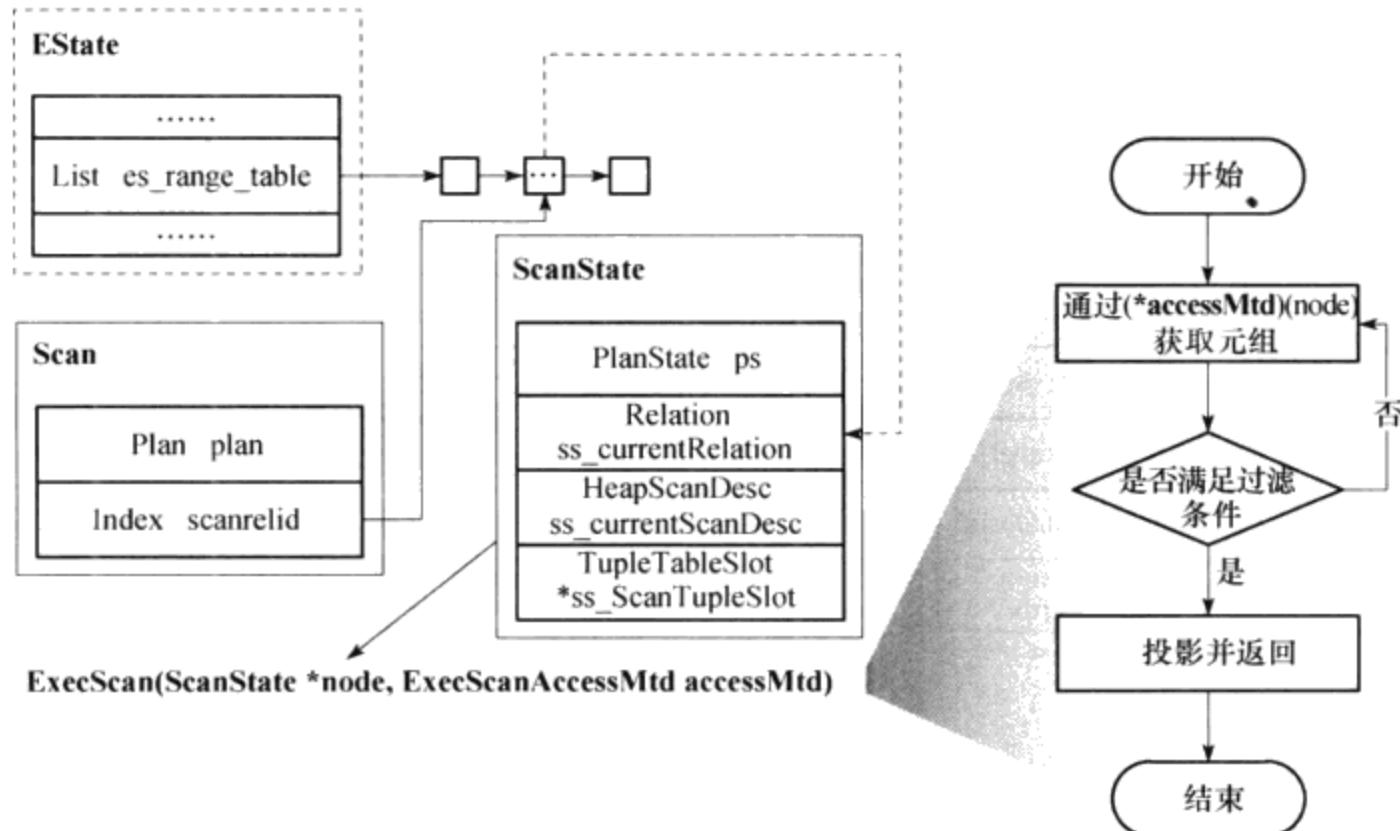


图 6-31 扫描节点公共函数及相关数据结构

1. SeqScan 节点

SeqScan 是最基本的扫描节点，它用于扫描物理表，完成没有索引辅助的顺序扫描过程。其计划节点 SeqScan 实际是 Scan 节点的一个别名，并未定义扩展属性。其执行状态节点 SeqScanState 也直接使用 ScanState。

SeqScan 节点的初始化由函数 ExecInitSeqScan 完成。该函数首先创建一个 SeqScanState 结构，将 SeqScan 节点链接在 SeqScanState 结构的 ps 字段中。然后调用 ExecInitExpr 对计划节点的目标属性和查询条件进行初始化，并将它们链接到 SeqScanState 相应的字段中。接下来还将为计划节点分配用于存储结果元组和扫描元组的数据结构。最后通过计划节点中 scanrelid 字段的信息获取被扫描对象的 RelationData 结构，并链接在 ss_currentRelation 字段中，同时利用该信息调用 heap_beginscan 初始化扫描描述符 ss_currentScanDesc。

SeqScan 节点的执行函数是 ExecSeqScan，在该函数中会调用 ExecScan 函数，并将 SeqNext 函数的指针作为 ExecScan 函数 AccessMtd 参数的值。SeqNext 函数将通过存储模块提供的函数 heap_get-

next 获取下一条元组并返回。ExecScan 在利用 SeqNext 获得一个元组之后，还将根据计划节点中的查询条件和投影要求对得到的元组进行条件检查和投影操作，最后将满足要求的结果元组返回。其他扫描节点的执行函数都采用类似的方式管理，即统一调用 ExecScan，但根据节点的类型给 ExecScan 的参数 AccessMtd 赋予不同的函数指针。

SeqScan 节点的清理过程由函数 ExecEndSeqScan 完成，在该函数中需要额外调用函数 heap_end-scan 来清理 ss_currentScanDesc 内的信息。

2. IndexScan 节点

如果选择条件涉及的属性上建立了索引，则生成的查询计划中涉及表的扫描时会使用 IndexScan 节点。该节点能够利用索引进行表的扫描以获取满足选择条件的元组。

IndexScan 节点的定义如图 6-32 所示。除了继承 Scan 节点定义的属性外，IndexScan 扩展定义了 indexid 属性（用于存储索引的 OID）、indexqual 属性（用于存储索引扫描的条件）、indexqualorig 属性（用于存储没有处理的原始扫描条件链表）[○]以及 indexorderdir 属性（用于存储扫描的方向）。

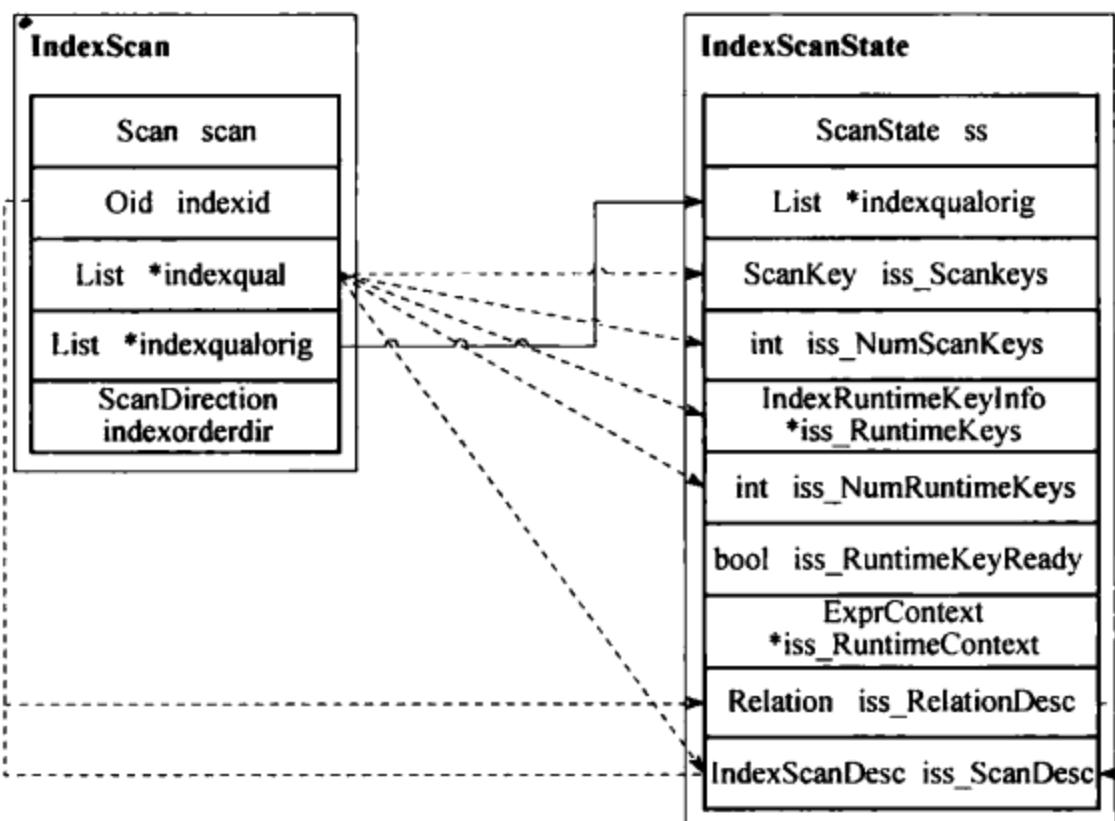


图 6-32 IndexScan 节点相关数据结构

IndexScan 节点的初始化过程由函数 ExecInitIndexScan 完成。该函数将构造 IndexScanState 节点，并使用 indexid 获取索引的 RelationData 结构存放于 iss_RelationDesc 字段中。同时，通过调用 ExecIndexBuildScanKeys 将 indexqual 中的索引扫描条件转换为扫描关键字（ScanKey，存储扫描满足的条件）以及运行时关键字计算结构（IndexRuntimeKeyInfo，执行时才能得到结果的表达式信息）分别存储在 iss_ScanKeys 和 iss_RuntimeKeys 这两个数组中。iss_NumScanKeys 和 iss_NumRuntimeKeys 则用于指示前面两个数组的长度，同时还要设置 iss_NumRuntimeKeys 为 false。最后将调用索引模块提供的 index_beginscan 初始化扫描描述符 iss_ScanDesc。而索引扫描未经特殊处理的原始约束条件链表

[○] 为了满足索引扫描的需要，indexqual 存储的条件表达式是由 indexqualorig 经过针对索引扫描处理后得到的条件表达式链表。

则用于构造 indexqualorig 字段^①。

IndexScan 节点的执行过程由 ExecIndexScan 函数完成，其执行过程同样由 ExecScan 统一管理，但对 IndexScan 节点将使用 IndexNext 函数来获取元组。ExecIndexScan 首先判断是否有 RuntimeKeys 且需要计算（iss_RuntimeKeyReady 为 false），如果存在则调用 ExecIndexReScan 函数计算所有的 iss_RuntimeKeys 表达式，并将其存储到关联的 iss_ScanKeys 中。接着调用 ExecScan 通过 IndexNext 获取元组，在 IndexNext 中将调用索引模块提供的 index_getnext 函数利用索引取得元组。

IndexScan 的清理过程由 EndIndexScan 函数完成，其中需要回收索引关系描述结构 iss_RelationDesc（调用 index_close）和索引扫描描述符 iss_ScanDesc（调用 index_endscan）。

3. BitmapIndexScan 节点

BitmapIndexScan 节点也是利用属性上的索引进行扫描操作，但是 BitmapIndexScan 得到的结果不是实际的元组，而是一个位图，该位图标记了满足条件的元组在页面中的偏移量。BitmapIndexScan 节点在第一次被执行时就将获取所有满足条件的元组并在位图中标记它们，而其上层节点中都会有特殊的扫描节点（例如下面将介绍的 BitmapHeapScan）使用该位图来获取实际的元组。

从图 6-33 可以看出，BitmapIndexScan 与 IndexScan 节点定义几乎相同，由于一次返回所有的元组，所以不需要记录扫描方向的 indexorderdir 字段。其执行状态节点 BitmapIndexScanState 与 IndexScanState 也很相似，但多出了表示索引关键字属性数组及其长度的字段。BitmapIndexScan 和 IndexScan 的执行过程也类似，只是在 BitmapIndexScan 处理过程中，初始化扫描状态使用 index_beginscan_bitmap 函数，该函数将调用 index_getbitmap 生成位图，并将其存放在执行状态记录节点的 biss_result 字段中。

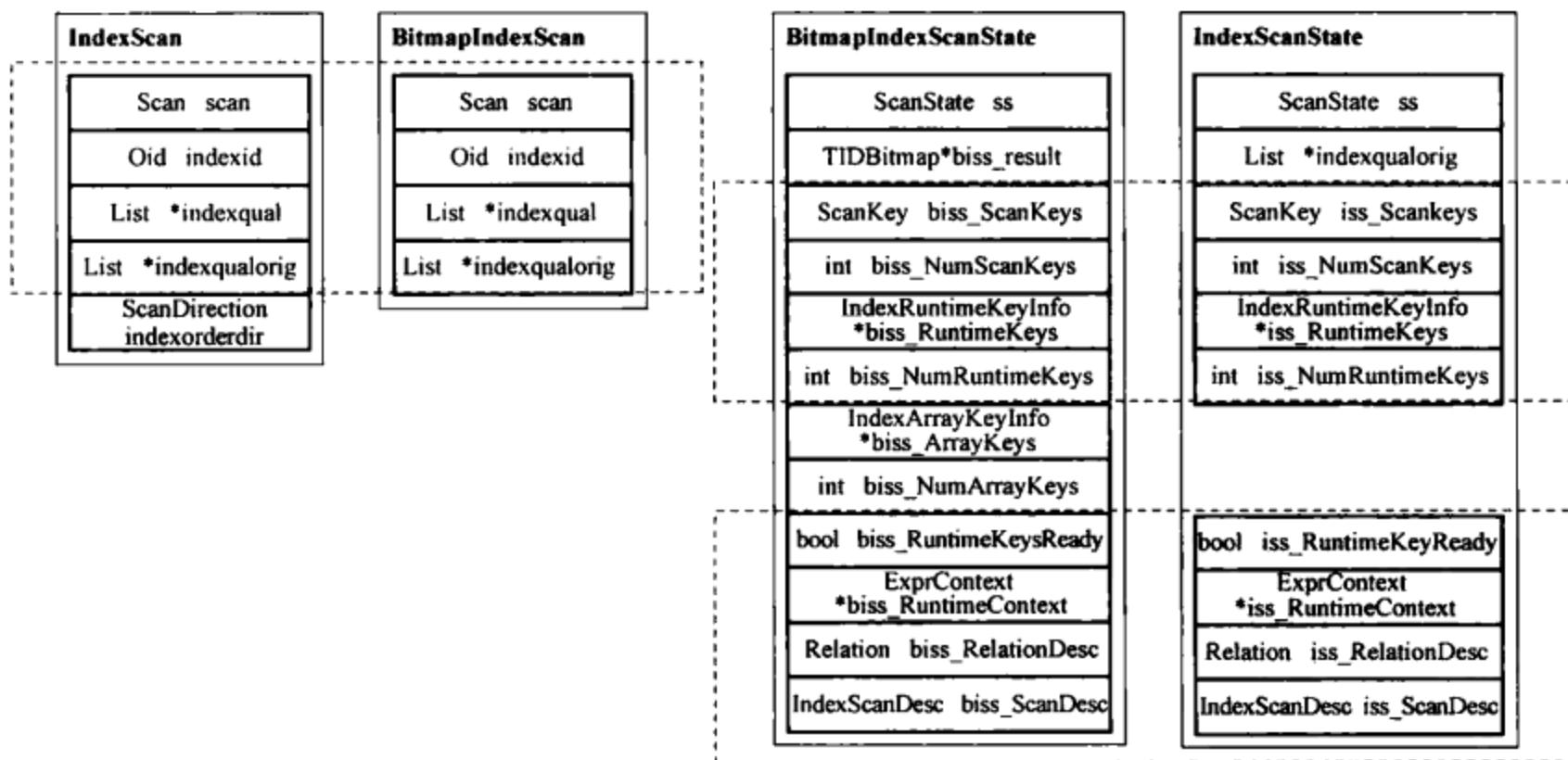


图 6-33 BitmapIndexScan 节点相关数据结构

^① 用于并发控制，执行过程发现当前处理元组被其他事务修改并已提交，需要检测该元组是否仍满足选择条件。

4. BitmapHeapScan

上一节介绍的 BitmapIndexScan 节点将输出位图而不是元组，为了根据位图获取实际的元组，PostgreSQL 提供了 BitmapHeapScan 节点从 BitmapIndexScan 输出的位图中获取元组。

BitmapHeapScan 节点定义如图 6-34 所示，该节点在 Scan 的基础上仅扩展了约束条件检查字段 (bitmapqualorig)，该字段与 IndexScan 节点的 indexqualorig 功能相同。当并发事务修改并提交了当前处理的元组时，需要重新扫描更新后的元组是否满足约束条件，而不是重新获取位图，因此将直接使用该表达式进行条件计算。BitmapHeapScan 有且仅有一个子节点（左子节点），显然这个左子节点必须是提供位图输出的计划节点。

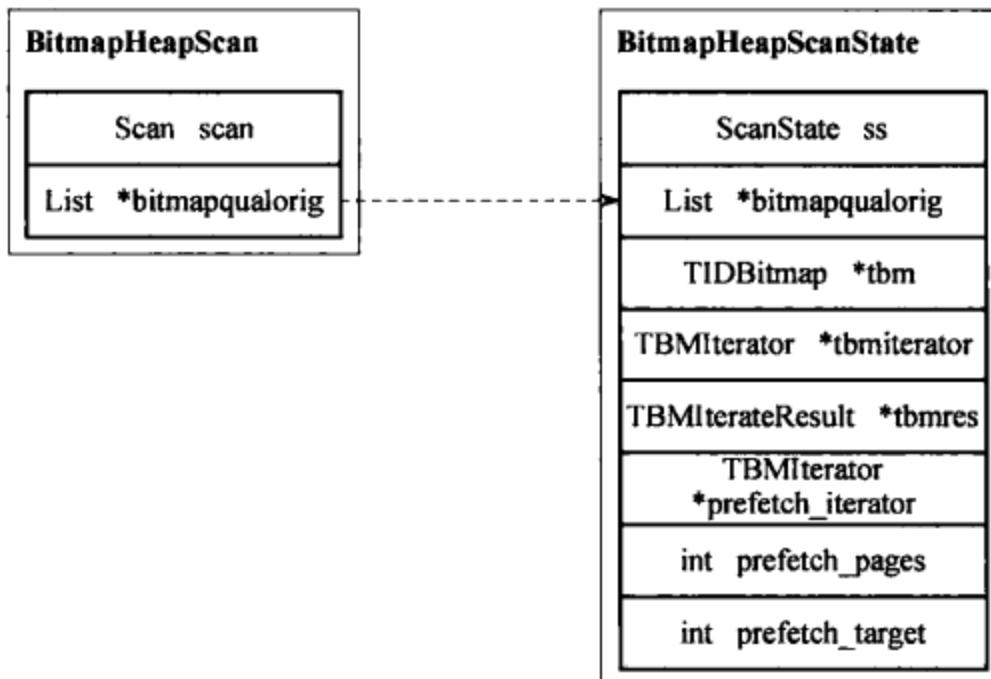


图 6-34 BitmapHeapScan 节点相关数据结构

初始化函数 ExecInitBitmapHeapScan 会根据节点中的 scanrelid 初始化扫描描述符 ss_currentScanDesc。其他的初始化设置在执行过程中进行。

执行函数 ExecBitmapHeapScan 会将 BitmapHeapNext 函数指针传递给 ExecScan，ExecScan 使用 BitmapHeapNext 获取元组。BitmapHeapNext 首先判断 BitmapHeapScanState 的 tbm（位图）是否为空，如果为空则调用 MultiExecProcNode 从左子节点获取位图，并调用 tbm_begin_iterate 初始化 tbmiterator。如果需要预取还要调用 tbm_begin_iterate 初始化 prefetch_iterator，并将 prefetch_pages 置 0，prefetch_target 设置为 -1。然后执行过程会利用 tbmiterator 遍历位图获取物理元组的偏移量，然后从对应的缓冲区中按照偏移量取出元组并返回。

清理过程 ExecEndBitmapHeapScan 需要调用左子节点的清理函数，然后清理 tbmiterator、prefetch_iterator 以及 tbm 位图，最后清理扫描描述符并关闭打开的表。

5. TidScan 节点

PostgreSQL 系统专用于标识元组物理位置的数据类型被称作 TID (Tuple Identifier)，一个 TID 由块号和块内偏移量组成，系统属性^① ctid 被定义为此种类型。在定义游标后，可以使用“UP-

^① PostgreSQL 系统中每个元组都会定义对用户隐藏的属性（例如 ctid），用于并发等系统功能，被称为系统属性。

DATE/DELETE…WHERE CURRENT OF…”语句对当前游标位置的元组进行修改/删除，此时生成的查询计划树仅包含一个 TidScan 节点，其扫描的对象是 TidScan 节点中保存的一个表达式链表，其中存储的表达式可以得到 ctid 的值，TidScan 节点将根据 ctid 值取得对应的元组。如图 6-35 所示，TidScan 节点只在 Scan 节点的基础上扩展了一个字段 tidquals 用于保存可以得到 ctid 的表达式链表。

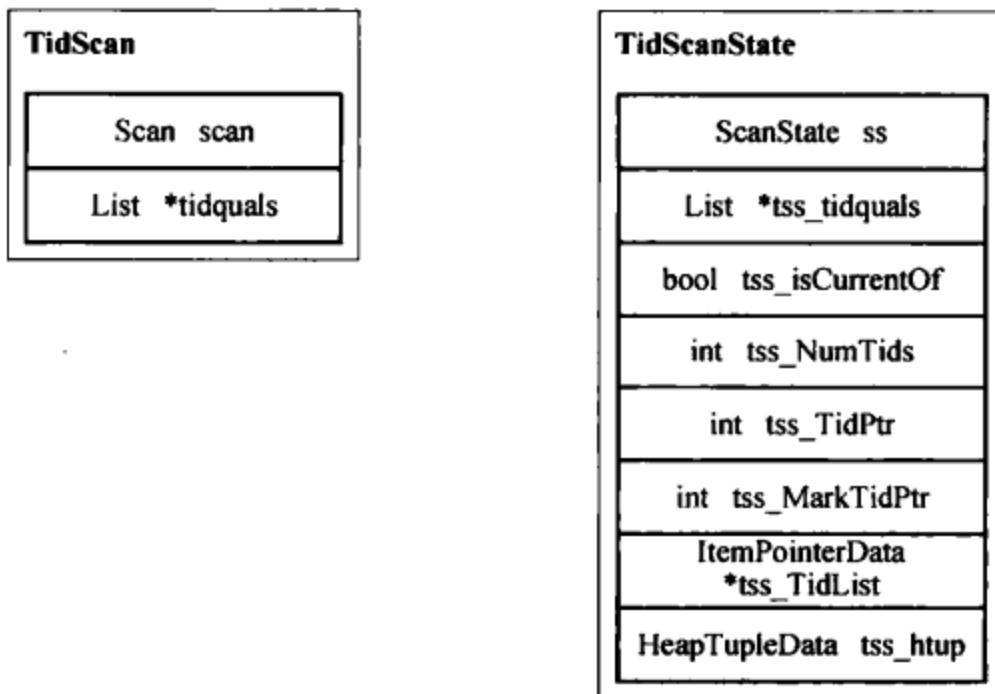


图 6-35 TidScan 节点相关数据结构

TidScan 节点的初始化函数 ExecInitTidScan 会根据 tidquals 初始化 TidScanState 中的 tss_tidquals 字段，然后调用 ExecInitExpr 初始化 tidquals 中的表达式，并根据节点中的 scanrelid 初始化扫描描述符 ss_currentScanDesc。

TidScan 节点的执行函数 (ExecTidScan) 也同样调用函数 ExecScan 来完成执行工作，其中传递给 ExecScan 函数的指针是 TidNext。函数 TidNext 首先需要通过计算 TidScanState 节点的 tss_tidquals 链表中的表达式来构造 tss_TidList 数组，该数组中存放的是一系列的 ctid，tss_NumTids 用于记录数组的长度，tss_TidPtr 用于记录当前处理的 ctid 在 tss_TidList 数组中的偏移量，初始值设为 -1。然后从 tss_TidList 中获取下一个 ctid 值，接着调用存储模块提供的 heap_fetch 根据该 ctid 获取元组并返回。出于并发的需要，当 TidScan 节点用于“CURRENT OF… (游标名)”语句时，获取的 ctid 可能已经被其他事务修改，需要获取此 ctid 对应元组的最新版本（利用 HOT 链，见 3.2.1 节），然后再调用 heap_fetch 进行获取。

清理过程 ExecEndTidScan 不需要特殊的操作，直接释放相关内存上下文和初始化时分配的空间。

6. SubqueryScan 节点

SubqueryScan 节点的作用是以另一个查询计划树（子计划）为扫描对象进行元组的扫描，其扫描过程最终被转换为子计划的执行。SubqueryScan 节点（数据结构见图 6-36）在 Scan 节点之上扩展定义了子计划的根节点指针（subplan 字段），而 subrtable 字段是查询编译器使用的结构，执行器运行时其值为空。

显然，SubqueryScan 节点的初始化过程（ExecInitSubqueryScan 函数）会使用 ExecInitNode 处理 SubqueryScan 的 subplan 字段指向的子计划树，并将子计划的 PlanState 树根节点指针赋值给 Sub-

queryScanState 的 subplan 字段。

SubqueryScan 节点的执行（ExecSubqueryScan 函数）则是通过将 SubqueryNext 传递给 ExecScan 函数处理来实现的。SubqueryNext 实际则是调用 ExecProcNode 处理 subplan 来获得元组。

在清理过程中，需要额外调用 ExecEndNode 来清理子计划。

7. FuncitonScan 节点

在 PostgreSQL 中，有一些函数可以返回元组的集合，为了能从这些函数的返回值中获取元组，PostgreSQL 定义了 FunctionScan 节点，其扫描对象为返回元组集的函数。如图 6-37 所示，FunctionScan 节点在 Scan 的基础上扩展定义了存储函数调用信息的 funcexpr 字段、函数结果的列名字段（funccolnames 字段）、列类型字段（funccoltypes 字段）以及对应的 typmod 字段（funccoltypmods 字段）。

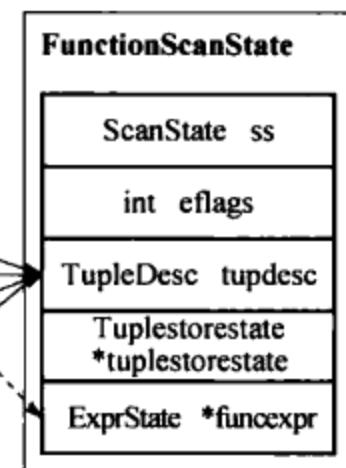
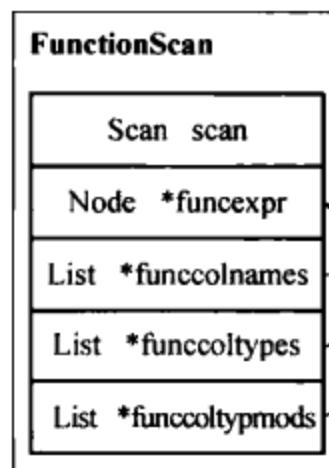
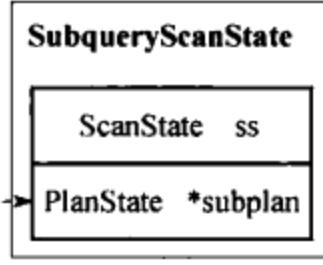
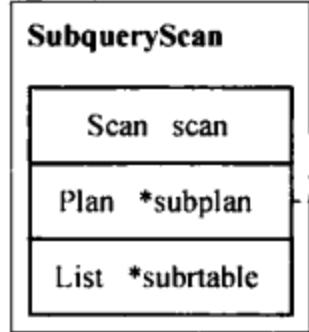


图 6-36 SubqueryScan 节点相关数据结构

图 6-37 FunctionScan 节点相关数据结构

FunctionScan 节点的初始化过程（ExecInitFunctionScan 函数）会初始化 FunctionScanState 结构，并根据 FunctionScan 中的 funcexpr 字段构造用于表达式计算的结构（存储在 funcexpr 中），还要构造函数返回元组的描述符存储在 tupdesc 中，此时用于存储函数结果集的 tuplestorestate 字段为 NULL。

在 FunctionScan 节点的执行过程（ExecFunctionScan 函数）中，将 FunctionNext 传递给 ExecScan 函数，FunctionNext 函数首先判断 tuplestorestate 是否为空（首次执行时为空），如果为空则执行函数生成所有结果集并存储在 tuplestorestate 中，此后每次执行节点将获取结果集中的一个元组。

FunctionScan 节点清理过程需要清理 tuplestorestate 结构。

8. ValuesScan 节点

在 PostgreSQL 中，可以使用 VALUES 子句给出一系列元组，ValuesScan 节点可以对 VALUES 子句给出的元组集合进行扫描（INSERT 语句中的 VALUES 子句除外）。如图 6-38 所示，ValuesScan 节点中的 values_lists 存储了 VALUES 子句中的表达式链表。

ValuesScan 节点的初始化过程（ExecInitValuesScan 函数）处理 values_lists 中的表达式生成 Values 表达式，并存储在 ValuesScanState 的 exprlists 数组中，array_len 记录数组长度，curr_idx 和 marked_idx 用于存储数组中的偏移量。同时还会分配内存上下文 rowcontext 用于表达式计算。

ValuesScan 节点执行过程（ExecValuesScan 函数）调用 ExecScan 实现，ExecScan 通过 ValuesNext 获取扫描元组，ValuesNext 则通过 curr_idx 从 exprlists 中获取需要处理的表达式，并计算出结果元组返回。

在 ValuesScan 节点清理过程（ExecEndValuesScan 函数）中需要释放内存上下文 rowcontext。

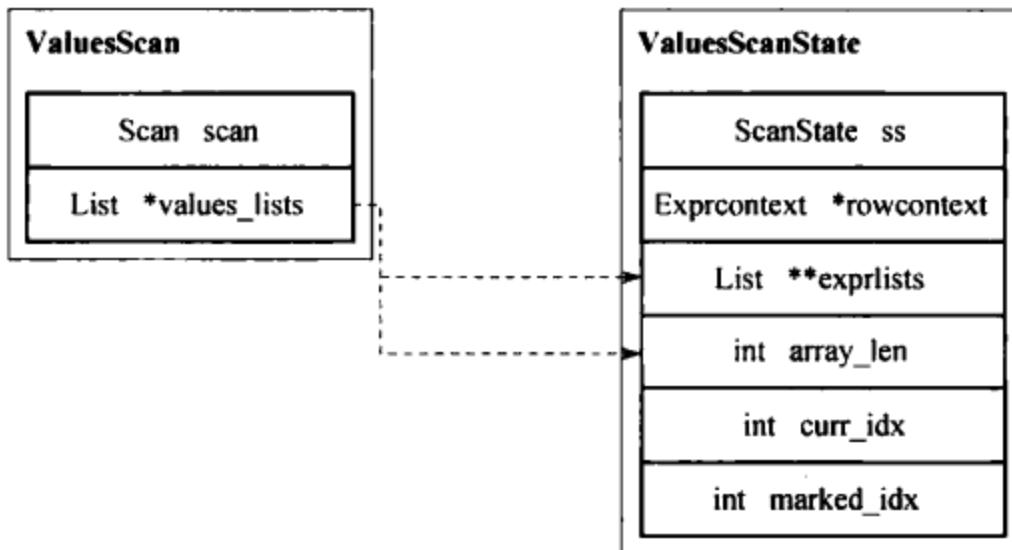


图 6-38 ValuesScan 节点相关数据结构

9. CteScan 节点

CteScan 节点用于扫描 SELECT 查询中用 WITH 子句定义的子查询。在 SELECT 查询中，允许使用 WITH 语句定义一个或多个命名的子查询，在主查询（Primary Query）中使用名称来引用该子查询，子查询在主查询中的效果类似于临时表或视图。CteScan 扫描的对象是一个子查询，这种子查询称为公共表表达式（Common Table Expression）。

在介绍控制节点 RecursiveUnion 时，我们曾使用 WITH 定义了递归子查询，如图 6-39 所示。图中查询计划的根节点为聚集函数计算节点（Aggregate，聚集节点，将在 6.4.3 节中介绍），聚集节点从 CteScan 节点中获取元组并进行聚集计算，而 CteScan 扫描的对象 t（t 是一个 CTE）被作为一个初始子计划（SubPlan）处理，而 t 中的元组则由 RecursiveUnion 为根节点的初始子计划构造。在初始化过程中，将把 CteScan 的子计划的执行状态树存放于执行器全局状态（Estate）的 es_subplanstates 链表中，并在 CteScan 中的 ctePlanId 存储其子计划在该链表中的偏移量，对于同一个子计划的 CteScan 的 ctePlanId 相同。PostgreSQL 在实现时，还为每个 CTE 在一个全局参数链表中分配了一个空间，其偏移量存储在 cteParam 中，对应同一个 CTE 的 CteScan 对应的偏移量也相同。CteScan 节点相关数据结构如图 6-40 所示。

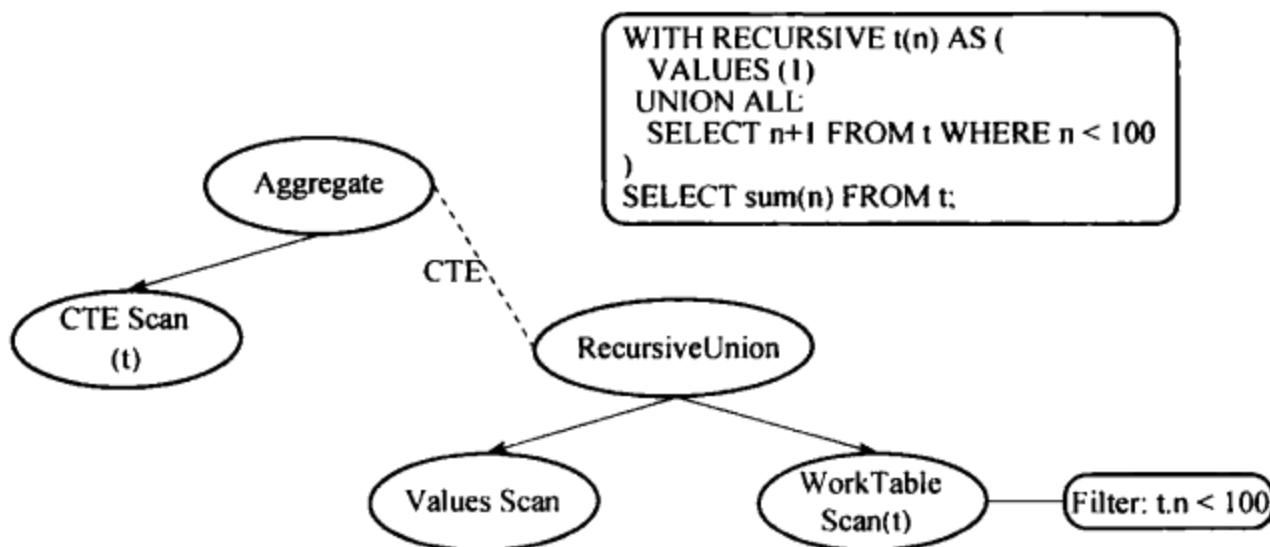


图 6-39 CteScan 示例

CteScan 节点的初始化过程（ExecInitCteScan 函数）将首先初始化 CteScanState 结构，通过 ctePlanId 在 es_subplanstates 中找到对应的子计划执行状态树，并存储在 CteScanState 的 cteplanstate 字段中。然后通过 cteParam 在执行器全局状态的 es_param_exec_vals 中获取参数结构 ParamExecData。若 ParamExecData 中 value 为 NULL，表示没有其他 CteScan 对此 CTE 初始化过存储结构，此时会初始化 CteScanState 的 cte_table 字段，并将 leader 和 ParamExecData 的 value 赋值为指向当前 CteScanState 的指针。若 ParamExecData 中的 value 不为 NULL，则将其值赋值给 leader，让其指向第一个 CteScan 创建的 CteScanState，而不为当前的 CteScan 初始化 cte_table。这样对应一个 CTE 全局只有一个元组缓存结构，所有使用该 CTE 的 CteScan 都会共享该缓存。

在执行 CteScan 节点时，将首先查看 cte_table 指向的缓存中是否缓存元组，如果有可直接获取，否则需要先执行 ctePlanId 指向的子计划获取元组。

CteScan 节点的清理过程需要清理元组缓存结构，但只需清理 leader 指向自身的 CteScanState。

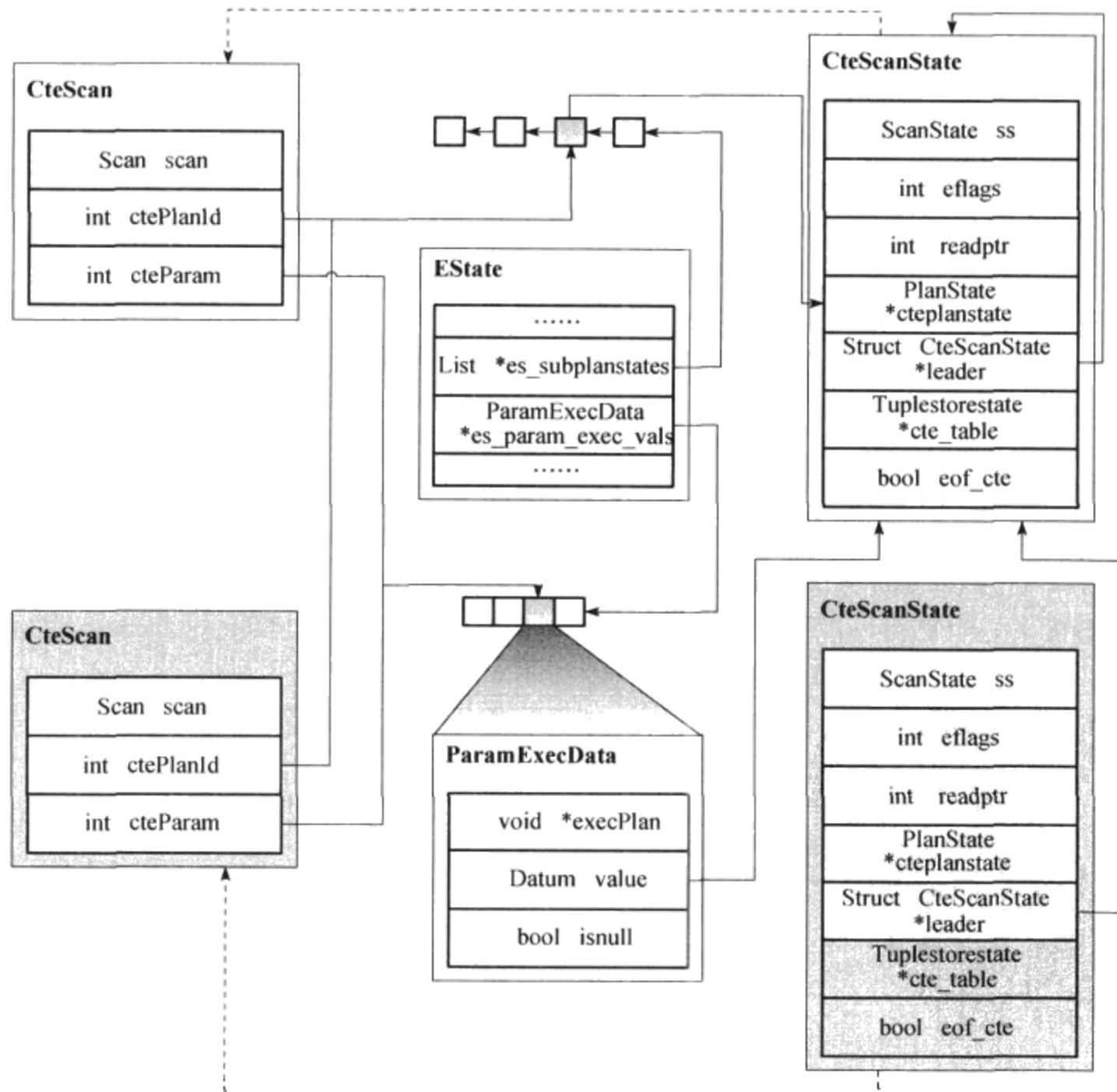


图 6-40 CteScan 节点相关数据结构

10. WorkTableScan 节点

WorkTableScan 会与 RecursiveUnion 共同完成递归合并子查询。RecursiveUnion 会缓存一次递归中的所有元组到 RecursiveUnionState 结构中，WorkTableScan 提供了对此缓存的扫描。

如图 6-41 所示，WorkTableScan 节点扩展定义了 wtParam 用于同 RecursiveUnion 节点间的通信，而 WorkTableScanState 节点的 rustate 字段中记录了 RecursiveUnionState 结构的指针，以便 WorkTableScan 在执行过程中可以从缓存结构中获取元组。WorkTableScan 节点的执行在介绍 RecursiveUnion 节点时已经提及，这里不再赘述。

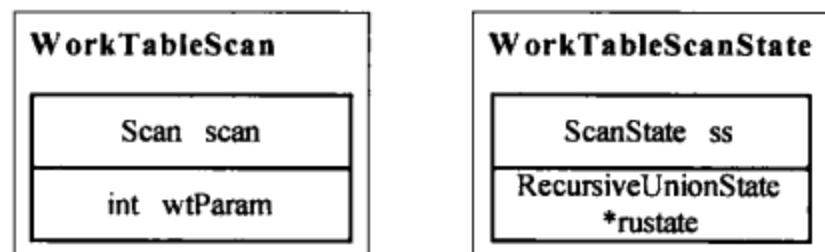


图 6-41 WorkTableScan 节点相关数据结构

6.4.3 物化节点

顾名思义，物化节点是一类可缓存元组的节点。在执行过程中，很多扩展的物理操作符需要首先获取所有的元组后才能进行操作（例如聚集函数操作、没有索引辅助的排序等），这时要用物化节点将元组缓存起来。表 6-4 列出了 PostgreSQL 8.4.1 中提供的物化节点。

表 6-4 物化节点

节点类型	文件	功能
Material	nodeMaterial.c	对子查询结果进行缓存
Sort	nodeSort.c	对底层节点返回的元组进行排序
Group	nodeGroup.c	对下层排序元组进行分组操作
Agg	nodeAgg.c	执行聚集函数
Unique	nodeUnique.c	执行去重操作
Hash	nodeHash.c	Hashjoin 辅助节点
SetOp	nodeSetOp.c	处理集合操作 Exist、Intersect 查询
Limit	nodeLimit.c	处理 Limit 子句
WindowAgg	nodeWindowAgg.c	处理窗口函数（新的聚集方式下的聚集函数）

物化节点需要有元组的缓存结构，以加快执行效率或实现特定功能（例如排序等）。物化节点的功能种类多样，实现过程也不尽相同，缓存的方式也有所不同，主要使用了 tuplestore 来进行缓存。

tuplestore 使用 Tuplestorestate 数据结构（见图 6-42）来存储相关信息和数据，其中 memtuples 字段指定了一块内存区域用于缓存数据，而当内存中缓存的元组达到一定的数量时会将元组写入 myfile 字段指定的临时文件。

tuplestore 通过 tuplestore_begin_heap 来构造和初始化，通过 tuplestore_end 进行释放。tuplestore_puttuple 函数用于将元组加入到 tuplestore，而 tuplestore_gettuple 函数则可以从 tuplestore 中获取元组。

在 tuplestore 的基础之上，PostgreSQL 还提供了 tuplesortstore（由数据结构 Tuplesortstate 实现），它在 tuplestore 内增加了排序功能：当获取所有的元组后，可对元组进行排序。如果没有使用临时文件，则使用快速排序，否则将使用归并法进行外排序。

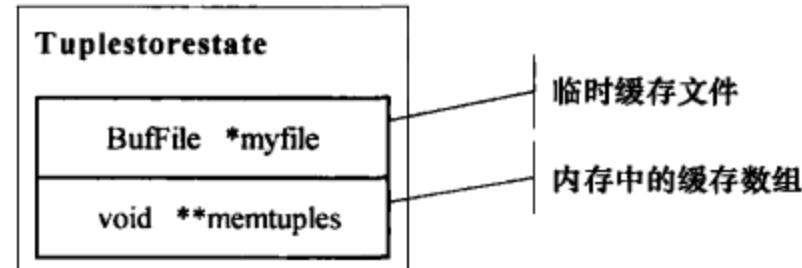


图 6-42 元组缓存结构 Tuplestorestate

1. Material 节点

Material 节点用于缓存子节点结果，对于需要重复多次扫描的子节点（特别是扫描结果每次都相同时）可以减少执行的代价。其实现方式在于将结果元组存储于特殊的数据结构 Tuplestorestate 中。

如图 6-43 所示，Material 节点并没有在 Plan 的基础上定义扩展属性，执行状态节点 MaterialState 扩展了 ScanState 节点的定义，增加了 tuplestorestate 字段用于缓存元组。eof_underlying 则表示下层节点已经扫描完毕，从而避免重复调用下层节点的执行过程，eflags 是一个状态变量，它表示当前节点是否需要支持反向扫描、标记扫描位置和重新扫描三种操作。

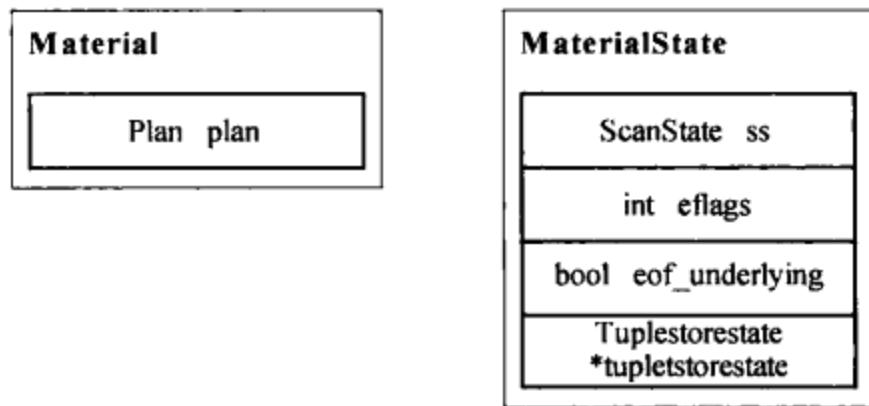


图 6-43 Material 节点相关数据结构

Material 节点的初始化过程（ExecInitMaterial 函数）主要是初始化 eflags 信息，并调用左子节点的初始化过程。

在 Material 节点执行过程（ExecMaterial 函数）中，首先判断是否已经初始化 tuplestorestate，如没有，则会调用 tuplestore_begin_heap 创建 Tuplestorestate 结构。然后把当前缓存中未返回的元组取出并返回，若不存在未返回的元组，则需要进一步判断下层节点是否扫描完毕（eof_underlying 为 true 表示扫描完毕）。如果下层节点没有扫描完成则会从下层节点获取元组放入缓存中，同时返回元组；如果下层节点已经完成扫描，则 Material 节点将返回空元组。

Material 节点的清理工作（ExecEndMaterial 函数）主要是对于元组缓存结构的清理，并调用左子节点的清理过程，最后释放 MaterialState 结构。

2. Sort 节点

Sort 节点（排序节点）用于对于下层节点的输出结果进行排序，该节点只有左子节点。排序节点先将下层节点返回的所有元组缓存起来，然后进行排序。由于缓存结果可能很多，因此不可避免地会用到临时文件进行存储，这种情况下 Sort 节点将使用外排序方法。

Sort 节点的定义如图 6-44 所示，其中保存了进行排序的属性个数（numCols）、用于排序的属性的属性号数组（sortcolIdx，长度为 numCols）和用于每个排序属性的排序操作符 OID 数组（sortOperators，长度也为 numCols）。另外，Sort 节点用一个布尔类型的字段 nullsFirst 记录是否将空值排在前面。

由于 Sort 节点仅对元组进行排序，不需要做投影和选择操作，因此在 Sort 节点的初始化过程（ExecInitSort 函数）中不会对 Plan 结构中的投影（targetlist）和选择（qual）链表进行初始化，只需构造 SortState 节点并调用下层节点的初始化过程。

在 Sort 节点的初次执行过程（ExecSort 函数）中首先会通过 tuplesort_begin_heap 对元组缓存结构初始化，然后循环执行下层节点获取元组，调用 tuplesort_puttupleslot 将每个获取的元组存放到缓存中。获取到下层节点的所有元组之后，调用 tuplesort_performsort 对缓存元组进行排序（该函数根据元组是否全部缓存在内存中来决定使用快速排序还是使用堆排序与归并排序）。后续对 Sort 节点的执行都将直接调用 tuplesort_gettupleslot 从缓存中返回一个元组。

Sort 节点的清理过程（ExecEndSort 函数）需要调用 tuplesort_end 对缓存结构进行清理回收，然

后调用下层节点的清理过程。

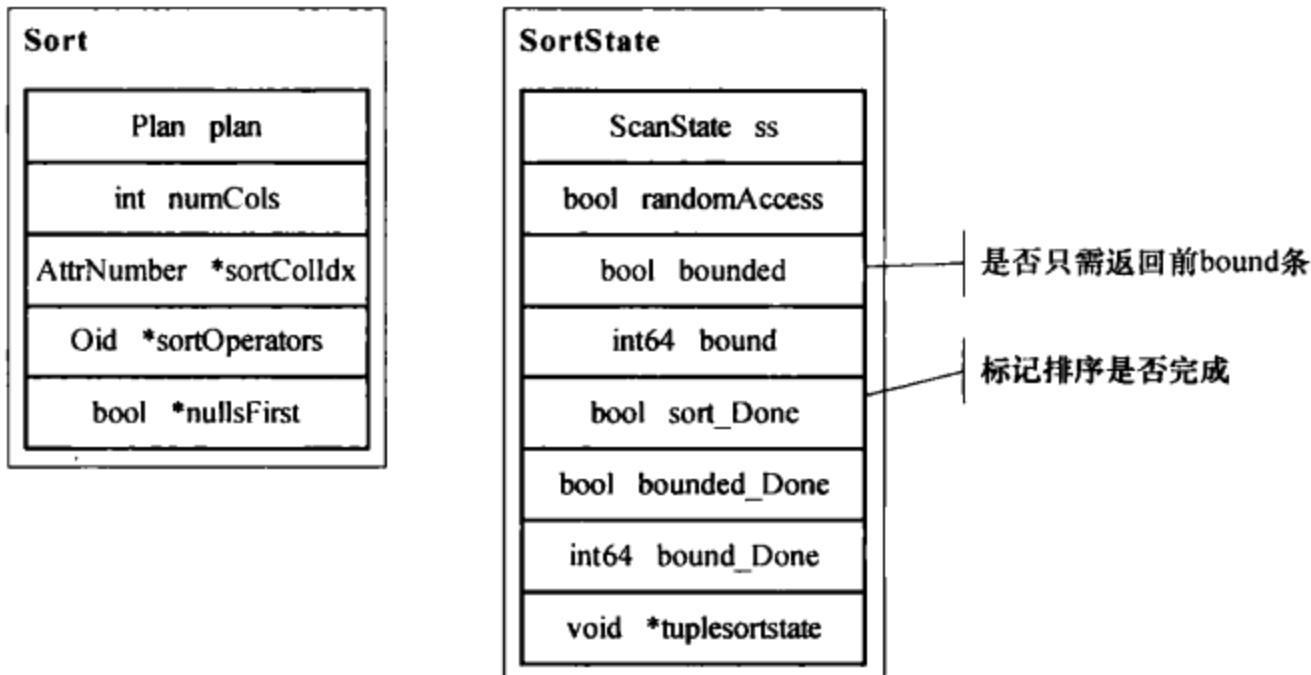


图 6-44 Sort 节点相关结构

3. Group 节点

Group 节点用于处理 GROUP BY 子句，将下层节点满足选择条件（HAVING 子句）的元组分组后，只返回该分组的第一个元组。该节点只有一个左子节点，且子节点必须返回在分组属性上已排好序的元组。

Group 节点的结构如图 6-45 所示，它在 Plan 的基础上扩展了以下几个字段：numCols 用于记录分组属性的个数，grpColIdx 数组记录了分组属性的属性号，grpOperations 数组记录了在分组属性上进行等值判断的操作符的 OID。

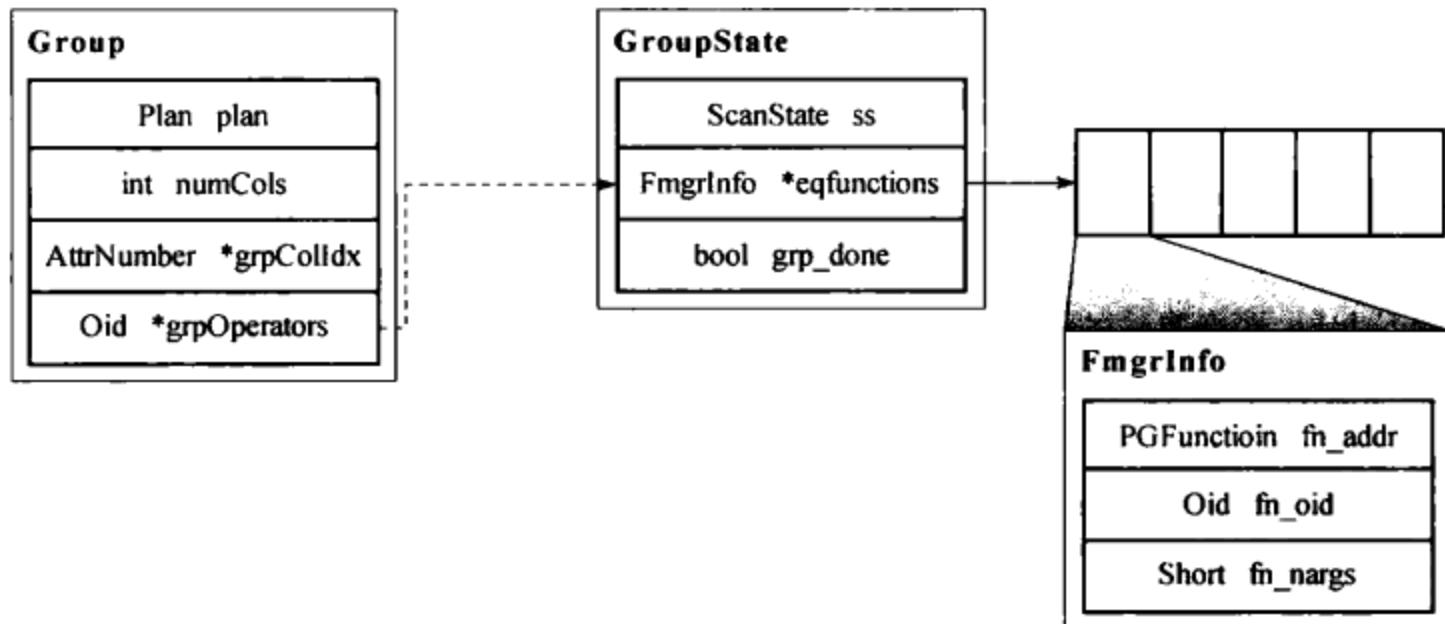


图 6-45 Group 节点的相关结构

Group 节点因为只有一个输入节点，因此初始化过程（ExecInitGroup 函数）只需调用左子节点的初始化过程。然后对每个属性上的比较操作符进行初始化，找到对应的操作符函数相关信息（用 FmgrInfo 结构保存，其中记录了函数指针 fn_addr、参数数量 fn_nargs 等信息）。

在执行过程中，由于左子节点返回的结果是按照分组属性排序的，因此只要发现连续的两个元组在分组属性上不等，即可判定前一个元组就是上一个分组的最后一个元组，而后一个元组是一个新分组的第一个元组。Group 节点的执行过程（ExecGroup 函数）如下：

- 1) 获取下层返回元组中符合 HAVING 子句条件（存储于 Plan.qual 选择条件链表中）的第一个元组，将其作为分组内的第一个元组，缓存该元组信息到状态节点的 ss_ScanTupleSlot 中，并输出该元组。
- 2) 依次获取组内的所有元组，直到获取到一个在分组属性上与当前分组不等的元组，表明当前分组结束；若下层节点返回空元组，表示分组操作完成，则会设置状态节点的 grp_done 字段为 true 并结束执行。
- 3) 扫描下一个满足 HAVING 条件的元组，缓存元组作为新分组的开始，并返回该元组。
- 4) 重复执行步骤 2、步骤 3。

由于执行过程中使用了状态节点中的 ss_ScanTupleSlot，在 Group 节点的清理过程（ExecEndGroup 函数）中需要调用 ExecClearTuple 进行元组的清理工作（参见 6.5.1 节的介绍）。

4. Agg 节点

Agg 节点用于执行含有聚集函数的 GROUP BY 操作，该节点能够实现三种执行策略：Plain（不分组的聚集计算）、Sorted（下层节点提供排好序的元组，类似 Group 的分组方法，然后进行聚集计算）、Hash（首先对下层节点提供的未排序元组进行分组，然后进行计算）。

Agg 的执行的一般方法是：首先初始化聚集计算的初始值，将其记录在中间结果中。然后，针对每一条输入元组使用迭代聚集函数进行迭代计算，得到新的中间结果。最后，如果有必要的话，使用结束函数进行处理。其伪代码如下：

```
transvalue = initcond
foreach input_tuple do
    transvalue = transfunc(transvalue, input_value(s))
result = finalfunc(transvalue)
```

其中，initcond 表示初始值，transvalue 为中间结果，input_tuple 是输入的元组，transfunc 是迭代聚集函数，finalfunc 是结束函数，result 是最终的结果。

Agg 节点的定义如图 6-46 所示，除了定义执行的策略类型（aggstrategy），还定义了聚集属性的数量（numCols）、聚集属性的属性号数组（grpColIdx）、分组函数数组（grpOperators）以及估计的分组数量（numGroups）。

在 Agg 节点的初始化过程（ExecInitAgg 函数）中，首先对目标属性 targetlist 和查询条件 qual 进行初始化，并找到其中的 Aggref 节点（用于表示聚集函数的表达式节点）。然后根据每个 Aggref 节点中存储的聚集函数信息进行初始化，为其构造一个 AggStatePerAgg 结构，其中存储了聚集函数信息运算相关的信息和内存上下文。如果有多个相同的 Aggref 节点，只会构造一个 AggStatePerAgg 结构（例如“SELECT sum(x)... HAVING sum (x) > 0”中的求和运算）。最后根据策略类型对相应状态信息进行初始化：

- 1) Plain（aggstrategy 值为 AGG_PLAIN）和 Sorted（aggstrategy 值为 AGG_SORTED）策略都不需复杂的分组过程，因此每次执行只需保存一个分组的中间结果，此时可使用 pergroup 结构数组保存一个分组的所有聚集函数的中间结果。由于 Sorted 方法与 Group 的分组方式一致，因此还需要 eq-

function 用于判断分组的分界线，并用 grp_firstTuple 字段缓存分组的第一个元组。

2) 对于 Hash (aggstrategy 值为 AGG_HASHED) 策略, Hash 表中不缓存真正的元组, 而只是使用元组计算得到的 Hash 值作为索引。Hash 表 (hashtable) 中只存储该分组的聚集函数中间值的数组 (由 AggStatePerGroup 结构存储), 而 hashslot 用于缓存 Hash 的元组中需要进行 Hash 的属性值, hash_needed 则对应于需要进行 Hash 的属性序号的链表。在已经获取了所有下层节点元组后, 从 Hash 表中依次获取分组的中间结果时会用到 hashiter。构造 Hash 函数信息从 Agg 节点的 grpOperators 获取, 信息存放于 hashfunctions 数组中。初始化时会对涉及的这些属性进行初始化。

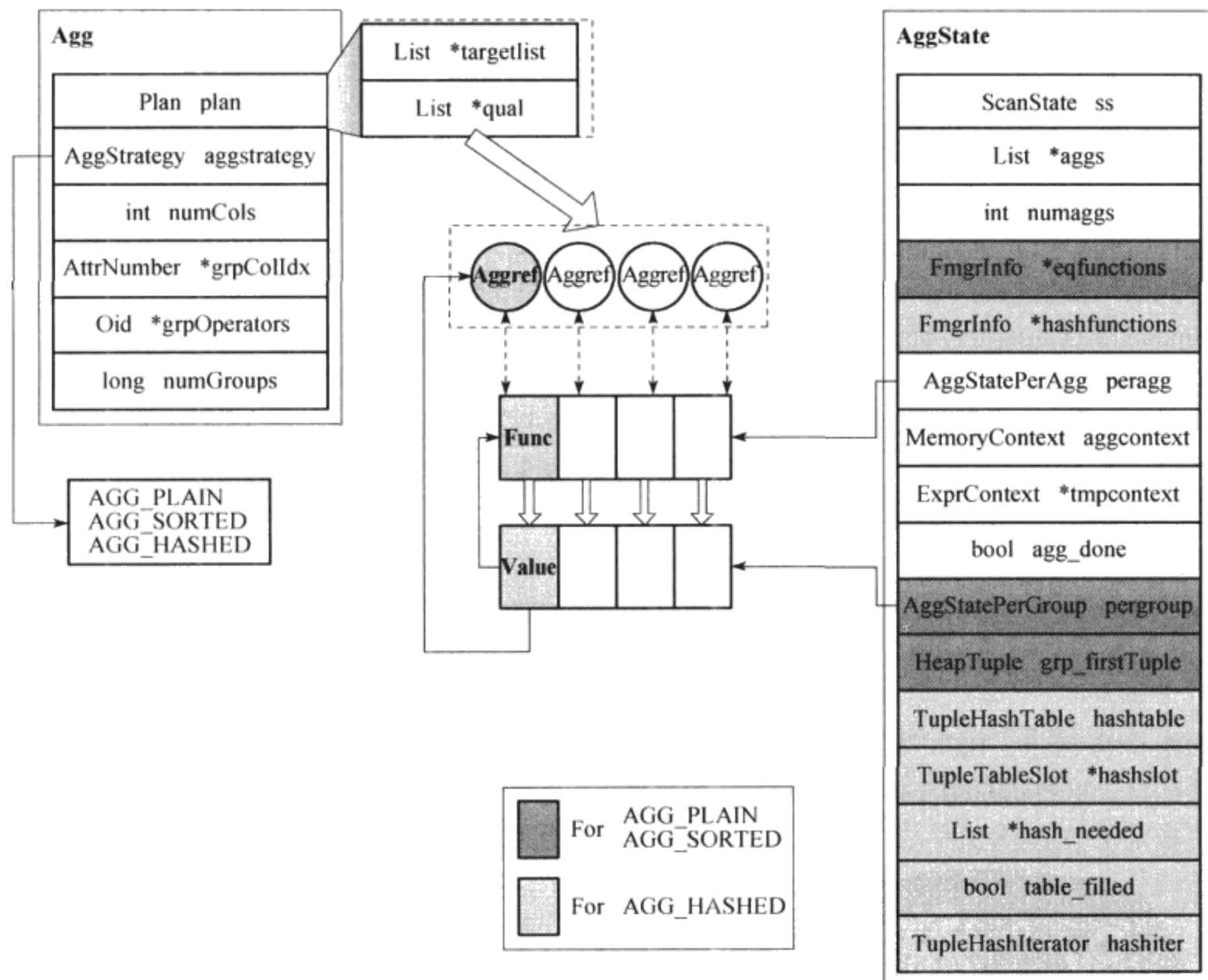


图 6-46 Agg 节点相关数据结构

图 6-46 展示了 AGG_PLAIN/AGG_SORTED 策略的执行状态。其中, peragg 用于存储聚集函数 (即 transfunc), pergroup 存储了当前所处理分组的状态, 获取的元组将与每一个 pergroup 值一并被 peragg 中的函数处理, 得到新的中间结果, 重新存储于相应的 pergroup 中。AGG_HASHED 过程则需要首先计算输入元组的 Hash 值, 然后从 Hash 表中获取对应分组的中间结果记录数组, 类似于 pergroup, 然后进行和以上两个策略相同的计算过程。

Agg 节点清理过程需要对使用到的内存上下文进行回收, 清理分配的 TupleTableSlot 结构, 并调用下层节点的清理过程。

5. Unique 节点

Unique 节点用于对下层节点返回的已排序元组进行去重操作。由于下层节点获取到的元组已经排序，因此在 Unique 节点的执行过程中只需要缓存上一个返回的元组，判断当前获取的元组是否和上一个元组在指定属性上重复。如果重复，则忽略当前元组并继续从下层节点获取元组；如果不重复，则输出当前元组并用它替换缓存中的元组。Unique 节点一般用于处理查询中的 DISTINCT 关键字，但这不是唯一的处理方式。如果要求去重的属性被“ORDER BY”子句引用时，一般会使用 Unique 节点进行处理（例如，“SELECT DISTINCT x… ORDER BY x” 中的属性 x）。

Unique 节点的定义如图 6-47 所示，其中 numCols 表示用于去重的属性数量，uniqColIdx 数组和 uniqOperators 数组分别存储了去重属性的属性号和对应的判断操作符。

在初始化过程中，会根据 uniqOperators 中的操作符列表初始化 UniqueState 结构中用于保存判断函数信息的 eqfunctions 字段。

Unique 节点的执行过程会首先从下层节点获取一个元组，如果是第一个元组，则直接输出。否则将使用 eqfunctions 中的函数判断当前元组与上次返回的元组在去重属性上是否相等，如果当前元组与上次返回的元组在去重属性上不相等则输出当前元组，如果相等则再从下层节点获取下一个元组进行同样的判断。

6. Hash 节点

Hash 节点作为 HashJoin 节点的辅助节点，共同完成 Hash 连接方法。Hash 节点利用将在 HashJoin 节点中介绍的 Hash 方法，将从左子节点获取的元组放入构造好的 Hash 表中。Hash 节点只有一个左子节点。

如图 6-48 所示，Hash 节点定义了 skew 方法（主要用于 HashJoin 节点，见 6.4.4 节）需要使用的相关信息。这些信息用于构造 Hash 表时，使用外连接表的元组统计信息来优化 Hash 表的组织结构，将最常用的 Hash 值单独存放，保持在内存中，以优化连接时 Hash 匹配的过程。这些信息包括了属性偏移量和属性的类型等。

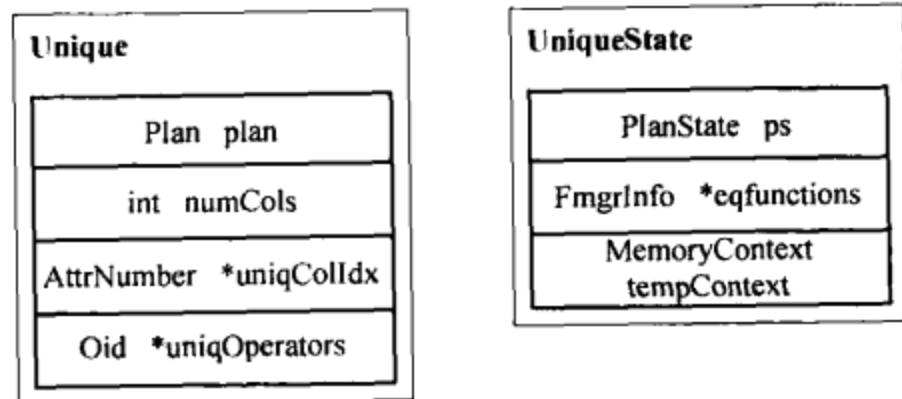


图 6-47 Unique 节点相关数据结构

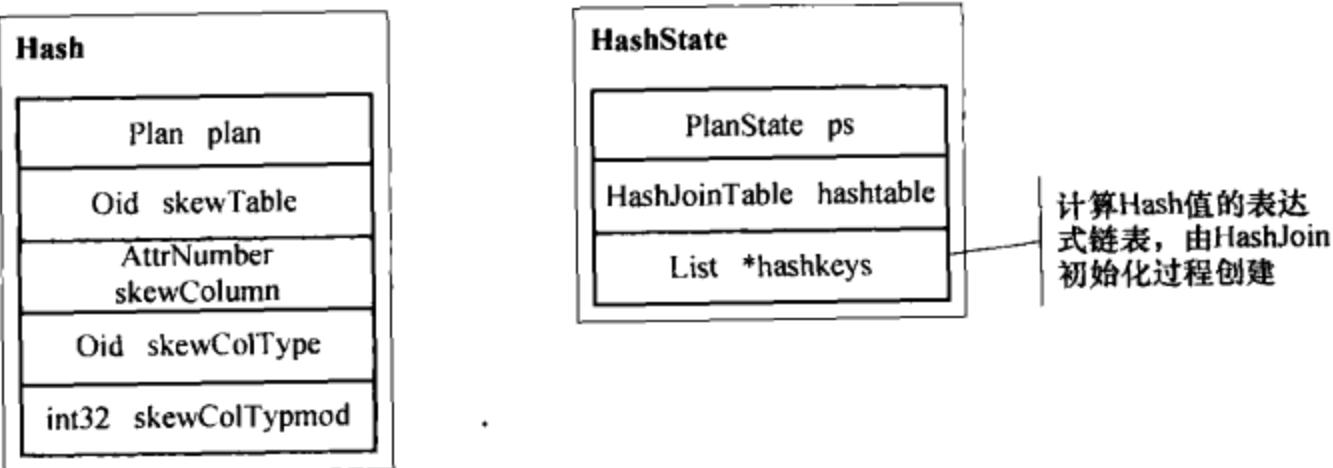


图 6-48 Hash 节点相关数据结构

Hash 节点的初始化过程仅构造 HashState 节点，并未构建 Hash 表和 hashkeys（将由 HashJoin 节点进行初始化）。

Hash 节点在执行时会从下层节点获取元组，并使用 hashkeys 中的表达式计算元组的 Hash 值，通过 Hash 值中的块号判断是否需要放入缓存文件中（Hash 表一次只能将一个块的内容放入内存中），然后将元组保存在 Hash 表相应的块中。Hash 节点的清理过程会调用下层左子节点的清理过程。

7. SetOp 节点

SetOp 节点用于处理集合操作，对应于 SQL 语句中的 EXCEPT、INTERSECT 两种集合操作，至于另一种集合操作 UNION，可直接由 Append 节点来实现。

一个 SetOp 节点只能处理一个集合操作（由两个集合参与），如果有多个集合操作则需要组合多个 SetOp 节点来实现。SetOp 节点仅有一个左子节点作为输入，其左子节点是一个 Append 节点或者是一个 Sort 节点（Sort 节点的子节点是一个 Append 节点），其低层的 Append 节点中只放置两个子计划用于表示参与集合操作的左集合和右集合。由于进行集合操作时需要区分元组来自哪一个子查询，因此在低层 Append 节点的投影操作时会为每个子计划输出元组增加一个子计划标记属性。

SetOp 有两种执行策略：排序（SETOP_SORTED）和 Hash（SETOP_HASHED）。排序策略下首先利用 Sort 节点对 Append 节点返回的元组集合进行排序，然后进行交、差的集合操作。Hash 策略下则通过 SetOp 节点提供 Hash 表对子查询元组进行 Hash 分组，然后进行集合操作。在 Hash 策略下，SetOp 节点的左子节点就是 Append 节点。在 SetOp 执行过程中需要附加子计划标记属性，这个属性在输入元组中的偏移位置被记录在 SetOp 节点的 flagColIdx 字段中。SetOp 节点的定义如图 6-49

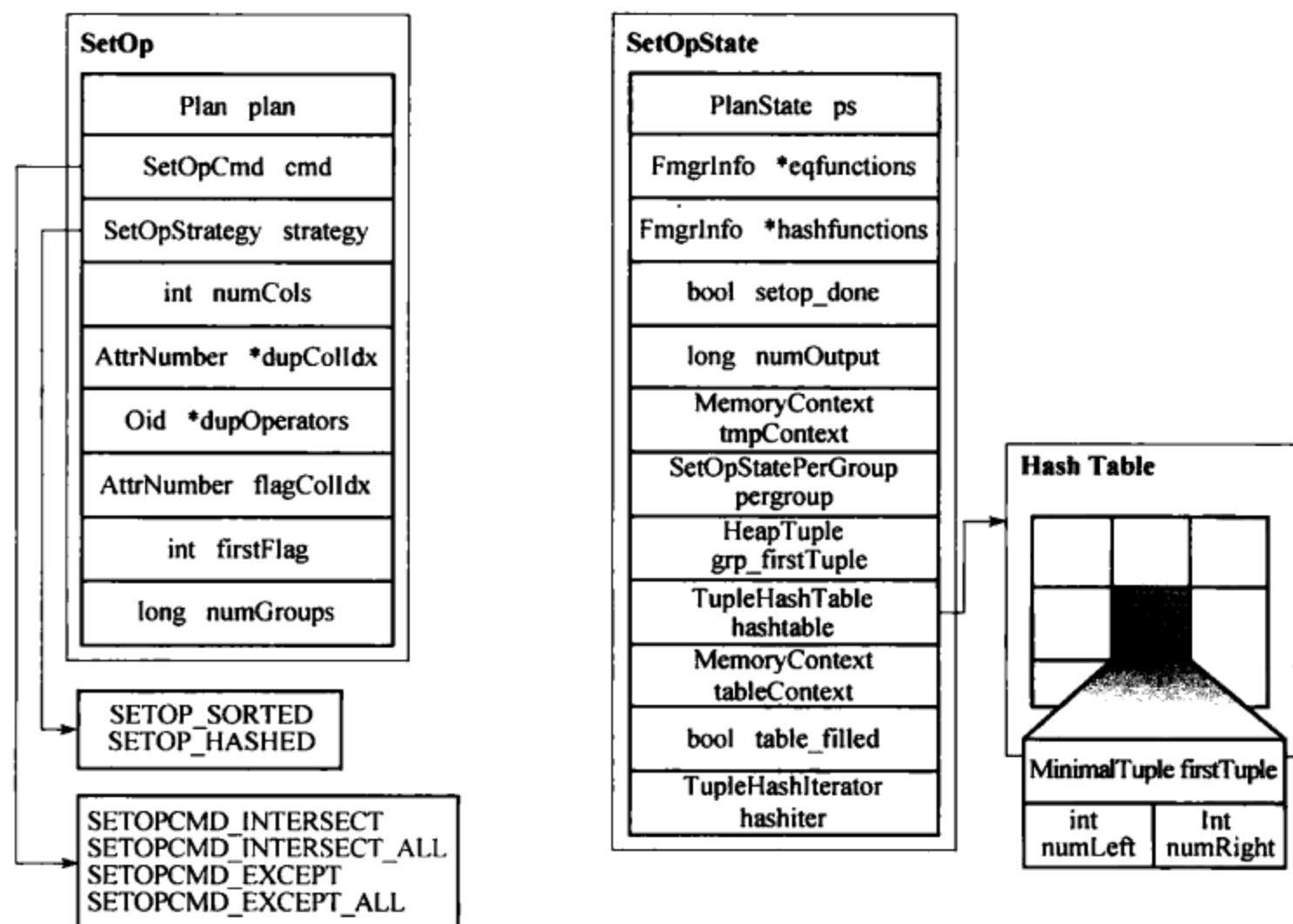


图 6-49 SetOp 相关数据结构

所示，其中除了标志属性的位置外，还有用于去重的属性号数组（dupColIdx）和相应的操作符 OID 数组（dupOperators）等信息，这两个数组用于判断元组是否相同。

SetOp 节点在执行时会对从左子节点中取到的每一个元组使用 numLeft 和 numRight 变量来记录该元组在左集合和右集合中出现的次数，然后根据集合操作命令的类型来确定该元组是否应该作为结果元组返回：

- 1) SETOPCMD_INTERSECT：对应 INTERSECT 操作，最后的结果集中不会有重复元组。其实现方式是，如果一个元组的 numLeft 和 numRight 都大于 0，则只输出一次该元组。

- 2) SETOPCMD_INTERSECT_ALL：对应于 INTERSECT ALL 操作，最后的结果集中允许有重复元组。其实现方式是，输出该元组的次数以 numLeft 和 numRight 的较小者为准。

- 3) SETOPCMD_EXCEPT：对应于 EXCEPT 操作，表示要求返回只出现在左集合中的元组，且结果集中不存在重复元组。如果一个元组的 numLeft 大于 0 而 numRight 等于 0，则输出一次该元组。

- 4) SETOPCMD_EXCEPT_ALL：对应于 EXCEPT ALL 操作，结果集中允许存在重复元组。如果一个元组的 numLeft 不小于 numRight，则输出 numLeft - numRight 次该元组。

SETOP_SORTED 策略的下层节点已经将需要执行的集合的操作合并在一起并进行了排序，SetOp 节点需要做的事情类似于 Group 节点：首先从下层节点获取一个分组的元组，仅缓存第一条元组，然后统计它在左右集合中的出现次数，然后通过集合操作的命令类型来决定如何返回元组。

SETOP_HASHED 策略的下层节点会依次返回两个子计划（对应于左右两个集合）中的元组。首先计算左集合中每个元组的 Hash 值，如果 Hash 表中没有，则插入 Hash 表中，否则，将其中 Hash 项内 numLeft 计数加 1。然后扫描右集合，计算 Hash 值后在 Hash 表中查找，找到则对 numRight 计数加 1，否则，不进行任何操作。扫描完成后，会依次从 Hash 表中获取元组和其对应的 numLeft 和 numRight，根据集合操作的命令类型来决定如何返回元组。

8. Limit 节点

Limit 节点主要用来处理 LIMIT/OFFSET 子句，它从下层节点的输出中挑选处于一定范围内的元组。该节点只有一个左子节点。

Limit 节点定义如图 6-50 所示，它在 Plan 的基础上扩展定义了 limitOffset 和 limitCount 两个表达式，用于计算偏移量和需要返回元组的数量。

在初始化过程中，会对 Limit 节点的 limitOffset 和 limitCount 两个表达式进行初始化，结果保存于 LimitState 节点的 limitOffset 和 limitCount 字段中。

在执行过程中，需要首先计算 LimitState 中的 limitOffset 和 limitCount 表达式，将结果保存于 offset 和 count 中，然后开始从下层节点获取元组，通过 position 记录已获取的元组数目，跳过前 offset 个元组，从 offset + 1 个元组开始返回，并在 offset + count 处直接

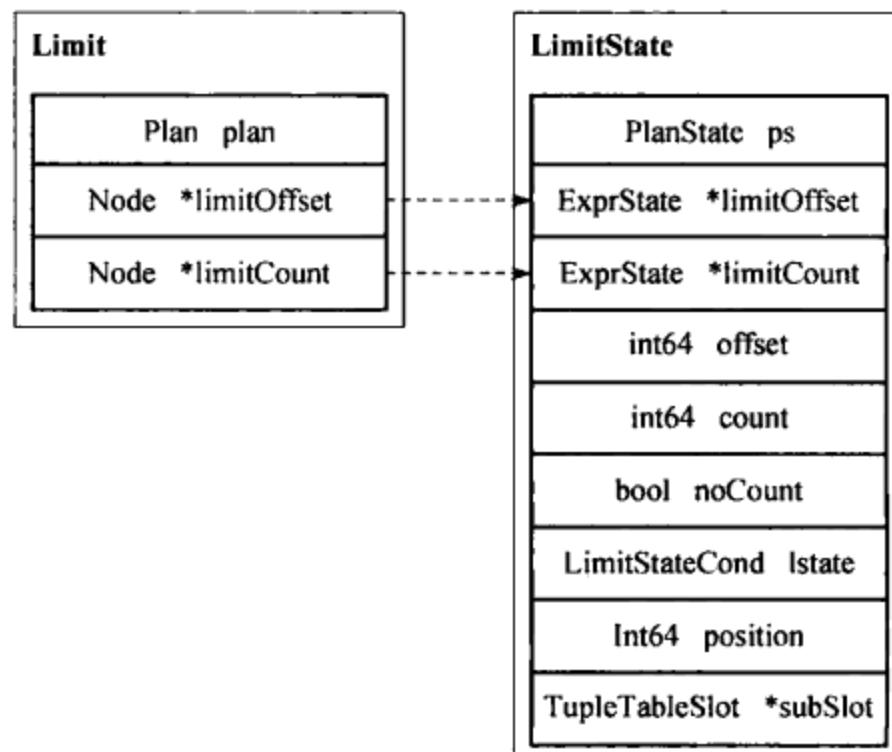


图 6-50 Limit 节点相关结构

返回空元组。

9. WindowAgg 节点

WindowAgg 节点用于处理窗口函数，窗口函数用于在与当前元组相关的一组元组上执行相关函数计算（包括在 GROUP BY 中使用的聚集函数）。窗口函数的使用类似于“SELECT…, avg (x) OVER (PARTITION BY y) FROM x;”的方式。

通常，GROUP BY 子句存在时，查询无法投影非分组属性，非分组属性只能出现在聚集函数中。在实际应用中可能需要投影非分组属性，例如，表 class 记录了班级 (cno)、民族 (nation)、人数 (num)，若需要统计班级中的各民族所占比例，使用 GROUP BY 就无法实现，如果使用窗口函数，可以使用如下语句：

```
SELECT cno, nation, num/sum(num) OVER (PARTITION BY cno) AS present FROM class;
```

上面的 SQL 语句中，将 sum 函数运用在与当前元组具有相同 cno 的一组元组上，并利用此聚合得到的总人数，计算出当前人数在总人数中的比例。这样能够实现更加灵活的函数操作，提供了更加丰富的查询方法。

窗口函数与在 GROUP 分组中进行聚集计算的不同主要体现在 SQL 语义上，在聚集函数的计算上两者是相似的。窗口函数需要处理与当前元组相关联的一组元组，在实现中要保留同一划分内的所有元组，计算得到聚集函数值后，再进行相关的投影操作。因此，在实现方法上，不论是 PARTITION 还是 GROUP 都需要先将元组在分组属性上进行划分，在分组过程和聚集函数计算的实现上是相似的。

WindowAgg 节点实现的功能类似于 Agg 节点，但不同点在于窗口函数不会造成同一分组中的元组被合并为一个，因此每个元组都可以生成一个结果元组，并可包含相关的聚集运算结果。WindowAgg 节点只有一个左子节点，用于提供已经在分组属性 (PARTITION BY 所指定的属性) 上排序的元组。

如图 6-51 所示，WindowAgg 节点定义中包含用于存储分区函数和所对应属性号的数组 partOperators 和 partColIdx、数组长度 partNumCols 以及在排序属性上判断是否相等的函数和相应的属性号数组 ordOperators 和 ordColIdx，数组长度用 ordNumCols 记录。

窗口函数也使用了 Agg 节点的聚集函数迭代计算方法，其实现方式是在 Agg 节

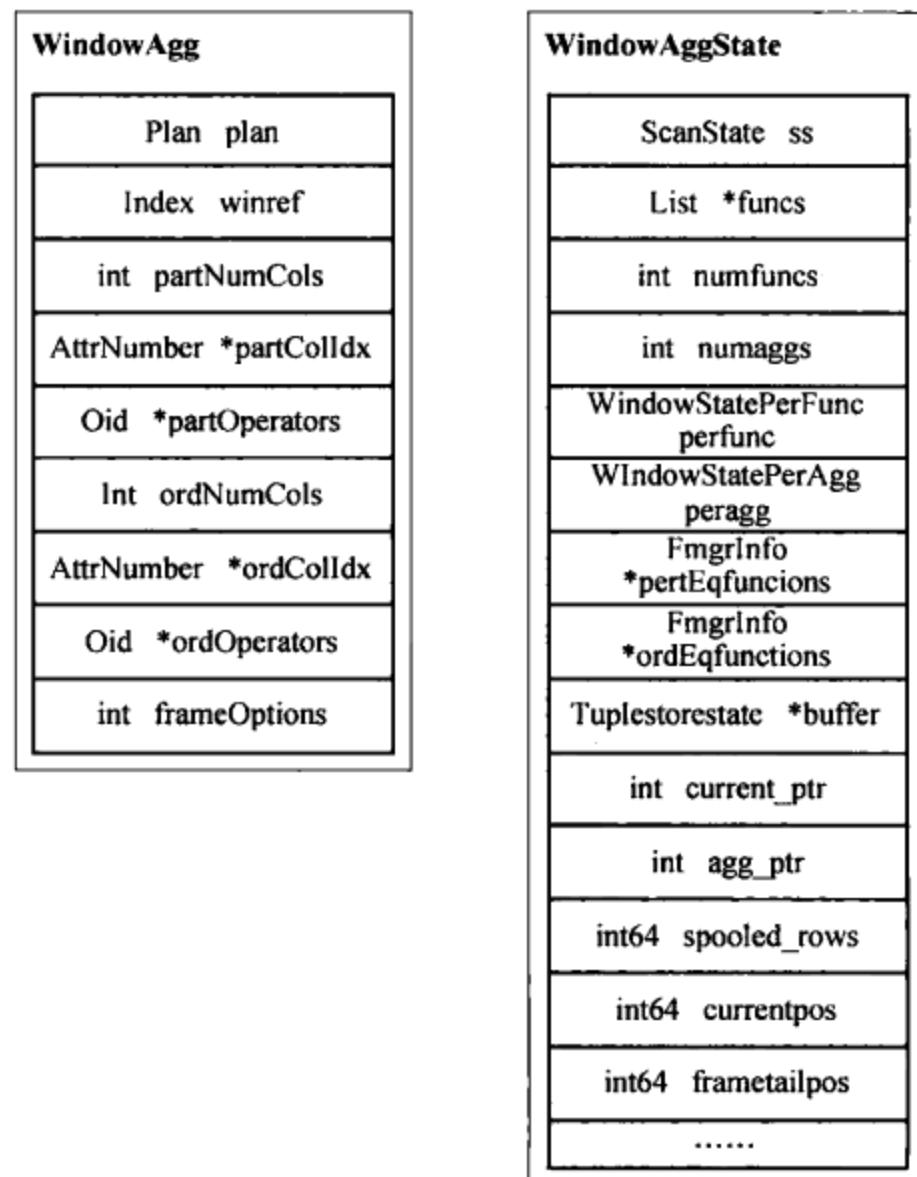


图 6-51 WindowAgg 节点相关数据结构

点的 AGG_SORTED 策略的执行中增加了缓存分组内所有元组的 tuplestore 结构，保存于状态节点中的 buffer 字段中。由于窗口函数不仅包含 Agg 中支持的聚集函数，还有新的窗口函数模式（能够支持随机的在分组内获取元组等），因此，在从 Plan 的 targetlist 中获取 WindowFunc，并将其信息初始化为 perfunc 指向的WindowStatePerFunc 数组后，需要将传统的聚集函数构造成新的WindowStatePerAgg 链表，以便在调用聚集函数处理过程中使用与 Agg 节点相同的方式。

初始化过程中，通过 WindowAgg 节点的 targetlist 构造 funcs 链表，其中包含窗口函数的表达式树。然后对 funcs 指向的表达式树进行初始化，构造函数相关调用信息存放于 perfunc 指向的数组中。对于传统的聚集函数，将在 peragg 中另外多初始化一个与 Agg 节点处理聚集函数时相同的数据结构，用于执行聚集计算。然后要对于分区判断函数和排序属性是否相等的操作函数进行初始化，分别保存于 perEqfunctions 和 ordEqfunctions 中。

在执行过程中，首先初始化一个 tuplestore 结构用于缓存元组，使用 perEqfunctions 判断是否在分区内：

- 1) 如果有传统的聚集类函数，获取分区内的所有元组缓存于 tuplestore 中。然后计算其聚集函数值，并将该值保存在对应函数信息WindowStatePerAgg 结构的 resultValue 中。依次获取分区每条元组，计算窗口函数，使用缓存的聚集函数结果进行投影。
- 2) 否则，直接扫描分区每条元组，计算该窗口函数，并对结果元组进行投影后返回。

6.4.4 连接节点

连接类型节点对应于关系代数中的连接操作，PostgreSQL 中定义了如下几种连接类型（以 T1 JOIN T2 为例）：

- 1) Inner Join：内连接，将 T1 的所有元组与 T2 中所有满足连接条件的元组进行连接操作。
- 2) Left Outer Join：左连接，在内连接的基础上，对于那些找不到可连接 T2 元组的 T1 元组，用一个空值元组与之连接。
- 3) Right Outer Join：右连接，在内连接的基础上，对于那些找不到可连接 T1 元组的 T2 元组，用一个空值元组与之连接。
- 4) Full Outer Join：全外连接，在内连接的基础上，对于那些找不到可连接 T2 元组的 T1 元组，以及那些找不到可连接 T1 元组的 T2 元组，都要用一个空值元组与之连接。
- 5) Semi Join：类似 IN 操作，当 T1 的一个元组在 T2 中能够找到一个满足连接条件的元组时，返回该 T1 元组，但并不与匹配的 T2 元组连接。
- 6) Anti Join：类型 NOT IN 操作，当 T1 的一个元组在 T2 中未找到满足连接条件的元组时，返回该 T1 元组与空元组的连接。

PostgreSQL 实现了三种连接操作，嵌套循环连接（Nest Loop）、归并连接（Merge Join）和 Hash 连接（Hash Join）。归并连接算法可以实现上述六种连接，而嵌套循环连接和 Hash 连接只能实现 Inner Join、Left Outer Join、Semi Join 和 Anti Join 四种连接。表 6-5 列出了 PostgreSQL 8.4.1 提供的连接节点。

如图 6-52 所示，连接节点有公共父类 Join，Join 继承了 Plan 的所有属性，并扩展定义了 join-

表 6-5 连接节点描述

节点类型	文件	描述
MergeJoin	nodeMergejoin.c	归并连接算法
NestLoop	nodeNestloop.c	循环连接算法
HashJoin	nodeHashjoin.c	Hybird Hash 连接算法

type 用以存储连接的类型，joinqual 用于存储连接的条件。对应的执行状态节点 JoinState 中定义了 jointype 存储连接类型，joinqual 存储连接条件初始化后的状态链表。

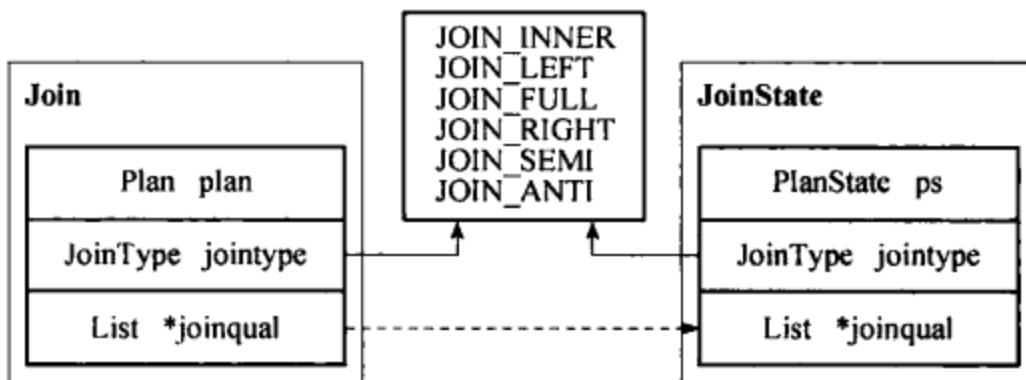


图 6-52 连接节点的公共数据结构

1. NestLoop 节点

NestLoop 节点实现了嵌套循环连接方法，能够进行 Inner Join、Left Outer Join、Semi Join 和 Anti Join 四种连接方式。图 6-53 展示了 NestLoop 节点及对应的状态节点 NestLoopState 的定义。NestLoop 节点并未对 Join 节点进行扩展，NestLoop 节点的初始化过程会将节点中连接条件（joinqual 字段）进行处理，转化为对应的状态节点 JoinState 中的 joinqual 链表。

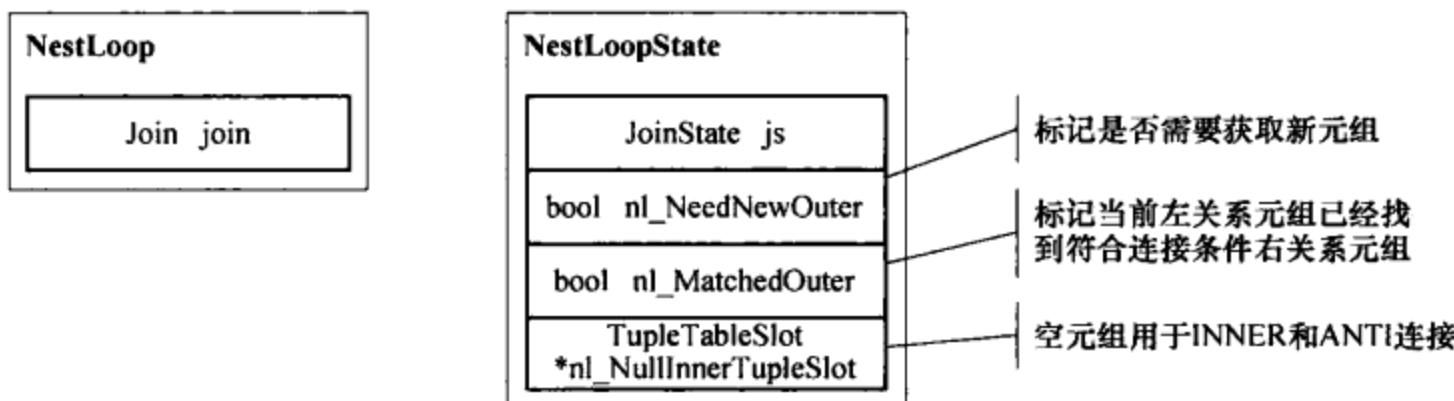


图 6-53 NestLoop 节点相关数据结构

循环嵌套连接的基本思想如下：

```

FOR each tuple s in S DO
    FOR each tuple r in R DO
        IF r and s join to make a tuple t THEN
            output t;
    
```

为了迭代实现此方法，NestLoopState 中定义了字段 nl_NeedNewOuter 和 nl_MatchedOuter。当元组处于内层循环时，nl_NeedNewOuter 为 false，内层循环结束时 nl_NeedNewOuter 设置为 true。为了能够处理 Left Outer Join 和 Anti Join，需要知道内层循环是否找到了满足连接条件的内层元组，此信息由 nl_MatchedOuter 记录，当内层循环找到符合条件的元组时将其标记为 true。

NestLoop 节点的初始化工作需要为 Left Outer Join 和 Anti Join 两种可能用到空元组的情况构造空元组，并存放在 NestLoopState 的 nl_NullInnerTupleSlot 中，还将进行如下两个操作：

- 1) 将 nl_NeedNewOuter 标记为 true，表示需要获取左子节点元组。
- 2) 将 nl_MatchedOuter 标记为 false，表示没有找到与当前左子节点元组匹配的右子节点元组。

NestLoop 执行过程需要循环执行如下操作：

1) 如果 nl_NeedNewOuter 为 true，则从左子节点获取元组，若获取的元组为 NULL 则返回空元组并结束执行过程。如果 nl_NeedNewOuter 为 false，则继续进行步骤 2。

2) 从右子节点获取元组，若为 NULL 表明内层扫描完成，设置 nl_NeedNewOuter 为 true，跳过步骤 3 继续循环。

3) 判断右子节点元组是否与当前左子节点元组符合连接条件，若符合则返回连接结果。

以上过程能够完成 Inner Join 的递归执行过程。但是为了支持其他几种连接则还需要如下两个特殊的处理：

1) 当找到符合连接条件的元组后将 nl_MatchedOuter 标记为 true。内层扫描完毕时，通过判断 nl_MatchedOuter 即可知道是否已经找到满足连接条件的元组，在处理 Left Outer Join 和 Anti Join 时需要进行与空元组的连接，然后将 nl_MatchedOuter 设置为 false。

2) 当找到满足匹配条件的元组后，对于 Semi Join 和 Anti Join 方法需要设置 nl_NeedNewOuter 为 true。区别在于 Anti Join 需要不满足连接条件才能返回，所以要跳过返回连接结果继续执行循环。

NestLoop 节点的清理过程没有特殊处理，只需递归调用左右子节点的清理过程。

2. MergeJoin 节点

MergeJoin 实现了对排序关系的归并连接算法，归并连接的输入都是已经排好序的。PostgreSQL 中 MergeJoin 算法实现的伪代码如下：

```

Join{
    get initial outer and inner tuples           初始化
    do forever{
        while(outer != inner){                  SKIP_TEST
            if(outer < inner)
                advance outer                 SKIPOUTER_ADVANCE
            else
                advance inner                 SKIPINNER_ADVANCE
        }
        mark inner position               SKIP_TEST
        do forever{
            while(outer == inner){
                join tuples                JOINTUPLES
                advance inner position     NEXTINNER
            }
            advance outer position       NEXTOUTER
            if (outer == mark)          TESTOUTER
            restore inner position to mark TESTOUTER
            else
                break      //return to top of outer loop
        }
    }
}

```

算法首先初始化左右子节点，然后执行以下操作（其中对于大小的比较都是指对连接属性值的比较）：

- 1) 扫描到第一个匹配的位置，如果左子节点（outer）较大，从右子节点（inner）中获取元组；如果右子节点较大，从左子节点中获取元组。
- 2) 标记右子节点当前的位置。
- 3) 循环执行左子节点 == 右子节点判断，若符合则连接元组，并获取下一条右子节点元组，否则退出循环执行步骤4。
- 4) 获取下一条左子节点元组。
- 5) 如果左子节点 == 标记处的右子节点（说明该条左子节点与上一条相等），需要将右子节点扫描位置回退到扫描位置，并返回步骤3；否则跳转到步骤1。

为了说明归并排序的连接算法，我们以 Inner Join 为例给出部分执行过程，如图 6-54 所示，两个 current 分别指向输入的当前元组，mark 用于标记扫描的位置。

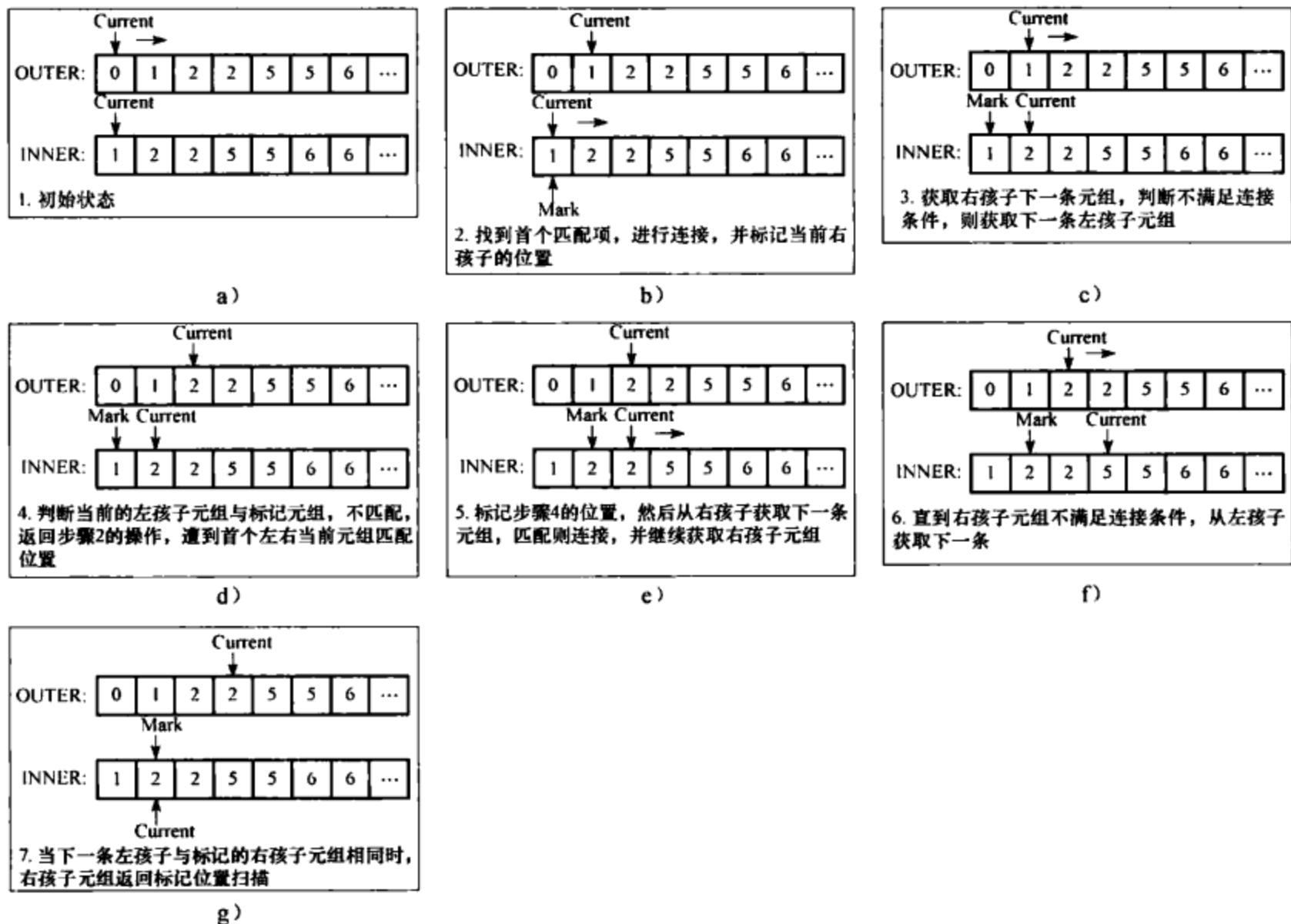


图 6-54 归并连接示例

- 1) 首先找到左右序列第一个匹配位置，图 6-54a 中 $\text{current}(\text{outer}) = 0$ 小于 $\text{current}(\text{inner})$ ，因此 outer 的 current 向后移动。
- 2) 如图 6-54b 所示，当找到匹配项后，则进行连接，使用 mark 标记当前 inner 的扫描位置，

并将 inner 的 current 向后移动。

3) 接着判断 $\text{current}(\text{outer}) = 1$ 小于 $\text{current}(\text{inner}) = 2$ ，则将 outer 的 current 向后移动，并判断 outer 是否与 mark 相同（这是为了发现 outer 的 current 与前一个相同的情况）。

4) 图 6-54d 显示 $\text{current}(\text{outer}) = 2$ 不等于 $\text{mark}(\text{inner}) = 1$ ，则继续扫描过程。

5) 判断两个 current 是否相同，发现 $\text{current}(\text{outer}) = 2$ 等于 $\text{current}(\text{inner}) = 2$ ，则进行连接，同样标记 inner 的当前位置，并将 inner 的 current 向后移动，如图 6-54e 所示。其中的 $\text{current}(\text{inner}) = 2$ 仍满足连接条件，因此连接完成后 inner 的 current 继续向后移动。

6) 如图 6-54f 所示， $\text{current}(\text{outer}) = 2$ 小于 $\text{current}(\text{inner}) = 5$ ，则将 outer 的 current 指针向后移动。

7) 此时判断 $\text{current}(\text{outer})$ 和 $\text{mark}(\text{inner})$ 相等，则将 inner 的 current 指向 mark 的位置，重新获取 inner 的元组进行匹配，如图 6-54g 所示。

8) 不断重复这样的匹配模式，直到 inner 或 outer 中的一方被扫描完毕，则表示连接完成。

MergeJoin 节点的定义如图 6-55 所示，该节点在 Join 的基础上扩展定义了 mergeclauses、mergeFamilies、mergeStrategies、mergeNullsFirst 字段。其中 mergeclauses 存储用于计算左右子节点元组是否匹配的表达式链表，mergeFamilies、mergeStrategies、mergeNullsFirst 均与表达式链表对应，表明其中每一个操作符的操作符类、执行的策略（ASC 或 DEC）以及空值排序策略。

在初始化过程中，会使用 MergeJoin 构造 MergeJoinState 结构，通过对于连接类型的判断来设置如下几个变量的值：

1) mj_FillOuter：为 true 表示不能忽略没有匹配项的左子节点元组，需要与空元组进行连接，在 LEFT JOIN、ANTI JOIN 和 FULL JOIN 时为 true。

2) mj_FillInner：为 true 表示不能忽略没有匹配项的右子节点元组，需要与空元组进行连接，在 RIGHT JOIN、FULL JOIN 时为 true。

3) mj_NullInnerTupleSlot：为右子节点元组生成的空元组，在 mj_FillOuter 为真时构造。

4) mj_NullOuterTupleSlot：为左子节点元组生成的空元组，在 mj_FillInner 为真时构造。

除此之外，需要将标记当前左（右）子节点元组是否找到能够连接的元组的变量 mj_MatchedOuter (mj_MatchedInner) 设置为 false，将存储左（右）子节点元组的字段 mj_OuterTupleSlot (mj_InnerTupleSlot) 设置为 NULL，并为 mj_MarkedTupleSlot 分配存储空间。

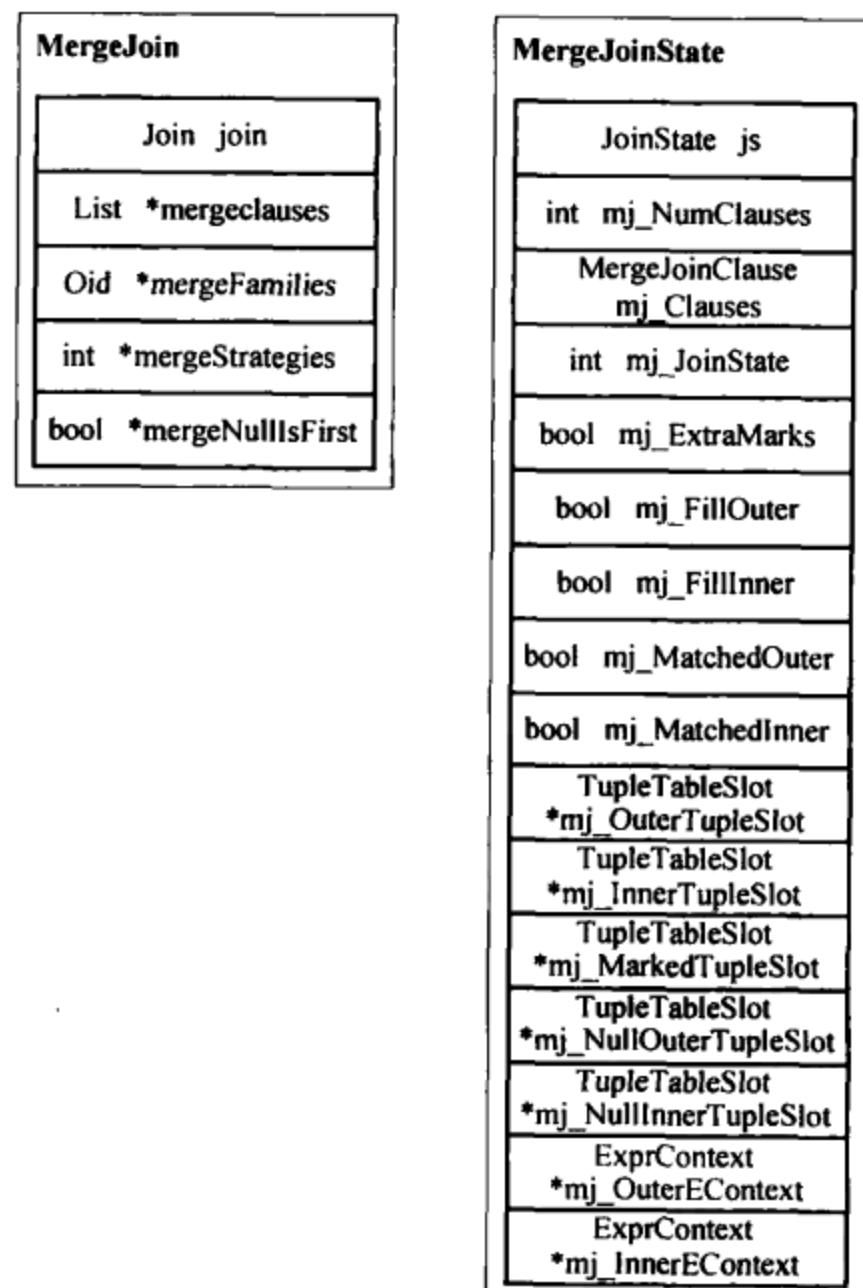


图 6-55 MergeJoin 操作相关数据结构

3. HashJoin 节点

HashJoin 节点实现了 Hash 连接算法，它能够实现 Inner Join、Left Outer Join、Semi Join 和 Anti Join 四种连接方式。

以下我们以表 R（左关系）与表 S（右关系）连接为例，说明 Hash 连接的实现过程。

1) 对一个表（例如 S）进行 Hash 时，其块和桶数量的确定和划分方法如下：

①首先对 S 分块（batch），估算存储 S 所占用的空间（inner_rel_bytes），Hash 所使用的内存空间被定义为 1 兆（hash_table_bytes），则分块的数量 $nbatch = \text{ceil}(\text{inner_rel_bytes}/\text{hash_table_bytes})^{\ominus}$ 。

②在内存中，每一个块又被划分为大小为 10 个元组的桶，因此，个数为 $nbucket = (\text{hash_table_bytes}/\text{tuplesize}) / 10$ ，其中 tuplesize 为元组大小的估计值。

③对于一个 Hash 值为 hashvalue 的元组，其所属的分块号为 $(\text{hashvalue}/nbucket) \% nbatch$ ，其对应的桶号为 $\text{hashvalue} \% nbucket$ 。

④在 PostgreSQL 实现中，为了能够使用位操作（位与和移位）实现取模和取余操作，将 nbatch 和 nbucket 取为不小于计算值的 2^n ，并使得 $2^{\log_2 nbuckets} = nbucket$ ，则块号的计算方法为 $(\text{hashvalue} >> \log_2 nbuckets) \& (nbatch - 1)$ ，桶号计算式为 $\text{hashvalue} \& (nbucket - 1)$ 。

2) 执行 HashJoin 的算法：

①顺序获取 S 中的所有元组，对每一条元组进行 Hash，并通过 Hash 结果获取块号和桶号。对于块号为 0 的元组，放入内存对应的桶内；否则放入为右关系每个块分别建立的临时文件中。此时，标记当前在内存中的块号 curbatch 为 0。

②从 R 中获取元组，进行 Hash，获取元组块号和桶号。当块号等于当前在内存中的块号时，直接扫描对应的桶，找寻满足条件的元组并进行连接；否则，将其存入为左关系每个块分别建立的临时文件中。执行过程，直到 R 被扫描完毕。

③从 S 的块号 curbatch + 1 对应的临时文件中读取所有存储的元组，将其 Hash 到相应的桶内，并将 curbatch 加 1。

④从 R 的块号 curbatch 对应的临时文件中依次读取所存储的元组，计算其桶号，并扫描桶中 S 的元组，寻找满足连接条件的元组进行连接。

⑤重复步骤 3 和 4，直到所有的块都被扫描为止。

为了实现上述过程，HashJoin 定义结构如图 6-56 所示，其中扩展定义了计算左右关系 Hash 值和 Hash 值比较的表达式。执行状态节点中存储了各种执行过程中用到的数据结构。和 NestLoop 节点类似，为了实现四种连接方法，HashJoinState 节点也定义了与 NestLoop 相同的几个属性。特别需要介绍的是，hj_CurBucketNo 用于标记当前放入内存的块号，CurHashValue 用于保存当前扫描的左子树元组计算得到的 Hash 值。此外，hj_CurBucketNo 用于标记另一个优化结构对应的桶号，针对于会生成多个块的右关系，当左关系比较大且无序时，PostgreSQL 在内存中分配了另一块内存空间，专门用于存储左关系在 Hash 属性上出现频率较高的 Hash 值所对应的右关系元组，每个桶对应一个 Hash 值，这样可以提高连接的效率。至于左关系中哪些 Hash 值的出现频率较高，可以从 pg_statistic 系统表中记录的统计信息中获取。这种方式被称为 skew 方法。

\ominus ceil 操作表示向上取整。

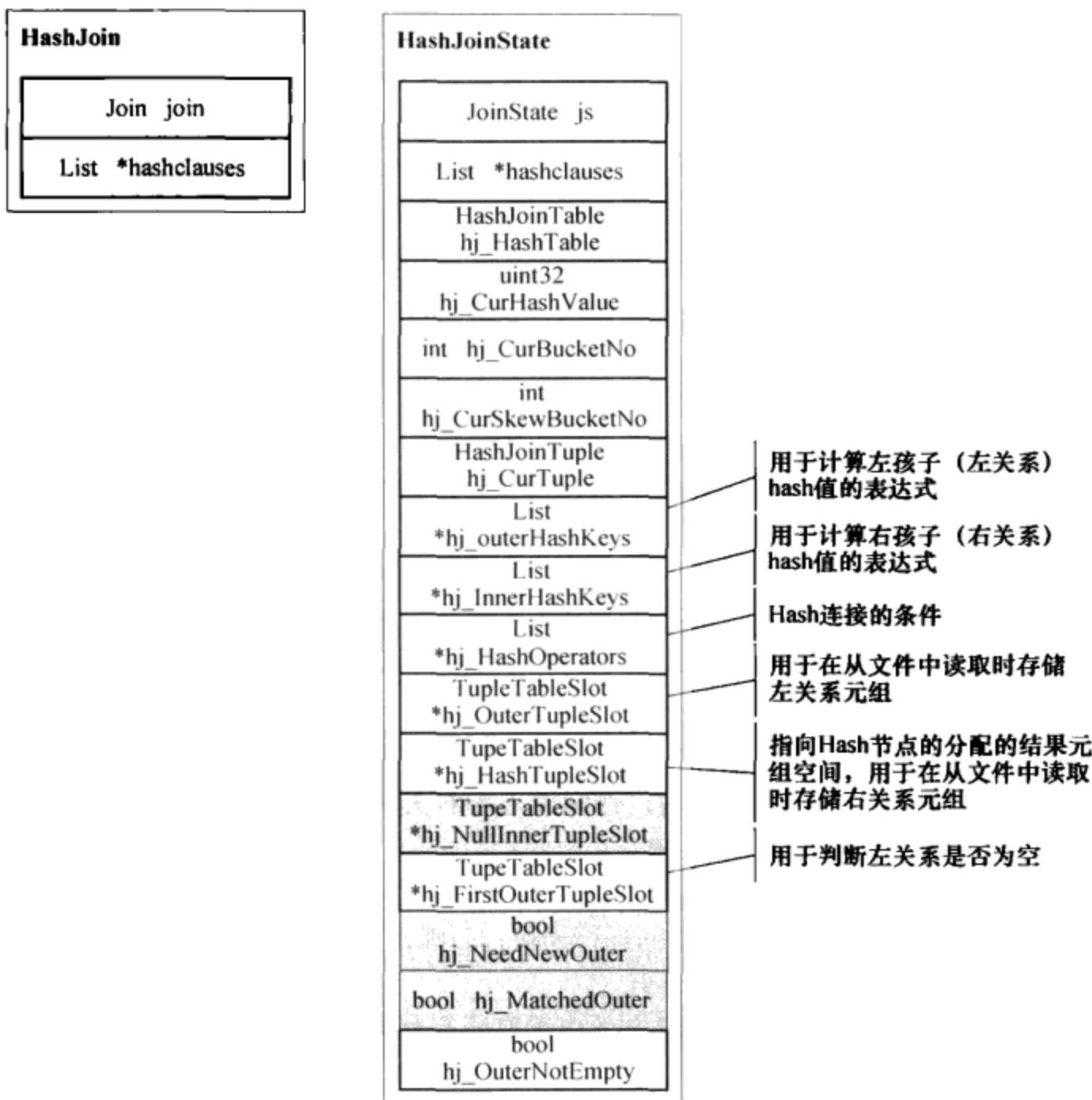


图 6-56 HashJoin 节点相关数据结构

HashJoin 是与 Hash 节点配合使用的，HashJoin 节点的右子节点一定是 Hash 节点。Hash 节点主要完成 HashJoin 算法的步骤 1 以及 Hash 表的管理。HashJoin 节点则负责处理 Hash 连接算法的其他步骤和功能。

6.5 其他子功能介绍

6.5.1 元组操作

PostgreSQL 使用元组存储所有信息，包括各种系统信息、数据等。因此，对于元组的处理就成为一个必不可少的操作。存储模块提供了很好的元组（HeapTuple）定义和操作接口，但该结构是面向物理元组的，结构解析和构造开销较大，不能满足执行器高效处理元组的需求。

执行器在执行过程中需要进行投影和属性选择判断，此时需要快速获取元组的数据。另外，

Hash、Material 等节点缓存元组时，要求元组体积更小，以节省存储空间。显然用于表示物理元组的结构 HeapTuple 已经无法满足这些需求。

为了节省元组占用的存储空间，PostgreSQL 中定义了 MinimalTuple 结构，去掉了 HeapTuple 结构中的事务相关信息。

为了能够统一地表示和处理各种形式的元组，执行器定义了数据结构 TupleTableSlot 结构和相关的处理函数，以便执行过程中按照需要使用相应形式的元组进行处理，并支持各种形式之间的转换。

执行器将物理元组存储于 TupleTableSlot 组成的数组中，该数组称为元组表 Tuple Table。在初始化过程中，统计所有计划节点的需求，为其分配足够数量的 TupleTableSlot，建立 Tuple Table，并存储于 Estate 结构的 es_tupleTable 字段中。在初始化每个节点时，节点会根据自身需求申请分配 TupleTableSlot 结构，用于存储节点的输出元组、扫描到的元组等。执行完成后在清理过程中会统一释放 Tuple Table 中的所有元组。

TupleTableSlot 的数据结构如图 6-57 所示。为了能够存储多种形式元组，TupleTableSlot 定义了 tts_tuple 用于存储 HeapTuple 形态的元组，用 tts_mintuple 存储 MinimalTuple 形态的元组，用 tts_values 指向元组的属性值数组（tts_nvalid 用于存储该数组的长度，tts_isnull 数组标记对应的属性是否

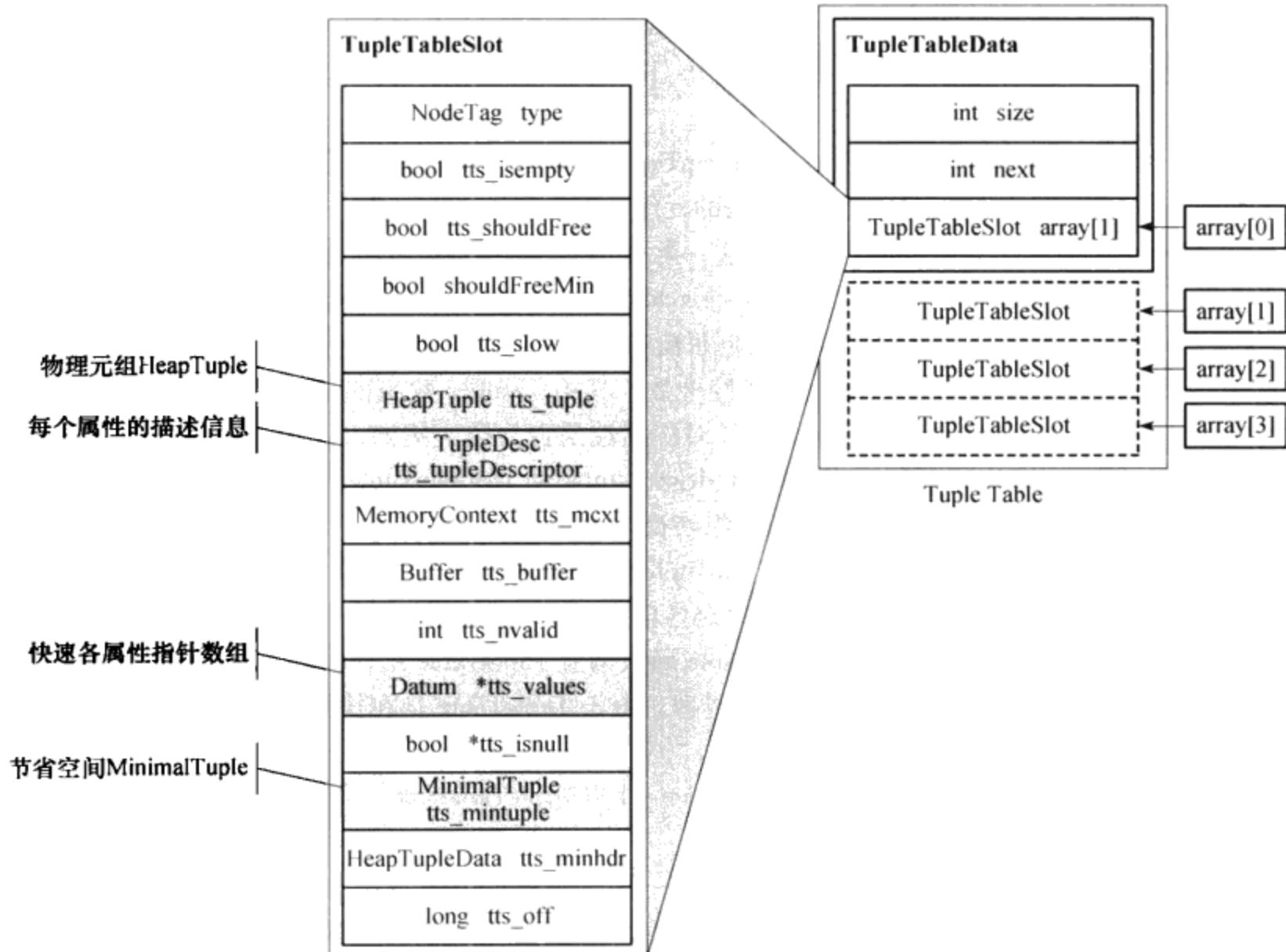


图 6-57 TupleTableSlot 相关数据结构

为空)。其中还定义了几个控制变量和资源指针：

- `tts_isempty`: 标记该 `TupleTableSlot` 结构是否为空, 当为 `true` 时表示没有存储任何的数据。
- `tts_shouldFree`: 标记是否需要释放 `HeapTuple` 指向的空间。
- `tts_shouldFreeMin`: 标记是否需要释放 `MinimalTuple` 指向的空间。
- `tts_slow` 和 `tts_off`: 执行过程中并不是一次性将所有属性从 `HeapTuple` 中解析出来, 只会解析出属性号低于给定值的所有属性, `tts_nvalid` 则记录了已经解析出的属性数。如果解析的属性中存在变长属性, `tts_slow` 将被设置为 `true`。`tts_off` 则标记当前解析到的偏移位置, 如果当前所需属性号大于已经解析的属性数 `tts_nvalid`, 则会使用 `tts_off` 继续进行属性的解析。
- `tts_mcxt`: 用于记录当前 `TupleTableSlot` 所在的内存上下文。
- `tts_buffer`: 若该元组是存储于磁盘的物理元组, 指向该元组所在的缓冲区。
- `tts_minhdr`: 与 `tts_mintuple` 配合使用, 使得 `MinimalTuple` 可以像 `HeapTuple` 一样进行各种操作。

执行器使用数据结构 `TupleTableData` 来组织 Tuple Table。`TupleTableData` 中记录了 Tuple Table 中存放的元组的个数 (`size` 字段) 和下一个空闲的 `TupleTableSlot` 结构的偏移 (`next` 字段)。`TupleTableData` 的 `array` 字段指向一个变长的 `TupleTableSlot` 数组。该数组的实际长度通过 `size` 字段来记录。

`TupleTableSlot` 实际管理了如下四种情况下的元组：

1) 在磁盘缓冲区中的物理元组 (`HeapTuple` 结构): 此时 `tts_buffer` 不为 `NULL`, `HeapTuple` 结构的元组存储于 `tts_tuple` 中, `tts_tuple` 则存放在 `tts_buffer` 指向的缓冲区中。例如从磁盘获取的元组即为此种情况。

2) 在申请的内存中构造的物理元组 (`HeapTuple` 结构): 元组将会 (但是还未被) 存储在磁盘上, 此时 `tts_buffer` 为 `NULL`, `HeapTuple` 结构的元组被存储于 `tts_tuple` 中。

3) 精简的物理元组 (`Minimal Tuple`): 这一类元组存储于特别申请的内存空间中。例如, 物化节点中需要缓存元组到文件中, 此时会调用 `ExecFetchSlotMinimalTuple` 获取 `Minimal` 状态的元组结构。

4) 虚拟元组 (`Virtual Tuple`): 虚拟元组并没有一个完整的数据结构, 它通过两个数组 `values` 和 `isnull` 来表示, 数组的长度等于元组的属性个数。其中, `values` 数组的每一个元素记录了对应属性的属性值, 而 `isnull` 数组的每个元素标记了对应属性是否为空。例如, 计划节点中将元组进行投影后的结果元组就处于此种状态。使用 `heap_deform_tuple` 可将 `HeapTuple` 形式的元组解构成虚拟元组, 同时也可使用 `heap_form_tuple` 将虚拟元组转换为 `HeapTuple` 结构。

前两种状态的元组的存储结构都使用 `HeapTuple` 结构存储物理元组, 两者的唯一不同在于是否存在于磁盘页面中。精简的物理元组是没有系统属性信息 (事务信息等) 的, 这一点和虚拟元组相同。

在 `InitPlan` 中, 会调用 `ExecCreateTupleTable` 函数构造 `TupleTable` 结构, 而 `TupleTable` 将在 `EndPlan` 函数中调用 `ExecDropTupleTable` 函数进行清理。由于 `TupleTable` 是在执行器初始化时分配, 在清理过程时回收, 因此其分配策略相对简单:

1) 构造时分配 `sizeof(TupleTableData) + (size - 1)(TupleTableSlot)` 大小的内存空间, `size` 表示 `TupleTableSlot` 数量。

2) 在初始化节点过程中, 当节点需要申请一个 `TupleTableSlot` 时, 会从 `next` 获取 `array[next]` 处的 `TupleTableSlot`, 并将 `next` 增加 1。该功能由函数 `ExecAllocTableSlot` 实现。

3) 回收 `TupleTable` 结构时会调用 `TupleTableSlot` 的清理函数 `ExecClearTuple` 处理每一个 `TupleTableSlot` 结构。

在执行过程中，如果计划节点从磁盘获取 HeapTuple 元组，需要使用函数 ExecStoreTuple 将 HeapTuple 存放于 TupleTableSlot 中，并设置相应的状态标志。执行投影操作需要将属性值和是否为空的标记写入 TupleTableSlot 的相应字段中，并使用 ExecStoreVirtualTuple 标记元组包含数据。INSERT 操作时需要先将构造的 TupleTableSlot 通过 ExecMaterializeSlot 构造成 HeapTuple 结构，然后存入物理介质。

6.5.2 表达式计算

处理 SQL 语句中的函数调用、计算式和条件表达式时需要用到表达式计算。PostgreSQL 系统中实现了表达式计算子系统，用于表示和执行 SQL 语句中的各种表达式。

如图 6-58 所示，在 PostgreSQL 系统中，表达式的表示方式与查询计划树的计划节点类似，用各种表达式计划节点（下称表达式节点）来完成相应的操作。表达式继承层次中的公共根类为 Expr 节点，其他表达式节点都继承 Expr 节点。表达式状态的公共根类为 ExprState，定义了类型（type）、辅助表达式节点指针（expr）以及用于实现该节点操作的函数指针（evalfunc）。

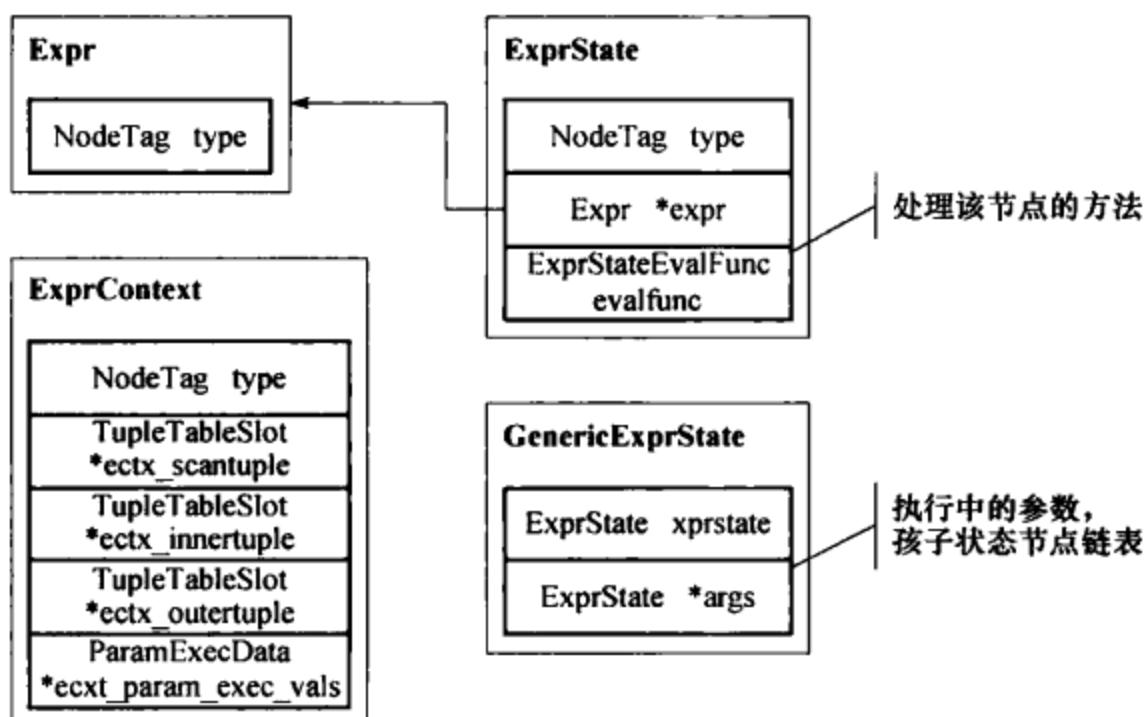


图 6-58 表达式节点相关数据结构

表达式节点与计划节点的不同在于，并不是所有的表达式节点都需要定义特定的状态类型。PostgreSQL 定义了一个 GenericExprState 基本的状态节点，没有特殊需要的表达式节点都用它来记录执行中的状态。GenericExprState 只在 ExprState 节点上扩展定义了子节点链表 args，这也是该操作的参数来源。ExprContext（表达式上下文）则充当了 Estate 的角色，表达式计算过程的参数以及表达式所使用的内存上下文等会存放在此结构中。每个计划节点都会有表达式计算，因此在 PlanState 中定义了 ps_ExprContext 用于存储该计划节点的表达式上下文。

此外，表达式节点并没有采用计划节点的方式去限制子节点的数量，而是使用类似于 Append 节点组织子计划的方式来实现不限个数的子节点方式：将子节点指针组织为链表（args）。

表达式的计算过程也分为三个部分：初始化、执行和清理。初始化过程使用统一接口 ExecInitExpr，根据表达式节点类型进行相应的处理。执行过程使用统一接口宏 ExecEvalExpr^①，执行过程

^① 通过表达式状态节点中 evalfunc 字段指向的函数来处理当前节点的调用过程实现。

类似于计划节点的递归方式。清理工作会通过释放内存上下文来实现。

初始化过程会为每个表达式节点生成对应的执行状态节点，并为状态节点的 evalfunc 字段初始化执行函数（用于执行该表达式节点功能的函数），然后调用 ExecInitExpr 对子节点进行初始化，并将子节点的执行状态节点组织成链表链存放于当前状态节点的 args 字段中。

执行过程 ExecEvalExpr 是一个宏定义，由于为每个表达式节点都定义的处理函数具有统一的参数和返回，因此 ExecEvalExpr 的定义比较简单，如下所示：

```
((* (expr) -> evalfunc) (expr, econtext, isNull, isDone))
```

这里将以执行可优化语句中涉及的投影和选择运算上的表达式为例，介绍表达式计算的相关内容。查询编译器为了尽早进行选择和投影运算，在生成查询计划时会将 SQL 语句中的选择条件尽量分配到下层计划节点，并在计划节点输出结果前投影需要的属性。因此，表达式广泛存在于计划节点的投影和选择列表中，用于进行属性值计算和条件选择，成为了计划节点处理过程的公共部分。

为了说明表达式表示的方式，我们仍以可优化语句执行中的查询语句为例。在图 6-21 的查询计划树中，有一个叶子节点是扫描 teacher 关系的 SeqScan 节点，图 6-59 显示了该节点上的投影和选择表达式。

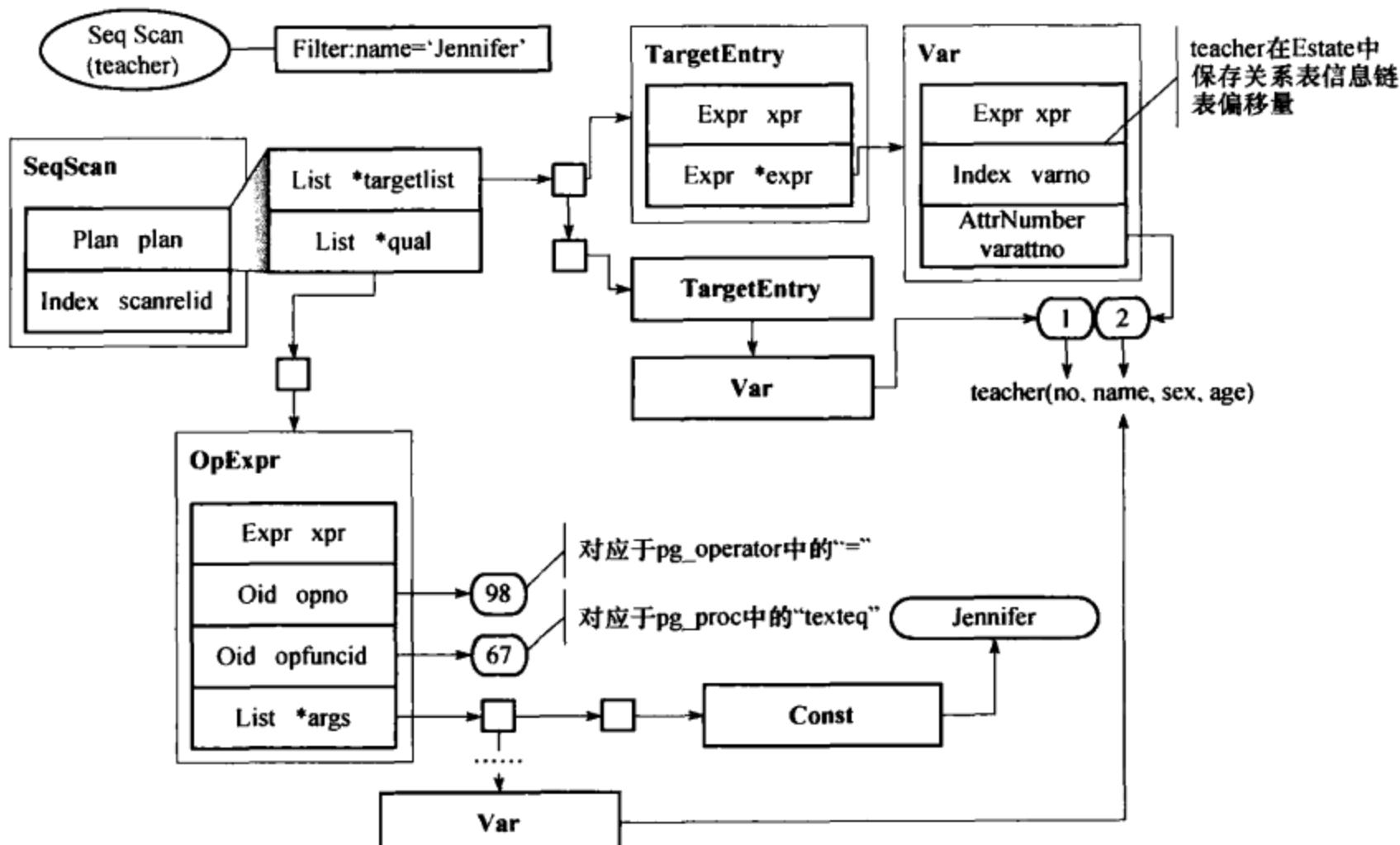


图 6-59 表达式的存储实例

在扫描节点的 targetlist（目标属性）字段中，使用 TargetEntry 结构链表表示需要从结果元组中投影出的属性：no 和 name，每一个目标属性用一个 TargetEntry 结构表示。其中，name 属性用于在这个 SeqScan 节点中选择元组，而 no 属性则用于上层连接节点。SeqScan 中的扫描条件记录在 qual

字段中，该字段也指向一个链表。“name = ‘Jennifer’”这个选择条件被表示成 qual 中的一个 OpExpr 类型的节点，该节点中的 opno 记录了操作符的 OID（即“=”），opfuncid 记录了 name 对应的数据类型所使用的等于操作符函数（即“texteq”）的 OID，而 args 字段中用一个链表记录了用于该操作符的参数（包括 name 和常量“Jennifer”）。OpExpr 的作用就是使用 texteq 函数将 teacher 的 name 属性值与常量“Jennifer”进行比较。

在计划节点初始化过程中，会调用 ExecInitExpr 对 targetlist 和 qual 中的表达式进行初始化，图 6-60 给出了初始化过程生成的表达式状态节点结构图。表达式状态节点中的 expr 指针是指向其关联的表达式节点的，图中标出了初始化后对应的状态节点类型以及赋值的操作函数（evalfunc）。执行过程中使用相应的操作函数，处理该状态节点，通过指向对应的 expr 指针获取其中保存的部分信息，并根据需求从 ExprContext 获取数据进行计算。

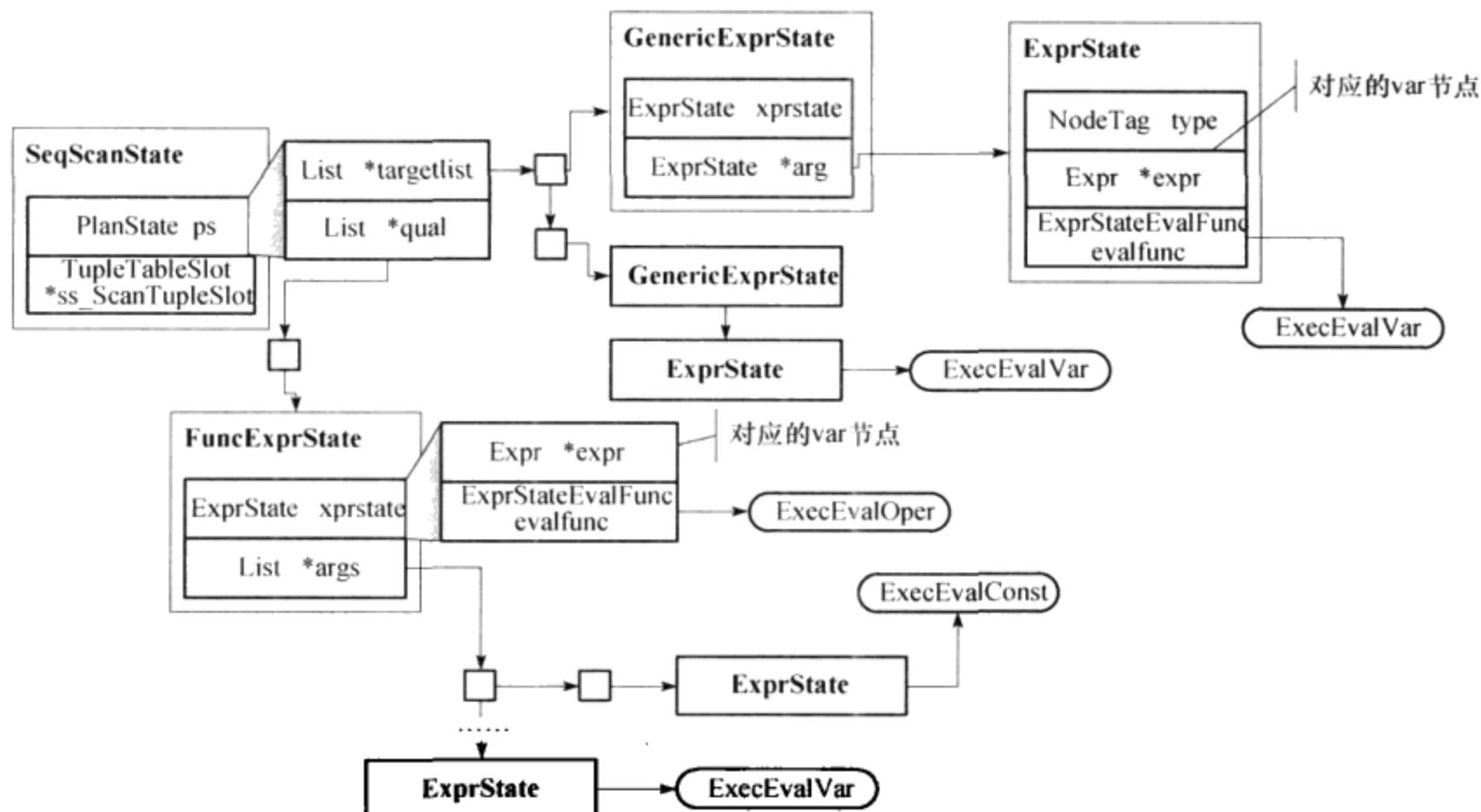


图 6-60 表达式状态结构图

PostgreSQL 中常用的表达式节点如表 6-6 所示。

表 6-6 常用表达式节点

节点类型	处理函数	功能
Var	ExecEvalVar	从元组中获取指定属性值
Const	ExecEvalConst	返回一个常量值
Aggregref	ExecEvalAggregref	返回已经计算好的聚集函数的值
WindowFunc	ExecEvalWindowFunc	返回已经计算好的窗口函数的值
FuncExpr	ExecEvalFunc	执行系统中定义的函数
OpExpr	ExecEvalOper	与 FuncExpr 类似，实现操作符运算，也是通过调用函数实现

6.5.3 投影操作

投影操作本身也是通过表达式计算来实现的。投影操作是一种属性过滤过程，该操作将对元组的属性进行精简，把那些在上层计划节点中不需要用到的属性从元组中去掉，从而构造一个精简版本的元组。投影操作中那些被保留下来的属性保存在查询计划的目标属性中，也被称为投影属性。

如图 6-61 所示，PostgreSQL 将投影属性组成了链表，每个属性使用一个 TargetEntry 结构进行存储，TargetEntry 也被定义为 Expr 的子类，并扩展定义了 expr 字段指向一个表达式树的根节点，此表达式树记录了需要映射的属性如何映射得到结果属性过程，执行过程也是通过处理这个表达式树实现。

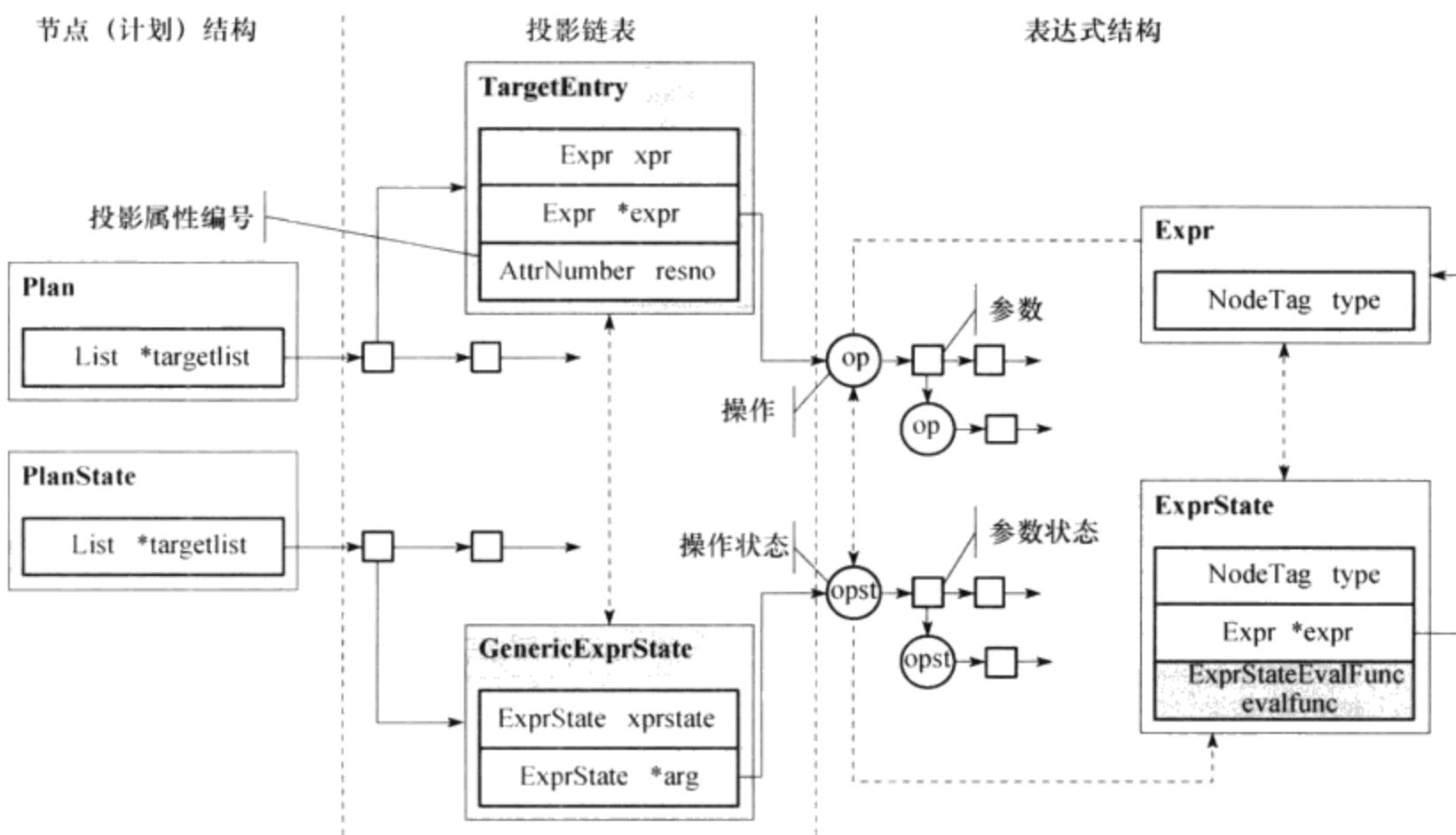


图 6-61 投影运算相关结构

执行过程前，首先会在计划节点初始化过程中为投影表达式链表生成对应的表达式状态链表。根据前面的介绍，对于任意 Expr 子类不一定会有特殊定义的 ExprState 子类来存储执行状态，TargetEntry 使用通用的表达式状态节点 GenericExprState，其中的 arg 就指向表达式状态树的根节点。表达式初始化过程会为表达式树生成相应状态树，然后构建投影运算的核心数据结构 ProjectionInfo。图 6-62 展示了 ProjectionInfo 的数据结构。

ExecBuildProjectionInfo 函数用于构造 ProjectionInfo，其参数包括初始化好的表达式状态链表 (targetlist)、表达式计算的上下文 (ExprContext) 以及投影结果输出的元组指针。

ExecBuildProjectionInfo 的执行流程如下：

- 1) 扫描输入的表达式状态链表获取一个表达式树，若已经没有表达式树需要处理则退出。否则，通过根节点 TargetEntry (表达式树都是以 TargetEntry 结构为根节点) 对应的 GenericExprState 获取 TargetEntry，从而获取其中保存的真正用于投影的表达式树根节点。

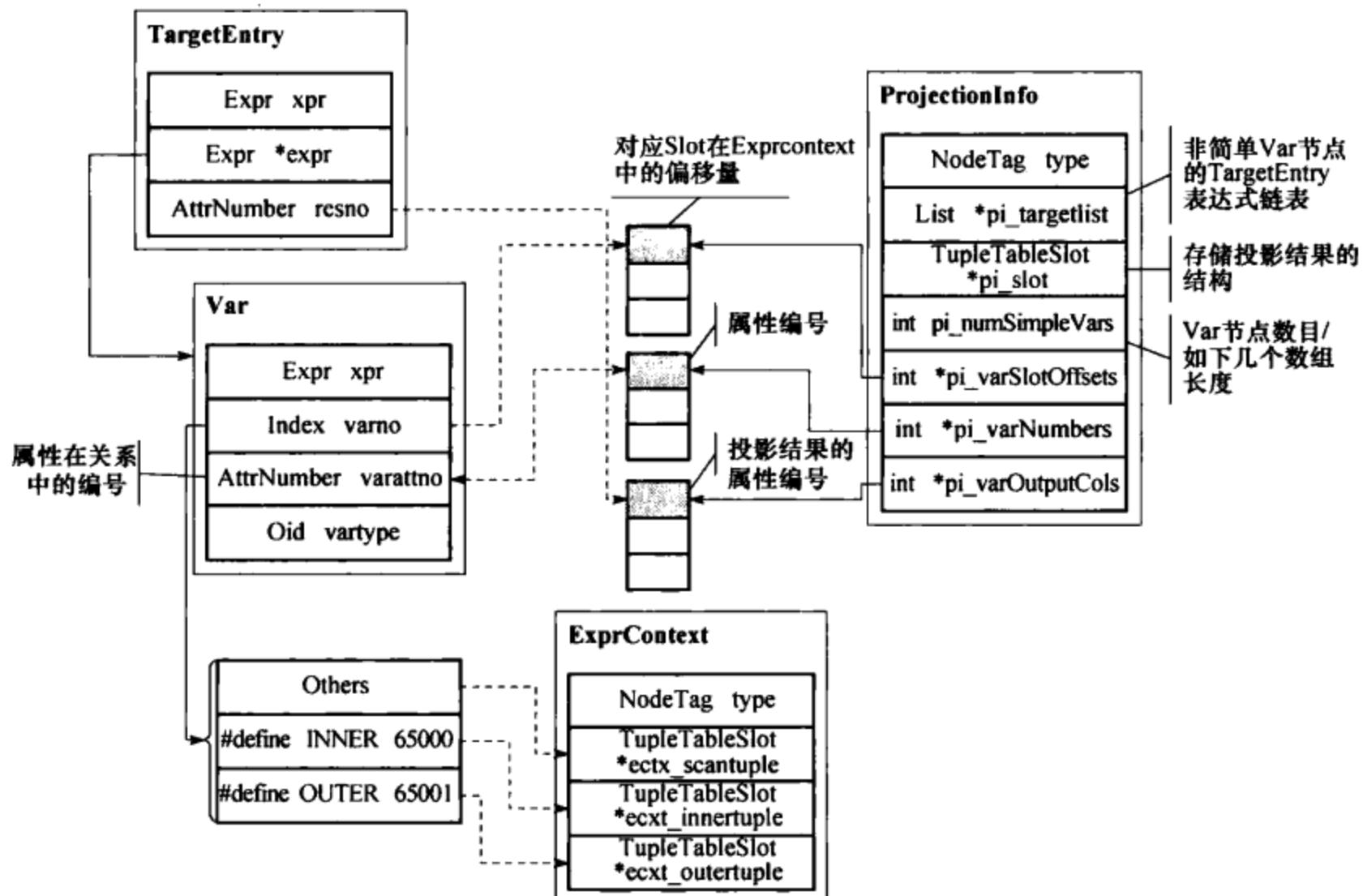


图 6-62 投影计算的 ProjectionInfo 结构

2) 判断 TargetEntry 中 expr 指向的表达式树根节点是否为 Var 节点。如果是 Var 节点，则该投影属性只是一个简单的属性映射，执行步骤 3；否则执行步骤 4。

3) 使用 Var 结构的 varno 获取原属性所在元组，将该元组在 ExprContext 结构中的偏移量存储在 pi_varSlotOffset [pi_numSimpleVars] 中，并将原属性在元组中的序号 (varattno) 存储在 pi_varNumbers [pi_numSimpleVars]，并将 TargetEntry 标记的输出属性在结果中的属性编号存储在 pi_varOutputCols [pi_numSimpleVars] 中，最后执行 pi_numSimpleVars 加 1 操作，执行步骤 1。

4) 将 TargetEntry 结构中表达式状态树放入 ProjectionInfo 结构的 pi_targetlist 链表中，继续执行步骤 1。

执行投影操作的函数为 ExecProject，使用 ProjectionInfo 结构作为参数进行投影运算。运算过程如下：

- 首先读取 ProjectionInfo 三个数组中存储的直接属性映射。按照执行的偏移获取原属性所在的元组，通过偏移量获取该属性，并通过目标属性的序号找到对应的新元组属性位置进行赋值。

- 对 targetlist 中复杂表达式树进行运算，将结果根据对应的 TargetEntry 的 resno 字段，赋值给结果元组的相应属性。

投影操作与表达式计算的不同之处在于对单一 Var 节点的特殊处理，这种处理方式有助于加快投影操作的速度。

6.6 小结

本章介绍了 PostgreSQL 中对于各种 SQL 语句的一般执行流程。对于用户输入的 SQL 语句，优化器将为可优化语句生成计划树，最终 Portal 会通过判断选择 Executor 来处理，而数据描述语句则在执行过程中由 Portal 使其统一进入 ProcessUtility 过程进行执行。

作为查询执行部分的入口，Portal 提供了对外的调用接口：PortalStart、PortalRun、PortalEnd，对内提供了执行流程和部件的选择。外层通过调用 Portal 接口，将计划器输出的执行计划传递给 Portal，Portal 通过对链表中操作的类型和链表长度等信息来决定选择怎样的执行过程，对于简单的查询语句直接调用 Executor，对于需要缓存输出直到执行完成的语句则需要为其增加缓存结构和输出过程，对于更为复杂的过程提供了更为通用的复杂处理流程。不论哪种执行过程，优化器生成的查询计划树由 Executor 来处理，而其他数据描述语句的功能由 ProcessUtility 来完成。

对于种类繁多的数据描述语句，每种都有一个 Stmt 类型的数据结构保存其语法分析的信息，通过对于 Stmt 类型的判断，ProcessUtility 会为其调用相应的语法解析和执行处理过程。

可优化语句则会提供 Plan 类的子类对象构成的计划树，每个计划节点数据结构共同继承于 Plan 节点，计划树的每种节点对应了一种物理操作，被分为四大类，除了对应于关系代数的各种操作外，还增加了扫描、物化、唯一等操作。每种节点根据所对应的操作的不同，执行器都为其实现了初始化、执行和清理流程。

Plan 只是一个查询计划，初始化过程中为这个计划里的每个节点生成状态节点并构造状态树，保存执行中的相关信息。最终通过 Executor 执行每种节点的执行函数来实现各种节点的物理操作，在执行中利用状态树中信息处理相关数据。各节点的执行过程普遍包含了投影和选择操作，以及对于下层节点的处理过程调用。

Executor 通过迭代的调用每个节点处理过程，从下层计划节点对应的执行流程中获取数据，并经过各种节点的处理流程，得到最终的输出结果。对于修改元组的相关操作则是在获取到元组后通过调用相关存储操作接口实现的。

执行过程需要使用存储和索引提供的接口存取数据，并涉及事务和权限控制等内容，这些会在相关章节给出更加详细的介绍。

习题

- 习题 6.1** 查询处理部分可划分为哪几部分，各部分的功能是什么？
- 习题 6.2** Portal 部分有哪几种执行策略，有何区别？
- 习题 6.3** PostgreSQL 中的物理操作节点分为几类，功能有何区别？
- 习题 6.4** 思考如果一个节点需要获取多个子节点的结果，PostgreSQL 中是如何实现的？

第 7 章

事务处理与并发控制

7.1 事务系统简介

事务是数据库操作的执行单位，一个事务提交之后或者完整执行，或者完全不执行。在 PostgreSQL 中，任何语句的执行都存在于事务环境中，任何语句开始执行之前事务就开始了。当语句执行完后，事务结束。通常情况下，我们提交的查询语句或者更新语句会被送入一个默认事务环境中执行。PostgreSQL 中负责管理事务运行的模块称为事务管理器，其结构如图 7-1 所示。

在了解事务处理的相关知识之前，我们先了解一下事务系统中各部分的作用。

1) 事务管理器是事务系统的中枢，它的实现是一个有限状态自动机（Finite State Machine），通过接受外部系统的命令或者信号（包括用户命令或者操作系统信息），并根据当前事务所处的状态，决定事务的下一步执行过程。7.1 节到 7.5 节将介绍事务管理器的相关内容。

2) 锁管理器实现了系统并发控制所需要的各种锁。在 PostgreSQL 系统中，事务执行的读阶段采用了多版本并发控制（MVCC），即对元组的读和写互不阻塞；而在事务执行的写阶段则需要由各种锁来保证事务的隔离级别。相关内容将在 7.6 节到 7.10 节中介绍。

3) 日志管理器用来记录事务执行的状态以及数据的变化过程，包括事务提交日志（CLOG）和事务日志（XLOG）。其中，CLOG 只用来记录事务执行的结果状态，而 XLOG 才是通常在“数据库原理”课程中所提及的日志，它记录了数据的变化过程并保持了一定的冗余数据。7.11 节将介绍日志管理器相关内容。

为了使读者更容易理解 PostgreSQL 的事务系统，我们需要先从整体上了解事务的执行流程，否则很可能陷入到事务执行中所涉及的锁操作、日志操作等细节之中而一头雾水。下面我们通过一个

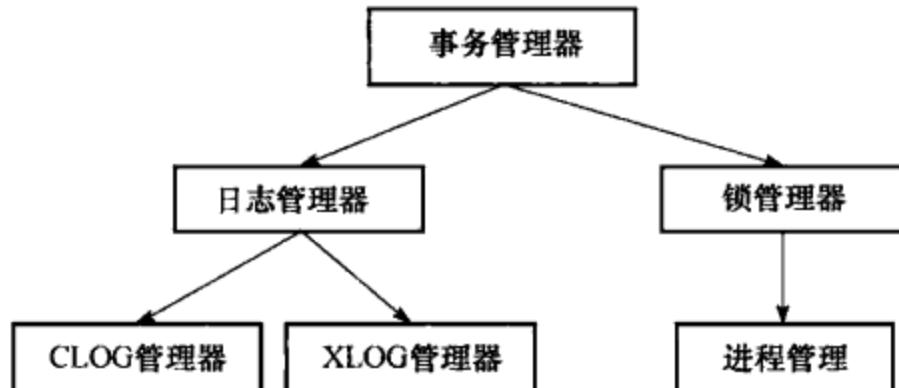


图 7-1 事务管理器与系统其他部分的联系

例子来再现 PostgreSQL 中事务的处理过程。

例 7.1 执行以下 SQL 命令：

```
BEGIN;
SELECT* FROM TABLE_A;
END;
```

PostgreSQL 在接收到上述语句之后，就会 Fork 出一个子进程 Postgres 来处理上述命令。对于这样一个由查询语句组成的简单事务，其实际执行过程如图 7-2 所示。首先会从总控进入事务块环境中，由于还没有事务信息，于是便调用底层事务处理函数启动一个事务，然后返回到上层，这时遇到“BEGIN”语句，于是改变事务块的状态，该状态标识后面在遇到“END”或者“ROLLBACK”命令之前的所有语句都属于同一个事务块。对于 SELECT 语句，由于它不改变事务块状态，所以直接进入事务底层，进入执行器执行该语句。最后遇到 END 语句，则调用事务底层函数提交事务并退出，然后结束事务块。总之，在 PostgreSQL 中，上层的事务块以及底层的事务共同构成了传统数据库中所提到的事务的概念。任何语句的执行总是先进入事务处理接口事务块中，然后调用事务底层函数处理具体的命令，最后返回事务块中。

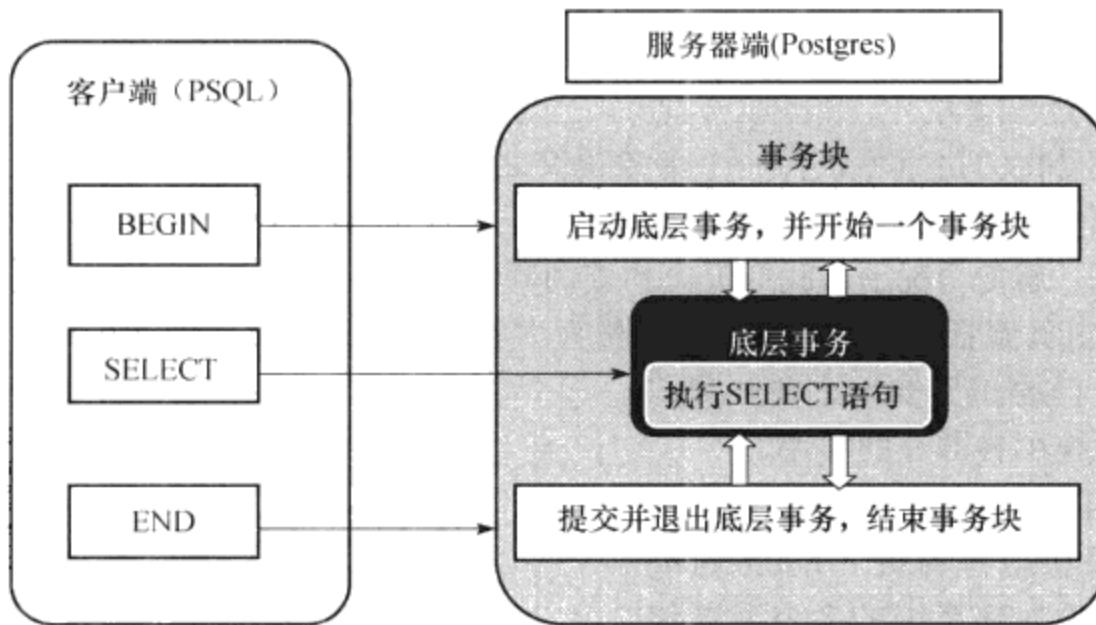


图 7-2 事务执行层次图

在下面两节中，我们将从事务系统的上层和下层来分别阐述事务的执行流程，读者可以利用调试工具来跟踪上面语句的执行过程，从而了解事务执行过程。

7.2 事务系统的上层

在数据库中引入事务是为了保证数据操作的 ACID 特性，而为了提高用户操作数据的灵活性，在 PostgreSQL 系统中，通过事务块来实现数据库内部底层事务和终端命令的交互，事务块对应着通常数据库理论中所提到的“事务”的概念。而 PostgreSQL 中的事务则体现了数据库理论中所提到的命令的概念，它的状态仅仅反映的是实际一个 SQL 语句（命令）的执行状态。为了在后文中加以区分，我们用带引号的“事务”表示通常数据库理论中谈到的事务概念；不带引号的事务表示 PostgreSQL 中与命令相对应的事务。在 PostgreSQL 中，一个事务块中包含多个事务，所以事务块的状态数量要比底层事务的状态数量多得多。

在 PostgreSQL 的事务处理层次中，位于事务系统上层的是事务块。PostgreSQL 执行一条 SQL 语句前会调用 StartTransactionCommand 函数，执行结束时会调用 CommitTransactionCommand 函数。如果命令执行失败，则会调用 AbortCurrentTransaction 函数。上述三个函数根据事务块的状态执行不同的操作，并调用不同的底层事务执行函数。

提示 PostgreSQL 将事务系统分成上层（事务块）和底层（事务）两个层次，通过分层的设计，在处理上层业务的时候可以屏蔽具体细节。

需要注意的是，上述三个函数只是进入事务系统上层的入口函数，并不处理具体事务。

7.2.1 事务块状态

在默认的情况下，PostgreSQL 中的一个事务处理一条 SQL 语句。事务块的设置使得一个事务能够处理多条 SQL 语句。在每一个 Postgres 进程中，自始至终都只存在一个事务块，当用户开始一个新的“事务”时，并不进入新的事务块，而是修改事务块的状态使之重新开始记录用户“事务”的状态。

事务块状态指的是整个事务系统的状态，它的变化反映了当前用户输入命令的变化过程和事务底层的执行状态的变化。比如，用户输入 BEGIN/ROLLBACK/END 语句会改变事务块状态，底层事务提交或者终止也会改变该状态。PostgreSQL 把事务块状态定义为一个枚举类型 TBlockState（数据结构 7.1）。

数据结构 7.1 TBlockState

```
typedef enum TBlockState
{
    /* 不在事务块中的状态 */
    TBLOCK_DEFAULT,           // 事务块的缺省状态
    TBLOCK_STARTED,           // 执行简单查询事务

    /* 处于事务块中的状态 */
    TBLOCK_BEGIN,              // 遇到事务的开始命令 BEGIN 时设置的状态
    TBLOCK_INPROGRESS,         // 表明事务块正在处理当中
    TBLOCK_END,                // 遇到事务结束命令 COMMIT/END 时设置的状态
    TBLOCK_ABORT,               // 事务出错，等待 ROLLBACK
    TBLOCK_ABORT_END,           // 事务出错，已接收 ROLLBACK
    TBLOCK_ABORT_PENDING,        // 事务处理中，已接收到 ROLLBACK
    TBLOCK_PREPARE,             // 事务处理中，已接收到 PREPARE(分布式事务处理)

    /* 子事务事务块状态，各状态的含义参考上面注释 */
    TBLOCK_SUBBEGIN,
    TBLOCK_SUBINPROGRESS,
    TBLOCK_SUBEND,
    TBLOCK_SUBABORT,
    TBLOCK_SUBABORT_END,
    TBLOCK_SUBABORT_PENDING,
    TBLOCK_SUBRESTART,
    TBLOCK_SUBABORT_RESTART
}TBlockState;
```

7.2.2 事务块操作

在 PostgreSQL 中，任何语句的执行都是通过事务块入口函数进入并执行的，执行完毕后，通过事务块出口函数退出，当系统在执行中遇到错误时，会通过事务块出错处理函数完成相关的出错处理。另外，用户也可以通过 BEGIN/END/ROLLBACK 等命令来改变事务块的状态，PostgreSQL 针对每一个能够改变事务块状态的命令提供了单独的函数，在本节中将阐述这方面内容。

1. 事务块基本操作

事务块基本操作函数包括 StartTransactionCommand、CommitTransactionCommand 和 AbortCurrentTransaction。其中，StartTransactionCommand 在每条语句执行前调用，CommitTransactionCommand 在每条语句执行后调用，AbortCurrentTransaction 则在系统遇到错误时调用。

(1) StartTransactionCommand 函数

在 PostgreSQL 中，每条语句执行前，都需要执行该函数，它的作用是通过判断当前事务块的状态进入不同的底层事务执行函数中。因此，它是进入事务处理模块的人口函数，直接在总控模块中被调用。

该函数根据当前事务块的状态决定下一步事务操作的方向。例如，假设当前事务块的状态为 TBLOCK_DEFAULT，那么调用该函数时说明当前系统中还不存在事务，于是调用底层事务操作函数启动一个事务，并设置事务块的状态为 TBLOCK_STARTED。

(2) CommitTransactionCommand 函数

在 PostgreSQL 中，每条 SQL 语句执行后，都需要执行该函数，它同 StartTransactionCommand 一样，也是根据当前事务块的状态决定下一步的事务操作，只不过它在每条语句执行结束后调用。

下面举例说明该函数的执行过程。假设当前事务块的状态为 TBLOCK_DEFAULT，调用该函数时说明此时系统没有任何事务信息，我们知道 PostgreSQL 中的任何语句都是在事务中执行的，而该函数又在执行语句后调用，说明现在处于一个非法状态，于是打印错误信息并返回。

(3) AbortCurrentTransaction 函数

在 PostgreSQL 中，当系统遇到错误时，会返回到调用点并执行该函数，其作用也是根据当前事务块的状态决定下一步的事务操作。只不过该函数在系统遇到错误时调用，其中可能进行事务出错退出操作或者事务清理操作。

假设当前事务块的状态为 TBLOCK_DEFAULT，如果系统执行遇到错误，那么会返回到 sigsetjmp 函数，然后调用该函数进行出错处理。如果底层事务的状态为 TRANS_DEFAULT，那么说明此时系统中没有任何事务信息，不需要进行事务出错处理，直接返回即可；若当前底层事务状态为 TRANS_START，说明此时底层事务已经存在，现在退出事务块需要调用 AbortTransaction 函数，再调用 CleanupTransaction 函数进行事务清理操作。

2. 事务块状态的改变

事务块状态改变函数主要包括 BeginTransactionBlock、EndTransactionBlock 和 UserAbortTransac-

tionBlock，它们的功能主要是根据用户输入的改变事务状态的命令以及系统本身的状态来改变事务块的状态。

前面已经提到，我们向 PostgreSQL 提交的任何一个 SQL 语句都被系统默认为一个事务块。如果要自定义一个事务块的边界，则要用 BEGIN 和 END 命令来实现。BEGIN 和 END 命令类似于一个大括号，在这个大括号中的 SQL 语句被系统认为是处于同一个事务块中的语句。

(1) 执行 BEGIN 命令

从例 7.1 中可以看到，执行 BEGIN 命令后，系统会进入一个事务块，这里状态的改变是调用函数 BeginTransactionBlock 来完成的。该函数判断当前事务块状态，如果事务块状态为 TBLOCK_STARTED，说明当前事务还没有进入到真正的事务块中，那么把事务块状态置为 TBLOCK_BEGIN，说明这是一个事务块的开始。如果当前事务块状态是 TBLOCK_INPROGRESS 或者 TBLOCK_SUBINPROGRESS，说明已经有一个事务块正在运行，这时系统会输出一个警告。对于 BeginTransactionBlock 函数来说，TBLOCK_END、TBLOCK_ABORT-END 等状态是无效的，一旦该函数遇到这些无效状态，PostgreSQL 会出现一个错误。

(2) 执行 END 命令

PostgreSQL 在接收到 END 命令时会调用函数 EndTransactionBlock，该函数指示一个事务块的结束。由于一个事务块的结束可能由 COMMIT 造成，也可能由 ROLLBACK 造成，所以 EndTransactionBlock 将根据返回值来确定系统的动作，返回 TRUE 表示函数成功 COMMIT，返回 FALSE 表示 ROLLBACK。

该函数的操作过程如下：

1) 判断当前事务块状态。如果事务正在执行当中，即事务块状态为 TBLOCK_INPROGRESS，则把事务块状态改为 TBLOCK_END，同时把返回值置为 TRUE；如果事务块状态为 TBLOCK_SUBINPROGRESS，说明当前处于一个活动的子事务内，递归地把子事务块状态置为 TBLOCK_SUBEND，最后把顶层事务块状态改为 TBLOCK_END，同时把返回值置为 TRUE。

2) 如果当前事务失败，即事务块状态为 TBLOCK_ABORT，则需要退出该事务块，把事务块状态置为 TBLOCK_ABORT-END，同时返回 FALSE；如果当前事务块状态是 TBLOCK_SUBABORT，即我们处于失败的子事务中，这时把整个事务看成是一个失败事务，递归地把子事务块状态置为 TBLOCK_SUBABORT-END，最后把顶层事务块状态置为 TBLOCK_ABORT-END，同时返回 FALSE。

(3) 执行 ROLLBACK 命令

函数 UserAbortTransactionBlock 用于执行用户提交的一个 ROLLBACK 命令。如果当前事务块状态是 TBLOCK_ABORT_PENDING，由于已经收到了用户发来的 ROLLBACK 命令，需要把事务块状态置为 TBLOCK_ABORT_PENDING 以退出该事务块。如果当前事务已经是一个失败事务，事务块状态为 TBLOCK_ABORT，则只需把该状态置为 TBLOCK_ABORT-END 即可。

图 7-3 展示了事务块几种主要状态之间的转换关系。

BeginTransactionBlock、EndTransactionBlock、UserAbortTransactionBlock 这三个函数都只改变事务块状态，并没有做实际的事务操作。

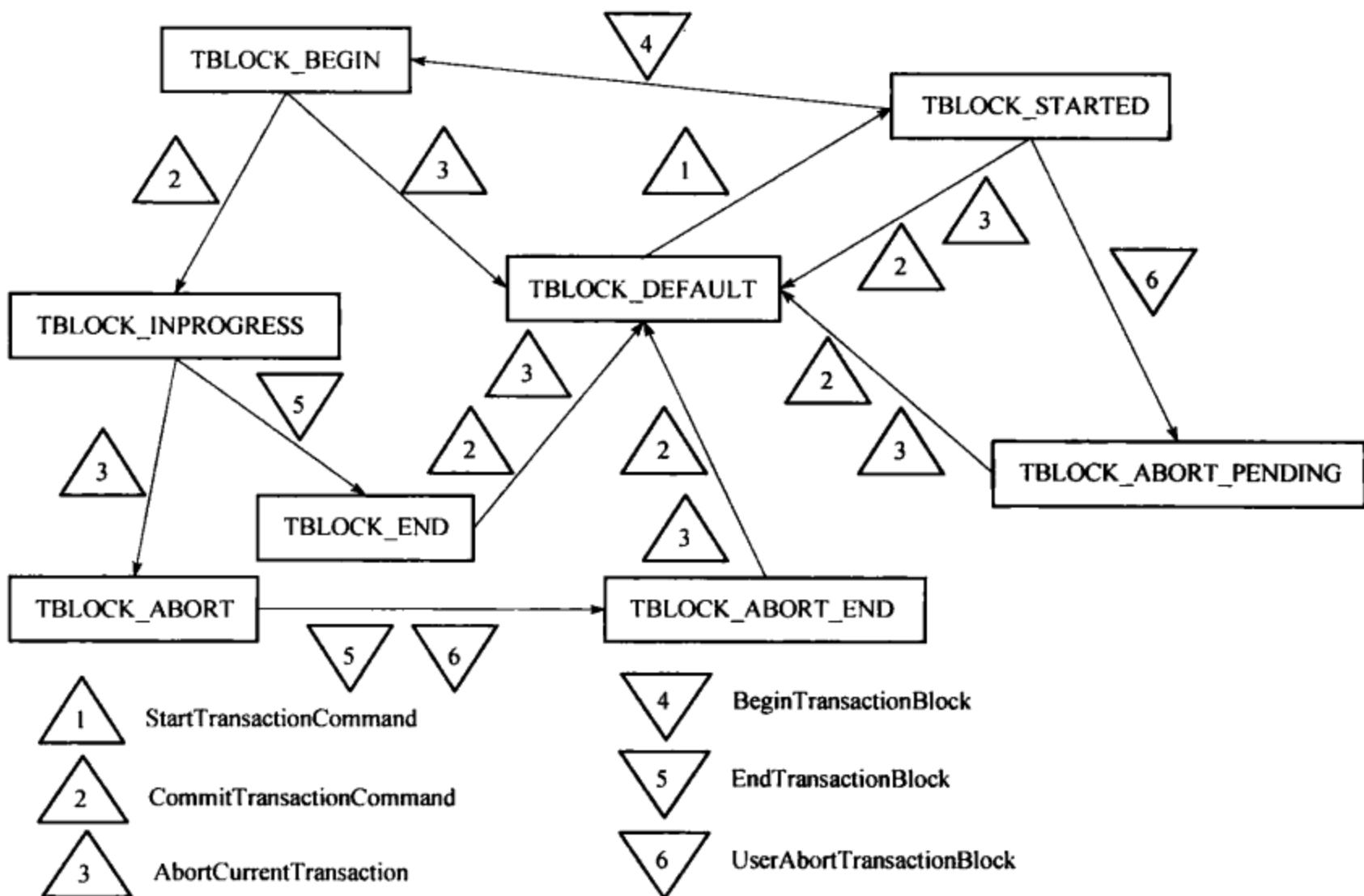


图 7-3 事务块主要状态之间的转换关系

7.3 事务系统的底层

上一节介绍了事务系统上层事务块的操作，主要涉及事务执行的入口和出口的相关处理，并不关注事务具体的执行细节。在本节，我们将详细介绍当事务进入到底层事务之后的具体处理过程，包括资源和锁的获取及释放、信号的处理、日志记录等。

7.3.1 事务状态

在 PostgreSQL 系统中，事务的状态定义在 TransState（数据结构 7.2）中。

数据结构 7.2 TransState

```
typedef enum TransState
{
    TRANS_DEFAULT,
    TRANS_START,
    TRANS_INPROGRESS,
    TRANS_COMMIT,
    TRANS_ABORT,
    TRANS_PREPARE
} TransState;
```

TransState 表示的是一种低层次的事务状态。而事务处理过程中存储事务状态相关数据的数据结构是 TransactionStateData（数据结构 7.3）。

数据结构 7.3 TransactionStateData

```

typedef struct TransactionStateData
{
    TransactionId      transactionId;           //当前事务 XID
    SubTransactionId   subTransactionId;         //子事务 ID
    char               *name;                   //保存点名称
    int                savepointLevel;          //保存点层数
    TransState          state;                  //低层次事务状态,例如 TRANS_INPROGRESS
    TBlockState         blockState;             //事务块状态
    int                nestingLevel;            //事务嵌套深度
    int                gucNestLevel;            //GUC 嵌套深度
    MemoryContext       curTransactionContext; //事务当前上下文
    ResourceOwner       curTransactionOwner;  //缓冲区等资源
    TransactionId*     childXids;              //提交的子事务链表
    int                nChildXids;              //提交的子事务个数
    int                maxChildXids;            //已分配的子事务存储空间
    Oid                prevUser;                //用户 OID
    bool               prevSecDefCxt;           //用户上下文环境设置标记
    bool               prevXactReadOnly;        //只读事务标记
    struct TransactionStateData* parent;       //指向父事务的指针
} TransactionStateData;

```

在 TransactionStateData 结构中，着重介绍以下几个字段的作用。

1) nestingLevel：代表当前事务所在的嵌套级别，其中顶层事务的嵌套级别为 1，每开启一个子事务（通过调用 PushTransaction 函数），该子事务的 nestingLevel 将在父事务的基础上加 1。

2) gucNestLevel：记录的是本事务 GUC（Grand Unified Configuration，全局统一配置）的嵌套深度。在 TopTransaction 中会分配一个 GUC 栈用来存储事务嵌套压栈后 GUC 变量的变化情况。当子事务提交时，其父事务变量会出栈，此时的 gucNestLevel 与 GUC 栈中的嵌套级别比较，用来确定 GUC 变量是否恢复，即当此时的 gucNestLevel 小于 GUC 栈的嵌套级别，表示当前的 GUC 变量需要恢复。顶层事务的 gucNestLevel 为 1，当函数 PushTransaction 把子事务压入栈的时候，首先会把全局变量 GUCNestLevel 加 1，然后再把该值赋给子事务的 gucNestLevel。当一个子事务中止时，在 gucNestLevel 大于等于当前事务嵌套级别（nestingLevel）的事务中定义的全局变量都会被丢弃。

3) curTransactionOwner：是指向 ResourceOwnerData（见 3.6 节相关内容）的指针，用于记录当前事务占有的资源。

4) prevUser：记录的是上一个 CurrentUserId 设置，因为事务执行过程中 CurrentUserId 可能改变。当开启子事务，父事务入栈时，PostgreSQL 会调用 GetUserIdAndContext 将当前 UserId 保存起来，当事务出栈时会调用 SetUserIdAndContext 恢复父事务曾经的 UserId。CurrentUserId 对应的是目前执

行环境下的用户名。例如，访问权限查询函数 `has_table_privilege` 查询用户是否有访问表的权限，该函数缺省的用户即 `CurrentUserId`。

5) `prevXactReadOnly`: 记录本事务上层事务的读写状态，顶层事务不使用这个字段，并把 `prevXactReadOnly` 置为 FALSE。在子事务执行 `CommitSubTransaction` 之后，需要恢复上层事务的读写状态标志。当前事务的读写状态标志通过全局变量 `XactReadOnly` 记录，每当子事务提交时，需要将子事务的 `prevXactReadOnly` 值赋值给全局变量 `XactReadOnly` 以恢复上层事务的读写标志状态。

7.3.2 事务操作函数

上一节中我们介绍了改变事务块状态的函数 `BeginTransactionBlock`、`EndTransactionBlock`、`UserAbortTransactionBlock`，它们并没有做实际的事务操作。实际的事务操作由 `StartTransaction`、`CommitTransaction`、`AbortTransaction`、`CleanupTransaction` 等函数完成。

1. 启动事务

启动事务由函数 `StartTransaction` 完成，该函数主要完成事务的初始化工作。函数的执行过程如下：

1) 将全局变量 `CurrentTransactionState` 指向 `TopTransactionStateData`，然后设置底层事务的状态为 `TRANS_START`（数据结构 7.4）。

2) 确保旧快照已经被释放，重置事务状态变量，并初始化事务内部计数器。

3) 使用 `AtStart_Memory` 和 `AtStart_ResourceOwner` 函数分别完成内存上下文和资源跟踪器的初始化。

4) 产生一个事务标识 XID 给当前事务。顶层事务通过 `VirtualTransactionID` 来识别，`VirtualTransactionID` 包含两个部分：`BackendId` 和 `LocalTransactionId`。其中 `BackendId` 就是运行此事务的进程 ID 号。`LocalTransactionId` 通过调用 `GetNextLocalTransactionId` 函数进行分配。`GetNextLocalTransactionId` 函数的实现非常简单，通过一个全局变量 `nextLocalTransactionId` 来产生 `LocalTransactionId`。分配完 `VirtualTransactionID` 之后，调用函数 `VirtualXactLockTableInsert` 锁住这个虚拟事务 ID，最后再把 `VirtualTransactionID` 中的 `LocalTransactionId` 插入到当前进程队列（`PGPROC`）中。

5) 设置时间戳。设置事务的开始时间戳（`xactStartTimestamp`）为当前命令开始的时间戳（`stmtStartTimestamp`），并初始化事务的结束时间戳为 0。

数据结构 7.4 TopTransactionStateData

```
static TransactionStateData TopTransactionStateData = {
    0,                      // 事务标识 ID
    0,                      // 子事务 ID
    NULL,                  // 检查点名称
    0,                      // 检查点嵌套层数
    TRANS_DEFAULT,          // 当前事务状态
    TBLOCK_DEFAULT,         // 所处事务块状态
    0,                      // 事务嵌套级别
    0,                      // GUC 嵌套级别
    NULL,                  // 当前事务的内存上下文
    NULL,                  // 当前事务的资源跟踪器
    NULL,                  // 已提交的子事务链表
    0,                      // 已提交的子事务个数
    0,                      // 分配的子事务 XID 存储空间
    InvalidOid,             // 父事务 CurrentUserId
    false,                  // 父事务 SecurityDefinerContext
    false,                  // 只读事务标记
    NULL                    // 指向父事务指针
};
```

6) 初始化当前事务状态变量。包括:

- 设置事务的嵌套级别为 1, 调用 StartTransaction 启动的一定是顶级事务, 子事务启动会调用 StartSubTransaction, 所以对于顶级事务其事务的嵌套深度为 1。
- 初始化当前事务的 GUC 嵌套深度为 1, 并且初始化已提交子事务列表为空。

7) 初始化 GUC、Cache 等结构。到这里, 启动事务所需要的操作已经基本完成。最后, 把事务状态设置为 TRANS_INPROGRESS。

2. 提交事务

提交事务由函数 CommitTransaction 完成, 其执行过程如下:

1) 检查事务状态: 确保当前事务状态为 TRANS_INPROGRESS。检查完事务状态之后进入预提交过程。

2) 触发所有延迟的触发器: 这里的“延迟”是指 SQL 中 After 语句触发器。After 表示该触发器的动作在其触发事件进行完毕之后才被执行。当所有延迟触发器执行完毕后, 关闭所有延迟触发器, 并检查是否有需要删除的表或对象, 如果有则将其删除。

3) 关闭大对象: 在低层次清除操作开始前, 调用函数 AtEOXact_LargeObject 关闭所有大对象, 并释放大对象内存上下文以防永久性的内存泄漏。如果 pendingNotifies 链表中还有本事务执行过程中产生的延迟通知, 就把这些通知发送出去。

4) 如果修改了 pg_database、pg_authid 或者 pg_auth_members, 则执行更新操作并通知 Postmaster 进程。这个过程由函数 AtEOXact_UpdatePasswordFile 完成, 该函数是事务真正提交前的最后一步。如果执行完 AtEOXact_UpdatePasswordFile 之后系统出现错误异常中断本事务, Postmaster 可以接收到这个消息并做善后处理。

5) 提交事务: 现在可以把当前事务状态置为 TRANS_COMMIT 了, 并执行 RecordTransactionCommit 函数, 用于将当前所提交的事务所产生的日志写回磁盘。返回最后处理完的当前事务或子事务的 XID, 如果当前事务没有 XID, 则返回 InvalidTransactionId。

6) 清理操作: 到目前为止, 事务已经“真正”提交了, 但还需进行一些提交后的清除工作。此时, 即使再出现什么错误也不能中止事务了, 所以在这一步释放的只能是一些非关键资源。资源释放的大致顺序如下: 先释放其他后台进程可见的资源(如文件、缓冲区锁等), 接着释放锁, 再释放进程的本地资源。之所以优先释放其他后台进程可见的资源, 后释放锁, 是为了保证等待本事务所持有锁的其他后台进程在获得锁的时候可以确定本事务已经完成了清除工作。因为本事务所在进程的本地资源与其他进程完全无关, 对其他后台进程来说是完全透明的, 所以留在最后释放。

7) 恢复初始状态: 设置事务状态为缺省状态 TRANS_DEFAULT, 然后准备接受下一事务。

3. 退出事务

退出事务由函数 AbortTransaction 完成, 该函数在系统遇到错误时调用, 属于系统终止事务退出时的行为。其执行流程如下:

- 1) 关中断, 因为进行事务退出后的清理工作时不能被别的事件打断。
- 2) 确保内存上下文和资源跟踪器有效。
- 3) 释放所有拥有的轻量锁。注意, 这里并不释放 Regular Lock (一般锁), 我们需要一直持有 Regular Lock 直到事务完成退出为止。
- 4) 检查事务状态。确保当前事务状态为 TRANS_INPROGRESS 或 TRANS_PREPARE。然后设置

事务状态为中止 TRANS_ABORT。

5) 终止执行事务。取消任何 AFTER 语句触发器，将 afterTriggers 列表置空，并置各个大对象描述符为空，同时清除延迟发送的通知信息。

6) 记录日志。调用 RecordTransactionAbort 函数，同时为其记录 XLOG 和 CLOG 日志。返回最后处理完的当前事务或子事务的 XID，如果当前事务没有 XID，则返回 InvalidTransactionId。

7) 清理资源。删除本事务创建的任何新表，处理 GUC 并重置服务器编程接口状态，中止后清除 oncommits 链表中的项，关闭所有仍然开启着的临时文件的 VFD 等。

8) 最后恢复中断。

4. 清理事务

清理事务的工作由函数 CleanupTransaction 完成，它在系统终止时调用，主要功能是释放事务所占用的内存资源。与 AbortTransaction 不同的是，该函数是终止事务退出时最后调用的函数，做最后的实际清理工作。AbortTransaction 函数并没有实际释放相关资源，只是切换一些资源的状态使其能够被其他事务获得。其执行流程如下：

- 1) 判断当前事务状态是否为 TRANS_ABORT 并释放 Portal 内存。
- 2) 释放 ResourceOwner、MemoryContext，并把事务状态恢复到默认情况。

7.3.3 简单查询事务执行过程实例

经过上面的描述，我们再来细化一下例 7.1 中 PostgreSQL 简单查询事务的实际操作过程（见表 7-1）。

表 7-1 简单查询事务执行过程

BEGIN	<ol style="list-style-type: none"> 1) 此时，事务块状态为 TBLOCK_DEFAULT。 2) 调用 StartTransactionCommand 函数，由于当前事务块状态为 TBLOCK_DEFAULT，于是调用 StartTransaction 函数，首先设置事务状态为 TRANS_START，然后完成有关内存、缓冲区、锁的处理后设置事务状态为 TRANS_INPROGRESS。最后设置事务块状态为 TBLOCK_STARTED。 3) 调用 ProcessUtility 函数处理 BEGIN 语句。 4) 遇到“BEGIN”，所以调用 BeginTransactionBlock 函数处理，它将设置事务块状态为 TBLOCK_BEGIN。 5) 最后调用 CommitTransactionCommand 函数，由于当前事务块状态为 TBLOCK_BEGIN，于是改变事务块状态为 TBLOCK_INPROGRESS，并准备读取下一条命令。
SELECT	<ol style="list-style-type: none"> 1) 调用 StartTransactionCommand 函数，由于当前事务块状态为 TBLOCK_INPROGRESS，表明正在继续该事务的命令，所以直接返回。 2) 调用 ProcessQuery 函数进入执行器处理 SELECT 语句。 3) 调用 EndTransactionBlock 函数，由于当前事务块状态为 TBLOCK_INPROGRESS，于是返回，并准备继续读下一个命令。
END	<ol style="list-style-type: none"> 1) 调用 StartTransactionCommand 函数，由于当前事务块状态为 TBLOCK_INPROGRESS，表明正在继续该事务的命令，所以直接返回。 2) 调用 ProcessUtility 函数处理 END 语句。 3) 调用 commitTransactionBlock 函数并设置事务块状态为 TBLOCK_END。 4) 最后调用 CommitTransactionCommand 函数，由于当前事务块状态为 TBLOCK_END，于是调用 CommitTransaction 函数提交事务，首先设置事务状态为 TRANS_COMMIT，然后真正提交事务，并清理相关资源，等着一切结束后，设置事务状态为 TRANS_DEFAULT，并返回。最后设置事务块状态为 TBLOCK_DEFAULT，整个事务执行结束。

7.4 事务保存点和子事务

事务保存点（SAVEPOINT）提供了一种机制，用于回滚部分事务。使用 SAVEPOINT 语句创建一个事务保存点，必要的时候可以用 ROLLBACK TRANSACTION TO 语句回滚到该保存点。

例 7.2 SAVEPOINT 的用法

```
BEGIN Transaction S;
insert into table1 values(1);
SAVEPOINT P1;
insert into table1 values(2);
SAVEPOINT P2;
insert into table1 values(3);
ROLLBACK TO savepoint P1;
COMMIT;
```

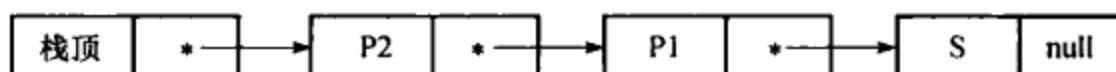
上述语句的执行结果是在表 table1 中插入数据 1。

7.4.1 保存点实现原理

PostgreSQL 在实现上把 SAVEPOINT 认为是定义子事务的标记。事务状态数据结构 TransactionStateData（数据结构 7.3）中有如下变量：

```
struct TransactionStateData*parent; /* 指向父事务的指针 */
```

parent 指针指向本事务的父事务，例 7.2 中事务的层次结构如下所示：



子事务的层次描述用一个链栈实现。栈顶元素拥有一个指向其父事务的指针。当启动一个新的子事务时，系统调用 PushTransaction 函数把描述该子事务的 TransactionState 结构变量压入栈中，这个变量可以标识该事务。相应的，PopTransaction 函数的功能是把栈顶事务弹出。

PushTransaction 函数为子事务创建一个 TransactionState 并压入事务状态堆栈中。在函数执行过程中，CurrentTrasactionState 会被切换到新创建的事务状态。PopTransaction 函数将当前事务状态弹出堆栈，把 CurrentTransactionState 切换到父事务的事务状态并转换资源所有者以及事务内存上下文。

保存点的作用主要就是实现事务的回滚，即从当前子事务回滚到该事务链中的某个祖先事务。在定义一个保存点时，需要用户输入保存点名称（例如 P1），从而标识当前子事务（参考数据结构 7.3）。在用户执行了一系列子事务之后回滚到该子事务时（ROLLBACK TO P1），会调用 RollbackToSavepoint 函数在事务链中从当前事务节点回溯查找该保存点名称，若找到，就直接把 CurrentTransactionState 指针移到查找到的事务节点处。由于 PostgreSQL 中 MemeoryContext 的存在，这里并不需要释放抛弃的子事务的相关内存资源，释放资源的操作在最后事务提交或者退出时通过 MemoryContext 来实现。

7.4.2 子事务

PostgreSQL 拥有一套与事务操作函数类似的子事务操作函数，包括 StartSubTransaction、Commit-

SubTransaction、CleanupSubTransaction 和 AbortSubTransaction 等。其实，在 PostgreSQL 中，子事务的主要作用是实现保存点，从而增强事务操作的灵活性。

在执行 PushTransaction 后会把当前事务块状态置为 TBLOCK_SUBBEGIN，在接下来为 SAVEPOINT 语句执行的函数 CommitTransactionCommand 中会调用 StartSubTransaction 函数并把 blockstate 设为 TBLOCK_SUBINPROGRESS。

在系统实现中，上层的事务在创建的过程中可以得到一个 XID，用于标识自己，这在 VACUUM 的过程中会用到。如果不对数据库进行写操作，就没有必要单独获得事务或其子事务的 XID。所以对于子事务，在它们调用 GetCurrentTransactionId 函数的时候才给它们分配 XID。

系统利用子事务 ID 标识每个子事务，最上层的子事务的值为 1，其子事务为 2，依此类推，0 用于表示无效的子事务 ID。

7.5 两阶段提交

PostgreSQL 使用两阶段提交（Two Phase Commit, 2PC）来支持分布式数据库的事务处理。

例 7.3 典型分布式数据库事务处理

A 在银行 1 的账户上有 \$100.

B 在银行 2 的账户上有 \$50.

A 转账 \$25 给 B.

如果“转账”的动作由于某些原因失败，A 和 B 的账户应该不会发生改变。但是由于银行 1 和银行 2 的数据库不同，简单直接地进行事务处理无法满足要求，因为“A 转账 \$25 给 B”的操作可以分解为三个步骤：

- 1) A 的账户减去 25。
- 2) A 发消息给 B，要求 B 的账户增加 25。
- 3) B 的账户增加 25。

现在假设步骤 1 和 2 都完成了，其中 A 所在的数据库已完成提交操作，这时步骤 3 由于某种原因执行失败，然后退出事务，这时候显然转账失败，即出现对账不平现象。假如此时 B 再发消息给 A，要求 A 进行减去 25 的操作，但由于某种原因，A 进行该操作时也失败了，怎么办呢？出现这种现象的根本原因在于“A 转账 \$25 给 B”操作不具有原子性，这时候就需要用到分布式数据库的事务处理。实现分布式事务处理的关键就是两阶段提交协议。

需要注意的是，PostgreSQL 只是支持分布式数据库中的两阶段提交协议，即给其提供了相关的操作接口，但并没有实现整个协议，两阶段协议的整个流程是由编程者在应用程序中保证的。所以下面叙述的两个过程是两阶段协议的内容，PostgreSQL 提供了最基本的操作接口。

7.5.1 预提交阶段

两阶段提交协议的第一个过程是预提交阶段，协议规定其执行步骤如下：

1) 对于分布式事务，某一个事务所在的数据库管理系统会被选择成为“协调者”。协调者在本地开始一个分布式事务，并且向其他数据库管理系统发送“Prepare”消息。发送消息时，会使用专门的事务 ID (GID) 来标识此分布式事务。这样数据库管理系统之间可以确定需要同步执行的事务。

2) 其他数据库管理系统接收到“Prepare”消息后，会试图开始一个本地事务以完成分布式事务的功能。它自行决定这个事务是提交还是终止，然后把它的决定发送给协调者。

3) 如果数据库决定提交上述的一个本地事务，它就进入“预提交”阶段。在此阶段，如果协调者没有发送终止的消息，它不能终止这个本地事务。

4) 如果数据库决定终止这个事务，它会向协调者发送取消的信息。然后由协调者进行全局性的取消动作。

在预提交阶段，各个分布式数据库系统将检查当前各自事务的情况，试图保证这个本地事务如果执行，就不会被终止。即使因系统出错而需要进行系统恢复，与这个事务相关的恢复操作也是应该 Redo 而不是 Undo。

在 PostgreSQL 系统中，当用户决定为当前事务做两阶段提交的准备时，利用命令 PREPARE TRANSACTION 可以完成。执行该命令时会调用 PrepareTransaction 函数，该函数的流程同 StartTransaction 一样，只不过其中插入了 StartPrepare 函数和 EndPrepare 函数，前者构建两阶段提交的相关记录的头部信息，后者完成记录数据信息的填充，并刷新到磁盘上。

7.5.2 全局提交阶段

两阶段提交协议的第二个过程是全局提交阶段，协议规定其步骤如下：

1) 如果协调者没有收到消息，就默认收到了“取消”消息。如果其他数据库返回给协调者的消息都是“Ready”，协调者将提交这个分布式事务，然后把“Commit”消息发送给其他数据库。如果协调者收到一条“取消”消息，则取消这个分布式事务，然后发送消息进行全局性的取消动作。

2) 本地数据库根据协调者的消息，对本地事务进行 COMMIT 或者 ABORT 操作。

在这个阶段，各个分布式数据库系统将进行事务的实质提交或者退出操作。

在 PostgreSQL 系统中，当用户准备提交一个早先为“两阶段提交”准备好的事务时，通过输入命令 COMMIT PREPARED 可以完成，执行该命令时会调用函数 FinishPreparedTransaction，它首先读取磁盘上先前预提交阶段所记录的信息，然后调用 RecordTransactionCommitPrepared 完成事务的最终提交操作。

如果用户决定取消一个先前为两阶段提交准备好的事务，则通过命令 ROLLBACK PREPARED 可以完成。执行该命令时它也会调用 FinishPreparedTransaction 函数，即首先读取磁盘上先前预提交阶段所记录的信息，然后会调用函数 RecordTransactionAbortPrepared 进行事务终止的操作。

两阶段提交协议是分布式数据库系统中保证分布式事务 ACID 特性的经典解决方案。PostgreSQL 为其提供了很好地操作接口，应用程序开发者利用该操作接口可以很好地实现该协议。

7.6 PostgreSQL 的并发控制

PostgreSQL 中存在多个会话试图同时访问同一数据的情况，并发控制的目标就是保证所有会话高效地访问，同时维护数据完整性。从本节开始到 7.10 节，将介绍 PostgreSQL 并发控制的实现机制。

PostgreSQL 为开发者提供了丰富的对数据并发访问进行管理的工具。在内部，PostgreSQL 利用多版本并发控制（MVCC）来维护数据的一致性。这就意味着当检索数据时，每个事务看到的都只

是一段时间之前（不同隔离级别下所看到的事务开始之前）的数据快照（一个数据库版本）。这样，如果对每个数据库会话进行事务隔离，就可以避免一个事务看到其他并发事务的更新而导致不一致的数据。

在 PostgreSQL 里也有表和行级别的锁定机制，因为 MVCC 并不能解决所有的并发控制情况，所以还需要使用传统数据库中的锁机制来保证事务的并发。另外，PostgreSQL 还提供了会话锁机制，利用它可以扩大锁的使用范围，即一次对某个对象加锁可以保证对于多个事务都有效。

SQL 标准考虑了三个必须在并行的事务之间避免的现象：

- 1) 脏读 (Dirty Reads)：一个事务读取了另一个未提交的并行事务写的数据。
- 2) 不可重复读 (Non-repeatable Reads)：一个事务重新读取前面读取过的数据，发现该数据已经被另一个已提交的事务修改过。
- 3) 幻读 (Phantom Read)：一个事务重新执行一个查询，返回一套符合查询条件的数据，发现这些数据因为其他最近提交的事务而发生了改变。

为了避免出现这三种现象，SQL 标准定义了 4 个事务隔离级别（见表 7-2）。

表 7-2 事务隔离级别表

隔离级别	脏读	不可重复读	幻读	隔离级别	脏读	不可重复读	幻读
读未提交	可能	可能	可能	可重复读	不可能	不可能	可能
读已提交	不可能	可能	可能	可串行化	不可能	不可能	不可能

在 PostgreSQL 里，可以请求四种可能的事务隔离级别的任意一种。但是在内部，实际上只有两种独立的隔离级别，分别对应读已提交和可串行化（由于实现 MVCC 的缘故）。如果选择了读未提交的级别，实际上用的是读已提交；若选择可重复读的级别，实际上用的是可串行化。所以实际的隔离级别可能比你选择的更严格。这是 SQL 标准允许的：四种隔离级别只定义了哪种现象不能发生，但是没有定义哪种现象一定发生。

下面介绍 PostgreSQL 中定义的两种隔离级别。

- 读已提交 (Read Committed)：这是 PostgreSQL 里的缺省隔离级别。当一个事务运行在这个隔离级别时，一个 SELECT 查询只能看到查询开始之前提交的数据而永远无法看到未提交的数据或者是在查询执行时其他并行的事务提交所做的改变。如果两个事务在对同一元组进行更新，第二个更新事务将等待第一个更新事务提交或者回滚。如果第一个更新回滚，那么它的作用将被忽略，而第二个更新者将继续更新最初发现的元组。如果第一个更新者提交，系统将重新计算查询搜索条件 (WHERE 子句)，如果元组符合条件，则第二个更新继续其操作，从该元组的已更新版本开始。
- 可串行化 (Serializable)：它提供最严格的事务隔离。这个级别模拟串行的事务执行，就好像事务一个接着一个地串行（而不是并行地）执行。如果两个事务在对同一个元组进行更新，可串行化的事务将等待第一个正在更新的事务提交或者回滚。如果第一个更新者回滚，那么它的影响将被忽略，这个可串行化的事务就可以在该元组上完成其更新操作。但是如果第一个更新者提交了，那么可串行化事务将回滚，从头开始重新进行整个事务。

在 PostgreSQL 系统中，事务的隔离级别所涉及的最小实体是元组，所以对元组的操作（例如读

取、插入、更新、删除)需要实施访问控制,它们是通过锁操作以及MVCC相关的操作来实现的,在本章后面的内容中将详细介绍。下面举例说明在PostgreSQL两种不同的隔离级别下事务的执行过程。假设有两个事务T1和T2并行执行,它们的执行情况如表7-3所示。

表7-3 两个事务的执行情况

事务T1	说明	事务T2
BEGIN	事务开始执行时,事务管理器为新事务分配XID;在共享内存中分配空间存放该事务结构;将新事务的XID返回给正在调用的进程	BEGIN
UPDATE C SET teacher = 'J. D. Ullman' WHERE CNO = '30001'	自动给cno为“30001”的记录加排他锁(Exclusive-Lock)	
SELECT * FROM C		SELECT * FROM C
END	事务提交时,事务管理器调用日志管理器,将日志中的状态改成“提交”,并将日志刷新到磁盘。同时调用锁管理器用于释放该事务所持有的所有的锁	
		UPDATE C SET cname = 'Database' WHERE CNO = '30001'

如果当前的隔离级别是读已提交,则事务T2中的“SELECT * FROM C”语句执行时看不到T1对数据库的修改,但“UPDATE C”语句可以看到。所以,在执行“UPDATE C”时,从已更新的数据版本开始。

如果当前的隔离级别是可串行化,则事务T2中的“SELECT * FROM C”语句执行时看不到T1对数据库的修改,“UPDATE C”语句也看不到T1对数据库的修改。那么可串行化事务T2将回滚,从头开始重新进行整个事务。

7.7 PostgreSQL中的三种锁

PostgreSQL实现并发控制的基本方法是使用锁来控制临界区互斥访问。后台进程对磁盘文件进行访问操作时,首先要获取锁。如果成功获得目标锁,则进入临界区执行磁盘读写访问,访问完成后退出临界区并释放锁;否则,进程睡眠直到被别的后台进程唤醒后重试。

PostgreSQL中定义了三种锁,分别是SpinLock、LWLock和RegularLock。

7.7.1 SpinLock

SpinLock是最底层的锁,使用互斥信号量实现,与操作系统和硬件环境联系紧密。SpinLock分为与机器相关的实现方法(定义在s_lock.c中)和与机器不相关的实现方法(定义在Spin.c中)。SpinLock的特点是:封锁时间很短,没有等待队列和死锁检测机制,事务结束时不能自动释放SpinLock。

作为一种最底层的锁,一般不直接使用SpinLock,而是利用它来实现其他锁(LWLock)。

毫无疑问,依赖于硬件的SpinLock机制肯定比不依赖于硬件的SpinLock机制速度快。因为不依

赖于硬件的 SpinLock 机制需要使用 PG 信号量来仿真 SpinLock。

如果机器拥有 TAS (test-and-set) 指令集，那么 PostgreSQL 会使用 `s_lock.h` 和 `s_lock.c` 中定义的 SpinLock 实现机制。如果机器没有 TAS 指令集，那么不依赖于硬件的 SpinLock 的实现定义在 `Spin.c` 中，需要用到 PostgreSQL 定义的信号量 PGSemaphore。

7.7.2 LWLock

LWLock（轻量级锁）主要提供对共享存储器的数据结构的互斥访问。LWLock 有两种锁模式，一种为排他模式，另一种为共享模式。轻量级锁不提供死锁检测，但轻量级锁管理器在 elog 恢复期间被自动释放，所以持有轻量级锁的期间调用 elog 发出错误消息不会出现轻量级锁未释放的问题。

LWLock 利用 SpinLock 实现，当没有锁的竞争时可以很快获得或释放 LWLock。当一个进程阻塞在一个轻量级锁上时，相当于它阻塞在一个信号量上，所以不会消耗 CPU 时间，等待的进程将会以先来后到的顺序被授予锁。

简单来说，LWLock 的特点是：有等待队列、无死锁检测、能自动释放锁。

1. LWLock 的数据结构

LWLock（数据结构 7.5）主要提供对共享存储器的数据结构的互斥访问。

数据结构 7.5 LWLock

```
typedef struct LwLock
{
    slock_t      mutex;           //保护 LwLock 和进程队列
    bool         releaseOK;
    char         exclusive;       //0 或 1,持有排他 LWLock 锁的后端数目
    int          shared;          //持有共享 LWLock 锁的后端数目(0..MaxBackends)
    PROC         *head;           //进程等待队列的头
    PROC         *tail;           //进程等待队列的尾,当头为空时,尾不确定
} LwLock;
```

可以看到，LWLock 的数据结构中定义了一个互斥量 `mutex`。使用 SpinLock 对此互斥量加/解锁即可实现对一个 LWLock 的互斥访问。

系统用一个全局数组 `LWLockArray` 管理所有的 LWLock。在 PostgreSQL 中，有 29 个 LWLock 被系统默认使用（如 `OidGenLock`、`WALInsertLock`、`CheckpointLock` 等，可参考 `LWLockId`）。用户还可以使用自定义 LWLock 来实现分段锁表。系统初始化 `LWLockArray` 时，为 `LWLockArray` 多分配了 2 个整型的空间，一个用来存储已分配的 LWLock 数量，另一个用来记录可分配的 LWLock 总数。

2. LWLock 的主要操作

(1) LWLock 的空间分配

Postmaster 启动后需要在共享内存空间中为 LWLock 分配空间，该操作定义在函数 `LWLockShmemSize` 中。分配空间时，需要计算出将要分配的 LWLock 的个数，计算过程由函数 `NumLWLocks` 给出，其计算公式如下所示：

```

numLocks = (int)NumFixedLWLocks + 2 * NBuffers + NUM_CLOG_BUFFERS
+ NUM_SUBTRANS_BUFFERS + NUM_MXACTOFFSET_BUFFERS
+ NUM_MXACTMEMBER_BUFFERS
+ Max(lock_addin_request, NUM_USER_DEFINED_LWLOCKS))

```

表 7-4 给出了公式中变量的含义。

表 7-4 NumLWLocks 中变量解释

变量	含义
NumFixedLWLocks	预先定义的 LWLock 数量
NBuffers	缓冲区个数，每一个缓冲区需要两个 LWLock
NUM_CLOG_BUFFERS	CLOG 缓冲区个数，每一个 CLOG 缓冲区需要一个 LWLock
NUM_SUBTRANS_BUFFERS	SubTrans 缓冲区个数，每个 SubTrans 缓冲区需要一个 LWLock
NUM_MXACTOFFSET_BUFFERS	MULTIXACT 中的每个 SLRU Offset Buffer 需要一个 LWLock（其中 Offset Buffer 存储的是 MULTIXACTID，系统设置用 8 个 Buffer。在 MVCC (7.10 节) 中，一个 MULTIXACT 标识当前元组涉及的多个处理过它的事务，这些事务 ID 共同组成了 MULTIXACT，系统使用 MULTIXACTID 来标识该 MULTIXACT，其中每个事务 ID 就是该 MULTIXACT 的成员）
NUM_MXACTMEMBER_BUFFERS	MULTIXACT 中的每个 SLRU Member Buffer 需要一个 LWLock（其中 Member Buffer 存储的是 MULTIXACT 的每个成员，在 MVCC (7.10 节) 中，由于一个元组中最多有两个事务 ID 存在 (Xmin 和 Xmax)，所以系统该 SLRU 成员个数是 SLRU Offset Buffer 的两倍，即 16 个 Buffer）
NUM_USER_DEFINED_LWLOCKS	留下一些额外的缓冲区供用户自定义使用

(2) LWLock 的创建

LWLock 的创建过程包括分配空间和初始化两个过程。其中分配空间的过程在 LWLockShmemSize 函数中给出，初始化即对 LWLock (数据结构 7.5) 结构的初始化，使它们处于“未上锁”的状态，并在 LWLockArray 队列的末尾初始化动态分配计数器。

(3) LWLock 的分配

LWLock 分配操作定义在函数 LWLockAssign 中，该函数的主要功能就是从共享内存中预定义的闲置 LWLock Array (表 7-4 中提到的 NUM_USER_DEFINED_LWLOCKS) 中得到一个 LWLock，同时 LWLockArray 前面的计数器会累加。基本过程是申请 SpinLock (访问 LWLock)，使用计数器得到闲置的 LWLock 数目，若没有闲置的 LWLock，输出错误；否则，修改计数器，释放 SpinLock，返回指向闲置 LWLock 的一个指针。

(4) LWLock 锁的获取

LWLock 的获取由函数 LWLockAcquire 定义，该函数试图以给定的轻量级锁 ID (LWLock ID) 和锁模式 (SHARE/EXCLUSIVE) 获得对应的轻量级锁。如果暂时不能获得，就进入睡眠状态直至该锁空闲。该函数执行流程如下：

- 1) 开中断，获取 SpinLock。
- 2) 检查锁的当前情况，如果空闲则获得锁，然后释放 SpinLock 退出。
- 3) 否则把自己加入等待队列中，释放 SpinLock，并一直等到被唤醒。
- 4) 重复上述步骤。

LWLock 还提供了另外一种获取 LWLock 的方式——使用 LWLockConditionalAcquire 函数，它与

上述函数的区别是若能获得此锁，则返回 TRUE，否者返回 FALSE。

(5) LWLock 锁的释放

LWLock 的释放由函数 LWLockRelease 完成，主要功能是释放指定 LockID 的锁。该函数的执行流程如下：

- 1) 获得该锁的 SpinLock。
- 2) 检查此次解锁是否唤醒其他等待进程。
- 3) 如果需要唤醒进程，遍历等待队列；如果遇到要求读锁的进程，从队列中删除，但保留一个指向它的指针。重复操作直到遇到要求写锁的进程。
- 4) 释放该 LWLock 的 SpinLock，把从队列中删除的进程唤醒。

另外，PostgreSQL 还提供能释放当前后端持有的所有 LWLock 锁的功能，该功能在函数 LWLockReleaseAll 中定义，主要在系统出现错误之后使用。其实现比较简单，就是重复地调用 LWLockRealease 函数进行 LWLock 的释放。

7.7.3 RegularLock

RegularLock 就是一般数据库事务管理中所指的锁，也简称为 Lock。RegularLock 由 LWLock 实现，其特点是：有等待队列，有死锁检测，能自动释放锁。

1. RegularLock 的锁方法类型

PostgreSQL 中使用了两种加 RegularLock 锁的方法：DEFAULT_LOCKMETHOD 和 USER_LOCKMETHOD，前者为默认锁方法，后者为用户锁方法。PostgreSQL 中一般将 DEFAULT_LOCKMETHOD 作为默认加锁方法。当然用户也可以定义自己的锁方法，比如建议锁（Advisory Lock）就是用户所创建的锁方法类型。

依据这两种不同的锁方法，可以生成两种不同类型的锁表。在 Postmaster 启动时，分配一块共享内存区作为 RegularLock 方法表并适当地初始化 RegularLock 方法表的各个字段，此外还将初始化锁方法数组和其他一些共享变量。后台进程启动后通过锁方法来引用 RegularLock。

用户锁方法应用于长期协作锁，即长事务。这种锁的用处在于提示应用程序某用户正在工作于某数据库元素上。因此，可以在一个数据库元组上加上用户锁，使得用户可以在长时间的检索和更新完该元组后再释放此锁。在该用户锁工作时，其他用户能检测到这个元组当前已经被另一应用级别用户加锁，其他用户在该元组上所能采取的操作将受到限制。

获取用户锁是无阻塞的，即进程或者得到用户锁，或者不能得到用户锁，不存在等待队列。因此如果一个进程持有了用户锁，则另一个进程就不能再获得它。当一个后台进程进程结束时，用户锁将自动被释放。默认锁和用户锁都是完整的，而且它们之间不互相干扰，即两种锁可以同时存在。

2. RegularLock 支持的锁模式

RegularLock 支持的锁模式有八种，按排他（Exclusive）级别从低到高分别是：

- 访问共享锁（AccessShareLock）：一个内部锁模式，进行查询（SELECT 操作）时自动施加在被查询的表上。
- 行共享锁（RowShareLock）：当语句中采用了 SELECT…FOR UPDATE 和 FOR SHARE 时将使用行共享锁对表加锁。
- 行排他锁（RowExclusiveLock）：使用 UPDATE、DELETE、INSERT 语句时将使用行排他锁对表加锁。

- 共享更新排他锁（ShareUpdateExclusiveLock）：使用 VACUUM（不带 FULL 选项）ANALYZE 或 CREATE INDEX CONCURRENTLY 语句时使用共享更新排他锁对表加锁。
- 共享锁（ShareLock）：使用不带 CONCURRENTLY 选项的 CREATE INDEX 语句请求时用共享锁对表加锁。
- 共享行排他锁（ShareRowExclusiveLock）：类似于排他锁，但是允许行共享。
- 排他锁（ExclusiveLock）：阻塞行共享和 SELECT... FOR UPDATE。
- 访问排他锁（AccessExclusiveLock）：被 ALTER TABLE、DROP TABLE 以及 VACUUM FULL 操作要求。

其中，排他模式的锁（ShareRowExclusiveLock、ExclusiveLock、AccessExclusiveLock）表示在事务执行期间阻止其他任何类型锁作用于此表；共享模式的锁（非排他模式的锁）表示允许其他用户同时共享此锁，但在事务执行期间阻止排他型锁的使用。排他模式和共享模式的锁都可以工作在以下授权级别上：

- Access：锁定整个表模式。
- Rows：仅锁定单独的元组。

访问共享锁是最低限制的锁，访问排他锁是限制最严格的锁。共享行排他型锁与排他锁相似，但它允许其他用户使用行共享锁。每种锁模式都有与之相冲突的锁模式，由锁冲突表（如表 7-5 所示）定义相关信息。

表 7-5 锁冲突表

编号	名称	用途	冲突关系
1	AccessShareLock	SELECT	8
2	RowShareLock	SELECT FOR UPDATE/FOR SHARE	7 8
3	RowExclusiveLock	INSERT/UPDATE/DELETE	5 6 7 8
4	ShareUpdateExclusiveLock	VACUUM	4 5 6 7 8
5	ShareLock	CREATE INDEX	3 4 6 7 8
6	ShareRowExclusiveLock	ROW SELECT...FOR UPDATE	3 4 5 6 7 8
7	ExclusiveLock	BLOCK ROW SHARE/SELECT...FOR UPDATE	2 3 4 5 6 7 8
8	AccessExclusiveLock	DROP CLASS/VACUUM FULL	1 2 3 4 5 6 7 8

在一般数据库系统中，使用 S（共享锁）、X（排他锁）、IS（意向共享锁）、IX（意向排他锁），以及 SIX（共享意向排他锁）五种锁即可完成加锁需要，由于 PostgreSQL 扩展了标准的 SQL 语句，所以对应于这些扩展的操作需要设计特定的锁完成加锁需求。下面详细解释以上几种在其他数据库中没有的锁。

1) RowShareLock 锁用于 SELECT FOR UPDATE/FOR SHARE 语句。下面详细说明这两种语句的作用以及加锁的冲突情况。

①FOR UPDATE 令那些被 SELECT 语句检索出来的行被锁住，就像要更新一样；从而避免它们在当前事务结束前被其他事务修改或者删除。也就是说，其他企图 UPDATE、DELETE、SELECT FOR UPDATE 这些行的事务将被阻塞，直到当前事务结束。同样，如果一个来自其他事务的 UPDATE、DELETE、SELECT FOR UPDATE 已经锁住了某个或某些选定的行，SELECT FOR UPDATE 将等到那些事务结束，并且将随后锁住并返回更新的行（或者不返回行，如果行已经被删除）。

②FOR SHARE 同 FOR UPDATE 类似，只不过它在每个检索出来的行上要求一个共享锁，而不

是一个排他锁。一个共享锁阻塞其他事务在这些行上执行 UPDATE、DELETE、SELECT FOR UPDATE，却不阻止他们执行 SELECT FOR SHARE。

2) ShareUpdateExclusiveLock 锁应用于 VACUUM 语句，VACUUM 回收已删除行所占据的存储空间。在一般的 PostgreSQL 操作里，那些已经 DELETE 的行或者被 UPDATE 过后过时的行并没有从它们所属的表中物理删除，在进行 VACUUM 之前它们仍然存在。因此有必要周期地运行 VACUUM，特别是在经常更新的表上。

3) ShareRowExclusiveLock 锁，目前 PostgreSQL 作为保留并没有使用。

4) ExclusiveLock 锁，除了一般意义上的排他锁，它还阻塞上面提到的 SELECT FOR UPDATE/FOR SHARE 语句。

总之，PostgreSQL 的 RegularLock 粒度可以是数据库的表、元组、页等，请求锁的主体既可以是事务，也可以是跨事务的会话。

3. RegularLock 的存储

共享存储器里有三种基本的锁结构：LOCK 结构、PROCLOCK 结构以及 LOCALLOCK 结构。LOCK 结构的存在是为了存储每个可锁的对象、可以保持锁或者请求锁。PROCLOCK 用于存储进程与锁之间的关系，通过该结构可以查询到当前进程阻塞了哪些进程，在死锁检测和消除中将频繁使用到该结构。另外，每个后台进程为当前进程中可锁对象和锁模式保留了没有共享的 LOCALLOCK 结构，它有助于快速处理加锁操作。对于某进程的加锁要求，首先在该进程的本地 LOCALLOCK 中查找当前进程是否已获得该锁，如果已获得，则直接对其引用计数加 1，授予该进程此次加锁请求。由于 LOCALLOCK 只对当前进程可见，所以对它的操作也不需要加 LWLock。

4. RegularLock 的数据结构

首先介绍一下锁模式和锁掩码。锁模式（LOCKMODE）是一个 1 ~ N 的整数，指示了一种锁类型，可表示前面提到的 8 种锁类型的一种。锁掩码（LOCKMASK）是一个位掩码，用于指示持有或请求的锁类型的集合，如 LockConflicts（数据结构 7.6）。

数据结构 7.6 LockConflicts

```
static const LOCKMASK LockConflicts[] = {
    0,
    /* AccessShareLock */
    (1 << AccessExclusiveLock),
    /* RowShareLock */
    (1 << ExclusiveLock)|(1 << AccessExclusiveLock),
    /* RowExclusiveLock */
    (1 << ShareLock)|(1 << ShareRowExclusiveLock)|
    (1 << ExclusiveLock)|(1 << AccessExclusiveLock),
    /* ShareUpdateExclusiveLock */
}
```

```

(1 << ShareUpdateExclusiveLock) |
(1 << ShareLock)|(1 << ShareRowExclusiveLock) |
(1 << ExclusiveLock)|(1 << AccessExclusiveLock),

/* ShareLock */
(1 << RowExclusiveLock)|(1 << ShareUpdateExclusiveLock) |
(1 << ShareRowExclusiveLock) |
(1 << ExclusiveLock)|(1 << AccessExclusiveLock),

/* ShareRowExclusiveLock */
(1 << RowExclusiveLock)|(1 << ShareUpdateExclusiveLock) |
(1 << ShareLock)|(1 << ShareRowExclusiveLock) |
(1 << ExclusiveLock)|(1 << AccessExclusiveLock),

/* ExclusiveLock */
(1 << RowShareLock) |
(1 << RowExclusiveLock)|(1 << ShareUpdateExclusiveLock) |
(1 << ShareLock)|(1 << ShareRowExclusiveLock) |
(1 << ExclusiveLock)|(1 << AccessExclusiveLock),

/* AccessExclusiveLock */
(1 << AccessShareLock)|(1 << RowShareLock) |
(1 << RowExclusiveLock)|(1 << ShareUpdateExclusiveLock) |
(1 << ShareLock)|(1 << ShareRowExclusiveLock) |
(1 << ExclusiveLock)|(1 << AccessExclusiveLock)
};


```

显然，系统中锁模式的个数不能比锁掩码的位数多。下面介绍几种 RegularLock 实现中所涉及的数据结构。

(1) 锁方法表

一个锁方法表的控制结构定义在 LockMethodData (数据结构 7.7) 中，它存在于共享内存中。

其中，numLockModes 指示在锁表上定义的锁类型数量，该数值必须小于或等于 MAX_LOCKMODES。如果 transactional 为 TRUE，说明这些锁会在事务结束后自动释放。conflictTab 显示锁类型冲突的位掩码数组，如果锁类型 i 和 j 冲突，那么 conflictTab[i] 的第 j 位为 1；锁模式的取值为 1 至 numLockModes，所以 conflictTab[0] 未使用。lockModeNames 用于 Debug，而 trace_flag 是指向本锁方法 GUC trace flag 的指针。

数据结构 7.7 LockMethodData

<pre> typedef struct LockMethodData { int numLockModes; bool transactional; const LOCKMASK *conflictTab; const char *const *lockModeNames; const bool *trace_flag; } LockMethodData; </pre>

(2) 加锁对象标识

LOCKTAG（数据结构 7.8）是在锁的 Hash 表格（hashtable）中寻找一个锁项所必须的码信息。一个 LOCKTAG 的值唯一标识一个可加锁对象。

数据结构 7.8 LOCKTAG

```
typedef struct LOCKTAG
{
    uint32      locktag_field1;
    uint32      locktag_field2;
    uint32      locktag_field3;
    uint16      locktag_field4;
    uint8       locktag_type;           // 锁类型标识
    uint8       locktag_lockmethodid;   // 锁方法 ID
} LOCKTAG;
```

LOCKTAG 结构中定义了一些 locktag_field，但没有具体指明这些变量的用途，因为针对不同的 LockTagType，这些 locktag_field 中保存的数据不一定相同。例如，LockTagType 为 LOCKTAG_TUPLE 时，要求的 locktag_field 包括 DB OID + REL OID + BlockNumber + OffsetNumber。如果 LockTagType 为 LOCKTAG_TRANSACTION，则只需要该事务的 XID。

(3) 加锁对象描述体

LOCK（数据结构 7.9）结构用于表示已经加锁的资源或是请求加锁的可锁资源。

数据结构 7.9 LOCK

```
typedef struct LOCK
{
    LOCKTAG          tag;           // 加锁对象的标识符 */
    LOCKMASK         grantMask;     // 当前在该对象上分配的所有锁类型的掩码 */
    LOCKMASK         waitMask;      // 当前在该对象上等待的所有锁类型的掩码 */
    SHM_QUEUE        procLocks;     // 与锁相关联的 PROCLOCK 对象队列 */
    PROC_QUEUE       waitProcs;     // 等待该对象上锁的进程等待队列 */
    int              requested[MAX_LOCKMODES]; // 记录每种模式的锁被要求的次数 */
    int              nRequested;    // 请求队列长度，即所有锁类型的锁请求的总的数量
    int              granted[MAX_LOCKMODES]; // 每一种锁类型上的已分配锁的数量 */
    int              nGranted;      // granted 数组中元素的个数
} LOCK;
```

其中，grantMask、waitMask 用于判断已持锁与请求锁是否冲突。

(4) 锁持有者信息描述体

对于同一可加锁对象，可能有几个不同事务持有或等待锁。我们需要为每一个这样的锁持有者（或即将持有者）存储持有者/等待者信息。这些信息保存在 PROCLOCK（数据结构 7.10）中。

数据结构 7.10 PROCLOCK

```

typedef struct PROCLOCK
{
    PROCLOCKTAG      tag;           // 可加锁对象的唯一标识
    LOCKMASK         holdMask;     // 显示该 PROCLOCK 所表示的已经分配的锁
    LOCKMASK         releaseMask;  // 显示该 PROCLOCK 已释放的锁
    /* 在 Lock 对象中有一个 PROCLOCK 链表,
     * 记录的是所有持有该锁的 PROCLOCK 对象,
     * 这里的 lockLink 就是记录本 PROCLOCK 对象在该链表中的位置 */
    SHM_QUEUE        lockLink;
    /* 每一个 PGPROC 对象中都有一个 PROCLOCK 链表,
     * 记录的是所有当前进程持有锁或即将持有锁的 PROCLOCK 对象,
     * 这里的 procLink 就是记录本 PROCLOCK 对象在该链表中的位置 */
    SHM_QUEUE        procLink;
} PROCLOCK;

```

在一个 PROCLOCK Hash 表中查找一个 PROCLOCK 项需要一个码信息 PROCLOCKTAG。PROCLOCKTAG 值为该对象唯一标识一个可加锁对象和一个持有者/等待者组合。

PROCLOCK 的所有者可能有两种：事务（通过运行它的后端的 PGPROC 和该事务的事务 ID 标识）或者会话（由后端 PGPROC 和事务 ID InvalidTransactionId 标识）。

当前会话 PROCLOCK 为用户锁和 VACUUM 持有的跨事务锁所使用。注意，一个单独的后端可以一次在多个不同事务下持有多个锁（包括会话锁）。我们把这样的锁都当成从不冲突（一个后端不会阻塞自己，即一个后端中不同事务所持有的各种锁相互间不冲突）。变量 holdMask 显示该 PROCLOCK 所表示的已经分配的锁。注意，如果一个进程当前正在等待锁，那么该进程和所等待的锁构成的 PROCLOCK 对象中的 holdMask 为 0。如果该 PROCLOCK 表明不是一个进程在等待锁，那么 holdMasks 为 0 的 PROCLOCK 对象将在方便的时候被回收。

每一个 PROCLOCK 对象被链接到关联的 LOCK 对象和所有者的 PGPROC 对象的列表中。注意，PROCLOCK 在创建的时候被马上插入到这些列表中，即使还没有锁被分配。一个正在等待分配锁的 PGPROC 也会被链接到该锁的进程等待队列中。

(5) 本地锁表

每一个后台进程也维护一个当前它感兴趣的（持有的或者希望持有的）每一个锁的本地 Hash 表。本地表 LocalLock（数据结构 7.11）记录了已经获取的锁的次数，这允许在不需要另外对共享内存进行访问的情况下对同一锁的多次加锁请求。同时也跟踪每一个资源所有者获得锁的数量，所以我们可以只释放属于某一资源所有者的锁。

数据结构 7.11 LocalLock

```

typedef struct LOCALLOCK
{
    LOCALLOCKTAG      tag;           // LocalLock 唯一标识
    LOCK              *lock;         // 关联 LOCK 结构
}

```

```

    PROCLOCK          *proclock;           //关联后端 PROCLOCK 结构
    uint32            hashcode;           //LocalTag 的 Hash 值
    int               nLocks;             //持锁数量
    int               numLockOwners;      //锁持有者数量
    int               maxLockOwners;     //已分配的存储空间
    LOCALLOCKOWNER   *lockOwners;        //锁持有者列表
} LOCALLOCK;

```

PostgreSQL 通过三张锁表联合管理系统的加锁情况。

- LockMethodLockHash：由锁方法 ID 到锁表数据结构的映射。
- LockMethodProcLockHash：锁方法对应的 PROCLOCK 对象的 Hash 表指针。
- LockMethodLocalHash：锁方法对应的 LOCALLOCK 对象的 Hash 表指针。

5. RegularLock 的主要操作

对于 RegularLock 的各种操作，相关函数定义在 lock.c 文件中。下面分别介绍这些操作，主要包括锁的空间计算、初始化、锁的申请、注销以及如何加锁、解锁等。

(1) RegularLock 的空间计算

锁空间计算操作由函数 LockShmemSize 实现，用于估计锁表所需共享内存大小。这些空间包括：

- “Lock Hash” 表所需空间。
- “ProcLock Hash” 表所需空间。
- 为安全考虑多分配 10% 的空间。

(2) RegularLock 的初始化

锁的初始化操作由函数 InitLocks 实现，该函数初始化锁管理器的数据结构，主要工作有初始化 LockMethodLockHash、LockMethodProcLockHash 和 LockMethodLocalHash 这三个 Hash 表。

(3) RegularLock 加锁

加锁操作在函数 LockAcquire 中定义，它有四个参数 LockTag、LockMode、SessionLock、dontWait：

- LockTag 是被锁对象的唯一标识。
- LockMode 指示要获得的锁模式（SHARE/EXCLUSIVE）。
- SessionLock 表示加锁的模式。如果为 TRUE，表示为会话加锁；如果为 FALSE，则为当前事务申请锁。
- DontWait 表示申请锁是否允许等待。如果为 TRUE，则在检查到无法获得锁之后不等待；如果为 FALSE，则可以等待。

该函数的返回值 LockAcquireResult（数据结构 7.12）表示加锁是否成功等结果信息。

数据结构 7.12 LockAcquireResult

```

typedef enum
{
    LOCKACQUIRE_NOT_AVAIL,      //不能获得锁,且不能等待
    LOCKACQUIRE_OK,             //成果获取到锁
    LOCKACQUIRE_ALREADY_HELD    //对已获得锁增加引用计数
} LockAcquireResult;

```

申请加锁的流程如下：

- 1) 首先用 LOCKTAG 和加锁模式得到具体的加锁类型，然后在本地表查找此加锁类型的信息。
- 2) 如果没有，则在本地表插入此信息；否则，分配空间以记录锁拥有者的信息。
- 3) 如果当前事务已经持有过此类型的锁，在本地表的计数器上加 1，然后直接退出。否则在全局的锁表（LockMethodLockHash）中查找这个锁。
- 4) 如果在“Lock Hash”中找不到，则在“Lock Hash”中插入一个新元素。然后在 ProcLock-Hash 表里也查找对应的 ProcLock。
- 5) 如果在 ProcLock Hash 表找不到，则插入该 ProcLock。
- 6) 检查这个类型的锁会不会与已加的锁发生冲突，如果不会，则加锁；否则，根据函数参数决定等待还是退出。如果退出，还需清除锁表中相应的元素以保持一致性。

(4) RegularLock 的释放

与加锁相对应的操作是解锁，RegularLock 的解锁操作定义在函数 LockRelease 中，该函数在本地锁表（LockMethodLocalHash）中查找锁标记为 LockTag 的锁，并释放该锁。如果 SessionLock 为 TRUE，则释放一个会话锁（SessionLock），否则，释放一个常规的事务锁。如果发现任何等待进程现在是可以被唤醒的，将请求的锁赋予它们并将其唤醒。

该函数的流程如下：

- 1) 用 LOCKTAG 和加锁模式得到具体的锁类型，然后在本地表查找此类型的锁的信息。
- 2) 找到此类型锁的拥有者，在它持有锁的计数器上减 1。如果它已经不再持有此锁，则删除这个拥有者的信息。
- 3) 如果这个类型的锁并没有真正释放，只是计数器减 1，直接退出。
- 4) 否则在全局的“Lock Hash”和“ProcLock Hash”表里查找此锁对应的 Lock 和 Proclock，调用 UnGrantLock 修改其信息。唤醒可以被唤醒的进程，并从“Local Hash”里移除该类型锁。

如果只需要释放一个本地锁，PostgreSQL 用函数 RemoveLocalLock 来完成。

另外，PostgreSQL 还提供了能一次释放当前进程所持有的指定锁方法的全部的锁。该操作在函数 LockReleaseAll 中定义，如果参数 AllLocks 为 FALSE，则释放除 SessionLock 之外所有的锁，否则释放所有的锁（包括 SessionLock）。

(5) RegularLock 的申请

该操作在函数 GrantLocalLock 中定义，如果当前进程已经获得所需的锁，那么只需要对当前该锁的引用计数加 1，即更新 LOCALLOCK 结构表示所需的锁已经被授予了。

(6) RegularLock 的注销

该操作为 RegularLock 的申请操作的逆过程，定义在函数 UnGrantLock 中。它取消对一个锁的分配，即更新 LOCK 和 PROCLOCK 数据结构以显示该锁不再由当前持有者持有或请求。如果在该锁上存在任何等待者则需要由 ProcLockWakeup 唤醒。

(7) 锁的冲突检测

该操作用来检测当前后端所请求的锁和已持有的锁是否冲突，它定义在函数 LockCheckConflicts 中。如果存在冲突，返回 STATUS_FOUND，否则返回 STATUS_OK。

一个进程所有锁之间并不冲突，即使是由不同事务 ID 所持有的（例如会话锁和事务锁并不冲突）。所以在决定请求的新锁是否和已持有锁之间是否冲突时，我们需要减去我们本身所持有的锁。

(8) 两阶段提交

两阶段提交中涉及以下三种关于锁的操作：

- 获取一个 Prepared 事务的锁。该操作定义在函数 lock_twophase_recover 中。该函数在数据库启动的时候被调用，这个重新获取锁的过程不会和其他事务有冲突。
- 分布式数据库预提交 COMMIT PREPARED 时释放该两阶段提交记录指示的锁。该函数定义在 lock_twophase_postcommit 中。
- 分布式数据库进行两阶段提交 ROLLBACK PREPARED 时释放锁操作，该操作定义在函数 lock_twophase_postabort 中。

(9) 锁的清理

该操作定义在函数 CleanUpLock 中，在释放锁之后执行，主要是清理 ProcLock 队列入口。

7.8 锁管理机制

上一节详细叙述了 PostgreSQL 中的锁的实现机制，包括底层 Spin 锁、LwLock 锁，以及最上层的 Regular 锁，那么这些锁是如何被系统其他模块调用的呢？在 PostgreSQL 系统中，结合其他模块的需要，针对不同的实体操作封装了相应的锁管理操作。锁的管理包括对不同粒度锁的操作，这些粒度包括表、内存页（目前只在索引中使用对内存页的加锁）、元组、事务 XID、事务 VXID 以及当前数据库中的一般对象。

7.8.1 表粒度的锁操作

表粒度的锁操作表示在数据库系统中加锁对象为一个表，其操作列举在表 7-6 中。

表 7-6 表粒度的锁操作

操作	实现函数	功能	说明
初始化锁	RelationInitLockInfo	初始化表描述符中的锁信息	如果要对一个表加锁，就是对一个 $\langle DBID, RELID \rangle$ 加锁，所以这里就是将这个二元组放到表结构的字段 (rd_lockInfo) 中。在创建任何表时都要调用该函数
加锁	LockRelation	对一个表加指定锁模式的锁	调用 RegularLock 模块中的 LockAcquire 函数完成。另外对已经打开的表获取一个额外的锁可以调用 ConditionalLockRelation 函数
解锁	UnlockRelation	释放对一个表所加的指定锁模式的锁	实际调用 RegularLock 模块中的 LockRelease 函数完成。与 LockRelation 成对使用
加会话锁	LockRelationIdForSession	在目标表上获取一个会话锁	会话锁是跨事务边界而存在的，所以应用范围更广
解会话锁	UnlockRelationForSession	释放对一个表所加的会话锁	与 LockRelationIdForSession 成对使用
加扩展锁	LockRelationForExtension	对一个表加指定模式的扩展锁	这个锁标签是用来实现表中新增页的互锁，需要这个锁是因为 BUFMGR/SMGR 中定义的 P_NEW 可能存在争用现象 (RaceCondition)。
解扩展锁	UnlockRelationForExtension	释放对一个表所加的扩展锁。	与 LockRelationForExtension 成对使用

7.8.2 页粒度的锁操作

页粒度的锁操作表示在数据库系统中加锁对象为一个页面，其主要操作如表 7-7 所示。

表 7-7 页粒度的锁操作

操作	实现函数	功能	说明
加锁	LockPage	获取一个页面级的锁	这个函数目前被一些索引访问方法用于对索引页面加锁。另外还有一个带条件的加锁函数 ConditionalLockPage，它只有在不阻塞当前进程的前提下才对一个页面加指定模式的锁，当且仅当获得锁时返回 TRUE
解锁	UnlockPage	释放对一个页面所加的锁	与 LockPage 成对使用

7.8.3 元组粒度的锁操作

元组粒度的锁表示在数据库系统中加锁对象为一个元组，其主要操作如表 7-8 所示。

表 7-8 元组粒度的锁操作

操作	实现函数	功能	说明
加锁	LockTuple	对一个元组加指定锁模式的锁	使用这种锁的时候要谨慎，因为不可能为每一个元组在共享的存储区提供一个单独的锁。另外，还提供一个带条件的加锁函数 ConditionalLockTuple，条件是在不阻塞进程的前提下，对一个元组加指定模式的锁，当且仅当获得锁时返回 TRUE
解锁	UnlockTuple	用于释放对一个元组的所加的锁	与函数 LockTuple 成对使用

7.8.4 事务粒度的锁操作

事务粒度的锁操作表示在数据库系统中加锁对象为一个事务，其操作如表 7-9 所示。

表 7-9 事务粒度的锁操作

操作	实现函数	功能	说明
加锁	XactLockTableInsert	获取一个锁以表明事务正在运行中	在事务或其子事务获取一个 XID 的时候会自动调用本函数，并且获取的锁会一直使用到事务的结束
解锁	XactLockTableDelete	释放指定事务上的锁	这里只能释放子事务所持有的锁，主事务锁在事务结束的时候被释放

7.8.5 一般对象的锁操作

数据库一般对象的锁操作表示在数据库系统中加锁对象为一般对象，其操作如表 7-10 所示。

表 7-10 一般对象的锁操作

操作	实现函数	功能	说明
加锁	LockDatabaseObject	实现对数据库一般对象的加锁	用于获得当前数据库的一般对象的一个锁，这种锁不能用在共享的对象上
解锁	UnlockDatabaseObject	实现对数据库一般对象的解锁	与函数 LockDatabaseObject 成对使用

7.9 死锁处理机制

在 PostgreSQL 中，事务可以按照任意顺序申请锁，所以必须考虑死锁的处理机制。PostgreSQL 对于死锁的预防分为两步：

1) 当进程请求加锁时，如果失败，会进入等待队列。如果在队列中已经存在一些进程要求本进程中已经持有的锁，那么为了尽量避免死锁，可以简单地把本进程插入到它们的前面。

2) 当一个锁被释放时，将尝试唤醒等待队列中的进程。如果其中的某个进程的要求与排在它前面但由于某些原因不能被唤醒的进程冲突，这个进程将不被唤醒。这么做可以保证互相冲突的加锁请求按照到达的先后被处理。

当然，仅仅依靠上述死锁预防方法还不能够完全消除死锁。PostgreSQL 还提供一套死锁检测机制。死锁检测的触发条件如图 7-4 所示。

PostgreSQL 使用等待图 (WFG) 来进行死锁检测。WFG 是一个有向图，图中的顶点表示申请加锁的进程，而图中的有向边表示依赖关系 (等待关系)，如图 7-5 所示。

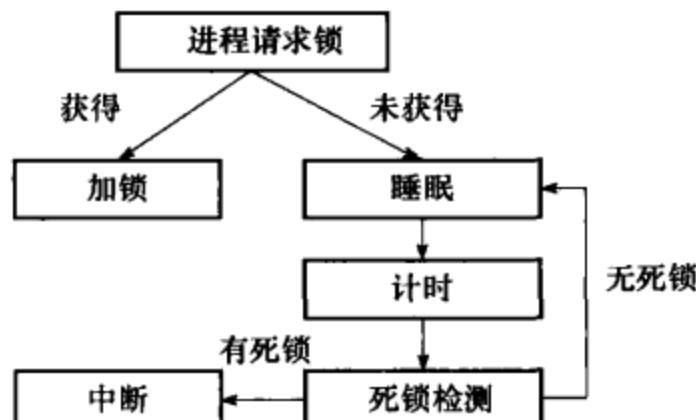


图 7-4 死锁检测的触发条件

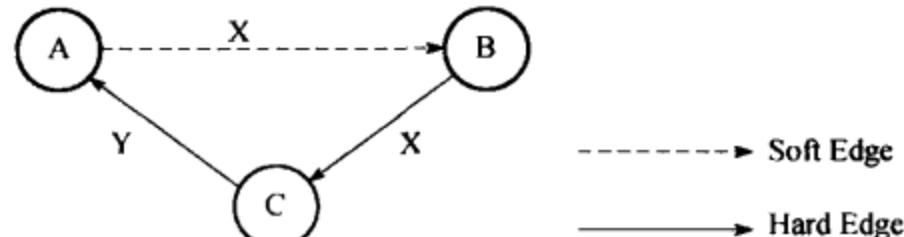


图 7-5 死锁检测等待图

如果进程 A 在等待进程 B (进程 A 申请的锁 a 与进程 B 持有的锁 b 冲突，或者进程 A 申请的锁 a 与进程 B 申请的锁 b 冲突并且进程 B 比进程 A 申请锁的时间更早)，那么从顶点 A 到顶点 B 就有一条有向边。系统出现死锁当且仅当 WFG 中出现环。

为了进一步讨论的需要，我们先介绍两个定义：Soft Edge 和 Hard Edge。

- 进程 A 和进程 B 都在某个锁的等待队列中，进程 A 在进程 B 的后面。如果进程 A 和进程 B 的加锁要求冲突，即进程 A 在等待进程 B，这时候顶点 A 到顶点 B 会有一条有向边。我们称这样的有向边为 Soft Edge。
- 如果进程 A 的加锁要求和进程 B 已持有的锁冲突，这时候顶点 A 到顶点 B 也有一条有向边。我们称这样的有向边为 Hard Edge。

如果 WFG 中的环是有 Soft Edge 的环，可以通过拓扑排序对队列进行重排，尝试消除死锁。

PostgreSQL 死锁检测与消除算法如下：

- 从每一个顶点开始出发，沿 WFG 中的有向边走，看能不能回到此顶点。如果找到一条路径满足此条件，则说明出现死锁。
- 如果此路径中没有 Soft Edge，直接终止这个事务。
- 否则记录所有出现的 Soft Edge。对于这个集合，递归地枚举它的所有子集，尝试进行调整。
- 对于上述每个子集，使用拓扑排序来排出一个调整方案，并逐一加以测试（一旦找到一个可行的调整方案，这个测试过程就结束了，PostgreSQL 并没有要求这个调整方案一定是最优的）。如果找不到一个可以消除死锁的方案，则死锁消除失败，然后终止这个事务。

假设有两个资源 X 和 Y。进程 C 持有 X 的共享锁，进程 B 请求 X 的排他锁，先进入等待队列；进程 A 请求 X 的共享锁，在进程 B 之后进入等待队列。进程 A 持有资源 Y 的排他锁，进程 C 请求 Y 的排他锁，进入等待队列，其 WFG 如图 7-5 所示。

PostgreSQL 对进程 A 的检测过程如下：

- 调用 DeadLockCheck 试图检测和消除死锁，DeadLockCheck 将调用 DeadLockCheckRecurse，递归进行死锁检测。
- DeadLockCheckRecurse 调用 TestConfiguration 测试当前的队列状态会不会发生死锁。TestConfiguration 调用 ExpandConstraints 进行约束性检查，若不满足，则死锁。其中可对 Soft Edge 进行调整，并调用 TopoSort 判断调整后的序列是否合法。
- 最后调用 FindLockCycle 判断是否出现环，如果出现并且不能调整，则死锁。

7.9.1 死锁处理相关数据结构

死锁处理方面的主要数据结构是 WFG 的定义，其中最重要的就是有向边（数据结构 7.13）。

数据结构 7.13 EDGE

```
typedef struct
{
    PGPROC      *waiter;      //等待的后端进程列表
    PGPROC      *blocker;     //阻塞的后端进程列表
    int          pred;        //拓扑排序时使用
    int          link;        //拓扑排序时使用
}EDGE;
```

在死锁消除时需要对一个锁的等待队列（数据结构 7.14）重新排序。

数据结构 7.14 WAIT_ORDER

```
typedef struct
{
    LOCK         *lock;        //这个等待队列所属的锁
    PGPROC      **procs;      //重新排序后的 PGPROC 数组
    int          nProcs;
}WAIT_ORDER;
```

已检测到的死锁环中每条边的信息保存在结构 DEADLOCK_INFO（数据结构 7.15）中。

数据结构 7.15 DEADLOCK_INFO

```
typedef struct
{
    LOCKTAG      locktag;          // 等待的加锁对象 ID
    LOCKMODE     lockmode;         // 正在等待的锁类型
    int          pid;             // 被阻塞后端的 PID
} DEADLOCK_INFO;
```

最后用一个枚举类型 DeadLockState（数据结构 7.16）表示 DeadLockCheck 函数的执行结果。

数据结构 7.16 DeadLockState

```
typedef enum
{
    DS_NOT_YET_CHECKED,           // 尚未执行死锁检测
    DS_NO_DEADLOCK,              // 没有检测到死锁
    DS_SOFT_DEADLOCK,            // 通过等待队列的重排序消除了死锁
    DS_HARD_DEADLOCK,             // 无法消除的死锁
    DS_BLOCKED_BY_AUTOVACUUM     // 没有死锁
} DeadLockState;
```

7.9.2 死锁处理相关操作

在 PostgreSQL 中，当进程不能获得锁进入等待队列时，就会触发死锁检测的操作。如果发现了“死锁”，进而尝试进行死锁的解除操作，死锁解除采用枚举的方法去尝试调整等待队列中进程的等待先后拓扑顺序，试图找到一种打破进程之间循环等待的状态。

(1) 死锁处理相关数据结构的初始化

每个服务进程启动的时候都需要初始化死锁检测器，即给死锁检测时需要用到的数据结构分配内存空间。

这个初始化过程由函数 InitDeadLockChecking 完成。每个服务进程都需要进行死锁检测器的初始化过程，因为服务进程不可能从 Postmaster 那里继承其工作空间的内容。另一方面，之所以在每个服务进程启动的时候就马上进行这个初始化过程，有两个原因：一是如果等到死锁检测器被启动的时候才进行初始化，可能已经没有空余的内存可以分配；二是死锁检测一般在 Signal Handler 中运行，而这部分的内存空间如果分配给死锁检测器的话不够安全。

(2) 死锁检测入口函数

死锁检测入口函数为 DeadLockCheck，用于在一个给定进程中检查死锁。一旦发现了死锁，该函数会重新排序锁的等待队列以消除这个死锁。如果这个死锁无法消除，则返回 DS_HEAD_DEADLOCK。该返回值提醒调用者中断正在处理的事务。

在执行该函数之前，调用者必须已经获取整个锁表结构的 LWLock。如果该函数执行失败，死锁详细信息将保存在 DeadlockDetails 中。

记录在 DeadlockDetails 中的死锁详细信息将由 DeadLockReport 函数打印出来。之所以将死锁报告和死锁检测函数分开实现，有两个原因：其一，死锁检测函数执行时需要持有整个锁表结构上的 LWLock，而在执行死锁报告函数的时候并不需要这些锁。如果在执行 DeadLockReport 的时候还持有这些 LWLock，会影响系统的并发性能。其二，DeadLockCheck 函数执行时处于 Signal Handler 关键内存中。

(3) 死锁检测调整函数

死锁调整定义在函数 DeadLockCheckRecurse 中，它递归地调整等待队列的排序，寻找无死锁的队列。其中最主要的操作就是检测调整后的等待队列是否满足“约束”条件。

如果等待队列中进程 A 一定得在进程 B 之前才能消除死锁，我们称这是一个约束。如果对等待队列重新排序后，进程 B 在进程 A 前面，那么称这个重新排序后的等待队列不满足约束条件。

如果等待队列的任何拓扑序列都无法消除死锁，则该函数返回 TRUE。如果发现可用的新等待序列，则保存这个新序列，并返回 FALSE。

(4) 等待图的环检测函数

该函数为 FindLockCycle，它调用 FindLockCycleRecurse 函数检测等待图中是否有死锁环。如果发现了死锁环，则通过调用参数返回一个 Soft Edge 的链表。

7.10 多版本并发控制

PostgreSQL 为开发者提供了丰富的管理数据并发访问的工具。在内部，PostgreSQL 利用多版本并发控制（MVCC，Multi-Version Concurrency Control）来维护数据的一致性。这就意味着当检索数据时，每个事务看到的只是一段时间之前的数据快照，而不是数据的当前状态。这样，如果对每个数据库会话进行事务隔离，就可以避免一个事务看到其他并发事务的更新而导致不一致的数据。

使用多版本并发控制的主要优点是对检索（读取）数据的锁请求与写数据的锁请求并不冲突，所以读不会阻塞写，而写也从不阻塞读。这就极大地提高了并发处理能力。

在讨论 MVCC 机制之前，先看一个多版本并发控制的例子，如表 7-11 所示。

表 7-11 多版本并发控制示例

事务 T1	事务 T2	表 A 的变化	事务 T2 的查询结果
BEGIN	SELECT * FROM A;	A	A
UPDATE		A→AA	
	SELECT * FROM A;		A
COMMIT			
	SELECT * FROM A;		AA

7.10.1 MVCC 相关数据结构

在 PostgreSQL 系统中，更新数据时并不是用新值覆盖旧值，而是在表中另开辟一片空间来存放新的元组，让新值与旧值同时存放在数据库中，通过设置一些参数，让系统识别它们。

PostgreSQL 的数据以元组的形式存放在表中，同时存储该元组相关的描述信息，这些信息存储在 HeapTupleFields（数据结构 7.17）中。

数据结构 7.17 HeapTupleFields

```

typedef struct HeapTupleFields
{
    TransactionId t_xmin;           // 创建此 tuple 的 XID
    TransactionId t_xmax;           // 删除此 Tuple 的 XID
    union
    {
        CommandId t_cid;           /* 创建或者删除 Tuple 的 command ID(即 Cmin 和 Cmax), 也可能两
                                      者都保存 */
        TransactionId t_xvac;       // 清理操作的事务 ID
    }t_field3;
}HeapTupleFields;

```

从 HeapTupleFields 结构中可以看出，Xmin、Cmin、Xmax、Cmax 和 Xvac 这五个字段被保存在三个物理空间中，Xmin 和 Xmax 分别有独立的字段存储，而 Cmin、Cmax 和 Xvac 共享一个字段。因为 Cmin 和 Cmax 都只在创建或删除事务的运行期间内有效，而一般情况下，在同一个事务中创建并删除同一个元组的概率比较低，Cmin 和 Cmax 同时存在有效值的概率也就很小，所以把它们放在同一个物理空间中。如果一个事务确实创建并删除了同一个元组，则我们使用一个 Combo Command ID 来保存 Cmin 和 Cmax。Combo Command ID 保存的是一个到实际 Cmin 和 Cmax 的映射。另一方面，Xvac 字段只由 VACUUM FULL 命令设置，此时并不需要 Cmin 和 Cmax。

在 PostgreSQL 中，Combo Command ID 是 32 位的，所以可以映射到 2^{32} 个 {Cmin, Cmax} 组合。在最坏的情况下，每个命令把它之前所有命令创建的元组都删除掉，假设这些命令总共有 N 个，那么，系统需要的 Combo Command ID 有 $N * (N + 1)/2$ 个。在这种最坏的情况下，Combo Command ID 最多可以保存 92682 个命令。在实际应用中，用户在达到这个命令限制之前，就已经耗尽了内存和硬盘空间，所以，把 Combo Command ID 设置成 32 位是足够的。

除了存储元组的事务、命令控制信息外，还需要存储元组的相关控制信息。这些信息存储在元组的头部 HeapTupleHeaderData（数据结构 7.18）中。

数据结构 7.18 HeapTupleHeaderData

```

typedef struct HeapTupleHeaderData
{
    union
    {
        HeapTupleFields t_heap;
        DatumTupleFields t_datum;
    }t_choice;
    ItemPointerData t_ctid;           // 本元组或者新元组的当前 TID
    uint16          t_infomask2;      // 标志位(属性、标志位的数量)
    uint16          t_infomask;       // 标志位(元组的事务信息)
    uint8           t_hoff;           // 头部长度
    bits8          t_bits[1];        // 填充位,起标记作用
}HeapTupleHeaderData;

```

一个新元组被保存在磁盘中的时候，其 t_ctid 就被初始化为它自己的实际存储位置。如果这个元组被更新了，该元组的 t_ctid 字段就指向更新后的新元组。由此可以看出，一个元组是最新版本的，当且仅当它的 xmax 空或者它的 t_ctid 指向它自己。如果要找到某个元组的最新版本，只需遍历由 t_ctid 构成的链表即可。

在 HeapTupleHeaderData 中，还有一个重要的字段（t_infomask）用来表示当前元组的事务信息（数据结构 7.19）。

数据结构 7.19 t_infomask

#define HEAP_HASNULL	0x0001
#define HEAP_HASVARWIDTH	0x0002
#define HEAP_HASEXTERNAL	0x0004
#define HEAP_HASOID	0x0008
/* 以上四种标志不涉及元组可见性判断，用于表示元组本身的基本信息 */	
#define HEAP_COMBOCID	0x0020 /* 复合命令 */
#define HEAP_XMAX_EXCL_LOCK	0x0040 /* xmax 持有排他锁 */
#define HEAP_XMAX_SHARED_LOCK	0x0080 /* xmax 持有共享锁 */
/* 如果 LOCK 位被设置了，表示 xmax 只是锁住了该元组，没有完成删除 */	
#define HEAP_IS_LOCKED	(HEAP_XMAX_EXCL_LOCK \
	HEAP_XMAX_SHARED_LOCK)
#define HEAP_XMIN_COMMITTED	0x0100 /* t_xmin 已提交 */
#define HEAP_XMIN_INVALID	0x0200 /* t_xmin 无效/中断 */
#define HEAP_XMAX_COMMITTED	0x0400 /* t_xmax 已提交 */
#define HEAP_XMAX_INVALID	0x0800 /* t_xmax 无效/中断 */
#define HEAP_XMAX_IS_MULTI	0x1000 /* 组合事务 */
#define HEAP_UPDATED	0x2000 /* 更新后的新元组 */

MVCC 的基本原理如图 7-6 所示，有两个并发执行的事务 T1、T2，T1 将元组 C 更新为 C'，但 T1 没有提交，此时假如事务 T2 要对该元组进行查询，它会通过 C 和 C' 中的头信息中的 Xmin 和 Xmax 及 t_infomask 来判断出 C 为有效版本，而 C' 为无效版本。

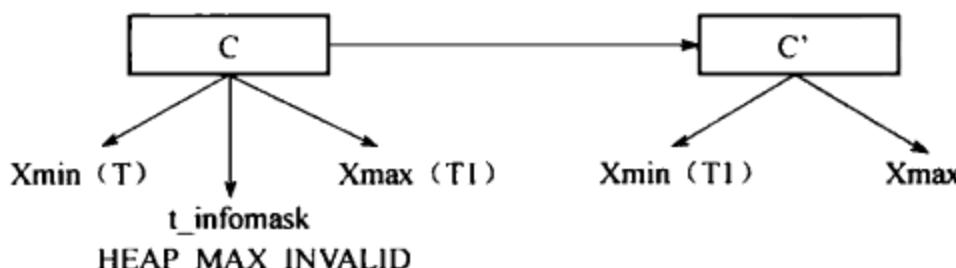


图 7-6 MVCC 基本原理

7.10.2 MVCC 相关操作

本节主要介绍 MVCC 操作所需要调用的相关函数，核心功能是用来判断元组的状态，包括元组的有效性、可见性、可更新性等。

(1) 判断元组对于自身信息是否有效

该操作定义在函数 `HeapTupleSatisfiesSelf` 中。如果函数返回值为 TRUE，则该元组是有效的。判断条件会考虑三个方面的因素：所有已提交的事务、当前事务前的所有命令、当前命令前的所有操作。判断过程如图 7-7 所示（图中其他判断条件出口则为不可见）。

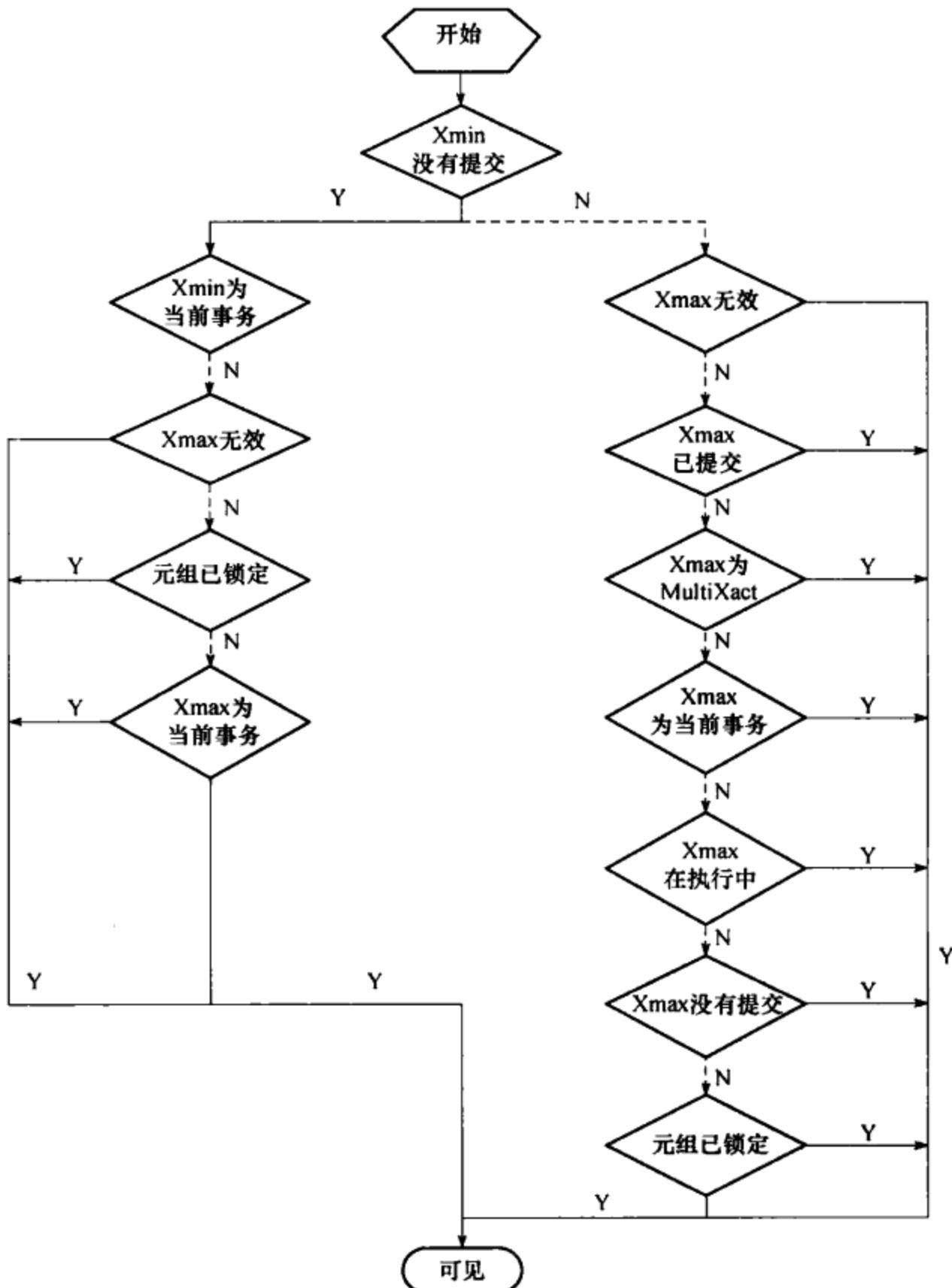


图 7-7 `HeapTupleSatisfiesSelf` 元组可见性判断图（局部）

(2) 判断元组对当前时刻是否有效

该操作定义在函数 `HeapTupleSatisfiesNow` 中。如果函数返回值为 TRUE，则该元组是有效的。该

函数和 HeapTupleSatisfiesSelf 类似，只不过它只考虑已经提交的事务和本事务当前命令之前的命令对元组造成的影响。

HeapTupleSatisfiesNow 函数根据图 7-8（图中其他判断条件出口则为不可见）所示逻辑进行元组有效性判断。

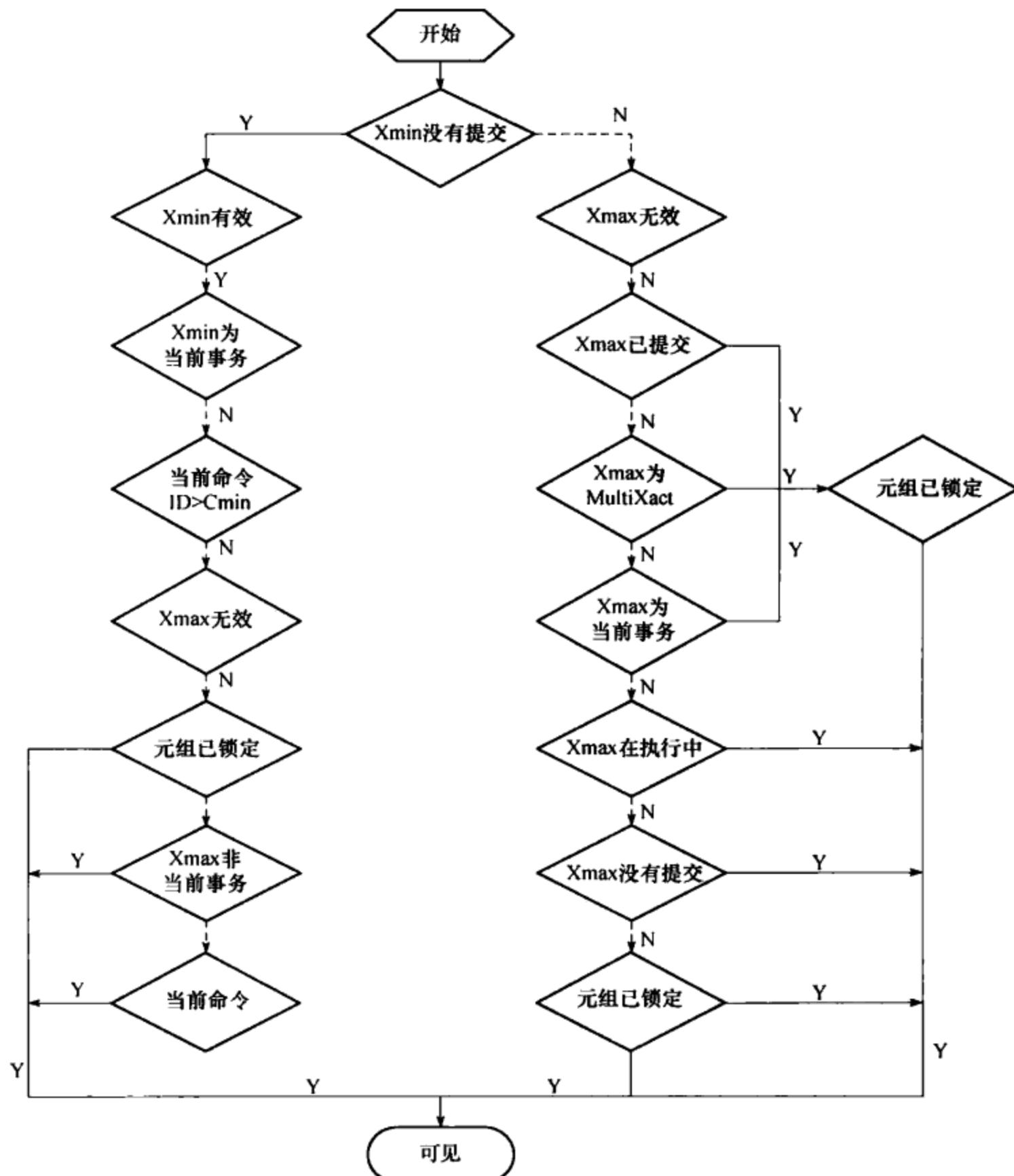


图 7-8 判断元组当前时刻是否有效流程图（局部）

(3) 判断当前元组是否已脏

如果该元组为脏，则需要立即从磁盘重新读取该元组，获得其有效版本。该操作定义在函数

HeapTupleSatisfiesDirty 中，它与元组自身有效性判断函数 HeapTupleSatisfiesSelf 所包含的事务的影响条件相同，即 HeapTupleSatisfiesSelf 返回 TRUE 的条件下 HeapTupleSatisfiesDirty 也返回 TRUE。

(4) 判断元组对 MVCC 某一版本是否有效

该操作定义在函数 HeapTupleSatisfiesMVCC 中，与 HeapTupleSatisfiesNow 判断流程基本相同，唯一不同的是当 Xmin 已经提交时，如果此时 Xid 仍然在 MVCC 快照中，那么此时该元组不可见。

(5) 判断元组是否可更新

该操作定义在函数 HeapTupleSatisfiesUpdate 中，与 HeapTupleSatisfiesNow 函数流程类似，但是它返回的信息更详细，返回结果定义在一个枚举类型 HTSU_Result (数据结构 7.20) 中。

其各个返回值的意义如表 7-12 所示。

表 7-12 元组是否可更新的返回结果表

返回值	代表的含义
HeapTupleMayBeUpdated	该元组可见，可以进行更新
HeapTupleInvisible	开始扫描时该元组不存在（可能是因为本事务后面的命令创建了这个元组）
HeapTupleSelfUpdated	该元组在当前扫描开始后被本事务更新了
HeapTupleUpdated	该元组被一个已提交的事务更新了
HeapTupleBeingUpdated	该元组正在被一个运行中的事务更新（非本事务）

该操作返回五种不同的结果，下面我们以返回 HeapTupleBeingUpdated 为例，说明其判断执行的过程（如图 7-9 所示）。

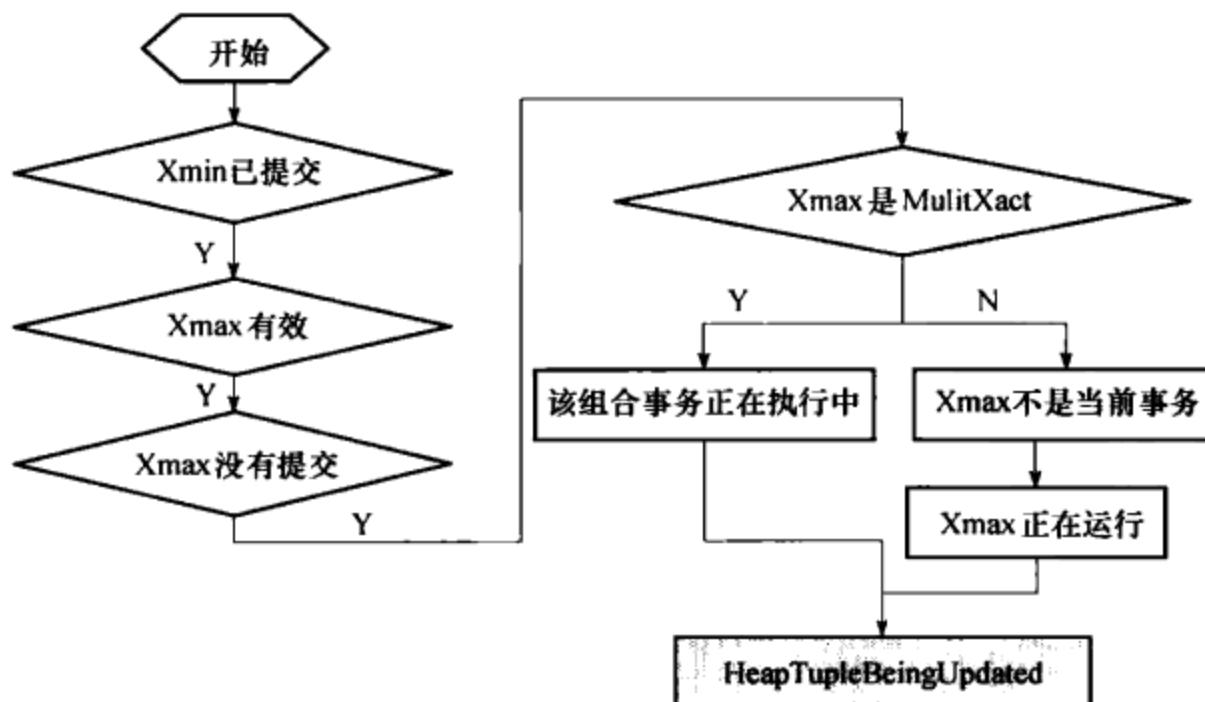


图 7-9 返回 HeapTupleBeingUpdated 流程图

数据结构 7.20 HTSU_Result

```

typedef enum
{
    HeapTupleMayBeUpdated,
    HeapTupleInvisible,
    HeapTupleSelfUpdated,
    HeapTupleUpdated,
    HeapTupleBeingUpdated
} HTSU_Result;
  
```

假设有事务 T 进行下列操作：

```
BEGIN
  INSERT INTO tableA tupleA; (Command ID = 1)
  UPDATE tableA tupleA; (Command ID = 2)
  .....
END
```

INSERT 语句执行完毕后，元组 A 的信息如下：

tupleA	Cmin = T	Cmax = NULL	Cmin = 1	Cmax = NULL
--------	----------	-------------	----------	-------------

处理 UPDATE 语句前，先检查数据版本的合法性，由于元组 A 中 Xmin 为当前事务，Xmax 无效，所以该元组对于更新语句可见。找到合法版本以后，再通过 HeapTupleSatisfiesUpdate 判断其是否可更新，返回结果 HeapTupleMayBeUpdated，因此可以对这个元组执行 UPDATE 操作。

7.10.3 MVCC 与快照

快照（SnapShot）记录了数据库当前某个时刻的活跃事务列表。通过快照，可以确定某个元组的版本对于当前快照是否可见。

快照的定义在数据结构 SnapshotData 中（如数据结构 7.21 所示）。

数据结构 7.21 SnapshotData

```
typedef struct SnapshotData
{
    SnapshotSatisfiesFunc satisfies;          // 函数指针
    TransactionId      xmin;                  // 所有 XID < xmin 对当前快照可见
    TransactionId      xmax;                  // 所有 XID >= xmax 对当前快照可见
    uint32              xcnt;                 // 当前活跃事务的长度
    TransactionId      *xip;                  // 当前活跃事务的链表
    int32               subxcnt;                // 当前活跃子事务个数
    TransactionId      *subxip;                // 当前活跃子事务的链表
    CommandId           curcid;                // 当前命令的序号
    uint32              active_count;          // 在活跃快照链表里的引用计数
    uint32              regd_count;            // 在已注册的快照链表里的引用计数
    bool                copied;                // 是否为共享内存中的全局快照
} SnapshotData;
```

SnapshotData 的结构对于事务的可见性判断提供了统一操作接口，其中 satisfies 函数指针指向 7.10.2 节中介绍的 MVCC 相关判断函数，通过该接口可以判断指定的事务 ID 对于当前快照是否可见。

在 PostgreSQL 系统中，默认有 7 种形式的快照，分别是：

```
SnapshotData SnapshotNowData = {HeapTupleSatisfiesNow};
SnapshotData SnapshotSelfData = {HeapTupleSatisfiesSelf};
SnapshotData SnapshotAnyData = {HeapTupleSatisfiesAny};
SnapshotData SnapshotToastData = {HeapTupleSatisfiesToast};
```

```

SnapshotData CurrentSnapshotData = {HeapTupleSatisfiesMVCC};
SnapshotData SecondarySnapshotData = {HeapTupleSatisfiesMVCC};
#define InitDirtySnapshot(snapshotdata) \
((snapshotdata).satisfies = HeapTupleSatisfiesDirty)

```

其中，不同的快照用于不同形式的可见性判断，需要注意的是如果当前事务的隔离级别是可串行化，那么只存在 CurrentSnapshotData，不存在 SecondarySnapshotData，而在读已提交的隔离级别，这两种快照都存在。

7.11 日志管理

日志是数据库系统必不可少的一部分，它以一种安全的方式记录数据库变更的历史。当系统出现故障后，数据库系统通过使用日志来重建对数据库所做更新的过程，以恢复数据库到一致状态，从而保证数据库的一致性和完整性。

PostgreSQL 采用的日志主要有 XLOG 和 CLOG，即事务日志和事务提交日志。XLOG 是一般日志记录，即通常意义上所认识的日志记录，它记录了事务对数据更新的过程和事务的最终状态。CLOG 在一般的数据库教材上是没有提及的，其实 CLOG 是 XLOG 的一种辅助形式，记录了事务的最终状态。因为每一条 XLOG 日志记录相对较大，如果需要通过日志判断一个事务的状态，那么使用 CLOG 比使用 XLOG 要高效得多。同时 CLOG 占用的空间也非常有限。此外，为支持嵌套事务，PostgreSQL 还引入了 SUBTRANS 日志记录——即子事务日志，记录每个事务的父事务的事务 ID，这样通过一个事务可以递归查找到其父事务，但是并不能通过一个事务查找到其子事务。同时，为了支持多版本并发控制，PostgreSQL 引入了组合事务 ID（MultiXactID），记录事务的组合关系，并维护从众多事务 ID 到 MultiXactID 的映射关系。

在 PostgreSQL 中，日志通过日志文件来存放。如果每个日志记录在创建时都被立即写到磁盘中，那么将增加大量 I/O 开销。因为通常向磁盘的写入是以块为单位进行的，而在大多数情况下，一个日志记录比一个块要小得多，为了降低写入日志带来的 I/O 开销，数据库系统在实现时往往设置了日志缓冲区，即先将日志记录写到主存中的日志缓冲区中，当日志缓冲区满了以后以块为单位向磁盘写出。

在数据库系统中，日志缓冲区通过日志管理器来管理。数据库系统并不直接操作磁盘日志文件，一般是通过日志缓冲区来使用日志文件，而缓冲区和磁盘之间的交互、同步则由日志管理器来完成。PostgreSQL 有四种日志管理器，即 XLOG（事务日志）、CLOG（事务提交日志）、SUBTRANS（子事务日志）以及 MULTIXACT（组合事务日志）。

PostgreSQL 通过同一种缓冲区来实现对 CLOG 日志、SUBTRANS 日志以及 MULTIXACT 日志的管理，即 SLRU 缓冲池——采用简单 LRU 算法作为页面置换算法的缓冲池。

如果要获取一个事务的状态，并不是直接通过日志管理器操作日志缓冲池，而是要通过事务日志接口例程进行操作。在事务日志接口例程中，一方面定义了可使用事务 ID 号的范围、可使用对象 ID 号的范围，另一方面提供了可设置和获取事务状态信息的接口例程。此外，为减少 I/O 次数，也使用了缓冲区机制，即建立一个单独的缓冲区，用于缓存最近获取的事务 ID 及其状态。而且为方便事务管理，建立了一个共享变量缓冲数据结构，存储下一可分配事务 ID、对象 ID 以及已分配对象数。

由于使用日志的资源有很多种，PostgreSQL 为对日志进行分类，使用了资源管理器（不是第 3

章中介绍的资源跟踪器) 的概念。资源管理器主要用于在日志系统中把各种需要记录日志的数据分类, 通过在日志中标识资源管理器号, 使系统在恢复或者读取日志记录时, 能够很方便地知道该日志记录的源数据属于哪一类, 从而通过资源管理器的方法可以准确地选择对应的方法。

7.11.1 SLRU 缓冲池

在介绍 SLRU 缓冲池之前, 先说明一下 PostgreSQL 对 CLOG 和 SUBTRANS 日志的物理存储组织方法。CLOG 和 SUBTRANS 日志在磁盘中由一个个小的物理文件组成, 每一个物理日志文件定义为一个段, 一个段由 32 个磁盘页面组成, 每一个磁盘页面的大小为 8KB。每一个段文件以段号来命名, 通过一个段号就可以找到对应的日志文件。进一步, 只要有日志的页面号, 就可以找到该页面所在的段文件及该页在段文件中的偏移, 从而将该日志页面调入到内存缓冲区中, 供日志管理器使用。也就是说, 通过二元组 $\langle \text{Segmentno}, \text{Pageno} \rangle$ 就可以定位日志页在哪个段文件中以及在该文件中的偏移位置。

SLRU 缓冲池由 8 个缓冲区组成。每一个缓冲区为一个页面, 对应一个磁盘块, 大小为 8KB。页面调度算法采用 SLRU 算法, 即简单的最近最少使用算法。因为缓冲区只有 8 个页面, 所以实现搜索时, 只需要使用简单的线性搜索算法。

写日志记录时, 先写入到对应日志页面所在的缓冲区, 再进一步写入到磁盘上的日志段文件。而读日志文件首先将要读的日志记录所在页面读入到缓冲池的某一个页面, 再在日志缓冲区中读具体的日志记录。

1. 缓冲池的并发控制

PostgreSQL 使用轻量级锁实现对 LRU 缓冲区的并发控制。其中使用一个控制锁 (用轻量级锁实现) 来保护整个数据缓冲区, 而对于每一个缓冲区, 当进行 I/O 同步时, 还使用一个 LWLock 来保护。一个进程在读或写一个页面缓冲区时并不持有控制锁, 而只是对当前所读入或写出的页面所在缓冲区加一个缓冲区锁。当改变一个缓冲区所对应页面或其状态时, 需要持有控制锁。但当一个写进程在持有一个缓冲区锁并将缓冲区状态由 DIRTY 改为 WRITE_IN_PROGRESS 时, 此时改变状态并不需要持有控制锁。

当缓冲区状态不是 EMPTY 或 CLEAN 时, 可能有进程正在对缓冲区进行 I/O 操作, 所以此时页面号码不能改变; 此时唯一可以进行状态改变的是在页脏了以后从状态 WRITE_IN_PROGRESS 转换为状态 DIRTY。

当一个缓冲区的任何状态转换动作涉及潜在的 I/O 操作时, 进程需要同时持有控制锁和缓冲区锁。请求锁的过程是先获取一个缓冲区锁再获取控制锁, 即不要在获取控制锁后再等待缓冲区锁。

如果一个进程试图读取一个页面, 则将其缓冲区状态标识为 READ_IN_PROGRESS, 然后释放控制锁, 再请求一个缓冲区锁, 在继续执行前重新检查缓冲区的状态 (目的在于防止其他人试图将相同的页面读取到另外一个缓冲区中)。

2. 缓冲池相关数据结构

缓冲池除了共享的数据缓冲区外, 还有一个非共享的控制结构 (数据结构 7.22), 提供缓冲池的控制信息。

数据结构 7.22 SlruCtlData

```

typedef struct SlruCtlData
{
    SlruShared      shared;           // 共享内存中共享的数据缓冲区指针
    bool            do_fsync;        /* 判断在写文件时是否进行同步写。默认对于 pg_clog
                                    /* 为真，对于 pg_subtrans 为假 */
    bool            (*PagePrecedes)(int, int); // 比较两个页面前后关系的方法接口
    char             Dir[64];         // 物理文件在磁盘的存储目录
}SlruCtlData;
typedef SlruCtlData *SlruCtl;

```

SLRU 缓冲池使用 SlruFlushData（数据结构 7.23）记录打开的物理磁盘文件。当进行刷写操作时，将对多个文件的更新数据一次性刷写到磁盘。

数据结构 7.23 SlruFlushData

```

typedef struct SlruFlushData
{
    int      num_files;           // 实际打开的文件个数
    int      fd[NUM_SLRU_BUFFERS]; // 缓冲区对应的已打开文件的 FD
    int      segno[NUM_SLRU_BUFFERS]; // 缓冲区对应页面所在的段号
}SlruFlushData;

```

3. 缓冲池主要操作

缓冲池主要的主要操作包括缓冲池的初始化、缓冲池页面的选择、页面的初始化、缓冲池页面的换入换出、缓冲池页面的删除等。

(1) 缓冲池的初始化

缓冲池的初始化定义在函数 SimpleLruInit 中，主要功能是在共享内存中初始化 SLRU 缓冲池，并分配内存空间供数据缓冲区使用。

在共享内存中记录了已存在的 SLRU 缓冲池，并提供了根据其名称找到该缓冲池共享数据缓冲区地址的索引，如 CLOG，其名称为“CLOG Ctl”；而对于 SUBTRANS，其名称为“SUBTRANS Ctl”。如果所需类型的缓冲池已经存在，那么直接使用该缓冲池，并不需要重新建立；如果不存在，那么在共享内存中申请空间并建立，其后需要在共享内存中用其名称注册该缓冲池。

其主要流程操作如下：

- 1) 先调用 ShmemInitStruct 函数在共享内存中查找指定名称的 SLRU 缓冲池。
- 2) 如果当前处于 Postmaster 运行之下，显然所需要的 SLRU 缓冲池还不存在，需要在共享内存中初始化该 SLRU 缓冲池，初始化过程包括申请 8 个缓冲区大小的空间，并设置缓冲池的数据缓冲区结构，设置每一个缓冲区状态为“空”，其初始化 LRU 值为 1，为每一个缓冲区申请一个轻量级锁，将轻量级锁作为缓冲池的控制锁。

- 3) 否则，在共享内存中一定能找到该名称的 SLRU 缓冲池。
- 4) 然后设置写操作是否进行同步。默认情况下为“真”，即需要进行写同步。最后，根据“Subdir”在缓冲池控制结构中设置数据文件的存储路径。

(2) 缓冲区的选择

该操作定义在函数 SlruSelectLRUPage 中，主要用于获取一个缓冲区。当需要一个缓冲区用于读入新页面时，需要调用这个函数获得一个可用的缓冲区，当需要创建并初始化一个缓冲区时，也需要调用本函数。

换入换出策略（LRU）就体现在该函数中，即每一个读入缓冲区的页面都有一个 LRU 值，其初始值为 0。当需要替换一个页面时，选取用于替换的页面的准则就是该页面的 LRU 值最大，即该页面最近最少使用。此外，由于对于日志记录，都是不断往日志文件中写入记录，所以对于当前日志文件中的页面号最大的页面，如果该页面在缓冲区中，那么并不把它替换出去，即不参与 LRU 替换，原因是该页面随时都需要被写入新数据，为提高效率，所以并不替换该页面。

其执行流程如下所示：

- 1) 首先逐个比较缓冲池中的缓冲区，如果所需要的页面已经在某一个缓冲区中，那么直接返回该缓冲区即可。
- 2) 如果能找到一个状态为“空”的缓冲区，那么直接返回该缓冲区以供使用。
- 3) 如果依然没有一个可用的缓冲区，那么在所有缓冲区中找到一个 LRU 值最大的缓冲区（需要注意的是，当前页面号最大的页面并不参与到这个 LRU 选择过程中）。
- 4) 然后判断找到的缓冲区的状态，如果为“干净”，那么直接返回该缓冲区；如果该缓冲区状态为“正在读入”，那么调用 SimpleLruReadPage 函数继续完成读入动作；如果缓冲区状态为“脏”或“正在写入”，那么调用 SimpleLruWritePage 函数继续完成写入动作。
- 5) 最后返回第一步，继续完成选择缓冲区的任务，直至选出一个可供重用的缓冲区为止。

(3) 缓冲池页面的初始化

该操作定义在函数 SimpleLruZeroPage 中，主要功能是在缓冲池中选择一个缓冲区供指定页面使用，将该页面初始化为全 0，并设置该页面标识值为当前最大值。在扩展日志文件时，该操作可以创建一个新的日志页面，方便后继的日志记录的写入，并最终写到磁盘。

主要流程如下所示：

- 1) 调用 SlruSelectLRUPage 函数以获取一个可用的缓冲区。
- 2) 然后将该缓冲区页面号设置为指定页面号，并设置其状态为“脏”，同时设置其为最近使用的页面，即该页面的 LRU 值初始化为 0，而其他页面的 LRU 值增一。
- 3) 最后初始化该缓冲区为全 0，并设置当前最大页面号为指定的页面号。最后返回所选择的缓冲区号。

(4) 缓冲池页面的读操作

该操作定义在函数 SimpleLruReadPage 中，主要作用是在缓冲池中找到所需要的页面。如果指定页面已经在缓冲池中，那么直接返回其所在的缓冲区；否则调用 SlruPhysicalReadPage 函数物理读入指定页面到某一缓冲区中，并返回该缓冲区号。

其执行过程如下：

- 1) 首先，调用 SlruSelectLRUPage 函数在缓冲池中寻找指定页面。

2) 如果指定页面在缓冲池中的某一缓冲区中且状态为非“空”，判断其状态是否为“正在读入”，如果不是，说明指定页面已在缓冲区中，所以可以直接使用，标识该页面为最近使用的页面，同时将缓冲池中的其他页面的 LRU 值增一，并返回找到的缓冲区。

3) 如果指定页面不在缓冲池中，需要物理读入到由 SlruSelectLRUPage 所找到的可用缓冲区中。接着，设置找到缓冲区的页面号为指定页面号，标识其状态为“正在读入”，并标识该缓冲区为最近使用。

4) 然后释放所持有的缓冲池控制锁，并请求缓冲区的排他锁。判断即将要写入的缓冲区的页面号是否仍为指定页面号且其状态仍为“正在读入”。如果不是，说明有其他进程已经打乱了之前所选定的缓冲区，必须返回到第一步重新执行。

5) 否则调用 SlruPhysicalReadPage 函数将指定页面读入到缓冲区中。然后，请求缓冲池的控制排他锁。判断是否成功读入指定页面，如果是，那么标识缓冲区状态为“干净”；否则，设置缓冲区状态为“空”。释放所持有的缓冲区锁，并标识指定页面为最近使用页面。最后返回页面写入的缓冲区号。如果没有成功读入，则报告出错信息。

另外，从物理文件中读入的操作在函数 SlruPhysicalReadPage 中实现，该函数的执行流程如下：

1) 首先，根据页面号获得该页面所在的物理段号以及段内偏移。因为物理段文件都是根据段号进行命名的，所以根据段号就可以很容易地找到所需打开文件的路径。

2) 然后根据物理文件路径打开该文件。检查该文件是否成功打开。在崩溃后重启的情况下，可能会接收到命令设置事务状态信息，而这些事务的日志处在已删除的段文件中，所以导致此时的文件不能成功打开，在这种情况下，并不返回失败，而是允许这种情况存在，只是将需要读入页面的缓冲区设置为全零，并成功返回。

3) 然后定位指定页面在物理段文件中的偏移，并调用 read 方法读入一个页面大小的数据到指定缓冲区中，此后关闭打开的文件。在这个过程中，可能在定位文件时出错、也可能在读入文件的时候出错，其应对措施都是关闭打开的文件，并返回失败。

4) 最后返回成功。

(5) 缓冲池页面的写操作

缓冲池页面的写操作定义在函数 SimpleLruWritePage 中，用于将一个指定页面由缓冲区写入到磁盘。调用本函数前要求持有缓冲池控制锁。

执行流程如下：

1) 首先，检查缓冲区状态是否为“脏”或“正在写入”，如果不是，说明不需要写入动作，直接返回。否则说明该页面有数据更新，需要写入到磁盘中。

2) 然后释放先前获得的缓冲池控制锁，并请求页面所在缓冲区的缓冲区排他锁。再次检查需要写入的页面仍在原来缓冲区中且其状态仍为“脏”或“正在写入”，如果不是，说明该页面可能被其他进程写入到磁盘或该页面被移动到其他缓冲区中，不再需要写操作，直接返回。

3) 否则，继续执行操作，将页面所在缓冲区状态修改为“正在写入”。接下来调用 SlruPhysicalWritePage 函数将页面写入到磁盘中。然后，检查页面写入动作是否成功完成，如果不是，检查缓冲区刷新结构 FData 中是否记录了已打开文件，如是，那么为安全起见，将 FData 中记录的所有打开文件关闭。重新申请缓冲池的控制锁；重新检查缓冲区状态，因为其他的进程有可能对该缓冲区重新进行写操作，所以其状态由“正在写入”改为“脏”；如果这种情况发生，那么在成功写入

页面到磁盘后无需修改缓冲区状态，否则，在将指定页面成功写入到磁盘后，需要将该页面缓冲区所在状态置为“干净”。

4) 最后，释放所持有的缓冲区排他锁，并检查刚才的写入动作是否成功完成，如果不成功，需要报告出错信息。

FData 的参数主要用于记录当前缓冲池所有打开的文件的文件描述符 (FD)，如果需要写入到磁盘的页面所在物理段文件已经打开，那么直接由 FData 中获得这些描述符；否则，需要重新打开该物理段文件（或重新创建），并把该文件的描述符放到 FData 中。

由于一开始涉及缓冲区状态相关的操作，因此对本函数调用前需持有缓冲池控制锁，在完成了缓冲区状态检查后即释放所持有的控制锁，接着申请缓冲区的排他锁用于写入操作，最后在函数 SimpleLruWritePage 退出前释放该缓冲区排他锁。而且，在退出前需要再次持有缓冲池的控制锁。

另外，物理写入文件的操作定义在函数 SlruPhysicalWritePage 中，该函数的执行流程如下：

- 1) 首先，由页面号转换获得指定页面在段文件中的段号以及段内偏移。
- 2) 然后检查要打开的段文件是否已经打开，如果是，那么其文件描述符会被放在 FData 中，所以不需要再次打开，直接由 FData 中获取即可；否则，根据文件名获得文件的路径，并打开该段文件。
- 3) 如果页面要写入的段文件并不存在，则需要创建一个新的段文件。然后，定位页面在段中的偏移，并通过系统调用将该页面由缓冲区写入到磁盘中。还要检查控制结构中的同步标志是否为“真”，如是，说明写操作需要进行同步写，要调用 pg_fsync 函数将页面同步写入到磁盘中。
- 4) 最后关闭打开的文件并返回。在这个过程中可能会遇到错误，需要设置错误位以及错误标识，方便上层调用函数知道出错原因并正确报错。

(6) 缓冲池页面的刷新

该操作定义在函数 SimpleLruFlush 中，用于在检查点或数据库关闭时，将缓冲池中的“脏”页面刷写到磁盘中，主要是通过调用 SimpleLruWritePage 函数完成写操作。

该函数的执行流程如下：

- 1) 首先，获取缓冲池的控制锁，并构造一个刷写数据结构 FData，以供后续的刷写操作记录打开的文件描述符。
- 2) 然后调用 SimpleLruWritePage 将缓冲池的各个缓冲区写回到磁盘中，并确保写后的每一个缓冲区的状态为“干净”或“空”。并释放所持有的缓冲池控制锁。
- 3) 然后，检查控制结构中的同步标识是否为“真”，如是，则将 FData 记录的打开的文件都同步写回到磁盘中，确保前面调用 SimpleLruWritePage 所执行的写操作都真正地物理写回到磁盘。
- 4) 最后逐个关闭打开的文件，并成功返回。在这个过程中的同步或关闭文件操作有可能会出错，则需要报告出错信息。

(7) 缓冲池页面的删除

该操作定义在函数 SimpleLruTruncate 中。由于部分日志文件随着检查点的建立而变得无用，因此可以删除。本函数会根据删除标记偏移量删除该页面之前的过时页面。由于物理存储是以段文件的形式存储的，因此执行删除操作时也是以段文件为单位执行的，即根据页号得到其所在的段文件，并删除此段号之前的所有段文件。

执行流程如下：

- 1) 首先，计算删除标记页面所在段的段号，并获取缓冲池的控制排他锁。

2) 检查删除标记页面是否合法, 如果该页面号比当前最大的页面号还要大, 说明不合法, 报错并释放缓冲池的控制锁后返回。

3) 检查缓冲池的每个缓冲区, 并确保删除标记页面之前的页面没有在缓冲区中存在或没有被使用。对于删除标记页面之前的页面, 如果在缓冲区中且其状态为“正在读入”, 则调用 SimpleLruReadPage 函数完成读入操作; 如果在缓冲区中且其状态为“脏”或“正在写入”, 那么调用 SimpleLruWritePage 函数完成写操作, 将页面写回到磁盘。

4) 返回到步骤 3, 直到删除标记页面之前的所有页面都不在缓冲区中为止。

5) 最后, 释放所持有的缓冲池控制锁, 并调用 SlruScanDirectory 函数真正执行删除操作, 将所有过时的段文件执行物理删除。

(8) 缓冲池页面相关物理文件的删除

该操作定义在函数 SlruScanDirectory 中, 主要为 SimpleLruTruncate 函数所调用, 以便从目录中删除过时的段文件。但是只在删除条件 (doDeletions) 为“真”时才真正执行删除动作。在执行删除动作时, 并不需要请求任何锁。

执行流程如下所示:

1) 首先, 根据参数计算需待删除页面所在的段, 因为真正执行物理删除时是以段文件为单位进行删除的, 所以需要计算待删除页面所在的段号, 并删除在该段之前的各段文件。

2) 然后由控制结构获得数据文件所在的目录, 取得目录下所有文件的文件信息。逐个比较各个文件的名称, 如果该文件的名称是以 4 位 16 进制数命名, 那么计算该段的第一个页面与参照页面的前后关系, 若位于参照页面之前, 则说明该段文件已经过时, 需要删除, 并置可删除文件的标识为“真”; 最后判断 doDeletions 是否为“真”, 如为“真”, 则真正执行删除; 否则, 不真正执行删除操作。在这个过程可能出错, 如果有错, 则报告出错信息。

3) 最后释放关联的目录, 并返回是否找到可删除文件的标识。

7.11.2 CLOG 日志管理器

CLOG 日志记录的是事务的最终状态。CLOG 日志管理器管理着 CLOG 日志缓冲池, 该日志缓冲池是基于 SLRU 缓冲池实现的。

1. CLOG 日志管理器相关数据结构

在 PostgreSQL 事务系统中一共定义了四种事务状态 (数据结构 7.24)。

数据结构 7.24 TRANSACTION_STATUS

#define TRANSACTION_STATUS_IN_PROGRESS	0x00 /* 事务正在执行 */
#define TRANSACTION_STATUS_COMMITTED	0x01 /* 事务已提交 */
#define TRANSACTION_STATUS_ABORTED	0x02 /* 事务被中止 */
#define TRANSACTION_STATUS_SUB_COMMITTED	0x03 /* 事务的子事务已提交 */

状态 TRANSACTION_STATUS_SUB_COMMITTED 是为嵌套事务而引入的(其他三种状态用于普通事务), 即对于嵌套事务, 是存在子事务的, 如果父事务没有提交, 而子事务已提交, 那么子事务的状态就设置为 TRANSACTION_STATUS_SUB_COMMITTED。

由于 CLOG 只记录事务的 4 个状态,因此记录一个事务状态信息的 CLOG 日志记录只需要 2 个比特,一个字节就可以存储 4 个事务的 CLOG 日志记录。在一个页面(大小为 8KB)中,可记录事务 CLOG 日志记录条数为 $8K * 8b / 2b = 32K = 2^{15}$ 。在一个由 32 个页面组成的段中,可记录的 CLOG 日志记录条数为 $32 * 2^{15} = 2^{20}$ 。

由于事务 ID 为 32 比特,因此可表示的事务数有 2^{32} 个,而 CLOG 日志文件是以段文件为单位的,这意味着可能存在的 CLOG 段的数量为 $2^{32} / 2^{20} = 2^{12}$ 。CLOG 日志段文件以段号命名,所以只需要 12 位即可,而实际上,PostgreSQL 以 4 位的 16 进制数来表示段号,并作为日志段文件的名称。CLOG 日志的数据位于 PGDATA/pg_clog 目录下。

通过一个 4 元组 $\langle \text{Segmentno}, \text{Pageno}, \text{Byte}, \text{Bindex} \rangle$ 就可以定位一条 CLOG 日志记录。其中 Segmentno 为段号,即实际的段文件名称; Pageno 为日志记录所在的段内的页偏移; Byte 为页面偏移; Bindex 为字节内偏移。

通过一个事务的事务 ID 就可以获得其日志记录对应的 4 元组,具体方法如下:

```
/* 根据事务 ID 获得对应事务所在页面 */
#define TransactionIdToPage(xid)
/* 根据事务 ID 获得事务在对应页面上的偏移 */
#define TransactionIdToPgIndex(xid)
/* 根据事务 ID 获得事务在对应页面上的字节偏移 */
#define TransactionIdToByte(xid)
/* 根据事务 ID 获得事务在记录其 CLOG 记录的字节上的位偏移 */
#define TransactionIdToBIndex(xid)
```

由于一个日志段由 32 个页面所组成,所以通过计算 CLOG 日志记录所在页面的页面号就可得到该日志记录所在的段。

CLOG 日志缓冲池就是一个 SLRU 缓冲池,在整个数据库系统中,CLOG 日志缓冲池只有一个,它注册在共享内存中,其名称为“CLOG Ctl”。PostgreSQL 定义了一个静态变量 ClogCtl,它是一个 SLRU 缓冲池控制结构,记录了 CLOG 的 SLRU 缓冲池的数据缓冲区在共享内存中的地址,以及对 CLOG 日志进行写操作的同步信息,默认要求 CLOG 日志的写操作是同步写操作。

2. CLOG 日志管理器主要操作

CLOG 日志管理器的操作主要包括对日志管理器的初始化、日志的读写操作、日志页面的初始化,以及日志管理器的新建、启动、关闭、扩展、删除等操作。

(1) 日志管理器的初始化

该操作定义在函数 CLOGShmemInit 中,用于在共享内存中初始化 CLOG 缓冲池。主要是设置 CLOG 共享内存控制信息中比较两个 CLOG 日志页面前后关系的函数信息,并调用 SimpleLruInit 函数在共享内存中初始化 CLOG 缓冲池,设置 CLOG 日志数据存储的磁盘路径,在共享内存中注册唯一的 CLOG 缓冲池,并设置 CLOG 日志管理器的控制结构,最后分配一个缓冲池的控制锁给 CLOG 日志管理器。

(2) CLOG 日志的写操作

该操作定义在函数 TransactionIdsetStatus 中,用于设置指定事务的状态,实际上就是在 CLOG 日志文件中写该事务的 CLOG 日志记录。此函数是一个底层的操作,并不直接提供给其他一般进程设

置事务状态，具体对事务状态更新由事务日志接口例程的函数来完成。

其主要流程如下：

- 首先判断设置的状态信息必须为“提交”、“中止”或“子事务提交”。写 CLOG 日志记录时，由于其初始化的日志页面为全 0，即状态“正在处理”，所以设置该状态是没有意义的。而且 CLOG 日志记录的是事务的最终状态，“正在处理”并不是一个事务的最终状态，所以设置一个事务的状态必须是“提交”、“中止”或“子事务提交”中的一种。由于涉及状态信息的修改，因此必须申请 CLOG 日志缓冲池的控制排他锁。

- 然后调用 SimpleLruReadPage 函数将指定事务的 CLOG 日志记录所在页面读入到 CLOG 日志缓冲池中的某一缓冲区中。如果指定页面已经在缓冲区中，则直接操作该缓冲区，而无需再次执行读入操作。然后，根据事务 ID 获得事务 CLOG 日志记录对应的四元组信息，进而设置该事务的 CLOG 日志中的状态信息。

- 最后置该条日志记录所在页面状态为“脏”，并释放之前持有的缓冲池控制锁。

(3) CLOG 日志的读操作

该操作定义在函数 TransactionIdSetStatus 中，用于从 CLOG 日志缓冲池中读取事务的 CLOG 日志记录，从而获取指定事务的最终状态信息。同 TransactionIdSetStatus 函数一样，此函数也是一个底层的操作，并不直接提供给其他一般进程获取事务状态，对事务状态的获取由事务日志接口例程的函数来完成。

其主要流程如下：

- 首先需要请求 CLOG 日志缓冲池的控制锁。
- 然后调用 SimpleLruReadPage 函数将指定事务 CLOG 日志记录所在的页面读入到缓冲区中，如果所需页面已经在缓冲池中，则直接使用。根据事务 ID 获得事务 CLOG 日志记录对应的四元组信息。由缓冲区中读出指定事务的 CLOG 日志记录，从而获得指定事务的最终状态。
- 最后释放所持有的 CLOG 日志缓冲池的控制锁并返回指定事务的最终状态。

(4) CLOG 日志页面的初始化

该操作定义在函数 ZeroCLOGPage 中，初始化指定页面为全 0；根据 writeXlog 设置是否要写一个 XLOG 日志记录。如果是为扩展 CLOG 日志页面创建一个新的 CLOG 日志页面并初始化该页面，则需要创建一条 XLOG 日志记录，记录所扩展 CLOG 日志页面的页面号。

基本流程如下：

- 首先调用 SimpleLruZeroPage 函数在缓冲池中选择一个可用缓冲区，并初始化该缓冲区为全 0。
- 然后判断是否需要写 XLOG 日志，如果 writeXlog 为 TRUE，则调用 WriteZeroPageXlogRec 函数将新创建的 CLOG 页面号记录到一条 XLOG 记录中；否则无需记录 XLOG 日志记录，直接返回。

(5) CLOG 日志段的创建

该操作定义在函数 BootStrapCLOG 中，在系统初始化时需要调用本函数，以创建第一个 CLOG 日志段。

执行基本流程如下：

- 首先获取 CLOG 日志缓冲池的控制锁。
- 然后在缓冲池中选择一个可用缓冲区供初始 CLOG 日志页面（页面号为 0）使用，调用 Ze-

roCLOGPage 函数将该缓冲区初始化为全 0。然后，调用 SimpleLruWritePage 函数将这第一个 CLOG 日志页面由缓冲区写到磁盘中。

3) 最后，设置该页面的状态为“干净”，并释放所持有的缓冲池控制锁。

(6) CLOG 日志的启动

该操作定义在函数 StartupCLOG 中，基本功能是当 Postmaster 或一个单独的服务进程启动时，启动 CLOG 日志管理器。

基本流程如下：

1) 首先由共享内存中的变量缓冲区获得下一可分配的事务 ID，请求 CLOG 日志缓冲池的控制锁。

2) 然后根据下一可分配的事务 ID 获得该事务的 CLOG 日志记录所在的 CLOG 日志页面，并设置该页面为当前最大的 CLOG 日志页面。

3) 然后调用 SimpleLruReadPage 函数将该日志页面读入到某一缓冲区中；根据下一可分配的事务 ID 获得该事务在页面中的偏移位置，将该 CLOG 日志页面从该偏移位置之后都设置为全 0。

4) 最后，设置该 CLOG 日志页面状态为“脏”，并释放所持有的 CLOG 日志缓冲池控制锁。

(7) CLOG 日志的关闭

该操作定义在函数 ShutdownCLOG 中，当关闭 Postmaster 或单独的服务进程时，需要调用本函数关闭 CLOG 日志管理器。它的实现是通过调用 SimpleLruFlush 函数将当前的 CLOG 日志缓冲池中的“脏”页面写回到磁盘中而完成的。

(8) 创建检查点时 CLOG 日志的操作

该操作定义在函数 CheckPointCLOG 中，作用是执行一个 CLOG 日志检查点，在 XLOG 执行检查点时被调用。该操作完成的任务就是通过调用 SimpleLruFlush 函数将当前 CLOG 日志缓冲池中的“脏”页面写回到磁盘中。

(9) CLOG 日志的扩展

该操作定义在函数 ExtendCLOG 中，为新分配的事务 ID 创建 CLOG 日志空间。当一个事务的 CLOG 日志记录位于一个 CLOG 日志页面的首部时，若这个事务要写一条 CLOG 日志记录，就需要调用 ExtendCLOG 以创建一个新的 CLOG 日志页面供该条 CLOG 日志记录写入。

执行流程如下：

1) 首先根据事务 ID 获得该事务 CLOG 日志记录所在的 CLOG 日志页面号。

2) 然后请求 CLOG 日志缓冲池的控制锁。

3) 将该日志页面初始化为全 0，并写一条 XLOG 日志记录该页面号。

4) 最后释放所持有的 CLOG 日志缓冲池控制锁。

(10) CLOG 日志的删除

该操作定义在函数 TruncateCLOG 中。由于日志检查点的建立使得一些事务的日志变成过时的，为了节省空间，应该把这些过时的日志记录删除。由于 CLOG 日志物理存储是以段为单位，因此该操作以段为单位，删除过时的 CLOG 日志记录。

基本执行流程如下：

1) 首先根据指定的事务 ID 获得该事务的 CLOG 日志记录所在页面号。

2) 接着调用 SlruScanDirectory 函数，以该页面为参照页面，检查是否有过时的 CLOG 日志段需

要删除，如果没有则直接返回。

- 3) 然后请求执行一个检查点。
- 4) 最后调用 SimpleLruTruncate 函数将过时的 CLOG 日志段文件删除。

(11) 构建 CLOG 时创建 XLOG 日志记录

该操作定义在函数 WriteZeroPageXlogRec 中。当创建了一个新的 CLOG 日志页面时，需要调用本函数，创建一条“ZEROPAGE”的 XLOG 日志记录，以记录所创建 CLOG 日志页面的页面号，方便在系统崩溃或恢复过程中重建 CLOG 日志。

基本流程如下：

1) 构造一个 XLOG 日志记录，记录指定 CLOG 日志页面的页面号到该记录中，并标识该日志的资源管理器号为 RM_CLOG_ID。

2) 调用 XLogInsert 函数将该条 XLOG 日志记录插入到 XLOG 日志中。

(12) CLOG 日志的 REDO 操作

该操作定义在函数 clog_redo 中，CLOG 日志的 REDO 在执行 XLOG 的 REDO 操作时被调用，主要是作为 CLOG 资源管理器的例程。

由于在扩展或创建一个新的 CLOG 页面时都要创建一条 XLOG 日志记录以保存该 CLOG 页面的页面号，因此在执行 XLOG 的 REDO 操作时，如果碰到资源管理器号为 RM_CLOG_ID、类型为 ZEROPAGE 的 XLOG 日志记录，那么调用该操作。

执行流程如下：

- 1) 从 XLOG 日志记录中获得该 CLOG 日志页面的页面号。
- 2) 然后调用 ZeroCLOGPage 函数在 CLOG 日志缓冲池中选择一个可用的缓冲区供该页面使用，并设置该页面为全 0。
- 3) 调用 SimpleLruWritePage 函数将该 CLOG 日志页面写回到磁盘中。
- 4) 在 XLOG 的 REDO 过程中完成 CLOG 日志记录的重建过程。在这个过程中需要持有 CLOG 日志缓冲池的控制锁。

7.11.3 SUBTRANS 日志管理器

PostgreSQL 中有嵌套事务的概念，它基本思想是嵌套事务中存在一个事务树。从根开始，每个事务都可以建立更低层次的事务（子事务），子事务被嵌套在父节点的控制区域之内。为此 PostgreSQL 引入了 SUBTRANS 日志，记录每个事务的父事务 ID。

SUBTRANS 日志通过 SUBTRANS 日志管理器来管理。SUBTRANS 日志管理器管理着一个缓冲池，同 CLOG 日志管理器的缓冲池一样，也是基于 SLRU 缓冲池实现的。SUBTRANS 日志管理器是一个类似提交日志管理器的管理器，存储的是每一个事务的父事务 ID。它是嵌套事务实现的一个基础部分。一个主事务的父事务是非法事务 ID，每一个子事务都有一个直接的父事务。遍历事务树可以很容易地由一个子事务到父事务，但是反过来并不能实现。

SUBTRANS 日志的健壮性要求和 CLOG 日志是完全不同的，因为需要记录的只是当前打开事务的子事务信息。所以在系统崩溃或重启时并不需要保存数据。由于在系统崩溃时不需要保存数据，因此也不需要和 XLOG 进行交互，也没有相应的 REDO 函数。在数据库启动时，只要使当前活跃的子事务页面为全 0 就可以了。

1. SUBTRANS 日志管理器相关数据结构

子事务日志记录的是一个事务的父事务 ID。事务 ID 为 32 比特，所以一条子事务日志记录也为 32 比特。一个页面所能存储的子事务日志记录条数为： $8K * 8b / 32b = 2^{11}$ 。一个段所能存储的子事务日志记录条数为： $32 * 2^{11} = 2^{16}$ 。同样，子事务日志记录也是以段文件为单位，并以段号命名。子事务日志的数据位于 PGDATA/pg_subtrans 目录下。

通过一个三元组 $\langle \text{Segmentno}, \text{Pageno}, \text{PageIndex} \rangle$ 即可定位一条子事务日志记录。其中 Segmentno 为段号，即实际的段文件名称，Pageno 为日志记录所在的段内的页偏移，PageIndex 为页内偏移。通过一个事务的事务 ID 就可以获得日志记录对应的三元组。

```
/* 根据事务 ID 获取记录该事务的 SUBTRANS 日志页面 */
#define TransactionIdToPage(xid) ((xid) / (TransactionId))
SUBTRANS_XACTS_PER_PAGE

/* 根据事务 ID 获取记录该事务所在 SUBTRANS 日志页面的偏移 */
#define TransactionIdToEntry(xid) ((xid) % (TransactionId))
SUBTRANS_XACTS_PER_PAGE
```

同样，由于一个子事务日志段由 32 个页面组成，因此通过计算子事务日志记录所在页面的页面号就可得到该日志记录所在的段。

子事务日志缓冲池是一个 SLRU 缓冲池，在整个数据库系统中，子事务日志缓冲池只有一个，它在共享内存中是经过注册的，其名称为“SUBTRANS Ctl”。

子事务日志缓冲池定义了一个静态变量 SubTransCtl，它是一个 SLRU 缓冲池控制结构，记录了子事务的 SLRU 缓冲池的数据缓冲区在共享内存中的地址，以及对子事务日志进行写操作的同步信息，默认要求子事务日志的写操作是非同步写操作。

2. SUBTRANS 日志管理器主要操作

子事务（SUBTRANS）日志管理器的操作包括日志缓冲池的建立、初始化、日志的读写操作，以及日志管理器的新建、启动、关闭、扩展、删除等。由于大部分操作同 CLOG 日志管理器操作流程一样，这里不再赘述。我们详细探讨子事务日志的读写操作。

(1) 子事务日志的写操作

该操作定义在函数 SubTransSetParent 中，用于设置指定事务的父事务 ID。实际上就是在子事务日志中记录指定事务的 SUBTRANS 日志记录。

基本流程如下：

- 1) 首先请求 SUBTRANS 日志缓冲池的控制锁。
- 2) 然后调用 SimpleLruReadPage 函数将指定事务的 SUBTRANS 日志记录所在页面读入到某一缓冲区中。
- 3) 然后根据事务 ID 获取其对应的三元组信息，以定位该事务的 SUBTRANS 日志记录在日志页面中的位置。
- 4) 然后设置指定事务的父事务 ID，即写指定事务的 SUBTRANS 日志记录。
- 5) 最后，设置页面的状态为“脏”，并释放所持有的缓冲池控制锁。

(2) 子事务日志的读操作

该操作定义在函数 SubTransGetParent 中，根据指定事务的 SUBTRANS 日志记录，获取该事务的

父事务 ID。对于一个顶层事务来说其父事务 ID 为无效事务 ID。

基本流程如下：

- 1) 首先检查该事务 ID 是否正常，如果不正常则其父事务 ID 为无效事务 ID。
- 2) 然后请求 SUBTRANS 缓冲池控制锁，调用 SimpleLruReadPage 函数读取指定事务的 SUBTRANS 日志记录所在页面到缓冲池中。
- 3) 然后计算指定事务 ID 的 SUBTRANS 日志记录对应的三元组信息，并读取其 SUBTRANS 日志记录，从而获得指定事务的父事务 ID。
- 4) 最后释放所持有的缓冲池控制锁，并返回指定事务的父事务 ID。

(3) 子事务日志的回溯

该操作定义在 SubTransGetTopmostTransaction 中，主要功能是根据指定事务逐级向上回溯其所在的嵌套事务树，获得指定事务所在嵌套事务树的根事务 ID。

由于嵌套事务形成了一个事务树，因此只需要通过指定事务，在子事务日志中逐级向上回溯找寻其父事务，直到遇到一个事务的父事务 ID 为无效事务 ID，则说明该事务为所要寻找的根事务 ID。

7.11.4 MULTIXACT 日志管理器

MULTIXACT 日志是 PostgreSQL 系统用来记录组合事务 ID 的一种日志。由于 PostgreSQL 采用了多版本并发控制，因此同一个元组相关联的事务 ID 可能有多个，为了在加锁（行共享锁）的时候统一操作，PostgreSQL 将与该元组相关联的多个事务 ID 组合起来用一个 MultiXactID 代替来管理。同 CLOG、SUBTRANS 日志一样，MULTIXACT 日志也是利用 SLRU 缓冲池来实现。

1. MULTIXACT 日志管理器相关数据结构

MultiXactID 是一个多对一的映射关系，需要在事务 ID 数组中标记哪一段映射到一个 MultiXactID（如图 7-10 所示）。所以在映射的过程中需要存储两种信息，即需要标识一段事务 ID 的偏移量（Offset），还需要记录这段偏移量的大小（nMembers）。

由于需要对 MultiXactID 的分配进行维护，于是定义了数据结构 MultiXactStateData（如数据结构 7.25 所示）。

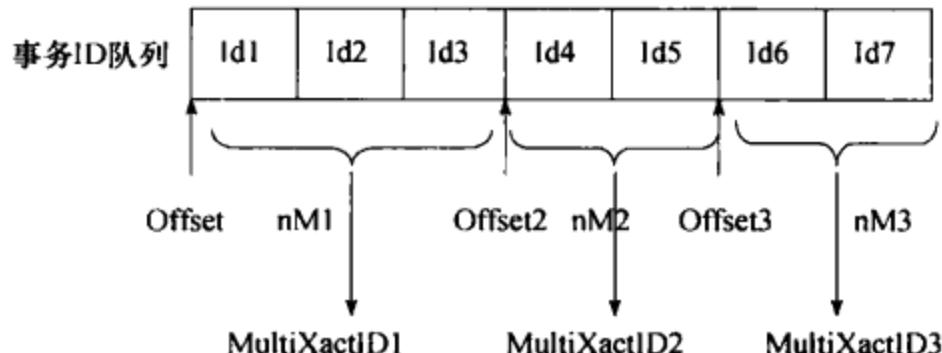


图 7-10 MultiXactID 组合关系

数据结构 7.25 MultiXactStateData

```
typedef struct
{
    MultiXactId     nextMXact;           //下一个可分配的 MultiXactID
    MultiXactOffset nextOffset;          //下一个对应 MultiXactID 的起始偏移量
    MultiXactId     lastTruncationPoint; //上一次删除位置
    MultiXactId     perBackendXactIds[1]; //MultiXactID 队列起始位置
}MultiXactStateData;
```

MultiXactStateData 用来管理和维护 MultiXactID，另外还需要定义一个统一的接口来操作 MultiXactID，即存储这种多对一的映射关系，PostgreSQL 定义 mXactCacheEnt（如数据结构 7.26 所示）来完成此工作。

MULTIXACT 日志存储在 PGDATA/pg_multixact/ 目录下，其中会有两个子目录 members 和 offset，分别存储 XactID 成员和偏移量。从一个 MultiXactID 映射到具体的存储位置是通过下面的变换来完成的：

```
/* 根据 MultiXactID 获得对应记录位于的页面 */
#define MultiXactIdToOffsetPage(xid) \
    ((xid) / (MultiXactOffset)MULTIXACT_OFFSETS_PER_PAGE)
/* 根据 MultiXactID 获得对应记录位于页面内的偏移量 */
#define MultiXactIdToOffsetEntry(xid) \
    ((xid) % (MultiXactOffset)MULTIXACT_OFFSETS_PER_PAGE)
/* 同理，以下为对于事务成员完成上述关系的转换 */
#define MXOffsetToMemberPage(xid) \
    ((xid) / (TransactionId)MULTIXACT_MEMBERS_PER_PAGE)
#define MXOffsetToMemberEntry(xid) \
    ((xid) % (TransactionId)MULTIXACT_MEMBERS_PER_PAGE)
```

MULTIXACT 日志的缓冲池用两个 SLRU 缓冲池来实现，分别是 MultiXactOffsetCtl 和 MultiXactMemberCtl，分别记录 Members 和 Offsets。在 Postmaster 启动后，就注册在共享内存中，管理全局的 MultiXactID。

2. MULTIXACT 日志管理器主要操作

同 CLOG 日志管理器一样，MultiXact 日志管理器的操作包括 MULTIXACT 日志系统的启动、关闭、检查点操作，以及 MultiXactID 的创建、扩增等。

1) MULTIXACT 的创建

该操作定义在函数 MultiXactIdCreate 中。MultiXactID 的创建指的是将两个 XID 组合成一个整体，结果返回一个 MultiXactID，用来代替这两个 XID。创建过程如下：

①首先在当前进程 MXactCache 中查找这两种 XID 的组合是不是已经生成了一个 MultiXactID，如果已经生成，直接返回该 MultiXactID 即可。

②调用 GetNewMultiXactId 函数从共享内存 MultiXactState 中获得一个新的 MultiXactID，并更新 MultiXactState 中的 nextMXact 和 nextOffset。

③向 XLOG 中写入当前操作的日志，记录当前分配得到的 MultiXactID 以及组合成该 MultiXactID 的 Xids 链表和数量。

④最后将该 MultiXactID 以及相关 Xids 更新到本地缓存中，返回 MultiXactID。

2) MULTIXACT 的扩增

该操作定义在函数 MultiXactIdExpand 中。我们知道，一个 MultiXactID 可对应多个 Xid，而上面的 MultiXactID 的创建只能完成两个 XID 到一个 MultiXactID 的映射过程，于是定义了 MultiXactI-

数据结构 7.26 mXactCacheEnt

```
typedef struct mXactCacheEnt
{
    struct mXactCacheEnt *next;
    MultiXactId          multi;
    int                  nxids;
    TransactionId        xids[1];
}mXactCacheEnt;
```

dExpand 来完成 MultiXactID 的扩增，即将一个 Xid 加入到该 MultiXactID 对应的 Xids 队列中。

①首先由 MultiXactID 找到对应的 Xids 队列和数量。

②检查该 MultiXactID 是否有对应成员存在，若不存在，创建一个仅由该 Xid 组成的 MultiXactID 返回即可。

③检查该 Xid 是否已在当前 MultiXactID 对应的 Xids 队列中，若存在，直接返回即可。

④将该 MultiXactID 对应的 Xids 链表取出，将新增的 Xid 放到其末尾，然后重新调用 MultiXactIdCreate 得到新的 MultiXactID，最后返回该值。

3) 检测 MULTIXACT 是否在运行中

该操作定义在函数 MultiXactIdIsRunning 中，其逻辑比较简单，当该 MultiXactID 中的任何一个成员属于当前事务 ID 或者正在运行中，则认为该 MultiXactID 也在运行中。

4) MULTIXACT 日志系统的启动

该操作在函数 StartupMultiXact 中定义，一般在 Postmaster 启动时通过调用 StartupXlog 进而调用 StartupMultiXact。主要功能是用来初始化管理 MultiXactID 的偏移量和成员的两个 SLRU 缓冲池。

5) MULTIXACT 日志系统的关闭

该操作在函数 ShutdownMultiXact 中定义，用来完成 MultiXactID 日志系统的关闭。通过调用 SimpleLruFlush 函数将管理 MultiXactID 的偏移量和成员的两个 SLRU 缓冲池刷新到磁盘即可。

6) MULTIXACT 日志系统检查点操作

在创建检查点时，MultiXactID 日志也需要刷新到磁盘，所以 CheckPointMultiXact 用来将管理 MultiXactID 的偏移量和成员的两个 SLRU 缓冲池刷新到磁盘。与 MULTIXACT 日志的关闭不同的是，在这个刷新过程中，允许其他进程向缓冲池中写日志。

7.11.5 XLOG 日志管理器

XLOG 是传统数据库理论中提到的事务日志，它详细地记录了服务进程对数据库的操作过程。XLOG 日志文件在内存中按页进行存放，每个页面大小为 8KB，每个页都有一个头部，头部信息之后才是 XLOG 日志记录。

每个 XLOG 文件都有一个 ID，但事实上它被分为一个个大小为 16MB 的 XLOG 段文件来存放。XLOG 文件号和段文件号可以用来唯一地确定这个段文件。确定日志文件内一个日志记录的地址时，只需用一个 XLOG 文件号和日志记录在该文件内的偏移量即可。

对于 XLOG 文件中每个段文件的第一个页面，其头部是一个长头部（Long Header）。一个 XLOG 文件头部是不是长头部，可以由其头部格式的标志位判断出来。

1. XLOG 日志管理器相关数据结构

XLOG 是数据库中的主要日志形式，一切操作的记录和恢复都是通过 XLOG 日志来最终完成。

(1) 日志页面头部信息

XLOG 日志分为许多的逻辑段文件，每一个段文件又分成许多个页面，每一个页面的大小为一个块的大小。对于每一个日志页面，需要在其头部写一个头部信息 XLogPageHeaderData（数据结构 7.27）。

数据结构 7.27 XLogPageHeaderData

```
struct XLogPageHeaderData
{
    uint16      xlp_magic;           //校验位
    uint16      xlp_info;            //标志位
    TimeLineID  xlp_tli;             //页面第一条记录的时间序列
    XLogRecPtr   xlp_pageaddr;        //XLOG 页面的首地址
}
```

其中，标志位 xlp_info 只使用最低两位“0”或者“1”，“0”表明该页的第一个 XLOG 记录接着上一个页的最后一个 XLOG 记录，这样的记录会接在一个 XLogContRecord 结构体之后；“1”表明该页的头部是一个长头部，即该页是该 XLOG 文件的第一页。

如果一个日志页面是一个逻辑段文件的第一个页面，那么在页面头部信息标志位会设置 XLP_LONG_HEADER 标志，那么将在原页面头部信息的基础上使用一个长的 XLOG 页面头部 XLogLongPageHeaderData（数据结构 7.28），主要是为更精确地定位文件服务。

数据结构 7.28 XLogLongPageHeaderData

```
typedef struct XLogLongPageHeaderData
{
    XLogPageHeaderData    std;          //头部信息
    uint64                xlp_sysid;     //pg_control 中的系统标识符
    uint32                xlp_seg_size;   //校验位,段的尺寸
    uint32                xlp_xlog_blkksz; //校验位,块的尺寸
}XLogLongPageHeaderData;
```

对于一个长的 XLOG 日志记录，可能在当前页面中没有足够空间来存储，系统允许将剩余的数据存储到下一个页面，即一个 XLOG 日志记录可以分开存储到不同的页面。但需要注意的是，一个 XLOG 日志记录的头部信息 XLogRecord 是不允许分开存储到两个不同页面的，如果一个页面的剩余空间不足已存储一个 XLOG 记录的头部，那么该页面的剩余空间将被舍弃，将这个新的 XLOG 记录存储到新的下一页。

当一个 XLOG 记录被分开存储到不同页面时，用 XLogContRecord（数据结构 7.29）来记录该页面除了本页面存储的内容外剩余的需存储数据信息。

数据结构 7.29 XLogContRecord

```
typedef struct XLogContRecord
{
    uint32      xl_rem_len;           //记录中剩余数据的总长度
}XLogContRecord;
```

(2) 日志记录控制信息

在 XLOG 中，每一个 XLOG 记录由 4 部分组成，如表 7-13 所示。

表 7-13 XLOG 日志记录结构

XLogRecord	XLOG 记录的头部信息
Data of RMGR	资源管理器的数据，长度 xl_len
Backup Block 1	备份数据块头部 BkpBlock + 块大小的备份数据
Backup Block 2	备份数据块头部 BkpBlock + 块大小的备份数据
Backup Block 3	备份数据块头部 BkpBlock + 块大小的备份数据

XLogRecord(数据结构 7.30)记录了 XLOG 的相关控制信息。另外,一个 XLOG 记录最多可以附 3 个备份块,每一个块对应一个磁盘块大小的数据,长度为 8KB。

数据结构 7.30 XLogRecord

```
typedef struct XLogRecord
{
    crc64           xl_crc;          //本记录的 CRC 校验码
    XLogRecPtr      xl_prev;         //在日志中的前一个记录
    TransactionId   xl_xid;          //事务 ID
    uint32          xl_tot_len;       //整条记录的总长度
    uint32          xl_len;           //资源管理器的数据长度
    uint8           xl_info;          //信息标志位
    RmgrId          xl_rmid;         //资源管理器
}XLogRecord;
```

其中,资源管理器号主要用于日志系统中,数据库系统把各种需要记录日志的数据分类,分配给它们对应的资源管理器号,系统在恢复或者读取日志记录时,能够很方便地知道该日志记录的源数据属于哪一类,结合信息标志位(xl_info)信息,就能知道数据库对源数据做的是哪种操作,从而迅速正确地调用对应的函数。共有 16 种资源,下面列举几类重要的资源:

- RM_XLOG_ID: 该条日志记录的是一个检查点的信息。
- RM_XACT_ID: 该条日志记录的是一个事务的提交或终止信息。
- RM_HEAP_ID: 该条日志记录的是对堆中元组进行修改的信息。
- RM_CLOG_ID: CLOG 中某一页的初始化。
- RM_BTREE_ID: 记录的是对 BTREE 进行修改。

其中,信息标志位(xl_info)的高四位由资源管理器使用,表示该日志是哪种类型的日志,低四位表示对应的块是否需要备份。高四位表示的信息主要包括以下几种:

- XLOG_XACT_COMMIT: 事务提交日志。
- XLOG_XACT_ABORT: 事务终止日志。
- XLOG_HEAP_INSERT: 插入元组日志。
- XLOG_HEAP_DELETE: 删除元组日志。
- XLOG_HEAP_UPDATE: 更新元组日志。

低四位里只用了三位:

- XLR_BKP_BLOCK_1: 表示日志记录(表 7-13)的第一个备份块。

- XLR_BKP_BLOCK_2：表示日志记录（表 7-13）的第二个备份块。
- XLR_BKP_BLOCK_3：表示日志记录（表 7-13）的第三个备份块。

(3) 日志记录数据信息

日志记录中的数据信息存储在结构 XLogRecData（数据结构 7.31）中。

数据结构 7.31 XLogRecData

```
typedef struct XLogRecData
{
    char             *data;           //资源管理器数据
    uint32          len;            //资源管理器数据的长度
    Buffer          buffer;         //该数据涉及的缓冲区
    bool             buffer_std;     //缓冲区存储标准
    struct XLogRecData *next;        //下一个结点指针
}XLogRecData;
```

XLOG 记录中的备份数据块的头部信息保存在 BkpBlock（数据结构 7.32）中。

数据结构 7.32 BkpBlock

```
typedef struct BkpBlock
{
    RelFileNode   node;           //表节点
    BlockNumber   block;          //块号
    uint16        hole_offset;    //空洞偏移量
    uint16        hole_length;    //空洞长度
}BkpBlock;
```

如果需要备份的块有空洞，则备份的时候只记录这个空洞的偏移量及其长度，但没有备份它，从而提高备份效率。

XLogRecPtr（数据结构 7.33）保存的是指向 XLOG 中某个位置的指针。该指针是 64 位的，以防止溢出。

数据结构 7.33 XLogRecPtr

```
typedef struct XLogRecPtr
{
    uint32        xlogid;         //日志文件号
    uint32        xrecoff;        //日志文件中的偏移量
}XLogRecPtr;
```

(4) XLOG 控制结构

在共享内存中用结构 XLogCtlData（数据结构 7.34）保存 XLOG。

数据结构 7.34 XLogCtlData

```

typedef struct XLogCtlData
{
    XLogCtlInsert      Insert;           // 插入一条日志后,最新的相关信息
    XLogwrtRqst        LogwrtRqst;       // 在该值之前的所有日志记录都已经刷新到磁盘
    XLogwrtResult      LogwrtResult;
    uint32              ckptXidEpoch;     // 检查点的创建时间
    TransactionId      ckptXid;
    XLogRecPtr          asyncCommitLSN;   // 最近异步提交日志序列号
    XLogCtlWrite        Write;
    char                *pages;           // 共享内存的缓冲区数组指针
    XLogRecPtr          *xblocks;          // 缓冲区内容对应的 XLOG 文件的内部指针
    Size                XLogCacheByte;    // XLOG 缓冲区的总大小
    int                 XLogCacheBlck;    // XLOG 最大缓冲区的下标
    TimeLineID          ThisTimeLineID;
    slock_t              info_lck;         // LogwrtRqst/LogwrtResult 的共享锁
} XLogCtlData;

```

其中, `LogwrtRqst` 指明需要写或同步写到的日志记录, 而此前位置的日志必须已经写或同步写好了。`LogwrtResult` 指示当前已经写或同步写结束的日志的字节位置。

(5) XLOG 日志的检查点策略

PostgreSQL 只使用 REDO 方法来进行数据库的恢复, 它不使用 UNDO 是因为其数据的多版本使得 UNDO 没有必要。但从本质上讲, PostgreSQL 日志并非是 REDO 日志, 而应该是 UNDO/REDO 日志, 其创建检查点时大体上服从 UNDO/REDO 日志创建检查点时要服从的规则, 只是做了一点改动。

PostgreSQL 建立检查点的规则如下:

- 标记当前日志文件中可用的日志记录点, 即逻辑上的检查点的开始位置。

这个值被作为检查点的 REDO 标记, 在这个过程中是不允许其他事务写日志的(事实上检查点的位置可能会在这之后, 因为在刷新数据缓冲区的过程中允许事务写日志)。

- 刷新所有被修改过的数据缓冲区, 在这个过程中允许其他事务写日志。
- 插入检查点日志记录, 并把该日志记录冲刷到日志文件中。

检查点的结构定义在 `CheckPoint` (数据结构 7.35) 中。

数据结构 7.35 CheckPoint

```

typedef struct CheckPoint
{
    XLogRecPtr          redo;
    TimeLineID          ThisTimeLineID;    // 目前的时间序列号
    uint32              nextXidEpoch;
    TransactionId       nextXid;          // 下一个可用的事务 ID
}

```

```

Oid          nextOid;           //下一个可用的 OID
MultiXactId nextMulti;
MultiXactOffset nextMultiOffset;
time_t       time;             //检查点的时间戳
}CheckPoint;

```

系统一般在正常关闭的情况下会自动地创建检查点，在日志的头信息中记录 XLOG_CHECKPOINT_SHUTDOWN，用户也可以手动的创建检查点，头信息是 XLOG_CHECKPOINT_ONLINE。检查点的日志与 COMMIT (ABORT) 的日志相似。

2. XLOG 日志管理器主要操作

XLOG 日志管理器的操作由以下几个部分组成：日志文件的建立、初始化、启动操作，日志的插入、归档、刷新、删除操作，以及日志文件的恢复、日志备份块的读取、日志的 REDO 操作等。下面分别介绍这些操作。

(1) XLOG 日志的启动

该操作定义在函数 BootStrapXLOG 中，功能是在系统引导时，初始化对 XLOG 的使用，即创建第一个 XLOG 日志文件，同时在 XLOG 中插入第一条检查点记录，并调用 BootStrapCLOG、BootStrapSUBTRANS 和 BootStrapMultiXact 对 CLOG、SUBTRANS、和 MULTIXACT 日志记录分别进行初始化。

(2) XLOG 日志文件的创建

该操作定义在 InstallXLogFileSegment 函数中，用于建立一个新的 XLOG 段文件，使其成为当前的段或一个将来的（比当前还要靠前）的段。如果成功建立文件，那么返回 TRUE，如果在最大范围内无法找到空段为新段文件使用，那么返回 FALSE。

(3) XLOG 日志的插入

该操作定义在函数 XLogInsert 中，根据一个 Rdata 链表以及相应的资源管理器信息向 XLOG 日志文件中插入一条 XLOG 记录。事务执行插入、删除、更新、提交、终止或者回滚命令时需要调用该函数。

其基本流程如下：

- 1) 首先判断缓冲区是否需要备份，若不需要备份（可能因为该 buffer 不是在最新的检查点后第一次被修改或该 buffer 已经被备份过了），记录数据在内存中的头信息和磁盘上的元组数据。若需要备份，则记录数据所在的文件号、块号以及磁盘上的数据块备份。

- 2) 然后根据 LogwrtRqst 和 LogwrtResult 来判断正在接受日志记录的缓冲区是否需要写回到磁盘，系统规定当缓冲区被写了一大半时就要写回磁盘。

- 3) 然后根据 XlogCtl 的 Insert 字段找到当前日志的插入点，开始写第一部分头信息，头信息这部分不能跨页存放。

- 4) 最后，返回可用来记录下一条日志的地址。

(4) 日志文件的归档

该操作定义在函数 XLogWrite 中，它的功能是按照 WriteRqst 指示的写请求将日志写/同步到磁盘日志文件中，调用本函数的时候需要持 WALWriteLock (日志写锁)。

(5) 日志文件的刷新

该操作定义在函数 XLogFlush 中，用于确保到达给定位置的所有 XLOG 数据都被刷写回磁盘，

调用该函数时并不需要持有日志写锁。

(6) 日志的打开操作

该操作定义在函数 XLogFileInit 中，用于创建一个新的 XLOG 段文件或者打开一个已存在的 XLOG 段文件。

其中，参数 “use_existent” 用于作为传入参数和返回参数。在传入时如果为 TRUE，那么标识允许使用已有的日志文件，否则任何已有的文件将被删除。作为返回参数时，如果为 TRUE，则说明使用一个已存在的文件。

如果参数 “use_lock” 为 TRUE，那么在将临时文件转移到目标位置时需要获取控制文件锁 ControlFileLock。除了在引导时期创建日志文件时不需要该锁外，在其他时候该锁都是必须的。但是调用者在调用本函数时不需要持有该锁。

整个过程如下所示：

- 1) 使用已有的段日志文件，如果该段文件存在的话，其实就是根据给定的段号、文件号打开该段文件。同时允许从给定段开始，在 XLOG 日志范围内寻找第一个为空的段号，若顺利找到，那么将所创建的临时文件作为一个将来的段日志文件，写到该位置；若不能找到，那么将删除所创建的临时文件，释放磁盘空间，因为不需要这个将来的段日志文件。

- 2) 如果该段文件不存在，那么在给定的段号位置创建一个新的段日志文件，并打开该文件。

- 3) 如果选择创建全新的段日志文件，此时不管给定的位置是否已经有文件存在，都将被删除，用新创建的文件替代。最后打开所创建的段日志文件。

(7) 日志文件的拷贝

该操作定义在 XLogFileCopy 中，通过拷贝一个已经存在的 XLOG 文件来创建一个新的 XLOG 段文件。该函数仅用于恢复阶段，所以没有加锁情况。

(8) 日志备份块的恢复

该操作定义在函数 RestoreBkpBlocks 中，主要功能是如果在一个 XLOG 记录中存在任何备份的块，那么将其恢复。

当一个备份块在 XLOG 中是有效的，我们无条件地将其恢复，即使数据库中的页面比它更新。这是为了防止数据库页面在崩溃过程中部分或者不正常地进行了写操作。

(9) 日志记录的读取操作

该操作定义在函数 ReadRecord 中，用于尝试读一条 XLOG 记录。如果 RecPtr 不是空的，尝试在该位置读一条记录。否则尝试在前面读入的最后一条记录之后读入一条记录。如果没有任何有效的记录，返回 NULL。

(10) 创建检查点。

该操作定义在函数 CreateCheckPoint 中，功能是建立不同的检查点。在 PostgreSQL 中，检查点有以下四种类型：

- 1) CHECKPOINT_IS_SHUTDOWN；数据库关闭时建立的检查点。
- 2) CHECKPOINT_END_OF_RECOVERY；数据库重启恢复完毕后建立的检查点。
- 3) CHECKPOINT_IMMEDIATE；尽可能快速刷新日志的检查点（ASAP）
- 4) CHECKPOINT_FORCE；即使 XLOG 为空，也要建立检查点。

创建检查点的基本步骤如下：

1) 首先检查要创建的检查点的类型是否是与系统关闭相关的检查点（例如 CHECKPOINT_IS_SHUTDOWN、CHECKPOINT_END_OF_RECOVERY），这时候需要去更新 ControlFile 文件信息。

2) 如果现在要创建的检查点不是 CHECKPOINT_IS_SHUTDOWN、CHECKPOINT_END_OF_RECOVERY 或 CHECKPOINT_FORCE，并且该记录插入的位置就是前一个检查点，此时为了节省空间，不再创建检查点，直接返回。

3) 然后检查此时是否还有事务处于提交阶段，若有则等待直到没有事务处于提交阶段。

4) 更新全局 ShmemVariableCache、ControlFile、MULTIXACT 等信息。

5) 调用 CheckPointGuts 将 CLOG、SUBTRANS、MULTIXACT、各种缓冲区以及两阶段提交相关数据刷新到磁盘。

6) 建立检查点日志，并刷新到磁盘。

7) 最后更新 ControlFile 信息，重置当前系统的最新检查点位置。

(11) 创建 RestartPoint

RestartPoint 是系统在恢复的过程中创建的日志恢复检查点。需要注意以下两点：

1) 如果此时系统不在恢复过程中，则无需创建 RestartPoint。

2) 如果当前恢复的位置已经包含在一个检查点之中，则不需要创建。

因为系统在恢复的过程中也可能崩溃，所以 RestartPoint 相当于日志恢复过程中的检查点。

3. XLOG 日志恢复策略

在 PostgreSQL 中，系统在崩溃后重新启动时会调用 StartupXlog 人口函数，该函数首先会扫描全局信息控制文件（global/pg_control）读取系统的控制信息，然后扫描 XLOG 日志目录的结构检测其是否完整，进而读取到最新的日志检查点记录，接下来根据日志记录序列的偏序关系检测到系统是否处于非正常状态下，若系统处于非正常状态，则触发恢复机制进行恢复。恢复完成后重新建立检查点并初始化 XlogCtl 控制信息，然后启动事务提交日志以及相关辅助日志模块。

日志中出现以下三种情况需要进行恢复操作。

1) 日志文件中扫描到 backup_label 文件。

2) 根据 ControlFile 记录的最新检查点读取不到日志记录。

3) 根据 ControlFile 记录的检查点与通过该记录找到的检查点日志中的 Redo 位置不一致。

在 PostgreSQL 系统中，日志建立的策略是采用改进的非静止检查点的 Redo 日志，恢复则是找到最近的合法检查点然后做 Redo 操作。恢复操作的具体步骤如下：

1) 首先更新控制信息到 ControlFile 中。

2) 初始化日志恢复时所用的资源管理器。

3) 从检查点日志记录的 REDO 位置开始往后读取日志记录。

4) 根据日志记录的资源管理器号选择对应的 RMGR，然后利用该 RMGR 做日志记录中所记录的操作过程（REDO 操作）。

5) 重复步骤 3 和 4 的过程，直至读取不到日志记录。

在上述恢复流程中，第 4 步的 REDO 操作会针对不同的日志类型做不同的恢复操作，下面介绍几种典型日志的恢复操作。

1) Database 类型的日志操作。这种类型日志中没有备份块，可能的操作有 Create/Drop。对于 Create 操作，首先强制刷新所有缓冲区，然后将日志中记录的源 DB 目录拷贝到新 DB 目录即可。对

于 Drop 操作，直接删除该数据库对应的缓冲区即可。

2) Heap 类型的 Redo 操作。首先根据日志序列号 (LSN) 以及日志记录找出是否有备份块，如果有则恢复备份块到 Page 中，并置页面为“脏”。然后根据日志标志位选择对应的操作，典型的有 INSERT/DELETE/UPDATE 操作。对于这几种操作，首先判断该日志记录是否有备份块，如果有则说明已恢复，直接返回即可；如果没有则需要读取日志重建 HeapTuple。

3) B-Tree 类型的 Redo 操作。Btree 是一种较复杂的索引结构，涉及的操作有叶子节点的插入操作，以及针对不同的位置（例如根节点、叶子节点、左右子树）节点的分割操作等，所以恢复时根据标志信息确定做哪种类型的 Redo 操作。

4) Xlog 类型的 Redo 操作。因为系统崩溃是不确定的，所以对于 Xlog 日志的操作也需要有日志记录。Xlog 类型的日志包括记录下一个可分配的 OID，设置检查点等操作。恢复的时候比较简单，将日志记录原信息拷贝出来即可。

当恢复过程完成后，重新建立检查点并重新初始化 XlogCtl 结构信息。然后继续调用 StartupCLOG 函数、StartupSUBTRANS 函数和 StartupMultiXact 函数完成事务提交日志及其他辅助日志模块的启动。如果根据系统日志检测到系统不需要恢复操作，那么将跳过恢复操作，然后完成事务提交日志等相关模块的初始化。

7.11.6 日志管理器总结

在 PostgreSQL 系统中，日志模块包括事务提交日志 CLOG 和数据日志 XLOG。其中，CLOG 是系统为整个事务管理流程所建立的日志，主要用于记录事务的状态，同时通过 SUBTRANS 日志记录事务的嵌套关系。而 XLOG 才是数据库日志的主体，记录数据库中所有数据操作的变化过程，因此，XLOG 日志管理器模块提供了日志操作的 API 供其他模块调用。

在 PostgreSQL 中，日志伴随着系统启动开始，始终在完成插入、刷新这样反复的过程，一旦系统崩溃还要完成数据恢复的功能。下面结合 PostgreSQL 总体控制模块及其他后台进程详细阐述日志模块在系统中的总体流程（如图 7-11 所示）。

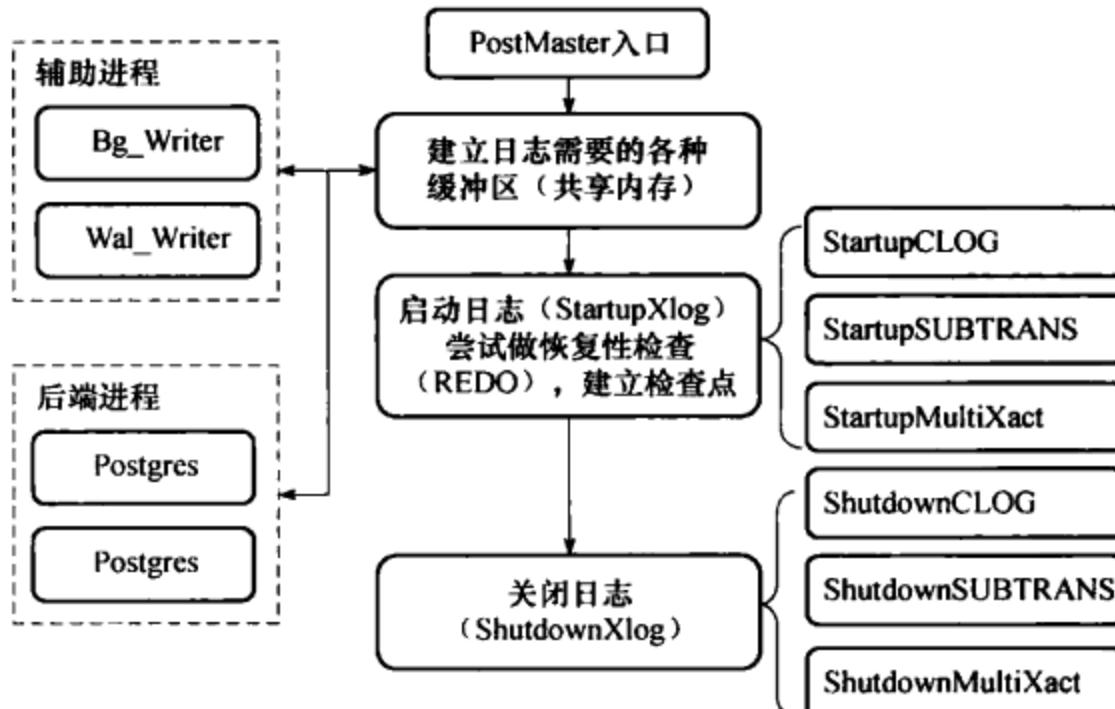


图 7-11 日志处理模块关系图

PostgreSQL 系统在启动时，首先会由 Postmaster 调用 InitCommunication 函数，进而调用 BaseInit 完成各种共享内存变量的分配和初始化，然后调用 StartupXlog 完成日志模块的启动。各种后端进程通过共享内存操作日志记录信息，其中 BgWriter 辅助建立检查点并刷新相关信息，WalWriter 辅助进程将在事务提交时设定的同步刷新点之前的日志刷写到磁盘，从而保证在写数据之前将日志信息刷新到磁盘的特性（WAL）。日志模块的启动和关闭是通过 StartupXLog 和 ShutupXlog 函数完成的。

7.12 小结

如果我们把整个数据库系统视为一个团队，那么事务系统在这个团队中扮演了“指挥官”的角色，它根据外部用户命令以及系统内部状态决定当前数据库系统中操作的执行方向。在 7.1 节到 7.5 节中，详细阐述了 PostgreSQL 事务管理器的操作，其核心功能就是根据当前状态和接收到的外部状态（用户命令）决定当前需要执行的操作，所以它承担了整个系统的决策功能，是“指挥官”的大脑。

在 7.6 节到 7.10 节中，详细阐述了 PostgreSQL 的并发控制机制，其核心功能是为了在保证数据一致性的前提下提高并发度，所以它承担了整个系统的调度协调功能，是“指挥官”的节拍器。

在 7.11 节中，详细阐述了 PostgreSQL 的日志管理机制，其核心功能是通过磁盘日志文件来记录数据库操作状态序列以及数据变化过程，所以它承担了整个系统的保障恢复功能。

总之，事务系统串联了整个数据库中各个不同的模块，事务系统的调度决策驱动了整个数据库系统的执行进程。

习题

习题 7.1 利用 BEGIN、END、SAVEPOINT、ROLLBACK 等命令控制一个事务的执行流程，并利用 GDB 跟踪其执行过程，然后解释你所看到的现象。

习题 7.2 死锁是怎样产生的？检测和预防的过程是怎样的？请设计一个会产生死锁的语句执行序列，在两个不同的终端中执行，并根据 7.9 节的阐述，给出你的解释。

习题 7.3 阐述 PostgreSQL 中并发控制的实现机制。

习题 7.4 结合数据库理论中介绍的 Undo、Redo、Undo/Redo 三种不同的日志类型，分析 PostgreSQL 的日志恢复方式。

第 8 章

数据库安全



数据库安全技术是指保护数据库以防止未授权用户窃取、篡改和破坏数据库中数据信息的技术。由于数据库系统中存储着大量重要的数据和各类敏感的信息，并且为持有不同权限的合法用户提供了数据共享服务，因此保护数据库的安全至关重要。数据库用户也最关心自身数据隐私的安全，而这种安全从根本上还是依赖于数据库系统提供的安全保障。

作为一种先进的对象关系型数据库管理系统，为了防止数据泄露，PostgreSQL 采取了一系列保护措施来保证数据库的安全性。本章将介绍和分析 PostgreSQL 中采用的安全技术。

8.1 PostgreSQL 安全简介

1994 年 4 月，美国国家计算机安全中心（NCSC）颁布了 TDI（即“可信计算机系统评估标准在数据库管理系统中的解释”），它将 TCSEC（可信计算机安全评价标准，由美国国防部于 1985 年颁布）扩展到了数据库领域。在 TCSEC 中，系统安全从四个方面划分为七个不同的安全级别，即 D、C1、C2、B1、B2、B3、A1，其中 A1 级别最高，D 级别最低。TDI 沿用了 TCSEC 的做法来描述每级的安全性。目前，国内外一些大型数据库系统都满足 C2 级以上的要求，即满足 C1 级别要求且具有审计机制。而 PostgreSQL 系统在 8.4.1 版本中依然缺乏明确意义上的审计机制，因此可认为其安全级别还是 C1 级。

PostgreSQL 首先通过用户标识和认证来验证访问数据库的用户身份，判断其是否为合法用户以及是否有权限访问数据库资源。然后通过基于角色的访问控制（Role Based Access Control，RBAC），并使用存取控制列表（ACL）方法控制访问请求和保护信息。

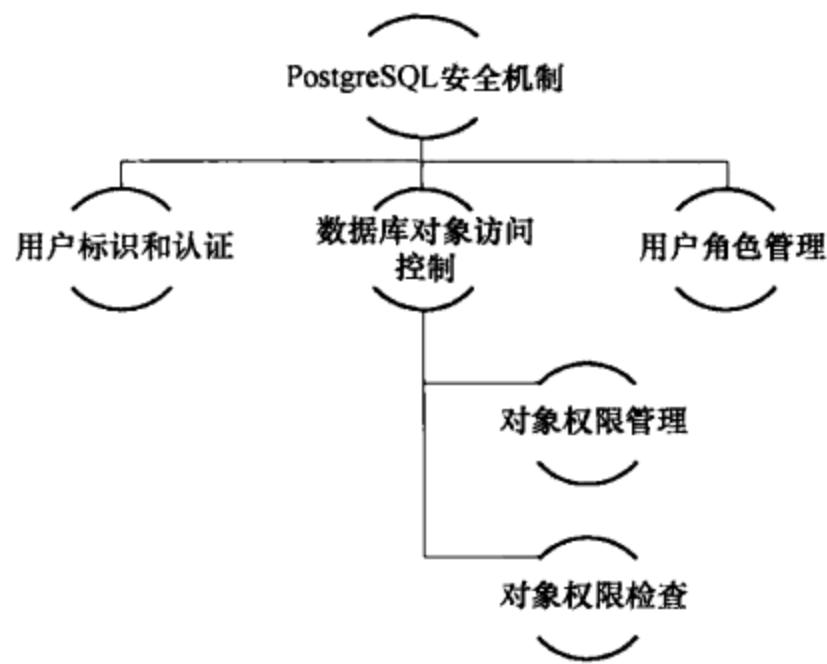


图 8-1 PostgreSQL 安全机制体系

(1) 用户标识和认证

PostgreSQL 通过用户标识和认证为系统提供最外层的安全保护措施。在客户端访问数据库资源之前，服务器首先需要通过身份认证模块来验证用户的合法身份，从而在数据库系统的前端之间建立安全的通信信道，防止非法用户连接数据库，保证只有合法的用户才能访问数据库资源。

(2) 用户角色管理

PostgreSQL 引入了角色概念进行权限管理。角色是权限的集合，管理员根据角色在系统中承担的职责分配具有不同权限的角色。用户权限的管理方式被简化成对角色的管理，这种改变不仅简化了权限管理过程，而且提高了权限管理的有效性和安全性。

(3) 数据库对象的访问控制

PostgreSQL 数据库对象包括表 (table)、序列 (sequence)、数据库 (database)、函数 (function)、语言 (language)、模式 (schema) 和表空间 (tablespace) 等。当用户或者进程在对数据库对象进行任何操作时，都必须进行相应的权限检查。同时，用户对象的权限发生变更时，需要对用户在对象上的操作权限进行维护和更新。

8.2 用户标识和认证

第2章曾介绍过，PostgreSQL 系统是一个基于客户端/服务器端 (C/S) 模式的数据库系统，每一个 PostgreSQL 会话都由后台服务进程 Postgres 和客户端进程组成。当客户端需要访问数据库时，先通过网络（通常是 TCP/IP 协议）将请求发送给服务器端的守护进程 Postmaster。Postmaster 会根据这个请求信息进行客户端认证。如果通过客户端认证，Postmaster 就为该客户端新建一个后台服务进程 Postgres，之后就由 Postgres 为客户端完成各种命令。因此，一次完整的客户端认证过程如图 8-2 所示：

- 1) 客户端和服务器端的 Postmaster 进程建立连接。
- 2) 客户端发送请求信息到守护进程 Postmaster。
- 3) Postmaster 根据请求信息检查配置文件 pg_hba.conf 是否允许该客户端连接，并把认证方式和必要信息发送到客户端。

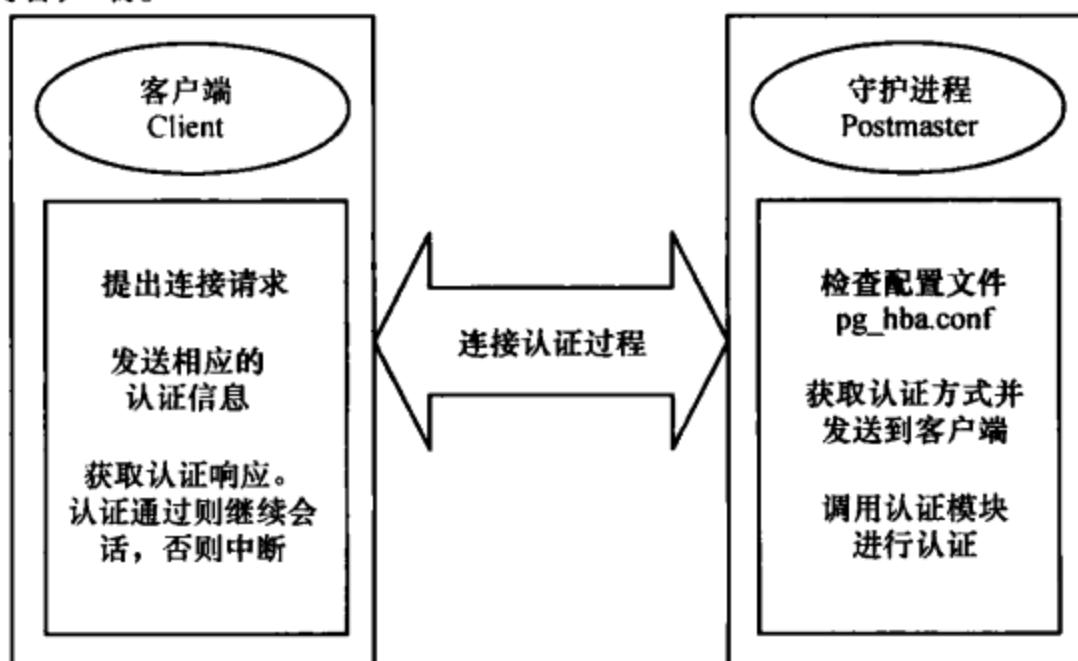


图 8-2 客户端连接认证过程

- 4) 客户端根据收到的不同认证方式，发送相应的认证信息给 Postmaster。
- 5) Postmaster 调用认证模块对客户端送来的认证信息进行认证。如果认证通过，初始化一个 Postgres 进程与客户端进行通信；否则拒绝继续会话，关闭连接。

上述的客户端认证过程实质上就是 PostgreSQL 系统进行用户标识和认证的过程。作为 PostgreSQL 系统的安全层面上的第一层保护，用户标识和认证机制可以防止非法用户进入到数据库内部，并通过丰富的认证方法加强认证的安全性和准确性。除了简单的口令认证之外，PostgreSQL 也支持操作系统认证、网络安全系统认证、加密协议认证等认证方式，在 8.2.2 节将对这些认证方法进行介绍。

8.2.1 客户端配置文件

从客户端认证的过程我们不难发现，配置文件 pg_hba.conf 起到了很重要的作用，客户端认证所需的配置信息都存放在其中。文件 pg_hba.conf（下文简称为 HBA 文件）存放在数据集簇的根目录下。HBA 文件是以文本文件的方式存储的，其中的内容是一组 HBA 记录的集合，文件中的每一行代表一条记录。其中，空白行及用井号 (#) 注释的行均被忽略。一条 HBA 记录由若干用空格和制表符分隔的字段组成。

HBA 记录中包括连接类型、客户端 IP 地址范围（与连接类型相关）、数据库名、用户名、匹配参数的连接所使用的认证方法。HBA 记录可以是下面七种格式之一：

local	database	user	auth-method	[auth-option]	
host	database	user	CIDR-address	auth-method	[auth-option]
hostssl	database	user	CIDR-address	auth-method	[auth-option]
hostnoss1	database	user	CIDR-address	auth-method	[auth-option]
host	database	user	IP-address	IP-mask	auth-method [auth-option]
hostssl	database	user	IP-address	IP-mask	auth-method [auth-option]
hostnoss1	database	user	IP-address	IP-mask	auth-method [auth-option]

其中各个字段的含义如下：

- 1) local：该类记录匹配企图通过 Unix 域套接字进行的连接。没有这种类型的记录，就不允许 Unix 域套接字的连接。
- 2) host：该类记录匹配企图通过 TCP/IP 进行的连接。host 类型的记录也可以匹配 SSL 和非 SSL 的连接请求。

注意 除非服务器带有合适的 listen_addresses 配置参数值启动，否则将不可能进行远程的 TCP/IP 连接，因为缺省的 listen_addresses 中只设置了本地回环地址 localhost，因此非本地回环地址的连接会被拒绝。

- 3) hostssl：该类记录匹配企图使用 TCP/IP 的 SSL 连接。但必须是使用 SSL 加密的连接。要使用这个选项，编译 PostgreSQL 的时候必须打开 SSL 支持。而且在服务器启动的时候必须打开 SSL 配置选项。
- 4) hostnoss1：该类记录与 hostssl 相反，它只匹配那些在 TCP/IP 上不使用 SSL 的连接请求。
- 5) database：声明记录所匹配的数据库名称。可以通过用逗号分隔的方法声明多个数据库，也

可以通过前缀@ 来声明一个包含数据库名的文件。这个字段除了实际的数据库名之外，还可以使用以下几种预定义的值：

- all：表明该记录匹配所有数据库。
- sameuser：表示该记录匹配和请求的用户名同名的数据库。
- samerole：表示请求的用户必须是一个与数据库同名的角色中的成员。

6) user：声明记录所匹配的数据库用户。这个字段可以是一个特定的用户名，或者是一个组名。当为一个组名时，组名前面应加上前缀“+”号。同时允许用“all”声明该记录匹配于所有用户。

注意 在 PostgreSQL 里，用户、角色和组没有真正的区别，+ 实际上只是意味着“匹配任何直接或者间接属于这个角色的成员”，而没有+记号的名字只匹配指定的角色。多个用户名可以通过用逗号分隔的方法声明。一个包含用户名的文件可以通过在文件名前面加前缀@ 来声明。

7) CIDR-address：声明可以匹配的客户端机器的 IP 地址范围。其格式由一个标准的点分十进制的 IP 地址（只能用数值而不能用域或主机名）和一个 CIDR 掩码长度组成。掩码长度表示客户 IP 地址必须匹配的高位二进制位数。在给出的 IP 地址里，该长度右侧的二进制位必须为零。IP 地址、“/” 和 CIDR 掩码长度之间不允许留有空白。例如，172.20.143.89/32 表示一个主机，172.20.143.0/24 表示一个小子网，10.6.0.0/16 表示一个大子网。如需声明单个主机，IPv4 地址情况下指定 CIDR 掩码为 32，IPv6 地址情况下指定为 128。以 IPv4 格式给出的 IP 地址会匹配拥有对应地址的 IPv6 连接，如 127.0.0.1 将匹配 IPv6 地址 ffff: 127.0.0.1。以 IPv6 格式给出的记录将只匹配 IPv6 连接，即使对应的地址在 IPv4-in-IPv6 范围内。注意，如果系统 C 库不支持 IPv6 地址，那么含有 IPv6 格式的记录将被拒绝。

8) IP-address、IP-mask：这两个字段的组合可以用于补充 CIDR-address 表示法。IP-mask 不是掩码的长度，而是一个完整的掩码。比如，255.0.0.0 等同于 IPv4 CIDR 的长度为 8 的掩码，而 255.255.255.255 等同于 CIDR 的长度为 32 的掩码。

9) auth-method：表示通过这条记录连接的时候使用的认证方法。下一节将对认证方法给出详细的介绍。可用的认证方法有：

- trust：无条件地允许连接。这个方法允许被该 HBA 记录匹配的客户端直接连入数据库。
- reject：无条件地拒绝连接。常用于“过滤”某些主机。
- md5：要求客户端提供一个 MD5 加密的口令进行认证。
- password：要求客户端提供一个未加密的口令进行认证。因为口令是以明文形式在网络上传递的，所以我们不应该在不安全的网络上使用这种方式。并且它通常还不能和线程化的客户端应用一起使用。
- gss：通过 GSSAPI 认证用户，只有在进行 TCP/IP 连接的时候才能使用。
- sspi：通过 SSPI 认证用户，只有在 Windows 下才可以使用。
- krb5：用 Kerberos V5 认证用户，只有在进行 TCP/IP 连接的时候才能使用。
- ident：获取客户端的操作系统用户名，通过映射文件 pg_ident.conf（记录了操作系统用户和数据库用户的对应关系）来判断操作系统用户名是否可以作为被允许的数据库用户名来连接数据库。

- ldap：使用 LDAP 进行认证。
 - cert：通过 SSL 客户端进行认证。
 - pam：使用操作系统提供的可插入认证模块服务（PAM）来认证。
- 10) auth-option：这个可选字段的含义取决于选择的认证方法。

客户端认证时数据库系统会对每个连接请求顺序地检查 HBA 文件里的记录，所以这些记录的顺序是非常关键的。通常，在数据库安全管理员设置 HBA 文件时，位置靠前的记录会采用比较严的连接匹配参数和比较松的认证方法，而靠后的记录会采用比较松的匹配参数和比较严的认证方法。

提示 一个用户要想成功连接到特定的数据库，不仅需要通过 HBA 文件的检查，还必须要有该数据库上的 CONNECT 权限。如果希望限制某些用户能够连接到某些数据库，赋予/撤销 CONNECT 权限通常比在 HBA 文件中设置规则简单。

8.2.2 认证方法

PostgreSQL 支持多种认证方法，包括信任认证、口令认证、kerberos 认证、基于身份（Ident）认证、LDAP 认证等。根据功能，这些认证方法可以归纳为以下几种认证模式：

- 1) 基于主机的认证：服务器端根据客户端的 IP 地址、用户名及要访问的数据库来查看配置文件从而判断用户的认证方式。如果是 trust 认证，则直接通过认证。
- 2) 口令认证：包括加密和非加密口令认证。合法的用户名和口令可以创建数据库连接，非法用户将被拒绝。
- 3) 操作系统认证：与操作系统集成完成用户认证。数据库中允许的操作系统用户可以不输入数据库用户名和密码直接与数据库连接。
- 4) 第三方认证：包括 kerberos 认证和 ident 认证。利用第三方认证服务器来进行认证。
- 5) SSL 加密：利用 OpenSSL 提供前后端进行安全连接的环境。

1. 信任（Trust）认证

若声明了 Trust 认证模式，PostgreSQL 无条件允许连接，即任何可以连接到服务器的客户端都可以通过任何数据库用户名（包括超级用户）连接，服务器端完全信任该连接而无需用户名和密码。

这个方法适用于那些服务器网络端口已经有足够保护措施的环境里。

2. 口令认证

PostgreSQL 的口令认证又分为明文口令（Password）认证和加密口令（MD5）认证。

明文口令认证要求客户端提供一个未加密的口令进行认证，传送过程较简单但安全性较差，一般不建议采用此认证方法。加密口令认证要求客户端提供一个经过 MD5 加密的口令进行认证，该口令在传送过程中使用了结合 salt（服务器发给客户端的随机数）的单向 MD5 加密，增强了安全性。

PostgreSQL 的数据库口令都存储在系统表 pg_authid 里，口令可以使用 SQL 语言命令 CREATE USER 和 ALTER USER 等管理（如 CREATE USER test WITH PASSWORD ‘1234’;），如果没有明确

设置口令，那么存储的口令为空并且该用户的口令认证总会失败。

3. Kerberos 认证

Kerberos 是一种适用于在公共网络上进行分布计算的工业标准的安全认证系统，Kerberos 只提供安全认证，但并不加密在网络上传输的查询和数据。

只有在进行 TCP/IP 连接的时候才能用 Kerberos V4 认证用户。Kerberos V5 是 Kerberos V4 的改良版，主要特点是不再依赖 DES 算法，同时增加了一些新特性。

PostgreSQL 可以支持 Kerberos V5，但 Kerberos 支持选项必须在编译的时候打开。PostgreSQL 运行时像一个普通的 Kerberos 服务，其名字是：servicename/hostname@realm。

4. GSSAPI 及 SSPI 认证

Kerberos 认证支持工业级安全 API：Windows 平台下的 SSPI 和 Unix、Linux 平台下的 GSSAPI，其中 GSSAPI 只有在 TCP/IP 连接时运行，SSPI 认证也只有当服务器与客户端均在 Windows 上时运行。

5. 基于身份 (Ident) 的认证

基于身份的认证方式将首先获取客户端的操作系统用户名，然后通过映射文件（记录了操作系统用户名和数据库用户的对应关系）来判断操作系统用户名是否可以作为允许的数据库用户名来连接数据库。判断客户端的用户名是非常关键的安全点，根据连接类型的不同，它的实现方法也略有不同。

(1) 基于 TCP/IP 的身份认证

该身份认证服务器可以判断是哪个用户从端口 X 连接到服务器的端口 Y。因为在建立物理连接后，服务器既可以知道 X，也可以知道 Y，因此它可以询问进行尝试连接的客户端的主机，并且理论上可以用这个方法判断发起连接的操作系统用户。

(2) 基于本地套接字的身份认证

身份认证也可以在使用 Unix 域套接字的系统上用于认证本地连接。这种情况下使用身份认证不会增加安全风险，实际上这也是在此类系统上使用本地连接时的首选方法。

(3) Ident 映射

当使用以身份为基础的认证时，在获取了初始化连接的操作系统用户名后，PostgreSQL 将判断它是否可以以它所请求的数据库用户的身份进行连接。这个判断是由 HBA 记录中 ident 字段后面的身份映射控制的。有一个预定义的身份映射是 sameuser，表示任何操作系统用户都可以以同名数据库用户进行连接。其他映射必须手工创建，非 sameuser 的身份映射定义在身份映射文件（缺省名 pg_ident.conf）里，该文件缺省存放在数据集簇的根目录下。身份映射文件包含下面通用的格式：

```
map-name ident-username database-username
```

map-name 是在 HBA 记录中引用这个映射的名称。另外两个字段分别说明哪个操作系统用户被允许以哪个数据库用户的身份进行连接。同一个 map-name 可以重复用于声明更多的用户映射。操作系统用户和数据库用户之间的映射是多对多的关系，即一个操作系统用户可以映射为多个数据库用户，一个数据库用户也可以映射多个操作系统用户。

6. LDAP 认证

LDAP 只用于验证用户名/口令对。因此，在使用 LDAP 进行认证之前，用户必须已经存在于数

据库里。用户可以在 HBA 文件的 `ldap` 关键字后面提供 LDAP 服务名，该参数的格式为：

```
ldap[ s ]://servername[:port]/base dn[ ;prefix[ ;suffix ]]
```

如果指定了 `ldaps` 而不是 `ldap`，那么 PostgreSQL 将使用安全传输层协议（Transport Layer Security, TLS）加密连接。需要注意的是，这仅仅加密 PostgreSQL 服务器与 LDAP 服务器之间的连接，而客户端与 PostgreSQL 服务器之间的连接并不受此影响。要使用 TLS 加密，需要在配置 PostgreSQL 之前先配置好 LDAP 库。

7. PAM 认证

这个认证类型操作起来类似 `password`，只不过它使用操作系统提供的可插入的认证模块服务（Pluggable Authentication Modules, PAM）作为认证机制。缺省的 PAM 服务名是“`postgresql`”。用户可以在 `pam` 关键字后面提供自己的可选服务名。

8. cert 认证

这种认证方法通过使用 SSL 客户认证进行认证，因此它仅应用于 SSL 连接。在这种认证方式下，客户端与服务器之间的数据是经过加密的，服务器端将要求客户端提供一个有效的证书。在连接建立后，客户端将向服务器端发出一个握手信息。作为回应，服务器会返回一个主密钥给客户端。在随后的连接过程中，客户端用服务器的公钥加密主密钥发送给服务器作为确认，服务器再用自己的私钥解密。这样，双方便有了一个共同的会话密钥，为后来的数据传输提供保密通信。

8.2.3 客户端认证

在 2.4.7 节中我们已经介绍过，Postmaster 接收到一个用户的连接请求之后会调用 `ClientAuthentication` 函数进行客户端认证。在进行认证之前，首先需要读取配置文件 `pg_hba.conf` 和 `pg_ident.conf`。前者用于决定认证的方法，后者主要用于初始化 Ident 认证的用户名映射。如果 HBA 文件加载失败，将直接报错。如果装载成功，则设置一个认证超时定时器，默认为 60 秒。超过时间将提示认证超时。

读取 HBA 文件的工作由 `load_hba` 函数实现，它的功能是按行解析 HBA 文件中的数据，并将其保存在一个节点类型为 `HbaLine`（数据结构 8.1）的链表结构里，该结构用于存储配置文件 HBA 中每个记录的各个属性，包括连接类型、数据库名、用户名、IP 地址、认证方法以及每种方法对应的变量等。函数 `load_hba` 是一个简单的按行解析文本文件的过程，其具体实现过程就不再分析。

数据结构 8.1 HbaLine

```
typedef struct
{
    int                     linenumber;           //记录行号
    ConnType               conntype;             //连接类型
    char                   *database;            //数据库名称
    char                   *role;                 //角色名,可以是用户或用户组
    struct sockaddr_storageb addr;              //IP 地址范围
}
```

```

    struct sockaddr_storage mask;           //子网掩码,和 IP 一起使用
    UserAuth auth_method;                 //认证方法
    char *usermap;                      //用来存储 ident 认证的用户名
    char *pamservice;                   //PAM 认证的服务名
    bool ldap tls;                     //表示在 LDAP 中是否使用 TLS
    char *ldapserver;                  //LDAP 服务器
    int ldapport;                      //LDAP 端口
    char *ldapbinddn;                  //LDAP 绑定的用户名
    char *ldapprefix;                  //LDAP 名称前缀
    char *ldapsuffix;                  //LDAP 名称后缀
    char *ldapbindpasswd;              //LDAP 绑定的密码
    char *ldapsearchattribute;         //表示 LDAP 查询的属性
    bool clientcert;                  //是否可以进行证书认证
    char *krb_server_hostname;          //krb 服务器名
    char *krb_realm;                   //krb 认证使用的域
    bool include_realm;                //表示是否把此认证信息加入到 krb 域中
} HbaLine;

```

客户端认证的过程通过调用 ClientAuthentication 函数完成，该函数只有一个类型为 Port 的参数（数据结构 8.2），Port 结构中存储着客户端的相关信息（如客户端主机名、端口、用户名、数据库名等）。Port 结构与客户端认证相关的部分字段参见数据结构 8.2。

数据结构 8.2 Port

```

typedef struct Port
{
    .....
    /* 认证过程中需要的信息 */
    HbaLine *hba;           //HbaLine 结构,一组记录
    char md5Salt[4];        //md5 随机数
    .....
} Port;

```

ClientAuthentication 函数的执行流程如图 8-3 所示。

此函数执行流程归纳如下：

- 1) 调用函数 hba_getauthmethod，检查客户端地址、所连接数据库、用户名在文件 HBA 中是否有能匹配的 HBA 记录。如果能找到匹配的 HBA 记录，则将 Port 结构中的相关认证方法的字段设置为 HBA 记录中的参数，同时返回状态值 STATUS_OK。理论上这个过程不会返回 STATUS_ERROR。
- 2) 如果在编译时选择了使用 SSL，在这里先要检查客户端是否已提供一个有效的证书（通过 Port 结构中 hba 字段的 clientcert 字段的值来判断）。

- 3) 根据不同的认证方法, 进行相应的认证过程。
- 4) 在认证过程中可能需要和客户端进行多次交互。最后返回值如果为 STATUS_OK, 则表示认证成功, 并将认证成功信息发送回客户端; 否则发送认证失败信息。

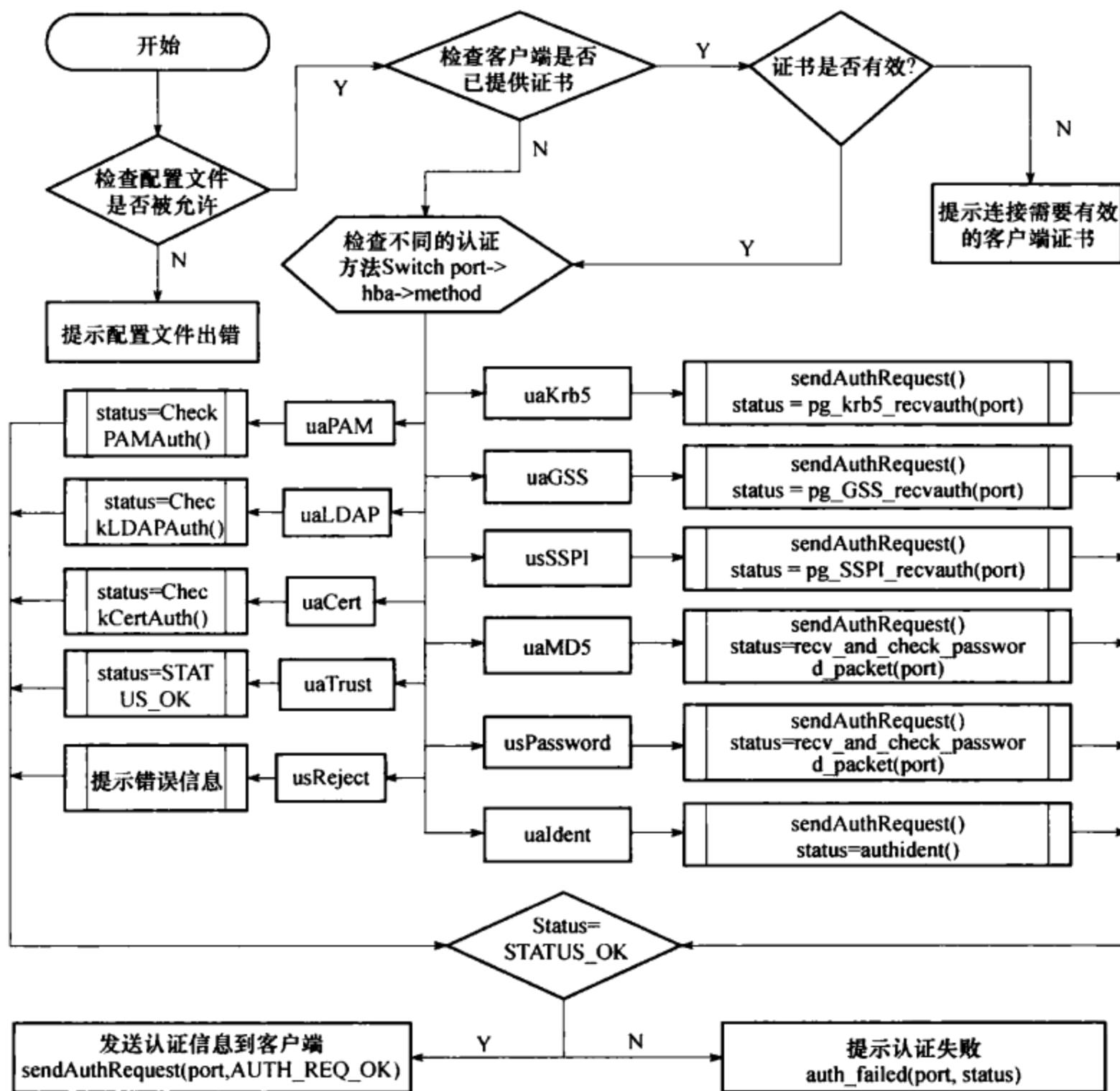


图 8-3 客户端认证流程图

在认证过程中, 服务器端要和客户端交互认证相关的信息, 有时不止一次。系统通过函数 `SendAuthRequest` 向客户端发送认证请求。此函数实现流程如下:

- 1) 初始化一个 buf。
- 2) 把请求信息转化成二进制加入 buf。
- 3) 判断是否为 MD5 认证, 是则加入 MD5 的随机数 Salt 信息。
- 4) 判断是否激活了 GSS 或 SSPI 认证, 是则将相关认证信息加入到 buf。
- 5) 调用 `pq_endmessage` 函数将 buf 内的信息发给客户端。

发送的请求信息中可能包括不同的认证信息，如表 8-1 所示。

PostgreSQL 提供多种不同的客户端认证方法，参见 8.2.2 节。这里以 MD5 认证为例描述 MD5 认证的过程：

1) 调用函数 `sendAuthRequest` 发送一个 MD5 类型的认证请求 `AUTH_REQ_MD5`。

2) 在 `sendAuthRequest` 中除了发送认证请求信息之外，还将调用 `pq_sendbytes` 函数将长度为 4 位的随机数 `md5Salt` 发送给客户端，并等待客户端回应。

3) 客户端将通过用户 `OID`、密码和 `md5Salt` 进行 MD5 加密，并将结果作为认证信息再次发送给服务器端。

4) 服务器端通过函数 `recv_password_packet` 获取客户端返回的结果。如果结果值为空，表明客户端没有发送密钥，返回 `STATUS_EOF`；如果结果值不为空，调用 `md5_crypt_verify`，获取从系统表 `pg_authid`（见表 8-2）中取出的用户的密码，同样将其和用户 `OID`、随机数 `md5Salt` 一起进行 MD5 加密。将服务器端计算的结果和客户端提交的结果进行比较，如果相等则返回 `STATUS_OK`。

5) 如果返回值为 `STATUS_OK`，则执行函数 `sendAuthRequest` 将认证成功信息发送到客户端，否则执行函数 `auth_failed` 将认证失败信息返回。

表 8-1 客户端认证请求类型

请求类型	值	描述
<code>AUTH_REQ_OK</code>	0	用户认证成功
<code>AUTH_REQ_KRB4</code>	1	KRB4 认证，已不再支持
<code>AUTH_REQ_KRB5</code>	2	KRB5 认证
<code>AUTH_REQ_PASSWORD</code>	3	明文加密认证请求
<code>AUTH_REQ_CRYPT</code>	4	密文密码认证请求，已不再支持
<code>AUTH_REQ_MD5</code>	5	MD5 加密认证请求
<code>AUTH_REQ_SCM_CREDS</code>	6	SCM 安全证书请求（Unix）
<code>AUTH_REQ_GSS</code>	7	GSS 认证请求
<code>AUTH_REQ_GSS_CONT</code>	8	GSS 二次认证请求
<code>AUTH_REQ_SSPI</code>	9	SSPI 认证请求

8.3 基于角色的权限管理

PostgreSQL 实现了基于角色的访问控制机制。角色将用户和权限联系起来，大大简化了对权限的授权管理。借助角色机制，当给一组权限相同的用户授权时，不需对这些用户逐一授权，只要把权限授予角色，再将角色授予这组用户即可。此后只要针对角色进行管理便可以了。同理，如果一组权限要改变，只需改变角色的权限，此时该角色所包含的所有成员的权限会被自动修改。本节主要描述 PostgreSQL 中角色的概念以及如何进行角色创建、删除等操作。

8.3.1 用户和角色

在 PostgreSQL 8.1 之前，用户和组都是独立类型的记录，但之后的版本就不再区分用户和组，而是用角色的概念统一了用户和组。一般来说，一个角色可以视为一个数据库用户，或者一组数据库用户。对于 PostgreSQL 来说，用户和角色是完全相同的两个对象，唯一不同的是在创建时角色的缺省值中没有 `LOGIN` 权限。换句话说，一个拥有 `LOGIN` 权限的角色可以认为是一个用户。因此，在后面的介绍中经常会用角色来代替用户进行说明。

我们可以通过将角色赋给一个用户使得该用户拥有这个角色中的所有权限。一个用户也可以属于不同的角色，拥有不同角色的权限。另外，角色还可以拥有数据库对象（例如表、模式、序列），并可以把这些对象上的权限赋予其他角色。角色用来控制哪些用户拥有访问哪些对象的哪些权限。8.4 节中将具体分析对象上的权限控制，这里我们先介绍权限的概念。

PostgreSQL 系统中的权限分为两种：系统权限和对象权限。系统权限是指系统规定用户使用数据库的权限（如连接数据库、创建数据库、创建用户等）。对象权限是指在表、序列、函数等数据库对象上执行特殊动作的权限，其权限类型有 SELECT、INSERT、UPDATE、DELETE、REFERENCES、TRIGGER、CREATE、CONNECT、TEMPORARY、EXECUTE 和 USAGE 等。例如角色 Johnny 想要更新表 Student 的权限，可以使用 SQL 语言命令 GRANT UPDATE ON Student TO Johnny 来授予相关权限。

每个角色都含有一组属性，这些属性定义了该角色的系统权限，以及与客户端认证系统的交互。角色属性包括：

- 登录 (LOGIN)：具有 LOGIN 属性的角色才可以用来连接数据库的用户。
- 超级用户 (SUPERUSER)：具有 SUPERUSER 属性的角色拥有系统最高权限。
- 创建数据库 (CREATEDB)：定义角色是否能创建数据库。
- 创建角色 (CREATEROLE)：定义角色是否可以创建新角色。
- 口令 (PASSWORD)：设置角色的口令。
- 继承 (INHERIT)：定义一个角色是否继承它所属角色的权限。
- 连接限制 (CONNECTION LIMIT)：声明该角色可以使用的并发连接数量。

这些角色属性在创建角色时可以使用 SQL 语言命令 CREATE ROLE 直接进行设定。例如，创建一个具有创建数据库和创建角色权限的角色时，可以使用 SQL 命令 CREATE ROLE admin WITH CREATEDB CREATEROLE。同样，这些角色属性在设定以后也可以使用 SQL 命令 ALTER ROLE 来修改。

提示 创建一个具有 CREATEDB 和 CREATEROLE 权限但不是超级用户的角色是一个很好的习惯。

你可以使用这个角色进行所有日常的数据库和角色管理。这个方法可以避免以超级用户身份操作时发生误操作而导致的严重后果。

当数据集簇创建完成后，数据库实例中就默认创建了一个数据库超级用户。超级用户与初始化数据集簇时使用的操作系统用户同名（一般习惯用 postgres 用户）。为了创建其他角色，首先必须使用这个初始用户（角色）连接。每个和数据库的连接都必须以一个角色身份进行，而这个角色决定了在该连接上的初始权限。例如，psql 程序使用 -U 命令行选项声明它代表的角色。图 8-4 描述了 PostgreSQL 的用户构建过程，即以超级用户为根节点，不断创建其他超级用户或其他角色等节点。系统中的角色从逻辑上构成了一棵以超级用户为根的树。

由图 8-4 可以得到，超级用户拥有任何权限，可以创建其他超级用户，也可以创建其他角色（如具有一定权限的用户或组），或者创建数据库对象等。

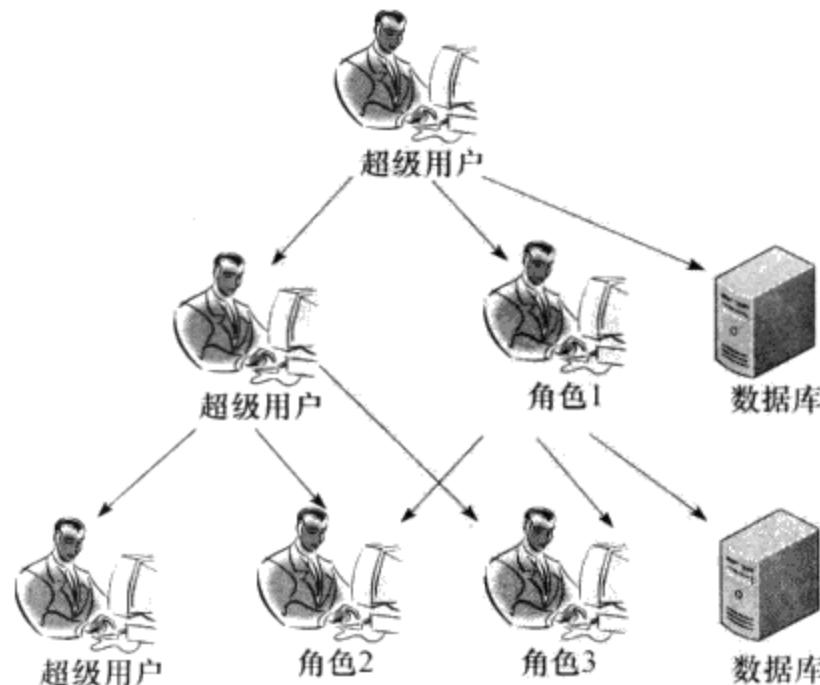


图 8-4 PostgreSQL 用户构建过程

8.3.2 角色相关的系统表

有关角色的属性信息可以在系统表 pg_authid（见表 8-2）中找到。另外，系统表 pg_roles 是系统表 pg_authid 公开可读部分的视图（口令域为空白），它提供了访问数据库角色有关信息的接口。

表 8-2 系统表 pg_authid

属性名	数据类型	注释
rolname	name	角色名称
rolsuper	bool	角色拥有超级用户权限
rolinherit	bool	角色自动继承其所属组角色的权限
rolcreaterole	bool	角色可以创建更多角色
rolcreatedb	bool	角色可以创建数据库
rolcatupdate	bool	角色可以直接更新系统表。如果这个字段没有设置为真，即使超级用户也不能这么做
rolecanlogin	bool	角色可以登录
rolconnlimit	int4	限制登录角色的最大并发连接数量。-1 表示没有限制
rolpassword	text	口令（可能是加密的），如果没有则为 NULL
rolvaliduntil	timestamptz	口令失效时间，如果没有失效期则为 NULL
rolconfig	text []	运行时配置变量的会话缺省

系统表 pg_auth_members（见表 8-3）存储了角色之间的成员关系，一般使用 GRANT/REVOKE 命令在该表上进行添加/删除角色成员操作。

表 8-3 系统表 pg_auth_members

属性名	数据类型	引用	注释
roleid	oid	pg_authid. oid	父角色的 OID
member	oid	pg_authid. oid	子角色的 OID
grantor	oid	pg_authid. oid	建立此父子角色关系的角色的 OID
admin_option	bool		如果 member 可以把 roleid 角色的成员关系赋予其他角色，则为真

PostgreSQL 对权限的管理主要是通过创建角色来包含权限，然后将其他角色作为该角色（即父角色）的成员加入，这样成员角色就可以获得父角色的权限。要创建一个角色，可以使用下面的 SQL 命令：

```
CREATE ROLE groupname;
```

一般情况下，角色是不具有 LOGIN 属性的，但是也可以进行设置。可以用 SQL 命令 GRANT 和 REVOKE 添加和撤销角色的成员关系。例如，下面的语句完成了如下的角色设置：

- 1) 创建角色 Johnny、Tom、David。
- 2) 将 Tom 授权给 Johnny，将 David 授权给 Tom。

```
CREATE ROLE Johnny LOGIN INHERIT;
CREATE ROLE Tom NOINHERIT;
CREATE ROLE David NOINHERIT;
GRANT Tom TO Johnny;
GRANT David TO Tom;
```

这个过程可以描述为：

1) 执行命令 CREATE ROLE 创建三个拥有不同权限的角色，即在系统表 pg_authid 中插入三个元组（分别代表 Johnny、Tom 和 David）：

rolname	rolsuper	rolinherit	rolcreaterole	rolcanlogin	...
Johnny	false	true	false	true	...
Tom	false	false	false	false	...
David	false	false	false	false	...

2) 执行命令 GRANT 赋予角色间的成员关系，即在存储角色成员关系的系统表 pg_auth_members 中插入两个元组：

roleid	member	...
Tom. oid	Johnny. oid	...
David. oid	Tom. oid	...

由于角色 Johnny 是 Tom 的角色成员，并且具有继承权限，所以当角色 Johnny 连接到数据库之后，该会话将立即拥有直接赋予 Johnny 的权限加上赋予角色 Tom 的权限。而赋予 David 的权限 Johnny 无法继承，因为 David 具有 NOINHERIT 属性，所以 Tom 无法继承其权限，自然 Johnny 也无法得到这些权限。

8.3.3 角色管理

PostgreSQL 中提供了对角色的各种管理操作，如角色的创建、修改、删除、授权、回收等。

1. 创建角色

如果要创建一个角色，可以使用 SQL 命令 CREATE ROLE。这个角色可以拥有数据库对象和数据库权限。这里的角色也可以是一个用户、一个组，或者两者兼有。其语法为：

```
CREATE ROLE name[[ WITH ]option[...]]
```

其中 option 可以是：

```
SUPERUSER | NOSUPERUSER | CREATEDB | NOCREATEDB  
| CREATEROLE | NOCREATEROLE | CREATEUSER | NOCREATEUSER  
| INHERIT | NOINHERIT | LOGIN | NOLOGIN | CONNECTION LIMIT connlimit  
| [ ENCRYPTED | UNENCRYPTED ] PASSWORD 'password'  
| VALID UNTIL 'timestamp' | IN ROLE rolename[,...]  
| IN GROUP rolename[,...] | ROLE rolename[,...]  
| ADMIN rolename[,...] | USER rolename[,...] | SYSID uid
```

执行该命令的用户必须具有 CREATEROLE 权限或者是超级用户，WITH 关键字可以省略。其中，SUPERUSER、CREATEDB、CREATEROLE、INHERIT、LOGIN、PASSWORD 等参数在角色属性中已经描述过其作用。这里主要描述其他一些 option 参数。如：

- VALID UNTIL：该子句设置角色的口令失效的时间戳。如果忽略了这个子句，那么口令将永远有效。

- IN ROLE：该子句列出一个或多个现有的角色，新角色将立即加入这些角色，成为它们的成员。请注意，没有任何选项可以把新角色添加为管理员；必须使用独立的 GRANT 命令来做这件事情。
- ROLE：该子句列出一个或多个现有的角色，它们将自动添加为这个新角色的成员。这个动作实际上就是把新角色作为一个“组”。
- ADMIN：该子句类似 ROLE，只是给出的角色被增加到新角色时带有 WITH ADMIN OPTION，即在 ADMIN 关键字后指定的角色具有将其他角色加入到将被创建的新角色的权利。创建角色通过函数 CreateRole 实现，该函数只有一个类型为 CreateRoleStmt（数据结构 8.3）的参数。

数据结构 8.3 CreateRoleStmt

```
typedef struct CreateRoleStmt
{
    NodeTag          type;
    RoleStmtType     stmt_type; // 将要创建的角色类型 ROLE/USER/GROUP
    char             *role;    // 角色名
    List             *options; // 角色属性列表, DefElem 结构体
} CreateRoleStmt;
```

其中，字段 stmt_type 是枚举类型 RoleStmtType，其取值包括 ROLESTM_ROLE、ROLESTM_USER、ROLESTM_GROUP，分别代表创建角色、创建用户、创建组用户。字段 options 用来存储角色的属性信息，它是一个 DefElem（数据结构 8.4）结构。

数据结构 8.4 DefElem

```
typedef struct DefElem
{
    NodeTag          type;
    char             *defnamespace; // 节点对应的命名空间
    char             *defname;     // 节点对应的角色属性名
    Node             *arg;        // 表示值或类型名
    DefElemAction    defaction;   // set/add/drop 等行为
} DefElem;
```

在 PostgreSQL 中，并没有明显区分角色和组之间的概念。当创建的角色类型为 ROLE 或者 GROUP 时，事实上都可以认为在创建组。如果创建的角色类型是用户，则该用户可以登录数据库，即其角色属性中的登录权限一定为真。

例如，要创建一个可以创建数据库，并且可以登录的带口令的角色 merry。其口令有效期到 2011 年 1 月 1 日，并且该用户可以创建数据库。对应的 SQL 语言命令为：

```
CREATE ROLE merry WITH LOGIN PASSWORD '12345'
CREATEDB VALID UNTIL '2011-01-01')
```

这个创建过程中的 options 如图 8-5 所示。LOGIN、CREATEDB、PASSWORD、VALID 选项都在 options 中对应有一个节点。

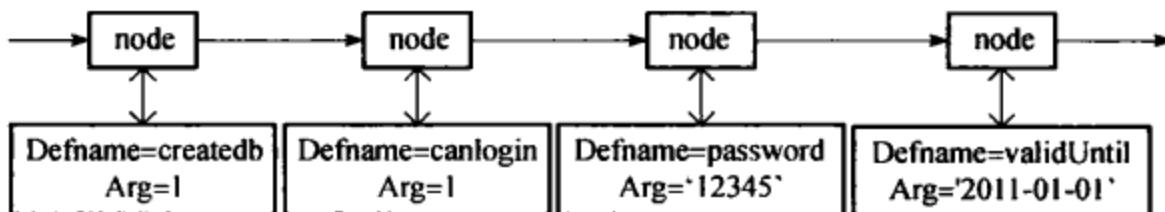


图 8-5 创建用户时的 options 实例图

同样，要创建一个拥有超级用户权限、可以创建数据库和管理角色、可以自动继承所属组权限的角色 admin 时，SQL 语言命令为：

```
CREATE ROLE admin WITH CREATEDB CREATEROLE
```

其对应的 options 如图 8-6 所示。

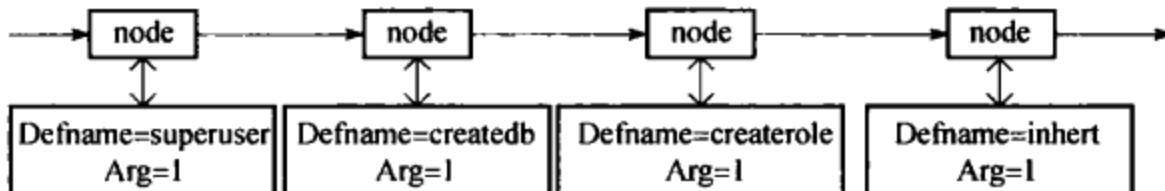


图 8-6 创建组时的 options 实例图

创建角色时，首先判断所要创建的角色类型。如果是创建用户，则设置其 canlogin 属性为 true，因为用户默认具有登录权限。而创建角色和创建组时，若角色属性即 options 参数没有声明的话，则 canlogin 默认为 false。在检查到所要创建的角色类型之后，开始循环提取角色的属性 options，并将提取的属性转换成对应的数据类型。再将转换后的角色属性值以及角色的信息一起构建一个 pg_authid 的元组，再将该元组写回系统表并更新索引。

由于角色在整个数据集簇中是全局定义的，因此角色命名必须唯一，不允许两个角色同名，创建角色时必须检查角色名是否已存在。CreateRole 的流程如图 8-7 所示。

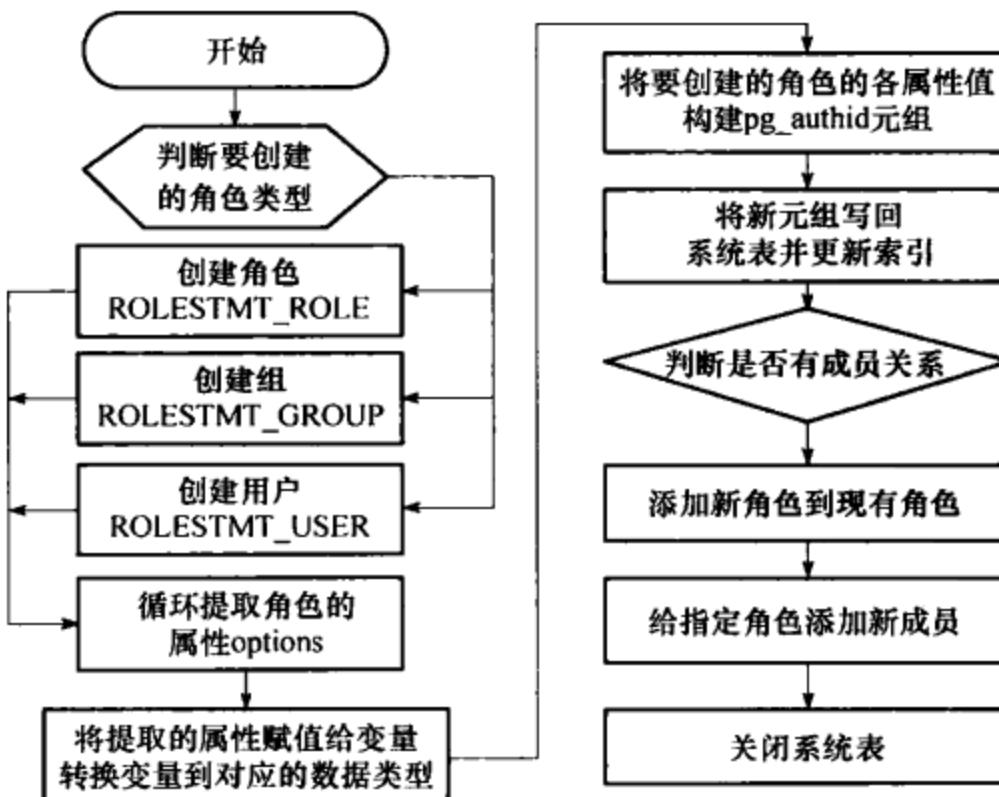


图 8-7 创建角色流程图

通常在创建角色时，可以指定新角色的成员或者将新角色加入到已有角色中，成为它们的成员。如 SQL 命令：

```
CREATE ROLE Johnny WITH CREATE ROLE INROLE Joe ROLE Tom
```

创建角色 Johnny，并将 Johnny 加入到角色 Joe 中成为 Joe 的成员，还将角色 Tom 加入到 Johnny 中成为 Johnny 的成员。

添加角色成员的过程通过调用函数 AddRoleMems 来实现，其原型函数如下：

```
AddRoleMems (
    const char *rolename,
    Oid roleid,
    List *memberNames,
    List *memberIds,
    Oid grantorId,
    bool admin_opt)
```

其中，参数 rolename 指定要被加入成员的角色名（父角色），roleid 指定要被加入成员的角色的 OID。memberNames 是要被添加进父角色里的成员名；memberIds 是要被添加进父角色里的成员 OID。grantorId 为完成 AddRoleMems 的角色的 OID。admin_opt 用于指示是否授予加入的成员加入其他成员的权限。其流程如图 8-8 所示。

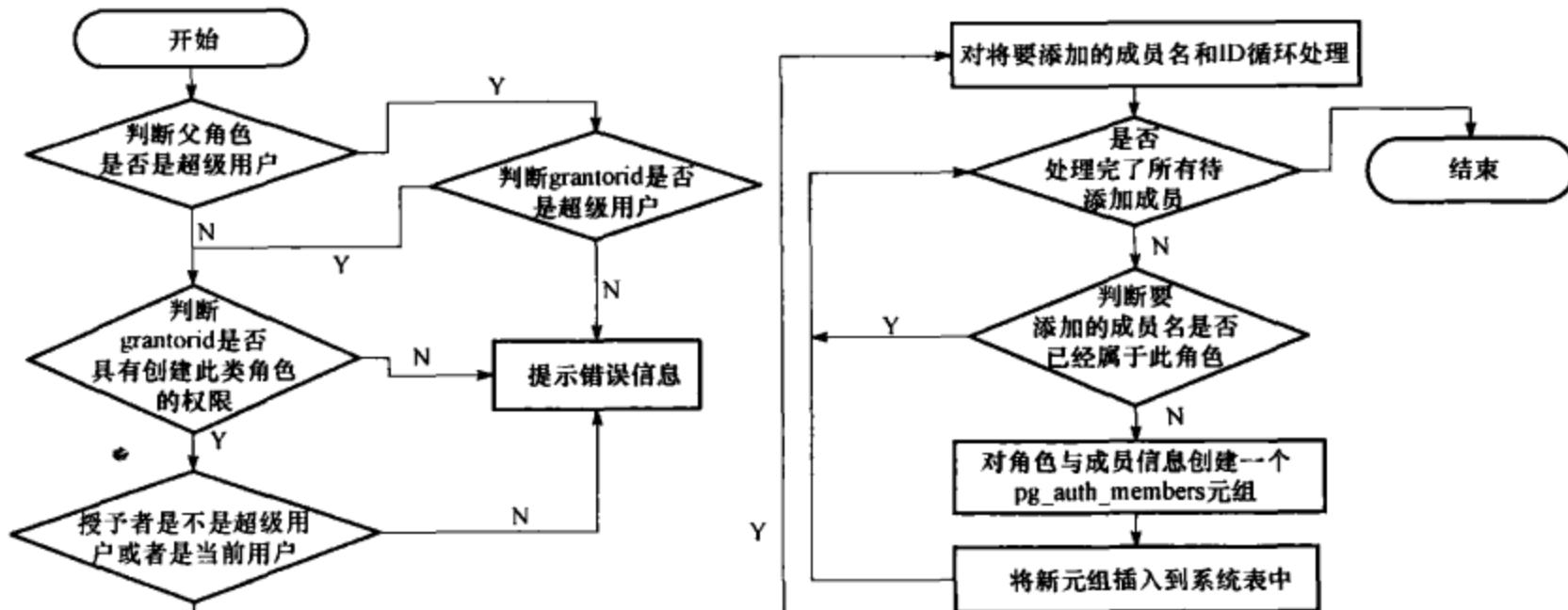


图 8-8 添加角色成员

2. 修改角色属性

如果要修改一个数据库角色，可以使用 SQL 语言命令 ALTER ROLE。该命令可以表示三种情况：修改角色属性、重命名角色或者为特定的配置变量修改角色的会话缺省值。其语法命令为：

- ALTER ROLE name[[WITH] option[…]]
 - ALTER ROLE name RENAME TO newname
 - ALTER ROLE name SET configuration_parameter{ TO | = } { value | DEFAULT }
- ALTER ROLE name RESET configuration_parameter

(1) 修改角色属性

超级用户可以使用该命令修改任何角色的任何属性，而拥有 CREATE ROLE 权限的角色也可以修改非超级用户角色的任何属性。但是只能给非超级用户角色设置。例如当需要改变角色 joe 的口令为“iter123”时，可以使用如下命令：

```
ALTER ROLE joe with password 'lter123'
```

角色属性的修改是通过调用函数 AlterRole 来实现的，该函数只有一个类型为 AlterRoleStmt（数据结构 8.5）的参数。

需要注意的是，AlterRole 函数也可以用来调整角色的成员（增加成员或者删除成员），因此结构体中的 action 字段值设置为 +1 和 -1 表示增加和删除成员关系。这个选项将在 GRANT 和 REVOKE 命令中具体描述。

修改用户角色属性时，首先循环判断 options，提取要修改的角色属性。然后查找系统表 pg_authid 判断是否已存在该角色，只有被修改的角色存在才可以修改其属性，否则提示出错。更改角色属性之前，必须检查是否有权限去更改原角色的属性。例如，只有超级用户才可以修改超级用户的角色属性。最后以被修改角色的 pg_authid 为基础新建一个元组，将要更改的属性更新到新元组中，再将新元组更新到关系表 pg_authid 中。修改角色属性的流程如图 8-9 所示。

数据结构 8.5 AlterRoleStmt

```
typedef struct AlterRoleStmt
{
    NodeTag      type;
    char         *role;           // 角色名称
    List         *options;        // 需要修改的角色属性列表
    int          action;          // +1 增加用户, -1 删除用户
} AlterRoleStmt;
```

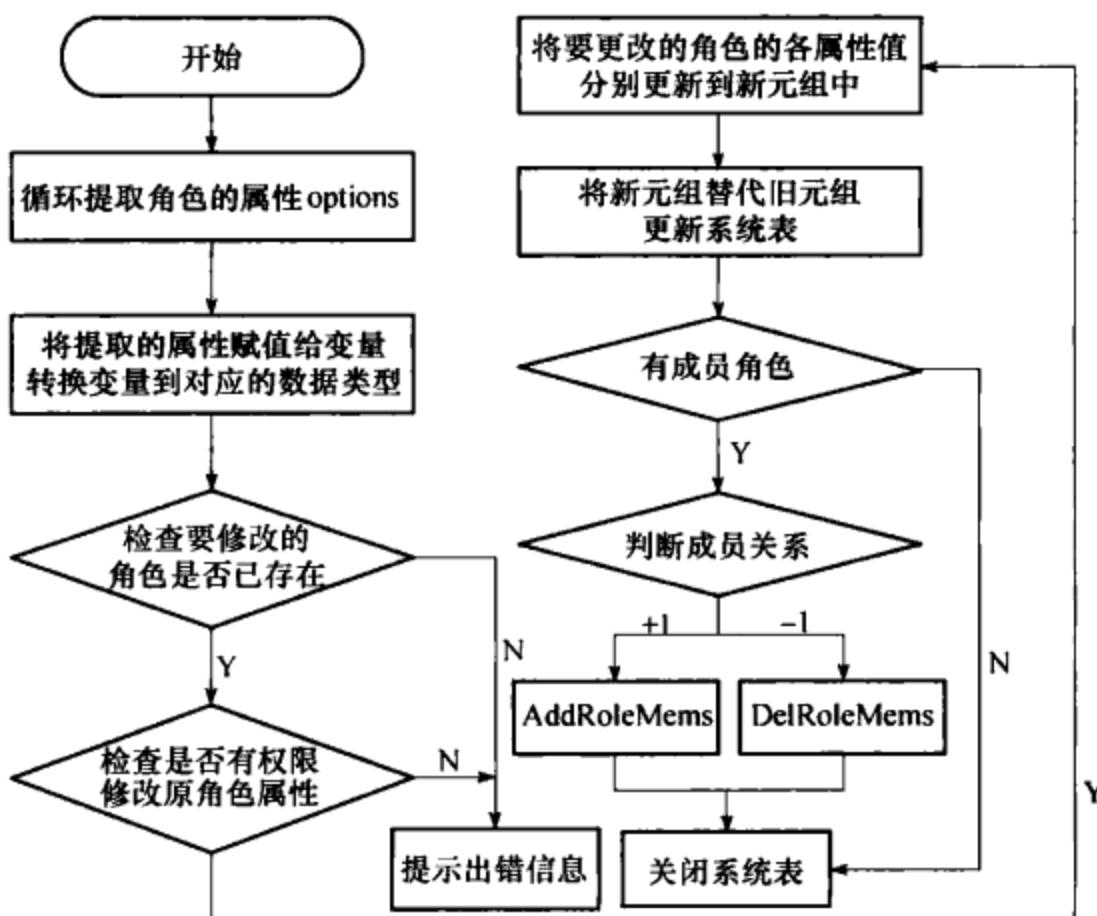


图 8-9 修改角色属性流程图

修改角色属性时，所有的属性都有可能被修改，修改命令中没有提到的属性则维持不变。若要改变角色的成员关系，即增加或删除成员，将在 GRANT/REVOKE 命令中介绍。

向角色中添加成员的过程在前面介绍 AddRoleMems 函数时进行了详细说明。而从角色中删除该成员，则是通过调用函数 DelRoleMems 来实现的，该函数原型为：

```
DelRoleMems(
    const char *rolename,
    Oid roleid,
    List *memberNames,
    List *memberIds,
    bool admin_opt);
```

其中参数 rolename 表示要从中删除成员的角色名，roleid 表示要从中删除成员的角色 OID。memberNames 表示要删除的角色成员列表，memberIds 表示要删除的角色成员的 OID 列表。admin_opt 表示是否仅仅只是撤销 memberIds 指定的成员的管理权限。删除角色成员的流程如图 8-10 所示。

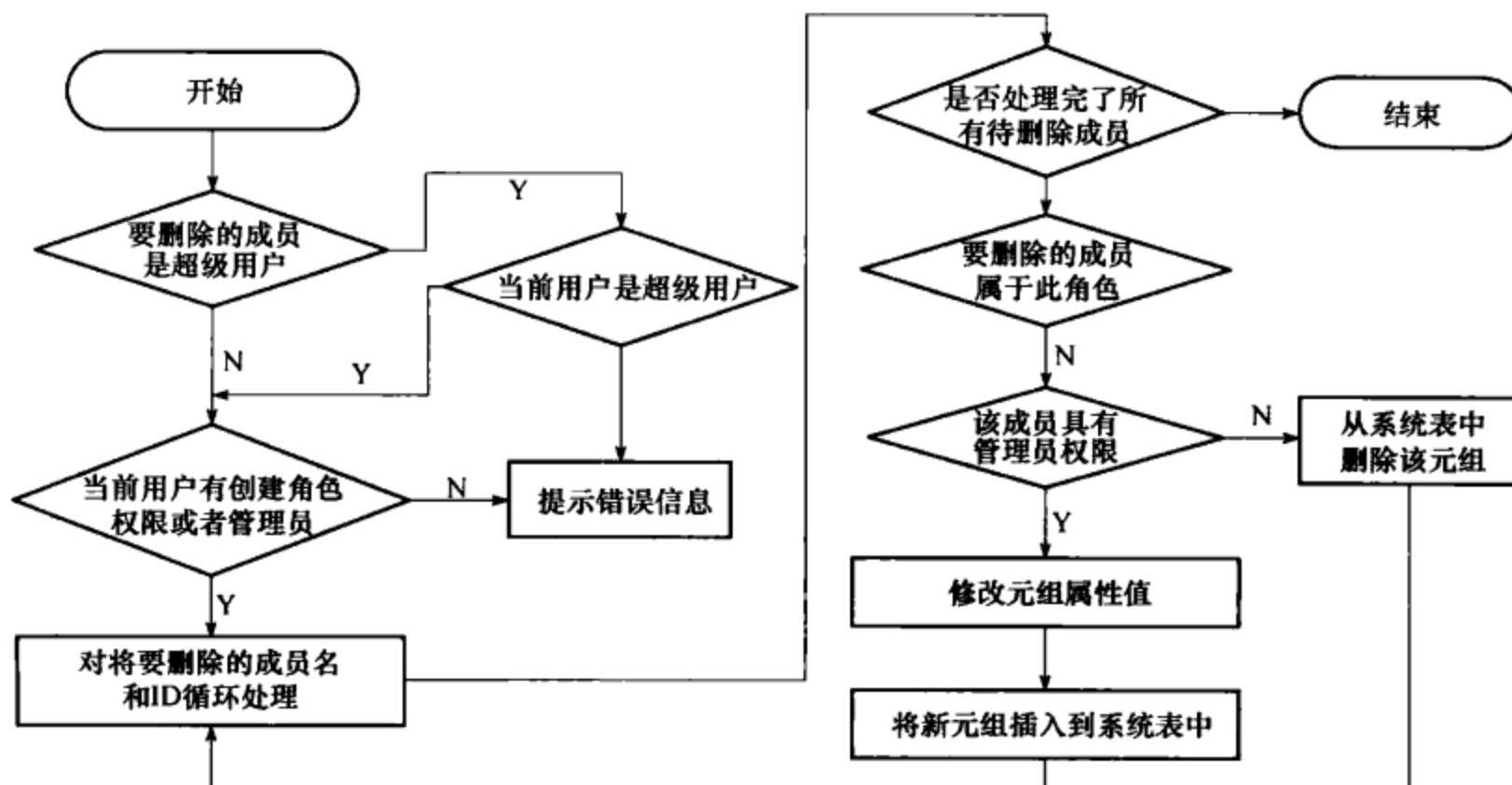


图 8-10 删除成员流程图

(2) 修改角色名称

通过修改角色的名称，拥有 CREATE ROLE 权限的角色可以给非超级用户进行重命名。当前会话的用户不能对自己重命名，若需要对当前会话的用户进行重命名，必须以另外一个用户的身份进行连接。这个过程通过调用函数 RenameRole 来用给定的名字重命名给定的角色。需要注意的是，MD5 加密口令使用角色名字作为加密的 salt。所以，如果口令用 MD5 加密，那么重命名某用户时会清空其口令。例如若需要修改角色名称 James（其密码经过 MD5 加密过）为 Cathy 时，可以使用 SQL 命令 ALTER ROLE James RENAME TO Cathy。此时该用户口令会被清空。其实现过程为：先打开存储角色信息的系统表 pg_authid 获取旧角色名，修改角色名属性，然后更新元组再写回系统表。重命名角色的流程如图 8-11 所示。

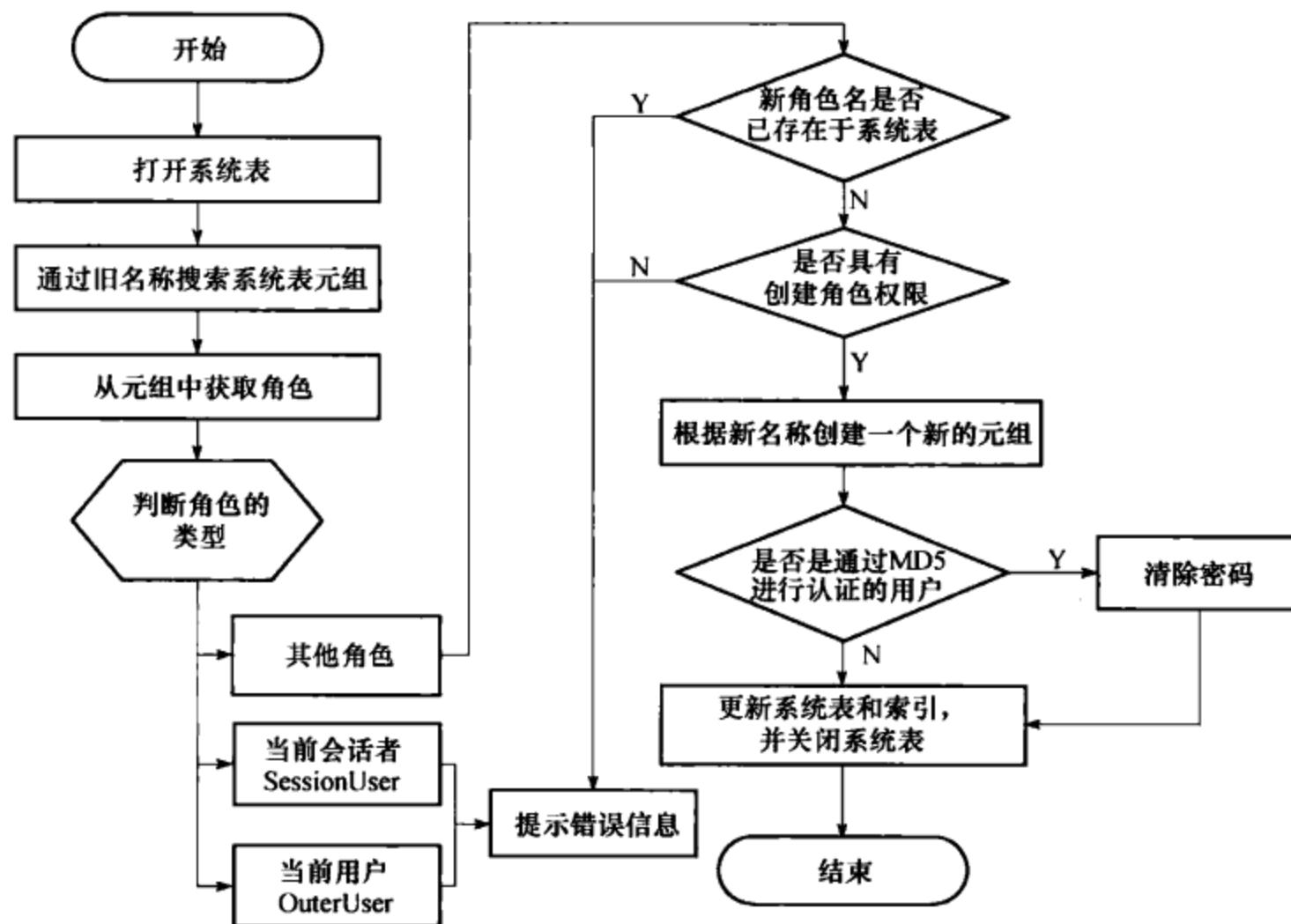


图 8-11 重命名角色流程图

(3) 修改一个角色的特定变量的会话缺省值

普通的角色可以改变自己的会话缺省值，超级用户可以修改任何人的会话缺省值，拥有 CREATE ROLE 权限的角色可以给非超级用户修改缺省值。在被修改角色随后开始一个新的会话之后，指定的数值会覆盖 postgresql.conf 文件里或者由命令行传入的参数值，而变成会话缺省值。例如，给角色 teacher 设置非缺省的 maintenance_work_mem 参数值 100000，可以使用 SQL 命令 ALTER ROLE worker_bee SET maintenance_work_mem = 100000。该过程通过函数 AlterRoleSet 来实现，该函数只有一个类型为 AlterRoleSetStmt 的参数（数据结构 8.6）。

数据结构 8.6 AlterRoleSetStmt

```

typedef struct AlterRoleSetStmt
{
    NodeTag          type;
    char             *role;      // 角色名
    VariableSetStmt *setstmt;   // set/reset 子命令
}AlterRoleSetStmt;

```

执行过程也很简单：先打开存储角色信息的系统表 pg_authid，按照参数中的角色名称找到对应的元组，然后修改其 rolconfig 属性，最后更新元组写回到系统表中。

3. 删除角色

如果要删除一个数据库角色，可以使用 SQL 命令 DROP ROLE。要注意的是，若要删除一个超级用户角色，则当前操作者也必须为超级用户；若要删除非超级用户角色，当前操作者必须拥有 CREATE ROLE 权限。该删除操作调用函数 DropRole 实现，该函数只有一个类型为 DropRoleStmt（数据结构 8.7）的参数。

数据结构 8.7 DropRoleStmt

```
typedef struct DropRoleStmt
{
    NodeTag      type;
    List          *roles;           // 要删除的角色列表
    bool          missing_ok;      // 判断角色是否存在
} DropRoleStmt;
```

首先要判断当前操作者是否有权限执行该操作。同时，对要删除的角色列表，还要检查待删除的角色是否存在：若 missing_ok 为 true 表示角色不存在。然后通过扫描系统表 pg_authid 和 pg_auth_members，获取所有涉及待删除角色的元组（角色）。注意，这里需要删除 pg_auth_members 中属于被删除角色的成员以及成员是被删除角色的元组。最后，删除该角色在 pg_shdescription 系统表中对应的注释信息。该函数流程如图 8-12 所示。

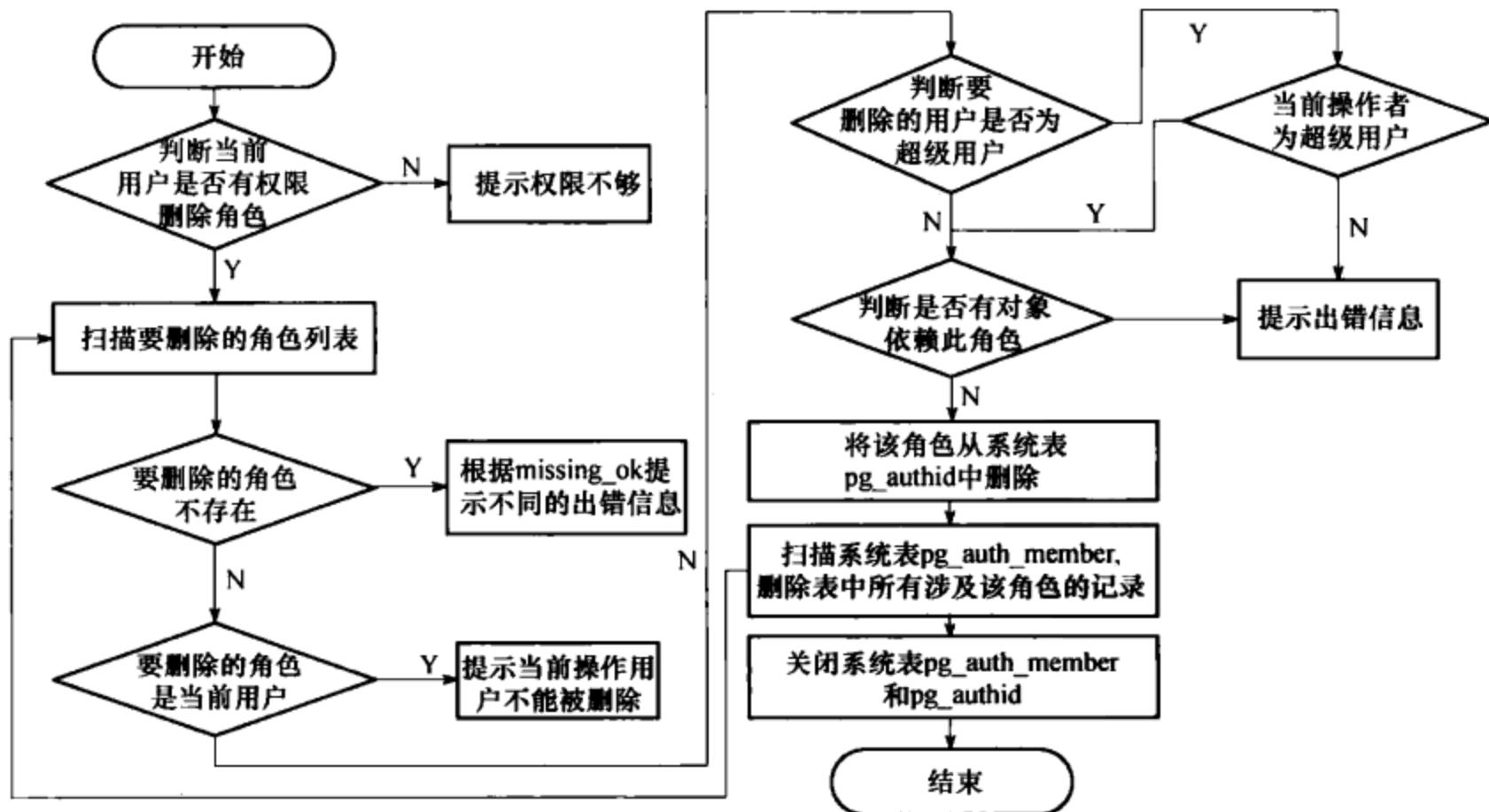


图 8-12 删除角色流程图

4. 授权和回收角色

如果要授权或回收角色的成员关系，可以使用 SQL 命令 GRANT/REVOKE。该命令有两个基本

变种：一种是赋予/撤销一个角色中的成员关系，而另一种是为角色授予在数据库对象（表、视图、序列、数据库、函数、过程语言、模式或者表空间）上的权限。尽管两个变种在很多方面都非常类似，但是由于功能不同，仍需分别描述。对于第一种，其命令的语法为：

```
GRANT role[,...] TO username[,...][WITH ADMIN OPTION]
```

此命令用于将一个或者多个角色作为成员加入到指定的角色中。如果声明了 WITH ADMIN OPTION，那么被加入的成员角色不仅可以将其他角色加入到父角色中，还可以从父角色中删除成员。如果没有该选项，普通用户不具有上述权限。当然，数据库超级用户可以给任何人赋予或者撤销任何角色的任何成员关系。拥有 CREATEROLE 权限的角色可以赋予或者撤销任何非超级用户角色的成员关系。此过程由函数 GrantRole 来实现，该函数只有一个类型为 GrantRoleStmt（数据结构 8.8）的参数。

数据结构 8.8 GrantRoleStmt

```
typedef struct GrantRoleStmt
{
    NodeTag          type;
    List             *granted_roles;      // 被授权或回收的角色集合
    List             *grantee_roles;     // 从 granted_roles 中增加或删除的角色集合
    bool             is_grant;           // true = 授权, false = 回收
    bool             admin_opt;          // 判断是否管理权限
    char             *grantor;            // 授权者, 默认为当前操作者
    DropBehavior     behavior;           // 回收行为(枚举类型)
}GrantRoleStmt;
```

字段 behavior 是一个枚举数据类型的 DropBehavior 结构，包括 DROP CASCADE 和 DROP RESTRICT。其中，DROP CASCADE 表示级联删除所有依赖于被删除角色的角色，DROP RESTRICT 表示拒绝删除那些有任何依赖角色存在的角色。具体描述可参见下一节。

授权角色时，grantee_roles 中的角色将被添加进 granted_roles 的角色中。此外，还要设置执行角色添加的授权者 grantor，该授权者应当是参数中给定的。若参数没有给出，则将授权者设置为执行当前授权操作的操作者。回收角色时，将 grantee_roles 中的角色从 granted_roles 中的角色里删除。此流程如图 8-13 所示。

5. 删除角色的数据库对象授权

如果要删除一个数据库角色所拥有的数据库对象授权，可以使用 SQL 命令

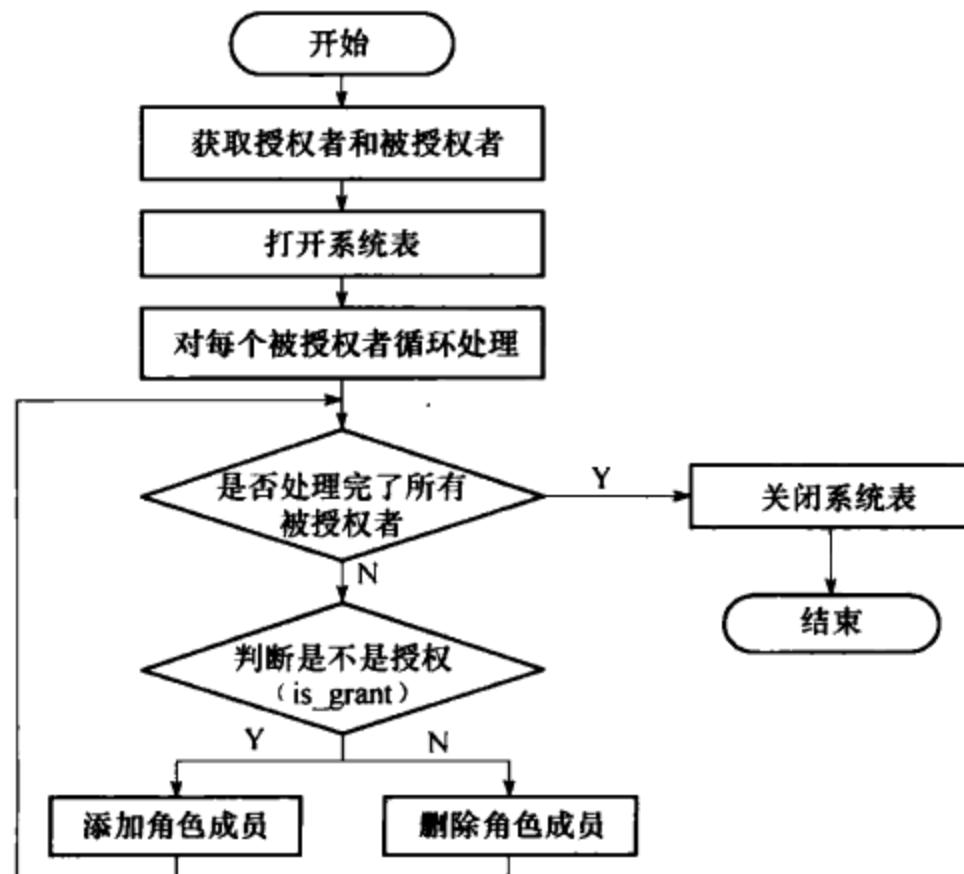


图 8-13 授权/回收角色流程图

`DROP OWNED`。执行该命令后，当前数据库中此角色在所拥有的对象上的权限都将被撤销。该过程由函数 `DropOwnedObjects` 来实现，该函数只有一个类型为 `DropOwnedStmt`（数据结构 8.9）的参数。

数据结构 8.9 DropOwnedStmt

```
typedef struct DropOwnedStmt
{
    NodeTag          type;
    List             *role;           // 角色列表
    DropBehavior     behavior;       // 授权行为
} DropOwnedStmt;
```

该函数在判断是否有权限对角色执行删除操作后，打开系统表 `pg_shdepend`。根据参数中的角色列表，检查每个角色的数据库对象权限并删除。该操作过程删除该角色上的所有数据库对象的权限，其流程如图 8-14 所示。

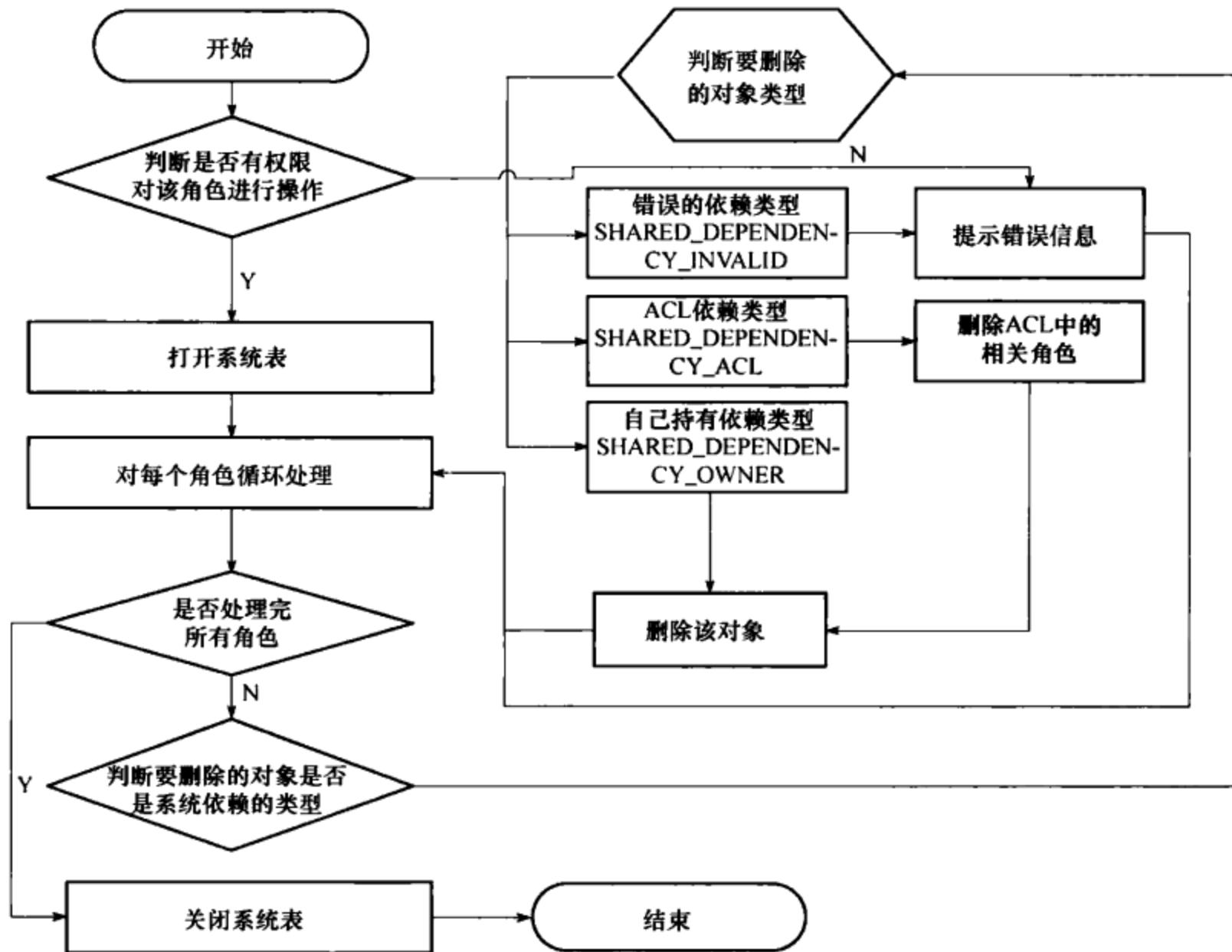


图 8-14 删除角色上的数据库对象

6. 修改数据库对象属主

如果要修改数据库对象的属主，则可以使用 SQL 命令 REASSIGN OWNED。该命令常用于在删除角色之前的准备工作。因为 REASSIGN OWNED 仅影响当前数据库中的对象，所以如果将被删除的角色在多个数据库中都拥有数据库对象，则需要在这些数据库中都执行一次该命令。其语法结构如下：

```
REASSIGN OWNED BY old_role[,...] TO new_role
```

old_role 是旧属主的角色名，new_role 表示将要成为这些对象属主的新角色的名字。该命令将所有旧属主拥有的数据库对象的属主更改为 new_role。该操作由函数 ReassignOwnedObjects 来实现，该函数仅有一个类型为 ReassignOwnedStmt（数据结构 8.10）的参数。

同样，该函数在执行操作之前需要检查是否有对该角色的操作权限。然后打开系统表 pg_shdepend。根据参数中的角色列表逐一将该角色拥有的所有对象的属主改为新角色。该操作过程流程如图 8-15 所示。

数据结构 8.10 ReassignOwnedStmt	
typedef struct ReassignOwnedStmt	
{	
NodeTag type;	
List *role; // 角色列表	
DropBehavior *newrole; // 授权行为	
}	ReassignOwnedStmt;

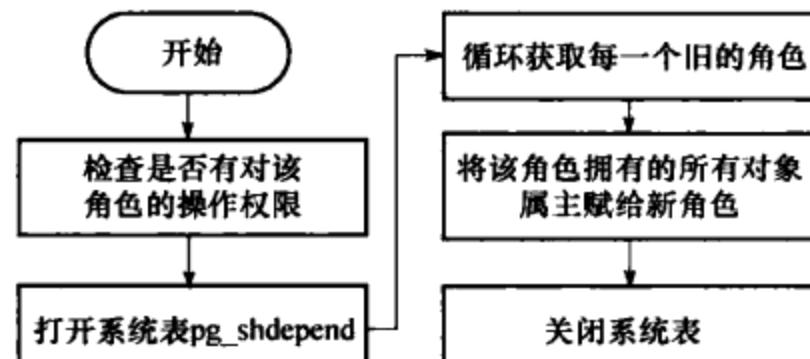


图 8-15 修改数据库对象属主

8.4 对象访问控制

在 PostgreSQL 中，建立安全的数据库连接是由用户标识和认证技术来共同实现，而建立连接后的安全访问保护则是基于角色的对象访问控制。前面已经阐述了角色的概念，本节则在角色概念的基础上进一步阐述对象访问控制机制。

数据库里的每个对象所拥有的权限信息经常会发生变化，比如授予对象的部分操作权限给其他用户，或者删除用户在对象上的操作权限，亦或是对用户在对象上的操作权限进行更新等。以上操作在 SQL 中体现为 GRANT 和 REVOKE 语句，且都涉及对象权限信息的动态管理，即权限控制中的对象权限管理。

为了保护数据安全，当用户要对某个数据库对象进行操作之前，必须检查用户在对象上的操作权限，仅当用户对此对象拥有进行合法操作的权限时，才允许用户对此对象执行相应操作。上述操作检查的过程被称为对象权限检查。

8.4.1 访问控制列表

访问控制列表（Access Control List，ACL）是对象权限管理和权限检查的基础，PostgreSQL 通过操作 ACL 实现对象的访问控制管理。在 PostgreSQL 中，每个数据库对象都具有 ACL，每个对象的 ACL 存储了此对象的所有授权信息。当用户访问对象时，只有它在对象的 ACL 中并且具有所需的权限时才能访问该对象。当用户要更新对象的权限时，只需更新 ACL 上的权限信息即可。

ACL 是存储控制项（Access Control Entry，ACE）的集合，其组织结构如图 8-16 所示。

每个 ACL 实际是一个由多个 AclItem（数据结构 8.11）构成的链表。每个 AclItem 对应一个 ACE。ACE 由数据库对象和授权列表构成，记录着可访问对象的用户或者执行单元（进程、存储过程等）。此外，ACE 中还记录了可在对象上进行权限（如读、写和执行等）操作的用户或者执行单元。PostgreSQL 中，ACE 由受权者、授权者以及权限位三部分组成。

其中，字段 `ai_privs` 是 `AclMode` 类型。`AclMode` 是一个 32 位的比特位，其高 16 位为权限选项位，低 16 位为该 ACE 中的操作权限位。每个操作权限占 1 个比特位，当该比特位的取值为 1 时，表示 ACE 中的 `ai_grantee` 对应的用户（受权者）具有此对象的相应操作权限。否则，表示用户没有相应权

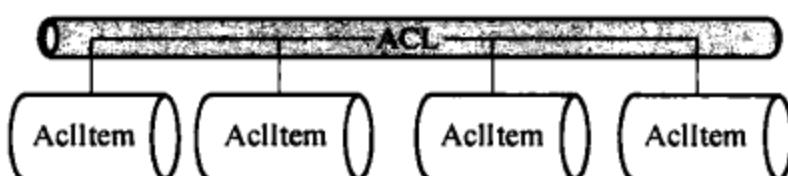


图 8-16 ACL 结构图

数据结构 8.11 AclItem

```
typedef struct AclItem
{
    Oid      ai_grantee;      //受权者的OID
    Oid      ai_grantor;     //授权者的OID
    AclMode ai_privs;        //权限位
}AclItem;
```

图 8-17 AclMode 的结构

在 AclMode 结构位图中，“ace type” 记录了 ACE 的被授权者（受权者）类型，它可以是用户、组或者 public。当它为 public 时，表示将该权限授权给数据库集群中的所有用户。“Grant options” 记录了各权限位对应的授出或者被转授选项。低 16 位分别记录了各个权限位的授予情况。从低到高的各个比特位代表的权限依次为 INSERT、SELECT、UPDATE、DELETE、TRUNCATE、REFERENCES、TRIGGER、EXECUTE、USAGE、CREATE、CREATE_TEMP、CONNECT 等，各权限在系统表中显示分别为 a、r、w、d、D、x、t、X、U、C、T、c 等。若当授权语句使用 ALL 时，则表示包含所有权限。

1. ACL 检查

在 PostgreSQL 中，主要通过访问控制列表来存储对象权限的相关信息。每个对象中都会有一个 ACL。对于指定的对象，用户可以查询该对象上是否存在某权限信息。对于不同的数据库对象（如数据库、表、语言、模式、命名空间、表空间等），都根据其不同的权限属性采用了不同的函数来实现。对于表来说，该检查过程可以用函数表示为：

has table privilege * * (PG FUNCTION ARGS)

注意，函数名中的最后两个星号分别代表参数中的用户信息和数据库对象信息，例如函数名为 `has_table_privilege_name_id` 表示参数中给出了用户名和表的 OID，该函数将检查给定的用户是否拥有在给定表上的指定权限。该表的操作权限检查函数基本流程如下：

1) 从 PG_FUNCTION_ARGS 中获得用户 OID、表 OID 和权限信息。

2) 执行 pg_class_aclcheck 获得用户在该表上的权限集，比较该权限集与操作所需权限集，如果返回的权限集合大于或等于操作所需权限集，则检查成功通过。否则失败。

其中，函数 pg_class_aclcheck 的执行过程为：如果用户要操作的是系统表，则禁止执行对此表的操作，即使用户是超级用户也是如此。取出被操作表的 ACL，将其一个副本，然后调用 aclmask 函数从这个副本中计算用户在该表上的权限集并返回。

2. ACL 更新

ACL 中存储着数据库对象的权限信息，对某对象进行权限授权/回收时，实际上就是在 ACL 上进行添加或删除指定的权限，同时更新 ACL 的过程。ACL 的更新很简单，将权限对应的标志位置为 0 或者 1 即可，这里主要对 ACL 更新中可能出现的循环授权情况进行介绍。

在执行授权操作时，可能会出现循环授权的情况。如图 8-18 所示，Johnny 将表 CheckTable 的 SELECT 权限赋予给 Tom；Tom 又将之转授给 David，最后 David 企图将同样的权限赋予给 Johnny。这一系列的授权操作构成了一个环，这个环是不允许出现的。

为了检测出循环授权的情况，在进行授权时会先执行函数 check_circularity，该函数将通过递归删除权限的方式来检查是否有环构成。其执行过程如下：

- 1) 获取旧的 ACL、要增加的授权信息以及对象的拥有者。
- 2) 创建一个 ACL 副本进行操作。
- 3) 递归删除所有被授予者所拥有的可再授予权限。
- 4) 检查授予者是否还有可再授予权限。如果没有，则授予者不能进行授权。

在上面的例子中，David 将权限赋予给 Johnny 之前将进行循环授权的检查：首先会收回 Johnny 所拥有的所有可再授予权限，此时 David 的可再授予权限将被回收（因为 David 的可再授予权限源自 Johnny），则 David 此时不能再给 Johnny 授权，这样就能检测出发生循环授权的问题。这种方法的前提是：当撤销权限时，将会调用函数 recursive_revoke 来递归地删除从这个权限授予出去的所有权限。在上面的例子中，当收回 Johnny 的权限时，由于 Tom 的权限是由 Johnny 授予，因此 Tom 的权限也被回收。同样地，David 的权限也被回收。这是对象上的权限级联删除过程。该函数的流程与 check_circularity 类似，也是递归删除所有被授予者所拥有的可再授予权限。不过在收回权限时不存在环状问题，因为在授权时就已经检测过。

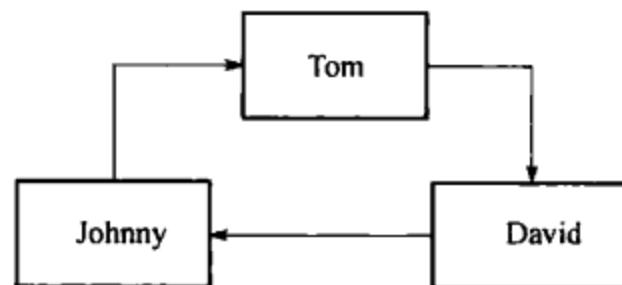


图 8-18 循环授权的例子

8.4.2 对象权限管理

对象权限管理主要是通过使用 SQL 命令 GRANT/REVOKE 授予或回收一个或多个角色在对象上的权限。

(1) 对象授权

对象授权语句 GRANT 的语法结构如下：

```

GRANT {{SELECT | INSERT | UPDATE | DELETE | REFERENCES | TRIGGER}
      [, ...] | ALL[PRIVILEGES]}ON[ TABLE]tablename[, ...]
      TO{username | GROUP groupname | PUBLIC}[ , ...]
  
```

[WITH GRANT OPTION]

此命令表示在数据库对象上给一个或者多个角色授予特定权限。关键字 PUBLIC 表示权限要赋予所有角色，包括那些以后可能创建的用户。PUBLIC 可以看做是一个预定义好的组，它总是包括所有角色。任何特定的角色的权限由三部分组成：直接赋予的权限、从所属的角色继承来的权限以及被赋予的 PUBLIC 权限。

如果声明 WITH GRANT OPTION，那么权限的被授予者可以将此权限再赋予他人，否则只能自身拥有被赋予的权限，这个选项称为“可再授予权限”。此选项不能赋予 PUBLIC。

就对象的所有者（通常就是创建者）而言，没有权限需要被赋予，因为所有者缺省就具有该对象上所有的权限。所有者出于安全考虑可以选择舍弃部分权限。不过，删除或修改对象的权限则属于创建者所固有，不能赋予或撤销。所有者同时隐式地拥有对象的可再授予权限。

(2) 对象权限回收

对象权限回收命令 REVOKE 的语法结构如下：

```
REVOKE[ GRANT OPTION FOR ]
  {{SELECT | INSERT | UPDATE | DELETE | REFERENCES | TRIGGER}
  [,...] | ALL[ PRIVILEGES ]}ON[ TABLE ]tablename[ ,... ]FROM username |
  GROUP groupname | PUBLIC }[ ,... ][ CASCADE | RESTRICT ]
```

REVOKE 命令用于撤销之前赋予角色的权限。关键字 PUBLIC 的意义和 GRANT 语句中相同。若指定了 GRANT OPTION FOR，则只是撤销角色对该权限的再授予能力，而不撤销权限本身。否则，权限（包括再授予权限）将被撤销。

如果用户持有某个权限，同时拥有再授予权限，并把权限授予了其他用户，那么其他用户所持有的此权限都被称为依赖性权限。如果被依赖用户持有的权限或者再授予权限被撤销，在声明了 CASCADE 关键字情况下依赖性权限也会被撤销，否则撤销动作会提示失败。

GRANT/REVOKE 命令都由函数 ExecuteGrantStmt 实现，该函数只有一个类型为 GrantStmt（数据结构 8.12）的参数。

数据结构 8.12 GrantStmt

```
typedef struct GrantStmt
{
    NodeTag      type;
    Bool         is_grant;           //true = 授权, false = 回收
    GrantObjectType objtype;        //被操作的对象类型
    List         *objects;          //被操作的对象集合
    List         *privileges;       //要操作的权限集合
    List         *grantees;         //受权者集合
    bool         grant_option;     //授出或者回收权限选项
    DropBehavior behavior;        //回收权限行为:restrict/cascade
    }GrantStmt;

}AclItem;
```

其中，字段 objtype 是一个 GrantObjectType 结构，记录被操作的对象类型。该结构为枚举类型，

包括列、表、视图、序列、数据库、模式等。字段 behavior 是一个 DropBehavior 结构，记录回收权限行为。此结构也为枚举类型，包括 CASCADE 和 RESTRICT。需要注意的是，在 PostgreSQL 8.4.1 中增加了对表或视图上的列权限设置的支持。

字段 privileges 是一个 AccessPriv 结构（数据结构 8.13），用于保存要操作的权限列表。

数据结构 8.13 AccessPriv

```
typedef struct AccessPriv
{
    NodeTag      type;
    char         *priv_name;      //权限名
    List         *cols;          //要设置权限的列的名字列表
}AccessPriv;
```

函数 ExecuteGrantStmt 先将 GrantStmt 结构转化成 InternalGrant 结构（数据结构 8.14），将命令的权限列表转化成内部的 AclMode 表示。当 privileges 取值为 NIL 时，表示授予或回收所有权限，此时 InternalGrant 的 all_privs 字段为 true，且 InternalGrant 中的 privileges 字段被设置为 ACL_NO_RIGHTS，也表示 ACL_ALL_RIGHTS。

数据结构 8.14 InternalGrant

```
typedef struct InternalGrant
{
    bool           is_grant;      //true = 授权, false = 回收
    GrantObjectType objtype;     //被操作的对象类型
    List           *objects;       //被操作的对象集合
    bool           all_privs;     //是否授予所有权限
    AclMode        privileges;    //以内部比特位形式表示的操作权限
    List           *col_privs;     //列权限
    List           *grantees;      //受权者集合
    bool           grant_option;  //授出或者回收权限选项
    DropBehavior   behavior;      //回收权限约束:restrict/cascade
}InternalGrant;
```

InternalGrant 中的 all_privs 和 privileges 只表示对象级的权限集合。列级的权限集合是用 col_privs 来表示的，它的值由 GrantStmt 中的链表 privileges 给出。GrantStmt 中的 grantees 链表和 InternalGrant 的 grantees 链表的节点类型是不一致的，前者是 PrivGrantee 结构，后者是 OID 结构，即 ExecuteGrantStmt 函数会把 PrivGrantee 链表转化为 OID 链表。若 PrivGrantee 结构中字段 rolename 值为空，会把 ACL_ID_PUBLIC 加入到 OID 链表中。

完成从 GrantStmt 到 InternalGrant 的转换之后，ExecuteGrantStmt 会调用函数 ExecGrantStmt_oids，该函数将根据对象类型调用相应用对象上的权限管理函数。ExecuteGrantStmt 的流程图如图 8-19 所示。

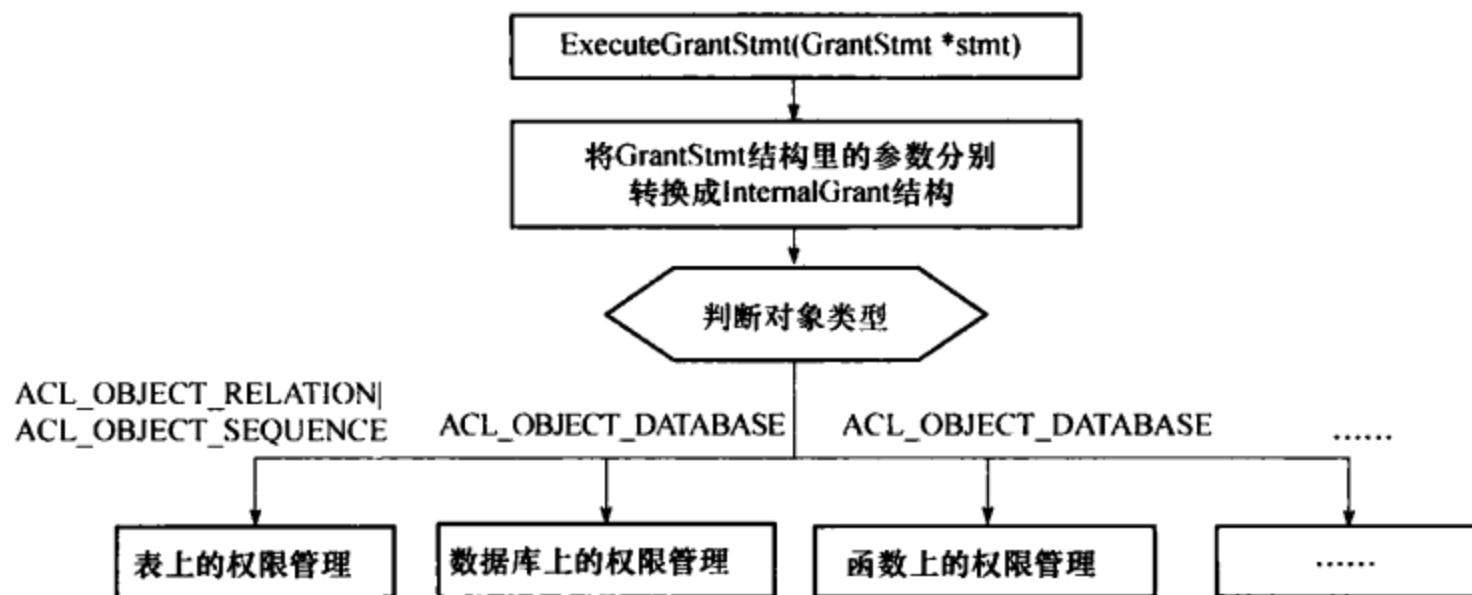


图 8-19 授权与回收权限功能入口函数流程

对于表、模式、数据库、函数等不同的数据库对象，它们的权限管理分别调用相应的函数来完成。此类算法的构成大同小异，核心过程是通过解析 GRANT/REVOKE 命令获取权限信息，然后和保存在 ACL 中的原授权信息共同计算出实际要授予/回收的权限，最后更新对象上的 ACL。本节以表为例描述对象上的权限管理过程。表上的权限管理函数是 ExecGrant_Relation，它的参数就是在 ExecuteGrantStmt 中转换得到的 InternalGrant 结构，其流程如图 8-20 所示。

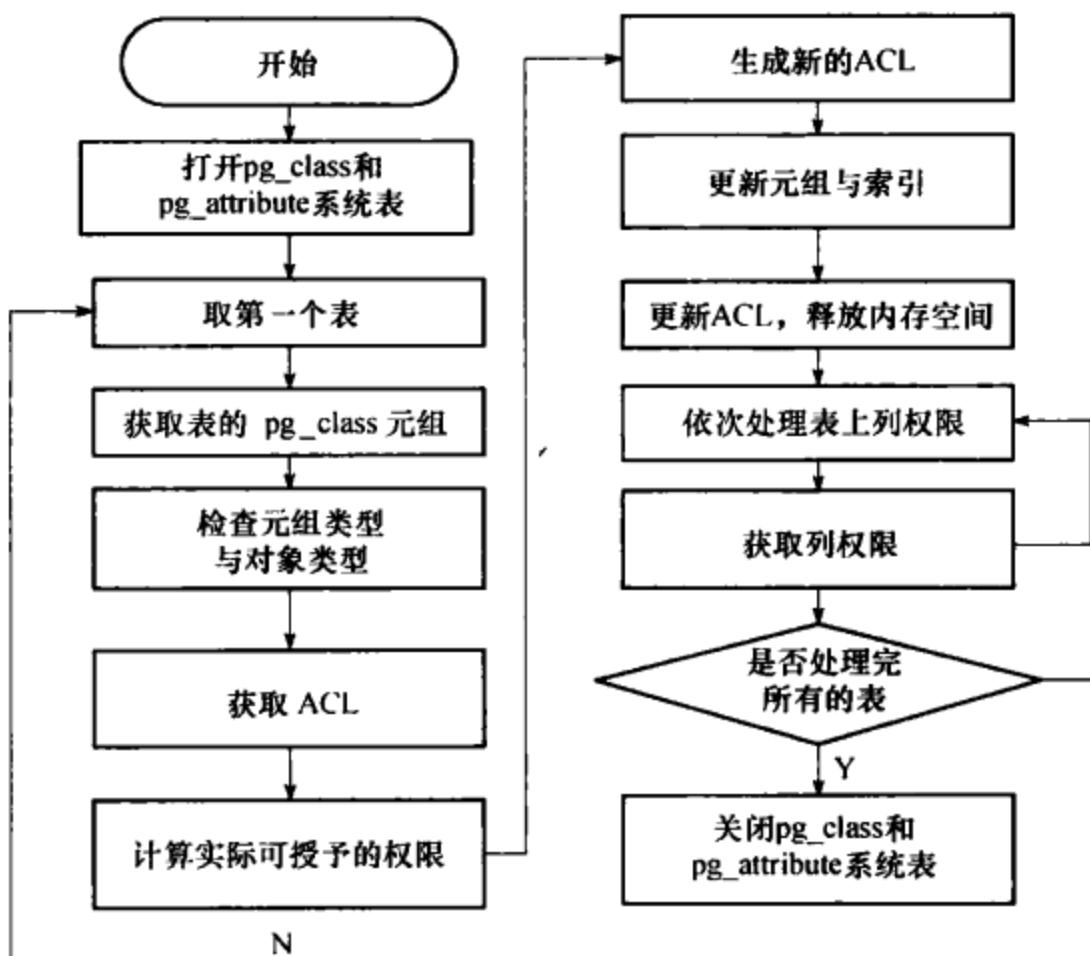


图 8-20 表的权限管理函数流程

获取旧 ACL 时，如果不存在旧的 ACL，则新建一个 ACL，并将默认的权限信息赋给该 ACL。根据对象的不同，初始的缺省权限含有部分可赋予 PUBLIC 的权限：对于表、模式、表空间没有公开访问权限；

对于数据库有 CONNECT 权限和创建 TEMP 表的权限；对于函数有 EXECUTE 权限；对于语言有 USAGE 权限。对象所有者可以随时回收以上权限。如果存在旧的 ACL，则将旧的 ACL 存储为一个副本。

生成新 ACL 时，如果是授予权限，则将命令中给出的要授予的权限添加到旧的 ACL 中；如果是回收权限，则将命令中给出的要被回收的权限从旧的 ACL 中删除。

列权限的处理通过调用函数 `expand_col_privileges` 实现，该函数将列级权限列表转化成数组形式表示。依次判断每一列的权限，调用函数 `ExecGrant_Attribute` 计算每一列实际可以授予/回收的权限，然后更新 ACL。每一列的 ACL 存放在该列对应的 `pg_attribute` 元组中。`ExecGrant_Attribute` 的处理过程与 `ExecGrant_Relation` 类似。需要注意的是，目前只有在表中才会有列级权限，所以 `ExecGrant_Attribute` 只能从 `ExecGrant_Relations` 处调用，而不能直接被 `ExecGrantStmt` 调用。

8.4.3 对象权限检查

在对数据库对象进行操作时，必须要对该对象上的权限进行检查。只有拥有该操作权限时，才可以执行该操作。例如，用户要查询表 `checktable`，可以执行命令：“`SELECT * FROM checktable`”。如果该用户在 `checktable` 上不具有查询权限，则不可以查询。通常数据库对象的拥有者具有对该对象进行一切操作的权限，超级用户也拥有对所有对象的全部操作权限，而非属主用户在对数据库对象操作之前需要进行权限检查。

下面也以表上的权限检查为例来描述权限检查函数的实现过程。表上的权限检查由函数 `ExecCheckRTEPerms` 实现，该函数只有一个类型为 `RangeTblEntry`（数据结构 8.15）的参数。

数据结构 8.15 RangeTblEntry

```
typedef struct RangeTblEntry
{
    NodeTag      type;
    RTEKind      rtekind;           // 对象类型
    .....
    AclMode      requiredPerms;    // 需要的访问权限(AclMode 类型)
    Oid          checkAsUser;      // 角色 ID
    Bitmapset   *selectedCols;    // 需要有 SELECT 权限的列
    Bitmapset   *modifiedCols;    // 需要有 INSERT/UPDATE 权限的列
} RangeTblEntry;
```

其中，字段 `requiredPerms` 表示需要的访问权限信息。`selectedCols` 表示在表的操作中需要有 SELECT 权限的列的集合，`modifiedCols` 则表示操作中需要 INSERT/UPDATE 权限的列的集合。该函数的基本流程如图 8-21 所示。

如果表级权限检查没有通过，我们可以进一步判断该用户是否在该表上具有列级权限。对列的操作检查有 `ACL_SELECT`、`ACL_INSERT` 和 `ACL_UPDATE` 三种。如果用户拥有列级权限，则检查通过，否则检查失败。例如，将表 `Student` 的第二列 `SELECT` 的权限赋予给用户 `Johnny`。当 `Johnny` 登录后，可以对该表第二列进行 `SELECT` 操作，但因为没有其他列的 `SELECT` 权限，所以不能对其他列进行 `SELECT` 操作。

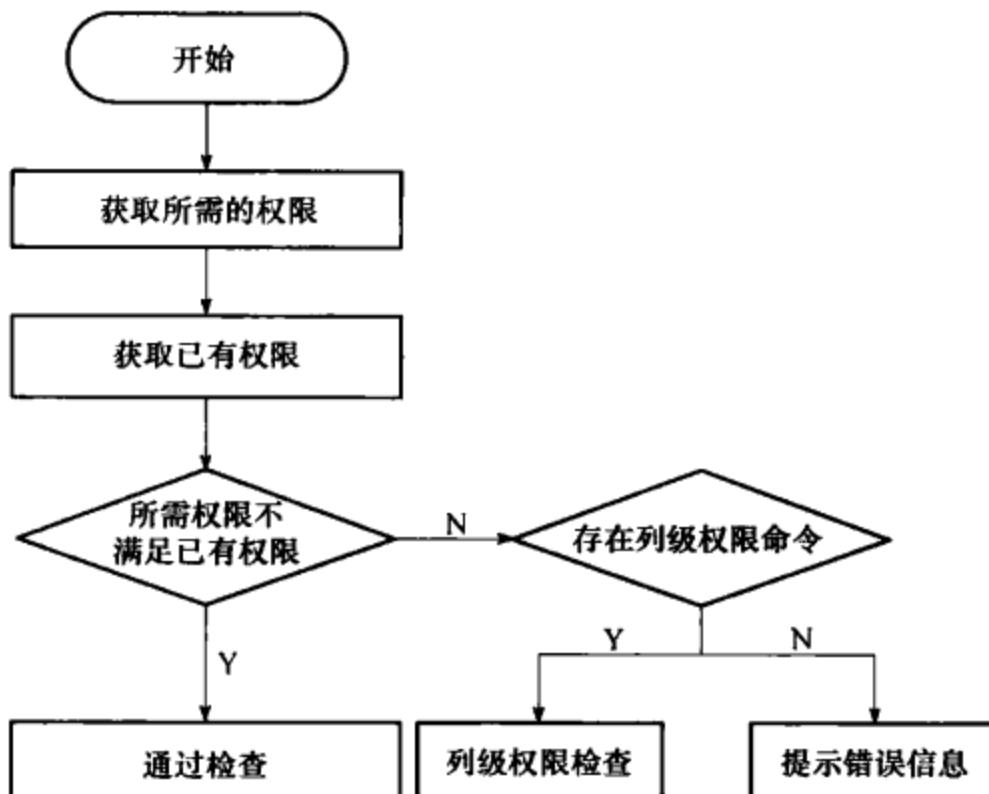


图 8-21 表操作权限检查过程

8.5 小结

PostgreSQL 通过用户标识和认证来防止非法用户访问数据库，并提供了多种不同的客户端认证方式。用户权限控制通过角色的概念来管理数据库权限。对象访问控制采用了比较完备的访问控制列表（ACL）技术进行权限的授予、回收和检查。

目前，国内外一些大型数据库系统都满足 C2 级以上要求，即满足 C1 级别要求并具有审计机制。而 PostgreSQL 数据库在 8.4.1 版本中依然缺乏明确意义上的审计机制，因此其安全级别还是为 C1 级。

习题

- 习题 8.1** PostgreSQL 通过哪些机制保证数据库的安全？
- 习题 8.2** 请描述客户端认证的过程，并采用不同的认证方式连接数据库。
- 习题 8.3** PostgreSQL 如何用角色的概念管理数据库访问权限？
- 习题 8.4** 假设有用户 A1 和 A2，A1 执行 SQL 语句：GRANT INSERT, DELETE ON TABLE TO A2，请描述在表 TABLE 上的权限管理过程。

附录 A

用Eclipse开发和调试PostgreSQL

对于那些习惯使用集成开发环境的读者来说，在 Linux 下开发和调试 PostgreSQL 可能会很不适应，这里我们将介绍一种用集成开发环境 Eclipse 开发和调试 PostgreSQL 的方法。

A. 1 安装 Eclipse

Eclipse 是著名的跨平台的开源集成开发环境（IDE），最初由 IBM 开发用来代替 Visual Age for Java。2001 年 11 月，IBM 将 Eclipse 贡献给了开源社区，现在由非盈利软件供应商 Eclipse 基金会管理。

Eclipse 最初是用来开发 Java 应用的，但通过 Eclipse 提供的插件机制已经将 Eclipse 支持的程序语言扩展到了 C/C++、PHP、Ruby、Python 等。Eclipse 中专门设立了一个分支用于提供对 C/C++ 语言的支持，称为 C/C++ 开发工具计划（CDT），它使用 GCC 作为编译器。包含 CDT 的 Eclipse 可以直接从 Eclipse 的官方网站 (www.eclipse.org) 下载。

笔者使用的操作系统是 Ubuntu 8.4，Eclipse 版本是 Helios（赫利俄斯，2010 年 6 月 23 日发布）。下载得到的是一个压缩包，如图 A-1 所示。



图 A-1 下载后的 Eclipse 压缩包

将压缩包解压后得到 eclipse 目录，其目录结构如图 A-2 所示。



图 A-2 Eclipse 目录结构

运行目录中的 `eclipse` 命令启动 Eclipse 开发环境，启动时首先要选择工作目录（`workspace`），如图 A-3 所示。



图 A-3 选择 workspace 界面

选定工作目录后即出现 Eclipse 的主界面，如图 A-4 所示。

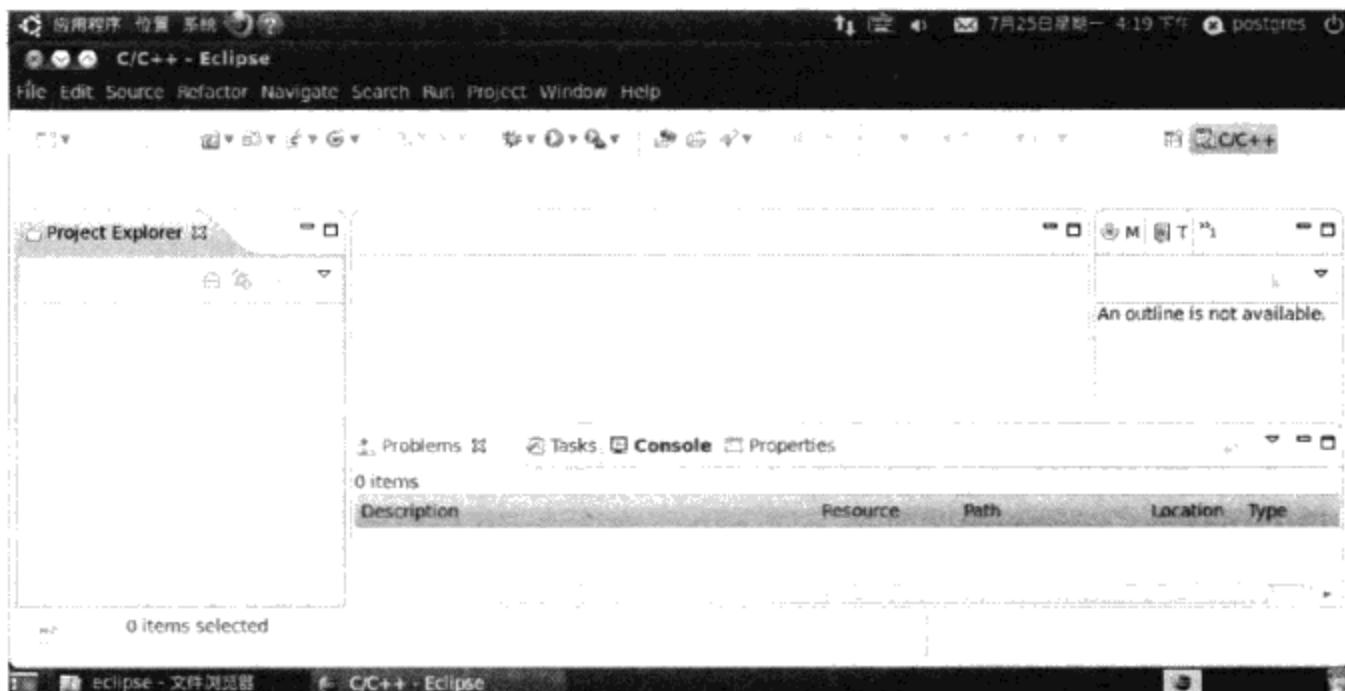


图 A-4 Eclipse 主界面

A.2 配置代码

配置 PostgreSQL 代码之前，先要安装下面几个组件：`libreadline5-dev`、`zlib1g-dev`、`Bison`、`Flex`，这些组件都可以通过 Ubuntu 的软件管理来安装，这里不再赘述。

在使用 Eclipse 导入 PostgreSQL 代码之前可以根据需要对 PostgreSQL 代码进行配置，例如笔者采用的是以下命令（在源代码目录下执行）：

```
. /configure --prefix = $HOME/project --enable-depend --enable-cassert --enable-debug
```

其中，“`--enable-debug`”是必需的，只有采用了这个参数才可能利用 Debug 工具对编译生成的 PostgreSQL 程序进行跟踪调试。

A.3 导入代码

在 Eclipse 主界面中，首先依次点击菜单 `File` 中的 `Import` 菜单项，在弹出的 `Import` 对话框（见

图 A-5) 中选择“Existing Code as Makefile Project”选项，点击“Next”按钮后进入“Import Existing Code”对话框（见图 A-6）。

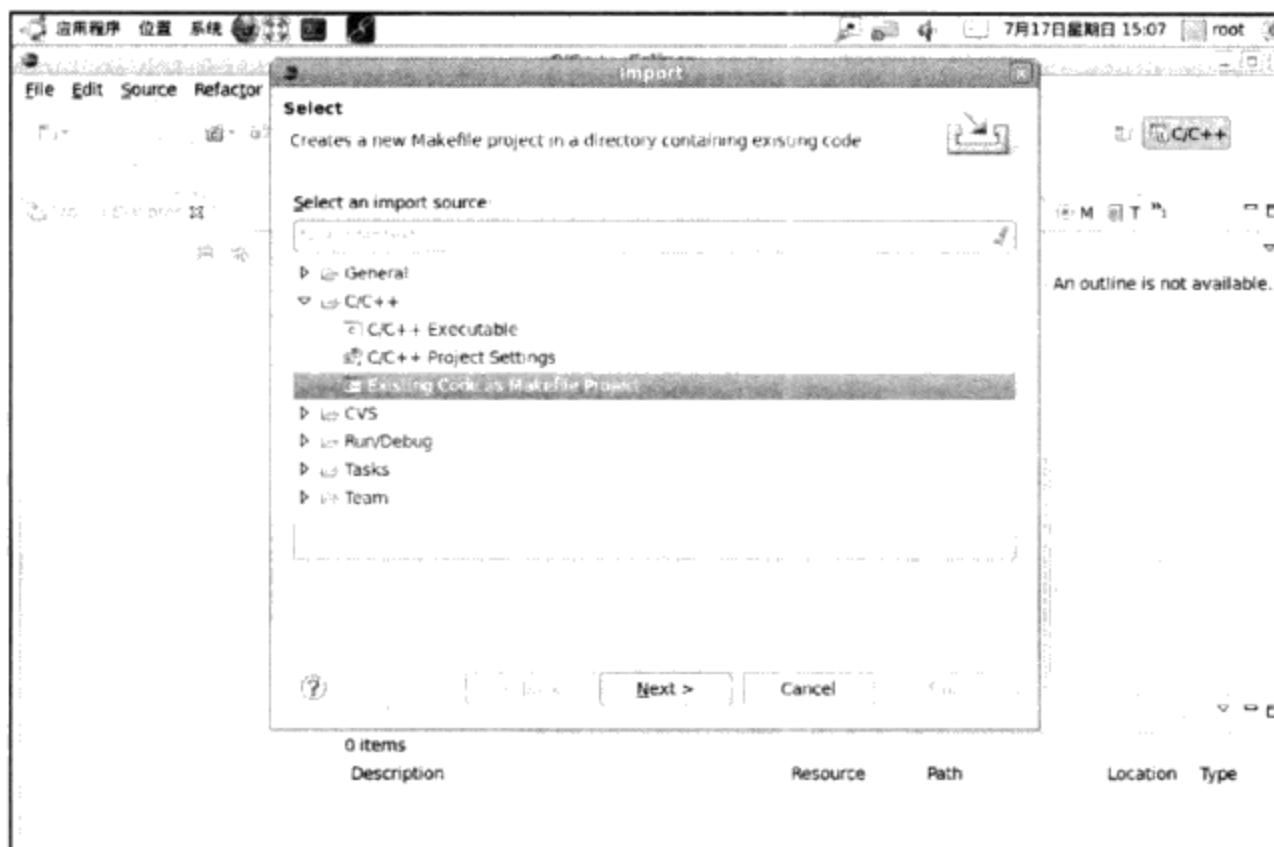


图 A-5 Import 对话框

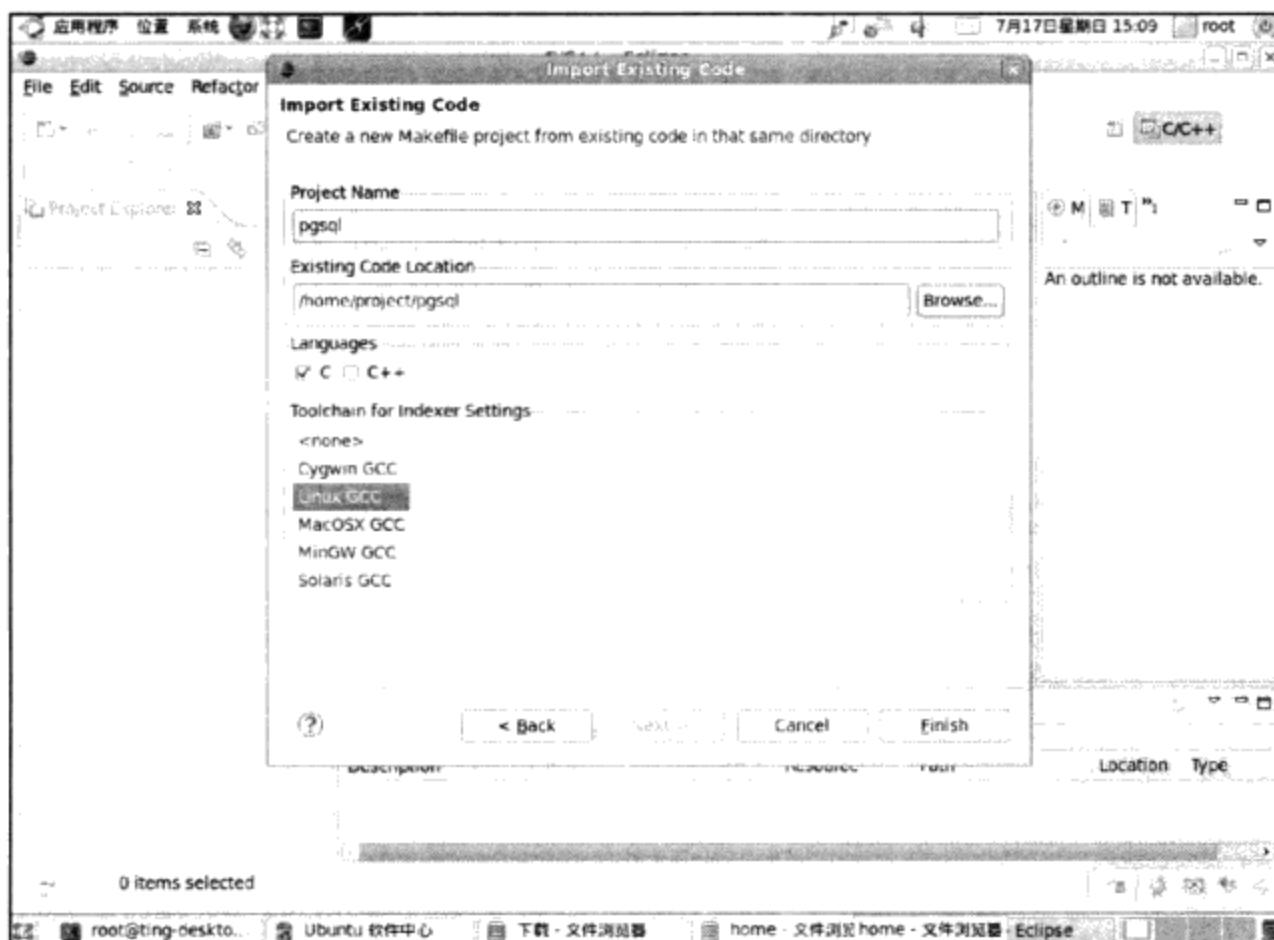


图 A-6 “Import Existing Code” 对话框

在“Import Existing Code”对话框中设置项目名称，将“Existing Code Location”指向解压后的 PostgreSQL 代码目录，注意一定要去掉“C++”复选框（PostgreSQL 是用 C 语言编写的）。

点击“Finish”按钮后，Eclipse 会花一定时间对代码进行编译，注意观察右下角的 Console 栏信息。直到提示如图 A-7 中的信息后才表示代码导入成功。

```
C-Build [pgsql]
labels -fno-strict-aliasing -fwrapv -g -fpic -L../../src/port -Wl,--as-needed -Wl,-rpath,/root/project/lib' --enable-new-dtags -L../../src/port -lgpgport -shared -o
autoinc.so autoinc.o
make[3]: 正在离开目录 '/home/project/postgresql/contrib/spi'
cp ../../contrib/spi/autoinc.so autoinc.so
make[2]: 正在离开目录 '/home/project/postgresql/src/test/regress'
make[1]: 正在离开目录 '/home/project/postgresql/src'
make -C config all
make[1]: 正在进入目录 '/home/project/postgresql/config'
make[1]: 没有什么可以做的为 'all'.
make[1]: 正在离开目录 '/home/project/postgresql/config'
All of PostgreSQL successfully made. Ready to install.
```

图 A-7 导入代码过程

导入完成后 Eclipse 界面如图 A-8 所示。

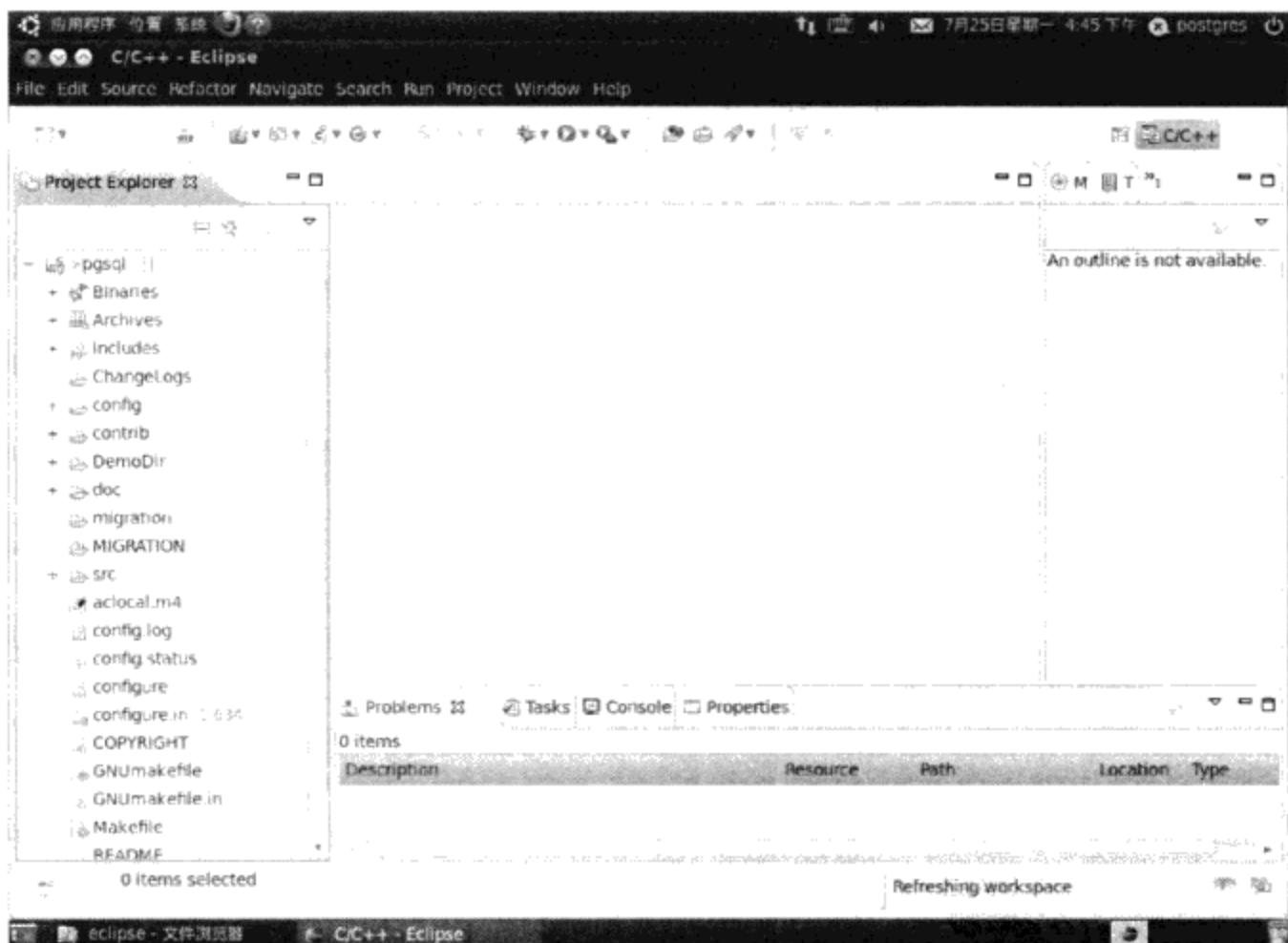


图 A-8 代码导入完成界面

A.4 开发与调试

A.4.1 新建“Make Target”

在调试 PostgreSQL 之前，首先需要将其安装到某个目录中，在 Eclipse 中是通过建立一个“Make Target”实现的。

在工程“pgsql”上单击右键，选择“New”菜单项将会出现如图 A-9 所示的菜单。



图 A-9 “Make Target”菜单项

选择“Make Target”菜单项后出现“Create Make Target”对话框（见图 A-10）。

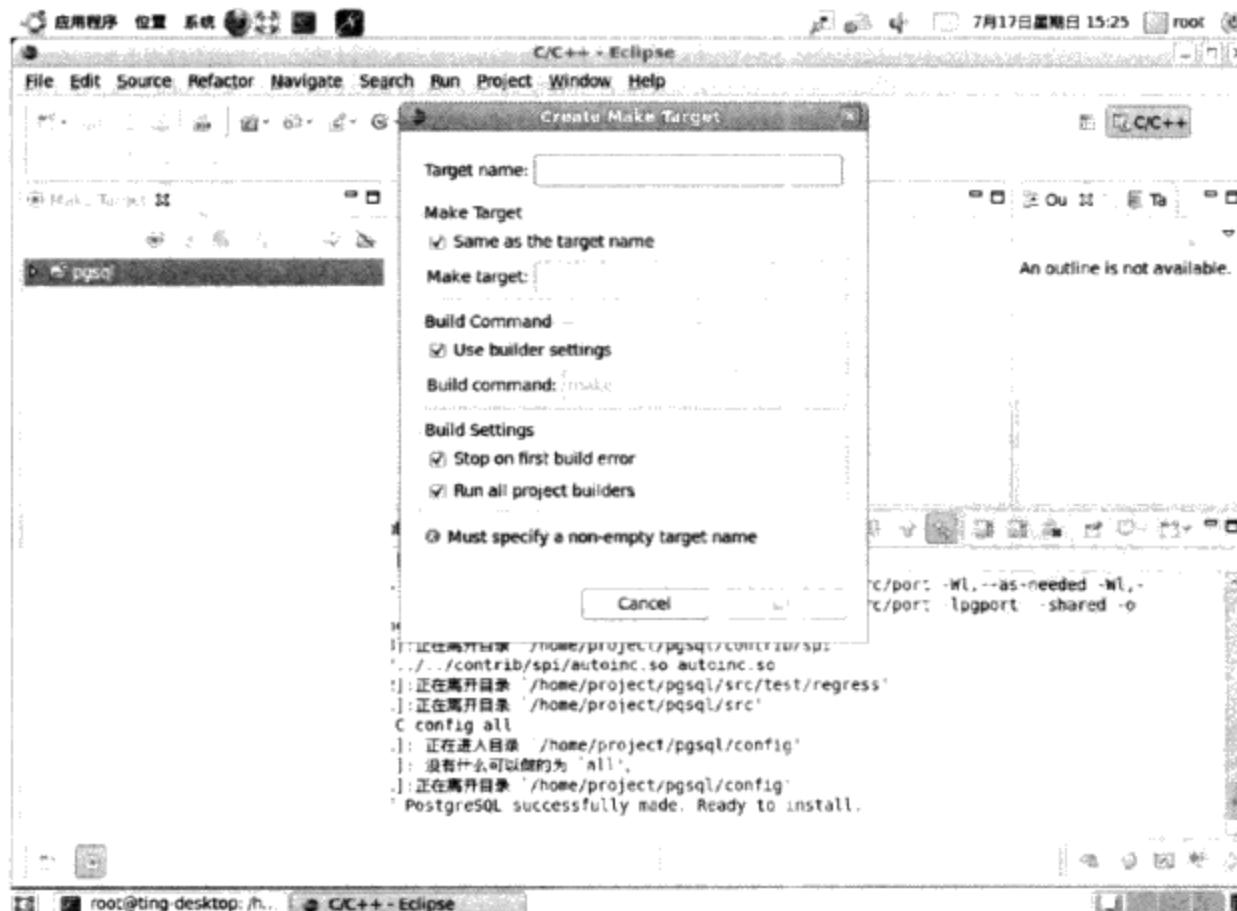


图 A-10 “Create Make Target”对话框

在“Target Name”文本框中输入 install，单击 OK 按钮创建“Make Target”，如果创建成功可以看到一个绿色的 install 选项（见图 A-11）。

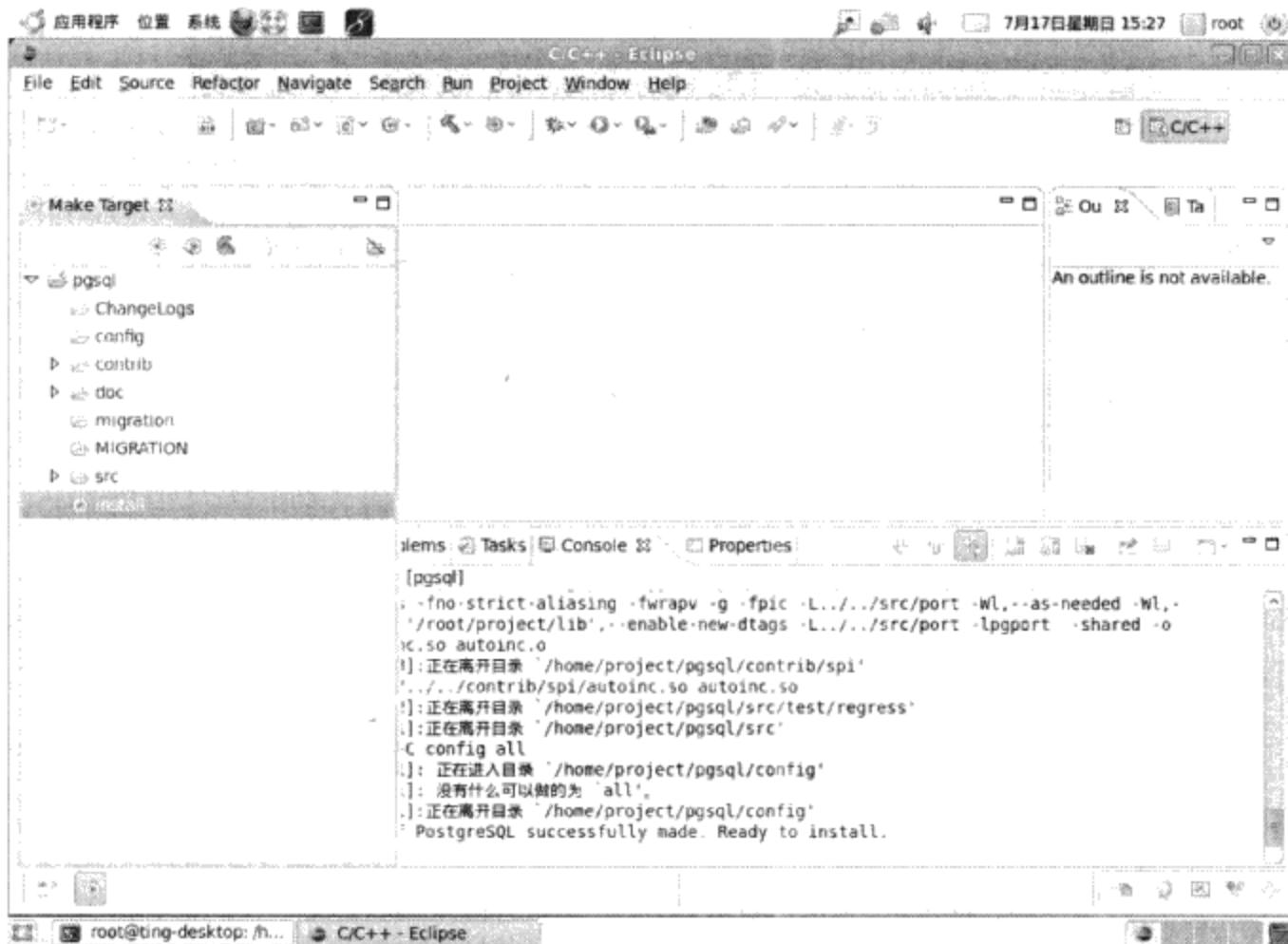


图 A-11 “Make Target” 创建成功

双击 install 项开始安装 PostgreSQL，当看到 Console 栏输出如图 A-12 所示的信息时表示安装完成。

```
C-Build [pgsql]
/bin/mkdir -p /root/project/lib/postgresql/pgxs/src/test/regress
/bin/sh ../../config/install-sh -c pg_regress '/root/project/lib/postgresql/pgxs/src/
test/regress/pg_regress'
make[2]:正在离开目录 '/home/project/pgsql/src/test/regress'
make[1]:正在离开目录 '/home/project/pgsql/src'
make -C config install
make[1]: 正在进入目录 '/home/project/pgsql/config'
/bin/mkdir -p /root/project/lib/postgresql/pgxs/config
/bin/sh ../../config/install-sh -c -m 755 ./install-sh '/root/project/lib/postgresql/pgxs/
config/install-sh'
make[1]:正在离开目录 '/home/project/pgsql/config'
PostgreSQL installation complete.
```

图 A-12 PostgreSQL 安装过程

A.4.2 尝试运行数据库

回到命令行的/home/project/pgsql 目录，输入如下指令初始化数据库：

```
export PATH=$HOME/project/bin:$PATH
export PGDATA=DemoDir
initdb
```

出现如图 A-13 所示的信息说明数据库启动成功。

```

initializing pg_authid ... ok
initializing dependencies ... ok
creating system views ... ok
loading system objects' descriptions ... ok
creating conversions ... ok
creating dictionaries ... ok
setting privileges on built-in objects ... ok
creating information schema ... ok
loading PL/pgSQL server-side language ... ok
vacuuming database template1 ... ok
copying template1 to template0 ... ok
copying template1 to postgres ... ok

WARNING: enabling "trust" authentication for local connections
You can change this by editing pg_hba.conf or using the -A option the
next time you run initdb.

Success. You can now start the database server using:

    postgres -D DemoDir
or
    pg_ctl -D DemoDir -l logfile start
postgres@ting-desktop:~/project/pgsql$ 

```

图 A-13 启动 PostgreSQL 数据库

回到 Eclipse，在工程上单击右键，在右键菜单中选择“Run as”中的“Run Configurations”菜单设置运行配置（见图 A-14）。



图 A-14 “Run Configurations” 对话框

单击“Run Configurations”对话框中的“New launch configuration”按钮进入到新建运行配置的界面（见图 A-15）。在“Main”页的“C/C++ Application”栏输入“src/backend/postgres”，在“Arguments”页中输入“-D DemoDir”。

最后点击“Run”按钮运行 postgres，出现图 A-16 所示的信息说明成功运行了。

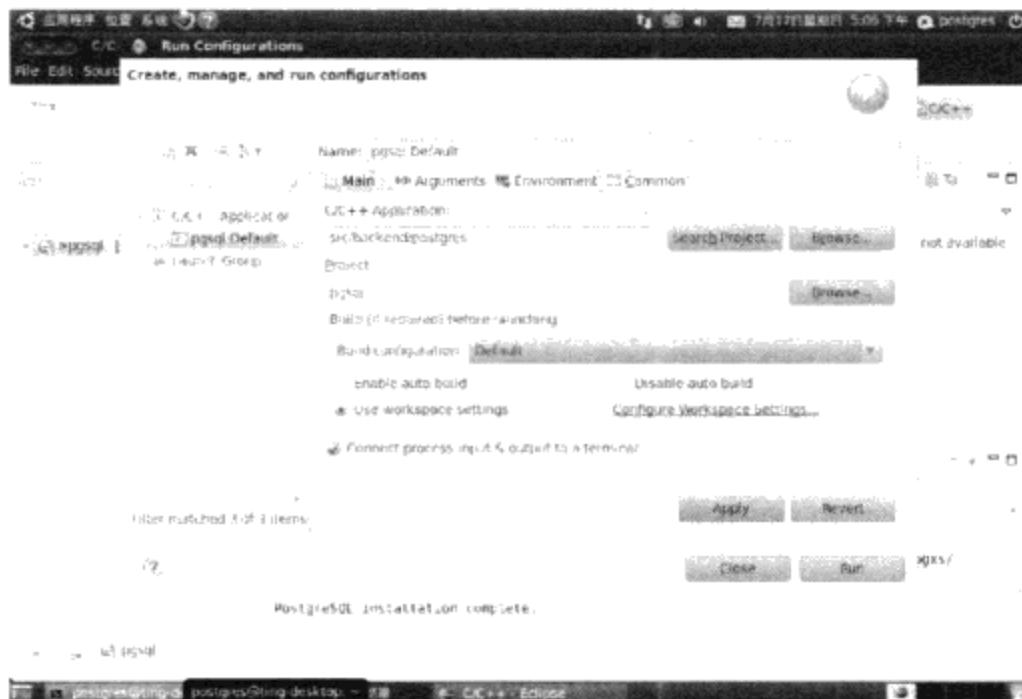


图 A-15 新建运行配置

```
pgsql Default [C/C++ Application] /home/postgres/project/pgsql/src/backend/postgres (11-7-17 下午5:08)
LOG: database system was shut down at 2011-07-17 17:00:10 CST
LOG: database system is ready to accept connections
LOG: autovacuum launcher started
```

图 A-16 在 Eclipse 中运行 PostgreSQL

A.4.3 调试数据库代码

运行成功后，一定要关闭原 Postmaster 进程才能进行调试，否则会提示“Postmaster 已经存在”的错误。

右键单击工程名，在右键菜单中依次选择“Debug As”、“Local C/C++ Application”（见图 A-17）。

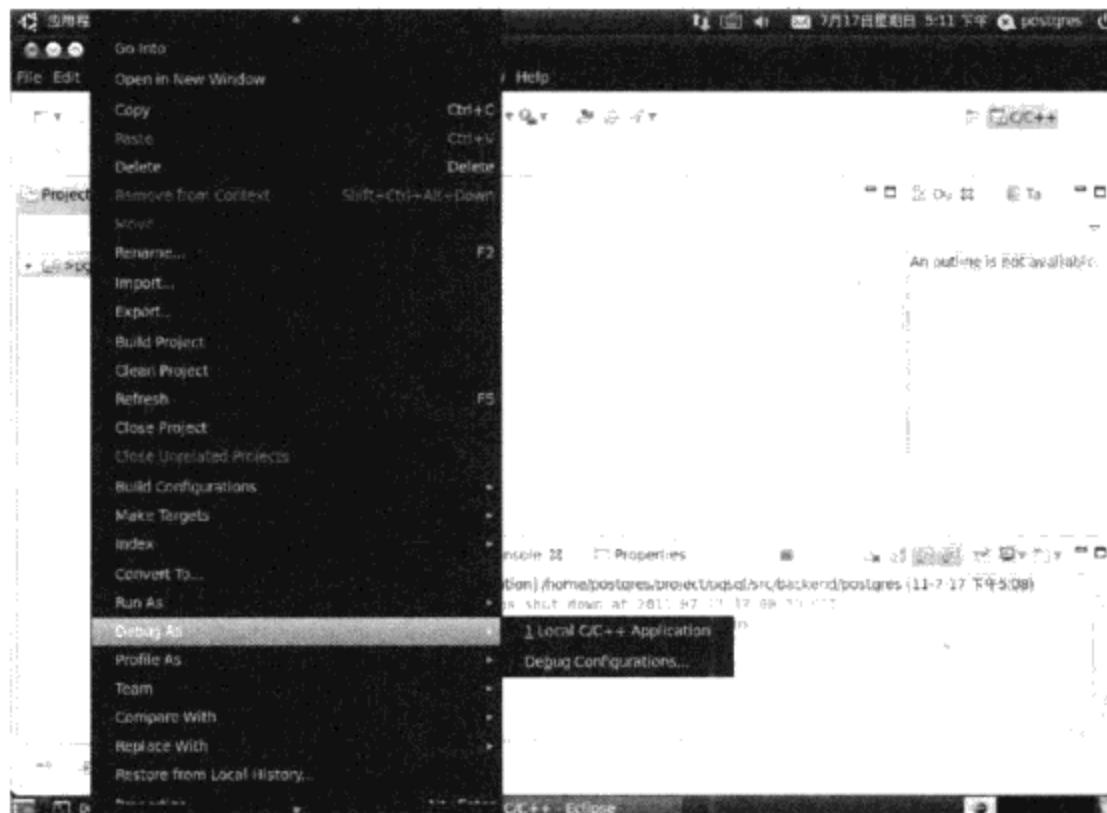


图 A-17 开始调试

然后选择“postgres”作为调试程序（见图 A-18）。

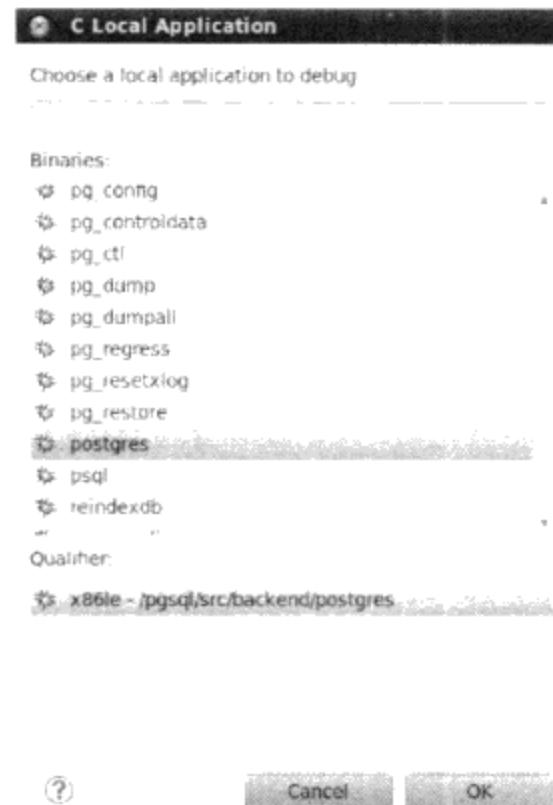


图 A-18 选择调试目标

出现图 A-19 所示界面说明进入调试模式成功。

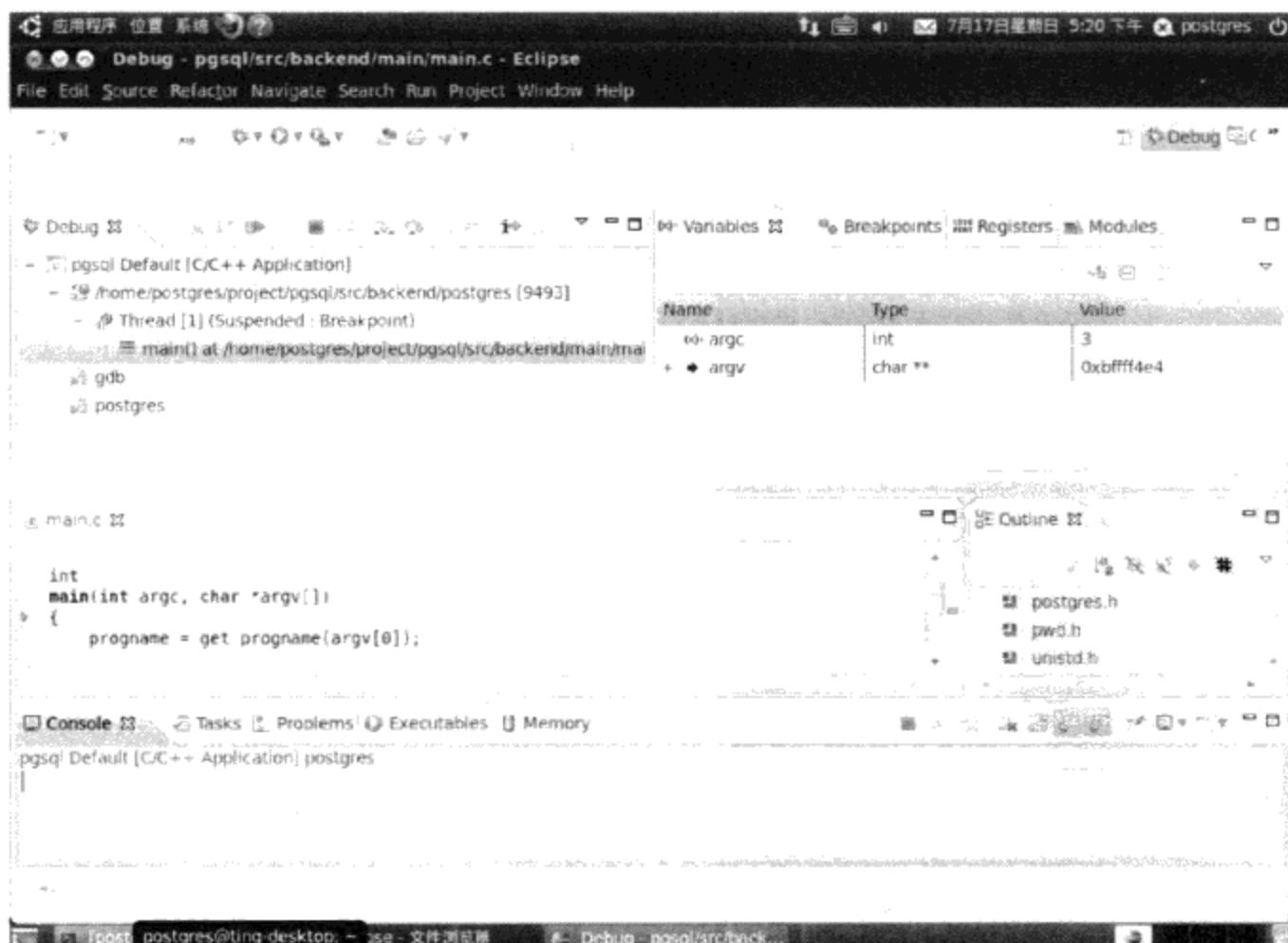


图 A-19 调试界面

因为程序会自动在 main 函数处停下，所以此时要继续点击工具栏的前进按钮让主进程运行完，否则将无法连接数据库。

运行完后会有和运行成功时一样的提示信息（见图 A-16）。

A. 4. 4 服务进程调试

先输入 psql 进行登录，由于是本机，可以直接登录数据库。此时画面会停住，这是由于 Eclipse 中使用了调试模式，因此要回到 Eclipse 里点击前进按钮直到其变灰为止，此时命令行开始有回应了（见图 A-20）。

```
postgres@ting-desktop:~/project/bin$ psql (9.1devel)
Type "help" for help.

postgres=#
```

图 A-20 psql 连接界面

输入“`select pg_backend_pid()`”查看进程号（见图 A-21）。

```
postgres@ting-desktop:~/project/bin$ psql (9.1devel)
Type "help" for help.

postgres=# select pg_backend_pid();
 pg_backend_pid
-----
 10783
(1 row)

postgres=#
```

图 A-21 查看服务进程号

回到 Eclipse 调试界面单击 “Run” -> “Debug configurations” -> “C\C++ Attach To Application”，然后选择创建，填入如图 A-22 所示的信息。

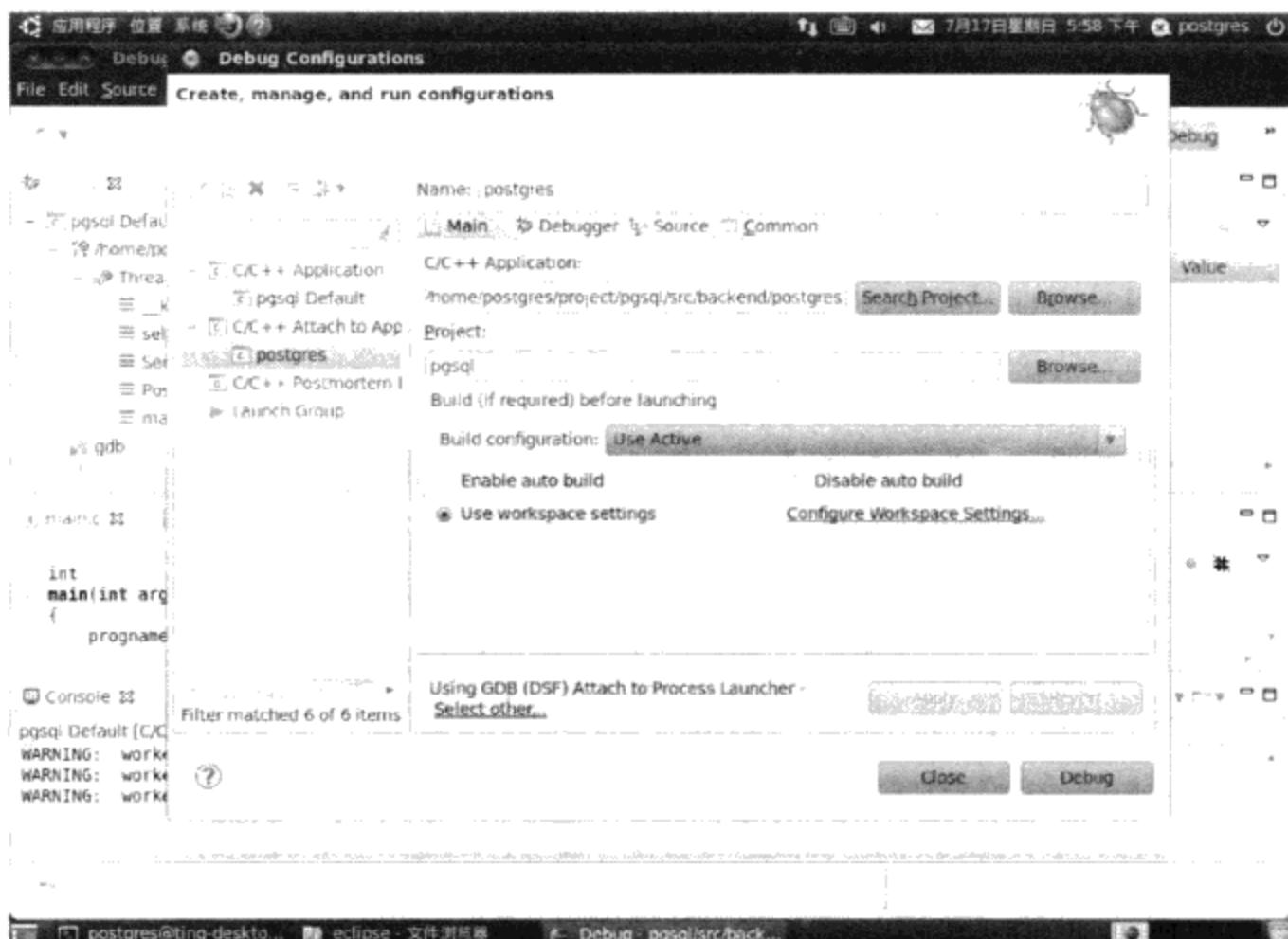


图 A-22 挂接调试进程

在“C/C++ Application”处依然填入 `postgre` 程序的地址，也可以是相对路径；“Build configuration”要选择“Use Active”，然后点击 Debug 按钮，会出现一个进程选择窗口，选择进程号对应的 `postgres` 进程（见图 A-23）。

单击 OK 按钮后，可以看到已经将调试工具挂上了这个 `postgre` 进程（见图 A-24）。



图 A-23 选择挂接进程

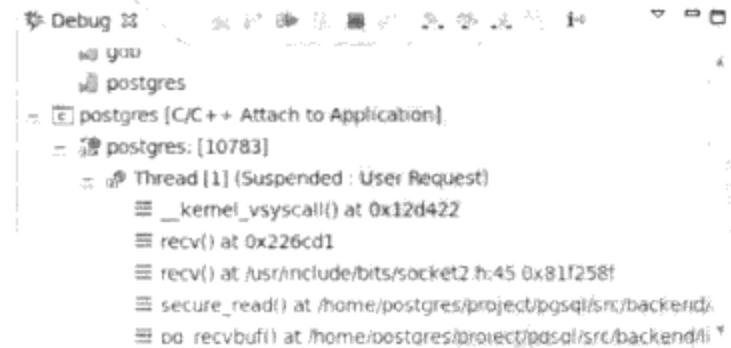


图 A-24 挂接成功

下面可以在代码中设置断点，比如在 `planner.c` 中设置了一个断点（见图 A-25）。

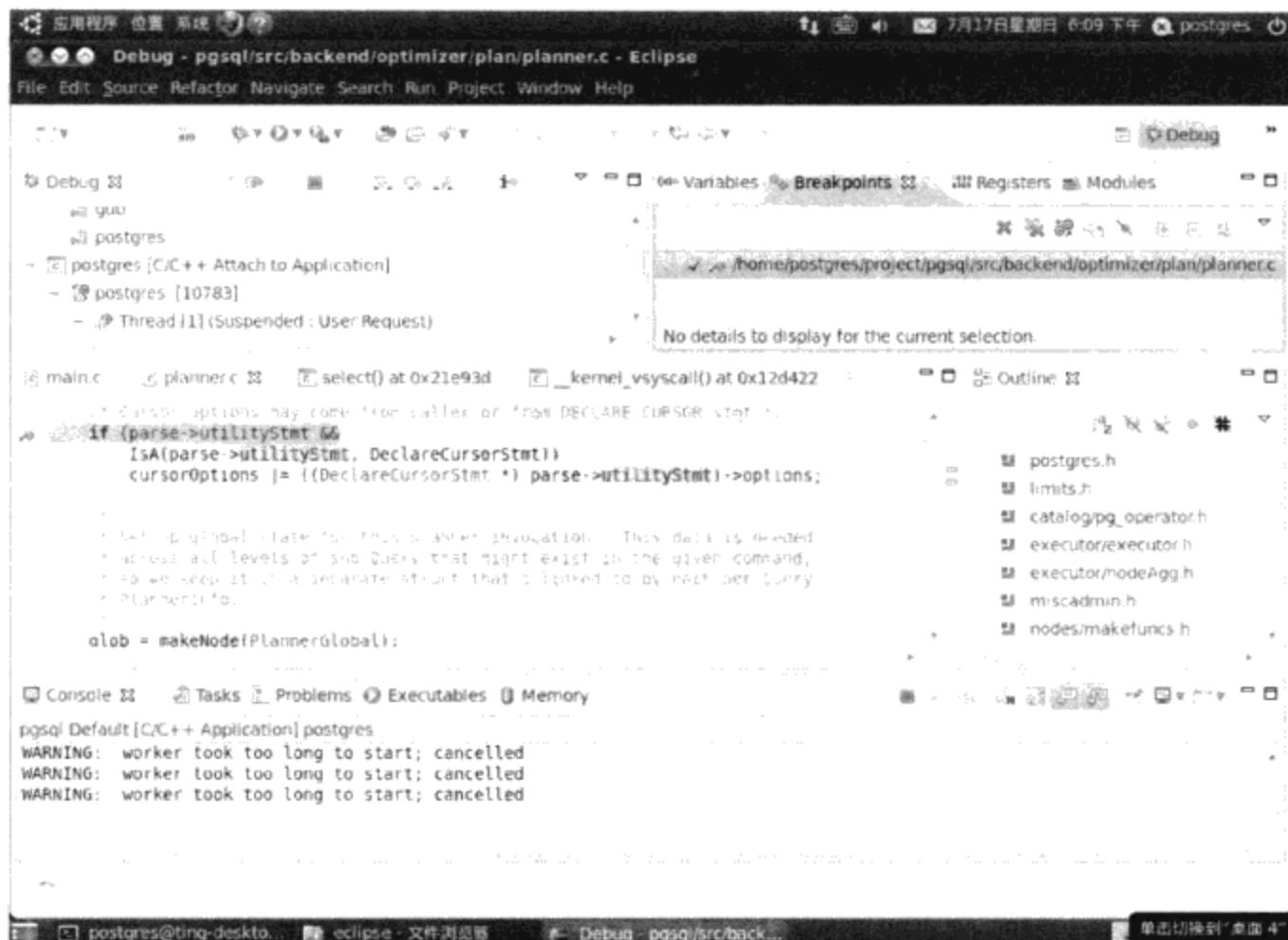


图 A-25 设置断点

然后回到命令行，输入一个 SQL 语句（见图 A-26）。

```
postgres@ting-desktop:~/project/bin$ psql
psql (9.1devel)
Type "help" for help.

postgres=# select pg_backend_pid();
 pg_backend_pid
-----
 10783
(1 row)

postgres=# insert into test values('2');

```

图 A-26 输入 SQL 语句

回到 Eclipse 界面按下前进按钮，如果断点设置正确，程序会停在断点处，然后我们可以通过右上角的 Variables 栏查看变量值（见图 A-27）。

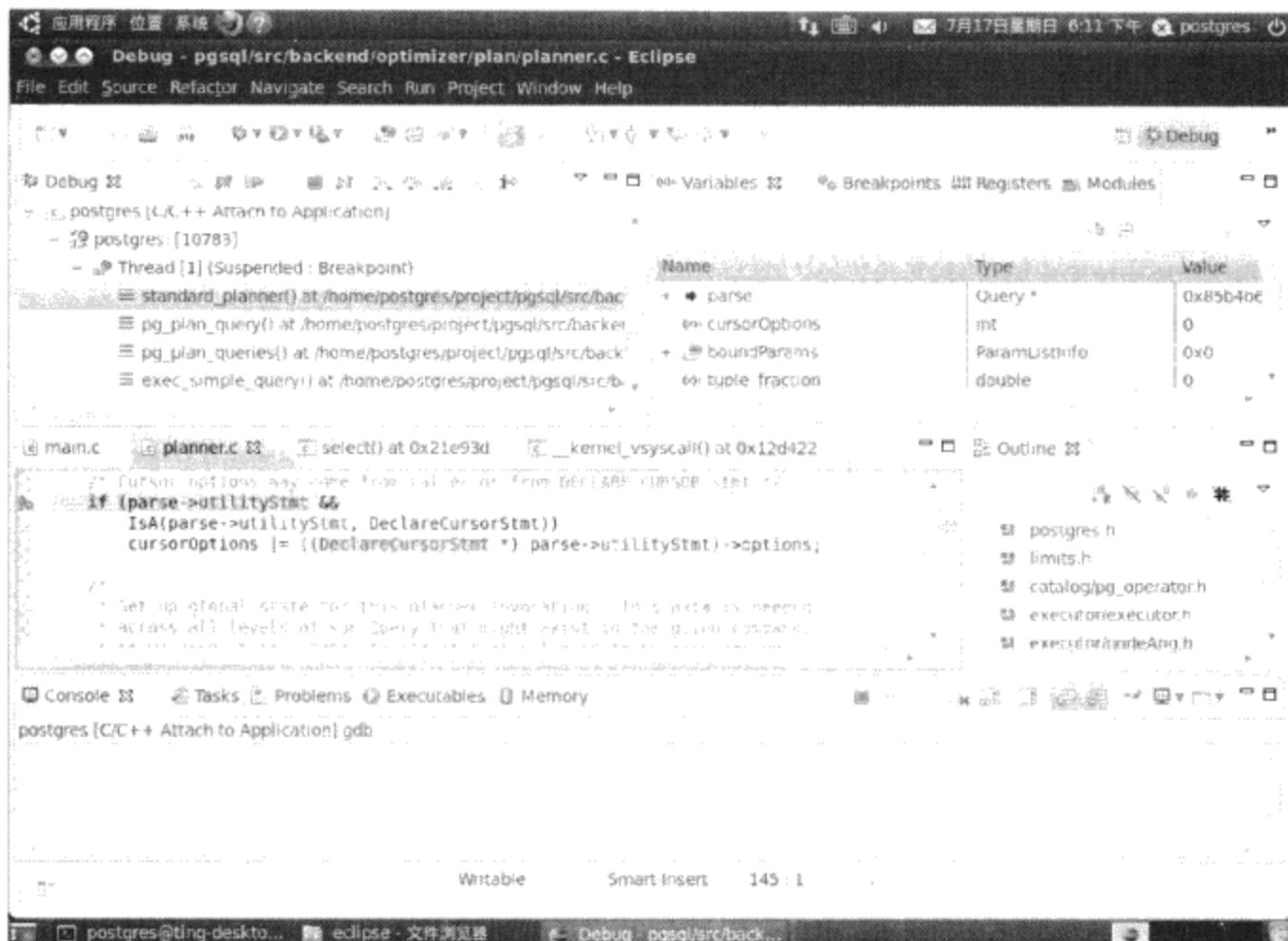


图 A-27 查看变量

到这里，我们就已经可以通过 Eclipse 来追踪 PostgreSQL 的执行过程了，在 Debug 工具栏上提供了“继续”、“停止”、“单步”等按钮，读者可以使用这些按钮控制程序的行进，从而帮助理解源代码。

A.5 小结

通过 Eclipse，我们可以构建一个集成开发环境来调试 PostgreSQL。除此之外，Eclipse 还提供了很多其他的功能，但限于篇幅我们不能一一列举，有需要的读者可以在 Eclipse 官方网站上找到相关的介绍。



一本打开的书，
一扇开启的门，
通向科学圣殿的阶梯，
托起一流人才的基石。

好书推荐



数据库系统基础教程（第2版）

作 者：(美)Jeffrey D.Ullman, Jennifer Widom
译 者：岳丽华等
中 文 版：7-111-10095-6
定 价：32.00
英 文 版：7-111-19344-4
定 价：59.00



数据库系统导论（第8版）

作 者：(美)C.J.Date
译 者：孟小峰 王珊等
中 文 版：ISBN 7-111-21333-8
定 价：75.00



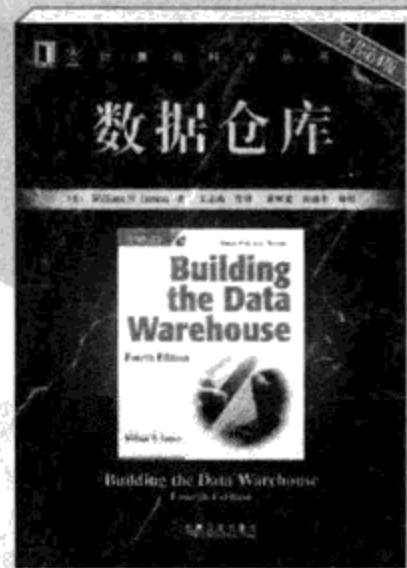
数据挖掘：概念与技术（第2版）

作 者：(加) Jiawei Han
译 者：范明等
中 文 版：7-111-20538-3
定 价：55.00
英 文 版：7-111-18828-4
定 价：79.00



数据挖掘：实用机器学习技术（第2版）

作 者：(新西兰) Ian H.Witten, Eibe Frank
译 者：董琳等
中 文 版：7-111-18205-7
定 价：48.00
英 文 版：7-111-17248-5
定 价：58.00



数据仓库（第4版）

作 者：(美) W.H. Inmon
译 者：王志海等
中 文 版：ISBN 7-111-19194-3
定 价：39.00



欲了解更多华章计算机图书出版动态，敬请您访问华章IT官方微博：<http://blog.csdn.net/hzbooks>
投稿服务热线:010-88379512 教材服务热线:010-88379061
读者服务热线:010-88379061 读者服务邮箱:tianchao@hzbook.com

[General Information]

书名 = PostgreSQL 数据库内核分析

作者 = 彭智勇 , 彭煜玮 编著

页数 = 447

出版社 = 机械工业出版社

出版日期 = 2012

SS号 = 12908238

DX号 =

URL = http://book2.duxiu.com/bookDetail.jsp?d

xNumber=&d=204018350E153ABB9B73FEBC82B56715