

# 阳光学院

## 本科毕业论文（设计）

题目：\_\_\_\_\_  
基于 JammDB 的数据库文件系统  
\_\_\_\_\_  
设计与实现  
\_\_\_\_\_  
学院：\_\_\_\_\_  
信息工程学院  
\_\_\_\_\_  
专业：\_\_\_\_\_  
电子信息工程  
\_\_\_\_\_  
年级：\_\_\_\_\_  
2021 级  
\_\_\_\_\_  
学号：\_\_\_\_\_  
2127380342  
\_\_\_\_\_  
姓名：\_\_\_\_\_  
卓堂越  
\_\_\_\_\_  
指导教师：\_\_\_\_\_  
李荣（中级工程师）  
\_\_\_\_\_

2025 年 4 月

# 基于 JammDB 的数据库文件系统设计与实现

## 摘要

传统文件系统在需要频繁访问大量小文件、频繁进行元数据更新，并且要求支持原子事务的使用场景中，常常面临读写效率低、操作一致性难以保障等问题。相比之下，数据库系统在数据结构组织、高效索引设计、并发处理和事务机制等方面具备成熟经验，这些优势为新型文件系统的设计提供了可借鉴的路径。

本研究以 Rust 编写的嵌入式键值 (key-value) 数据库 JammDB 为底层的存储引擎。该数据库具备良好的性能表现，原生支持强的事务操作和数据持久化等优点，并且采用 B+ 树作为索引结构，有效提升了如路径查找等类文件操作的响应效率。在 Linux 平台上设计并实现了一个 JammDB 数据库为底层的文件系统。该系统遵循 POSIX 文件接口规范，同时具备 FUSE3 用户态支持，能够在用户空间运行与测试。文件操作的事务性与一致性通过数据库的原子提交机制得以保障，性能方面则通过缓存优化、批量更新与多线程调度等手段进行了提升。

基于上述设计与实现，本研究编写自动化测试脚本集成了三种主流评测工具对其进行了功能和性能测试。结果显示，DBFS 在写入密集型与元数据操作场景中表现优于传统文件系统，在高并发读请求下也展现出良好的可扩展性。虽然在顺序读取方面尚存一定的差距，但是文件系统在响应延迟控制方面的任然表现出优异的成绩和展现出较强的应用潜力。

**关键词：**嵌入式键值数据库 数据库文件系统 fuse

# 目 录

第 1 章 绪论.....	1
1.1 研究背景.....	1
1.2 国内外研究现状.....	1
1.3 研究目标和意义.....	2
1.4 主要内容.....	2
第 2 章 相关理论与方法.....	4
2.1 数据库基础知识与 JAMMDB.....	4
2.2 LINUX 文件系统架构.....	4
2.3 数据库文件系统实现技术.....	5
2.4 RUST 编程语言.....	5
2.5 小结.....	6
第 3 章 DBFS 设计与实现.....	7
3.1 DBFS 的整体架构与 FUSE 适配设计.....	7
3.2 数据库引擎.....	8
3.3 DBFS 的接口设计与核心结构.....	9
3.4 元信息管理与目录树构建.....	12
3.5 文件数据存储.....	14
3.6 数据库接口导出.....	17
3.7 小结.....	19
第 4 章 DBFS 的 FUSE 实现与 LINUX 系统集成.....	20
4.1 FUSE 原理与工作流程.....	20
4.2 DBFS 的 FUSE 接口实现.....	21
4.3 FUSE 与 JAMMDB 的高效映射机制.....	23
4.4 性能优化措施.....	25
4.5 小结.....	27
第五章 文件系统性能对比分析.....	28
5.1 测试环境与方法论.....	28
5.2 POSIX 兼容性测试.....	29
5.3 元数据性能测试.....	30
5.4 FIO 性能测试.....	31
5.5 性能测试总结分析.....	33
结论.....	34

参考文献.....	35
附录.....	36
附录一 DBFS 核心代码片段 .....	36
附录二 核心数据结构.....	36
附录三 接口实现与数据结构.....	37
附录四 DBFS 测试图 .....	39
致谢.....	40

# 第 1 章 绪论

## 1.1 研究背景

计算机文件系统通过“文件”和树状目录的方式，将底层物理设备（如硬盘或光盘）上的数据块管理细节隐藏起来，为用户提供了简单直观的数据访问接口。用户只需记住文件所在的路径和名称，不用关心数据究竟存储在那些物理扇区；系统则自动承担空间分配与回收的工作，确保存储操作的透明和高效。

然而，随着数据规模的飞速增长以及应用场景的日益多样化，传统文件系统在处理海量小文件、频繁的元数据操作或需要事务支持的场景时，性能往往难以满足需求。相比之下，数据库技术凭借几十年积累的索引优化、并发控制和事务管理经验，提供了更加成熟的解决方案。

将数据库的优势引入文件系统，以其高效索引、可靠事务和出色并发能力来提升文件系统的整体表现，已成为一个颇具前景的研究方向。尤其是键值数据库，其数据模型简单、读写效率高，为文件系统提供了理想的后端存储方案。

## 1.2 国内外研究现状

数据库文件系统的研究最早可追溯至上世纪 90 年代。其中，微软的 WinFS 项目是早期将关系数据库与文件系统相结合的典型尝试之一。虽然该项目最终未能实现商业化，但其探索为后续研究提供了宝贵的经验和参考。此外，Oracle 推出的数据库文件系统 DBFS 允许将数据库表空间挂载为文件系统，主要面向企业级应用场景，这一实践进一步推动了该领域的发展。

开源社区，PostgreSQL 的 TableFS、MongoDB 的 GridFS 等项目也对不同的数据库作为文件系统后端的可行性进行了研究，扩大了该方向的研究范围，学术界也进行了研究，例如麻省理工学院（MIT）的 InversionFileSystem、卡内基梅隆大学（CMU）的 KVFS 等等，它们从各自的角度研究了数据库与文件系统的结合方式和技术路径。

与国外研究相比，国内在这方面的研究很少，基本上都是在一些应用场景上对文件系统进行了优化，比如分布式文件系统、对象存储等。而基于键值数据库实现一个完全 POSIX 兼容的文件系统，并且进行内核级实现的研究更是少之又少。

当下，已有研究存在一些关键短缺之处，其一，大多系统只在用户空间执行，所以它们的性能被限制了，其二，已有的系统对 POSIX 的兼容性不完全，很难与既有的应用程序匹配起来，其三，缺乏系统的性能考量与改良手段，于是这些系统在实际应用中的表现就无法得到全面的考察与改进。

### 1.3 研究目标和意义

本研究旨在设计并实现一种基于键值数据库的文件系统 DBFS，探索数据库技术在文件系统领域的应用潜力。具体研究目标包括：

设计一种基于 JammDB 键值数据库的文件系统架构，实现文件系统对象到键值对的高效映射，确保完整的 POSIX 兼容性与现有应用程序无缝集成，同时支持用户空间（FUSE）和内核空间两种实现方式，充分利用数据库事务机制提供强一致性保证和崩溃恢复能力，并通过系统化性能测试与优化提高 DBFS 在元数据操作和小文件处理方面的效率。

本研究具有重要的理论和实践意义。在理论层面，研究探索了数据库与文件系统这两种不同数据管理范式的融合模式，丰富了文件系统设计理论。通过将文件系统的层次结构映射到键值存储模型，研究提出了一种新型的数据组织方法，为解决传统文件系统在元数据处理和小文件存储方面的局限性提供了新思路。DBFS 的设计模型挑战了传统文件系统的设计范式，探索了一种将数据库事务特性与文件系统操作语义相结合的新方法，为文件系统理论研究提供了新的视角。

从实践方面来讲，DBFS 给一些特定应用场景带来了新的存储方案。对于元数据密集型的应用，像 Web 服务器，邮件系统，小文件存储场景，比如图片缩略图，日志文件等，DBFS 设计可以带来更有效的存储和访问机制。同时，DBFS 的事务支持为需要数据一致性保证的应用提供了额外的安全保障。DBFS 的实现填补了国内在基于键值数据库实现完整 POSIX 兼容文件系统的研究空白，其内核级实现更是克服了现有研究中多数系统仅在用户空间实现的局限性，为高性能数据库文件系统的实际应用奠定了基础。

而且，这个研究使用 Rust 语言完成的，Rust 语言有内存安全、并发方面的优点，给系统软件研发提供了一个例子，DBFS 采用双轨执行，就是用户空间和内核空间这两种方式，又给文件系统研发提供了一种协调开发速度和运行效率的新思路。

### 1.4 主要内容

本论文共分为六章，主要内容如下：

第一章为绪论，介绍研究背景、国内外研究现状、研究目标和意义，总结了本文的主要内容。

第二章详细阐述相关理论与方法，牵涉到数据库的一些基本常识、JammDB 特性、Linux 文件系统架构、数据库文件系统实现技术以及 Rust 编程语言的关键特性。

第三章着重论述 DBFS 的设计及达成，系统架构，数据组织与管理，事务与一致，缓存与性能，安全与权限，错误处理与恢复。

第四章主要介绍了 DBFS 在 FUSE 框架下的实现。

第五章进行文件系统性能对比分析，通过 POSIX 兼容性测试、元数据性能测试和 FIO 性能测试。

总结部分，总结了 DBFS 在设计理念、实现技术和性能特点方面的创新，同时客观分析了当前实现的局限性。

通过前面章节的研究探讨，本论文从各个方面都可以体现出利用键值数据库来做文件系统的各种设计和实现都是可行的并且有优势的地方，提出了一些在文件系统范畴里的更新的某些新想法和参考计划，DBFS 的成功完成既证明了数据库技术可以应用到文件系统里，也给解决传统文件系统在某些情况下遇到的性能问题提供了一种可能的解决方案。

## 第 2 章 相关理论与方法

### 2.1 数据库基础知识与 JammDB

数据库系统是用于组织、存储和管理数据的专业化软件系统。经过数十年的发展，数据库技术已经形成了完整的理论体系和技术架构。在文件系统设计领域，特别是基于数据库的文件系统开发中，数据库的核心概念与关键技术发挥着至关重要的作用。

事务是数据库管理系统执行过程中的一个逻辑单位，具有原子性、一致性、隔离性和持久性（ACID 特性）。原子性确保事务中的所有操作要么全部完成，要么全部不完成；一致性保证数据库从一个一致性状态转变为另一个一致性状态；隔离性确保多个事务并发执行时互不影响；持久性则确保已提交的事务永久保存<sup>[1]</sup>。在文件系统中，事务机制可以确保文件操作的原子性和一致性，例如创建文件时需要同时更新目录项和分配 inode 这样的复合操作<sup>[2]</sup>。

一般索引是提高查询效率的关键数据结构并且通常都是采用 B 树或 B+树实现。B+树作为一种多路平衡查找树，具有所有叶子节点位于同一层、非叶子节点只存储键值等特点，能够将查找时间复杂度控制在  $O(\log n)$ 。缓存机制通过将频繁访问的数据保存在内存中，显著减少磁盘 I/O，提高访问速度。常见的缓存策略包括 LRU 和 LFU 等<sup>[3]</sup>。

JammDB 数据库是数据库文件系统的核心存储引擎，这是一个用 Rust 编写的嵌入式键值数据库。它采用桶和键值对的两级存储结构，基于 B+树实现高效索引，并通过写前日志确保数据持久性<sup>[4]</sup>。JammDB 通过事务日志实现原子性，通过事务机制和约束检查保证一致性，支持锁机制和 MVCC 实现隔离性，采用 WAL 机制确保持久性。为了优化性能，JammDB 实现了内存缓存机制、批量操作支持、异步写入和数据压缩等特性<sup>[5]</sup>。

### 2.2 Linux 文件系统架构

Linux 属于单内核操作系统，它的文件系统架构的关键部件就是虚拟文件系统（VFS, VirtualFileSystem），VFS 在 Linux 内核里占据着非常重要的位置，它既给上层的应用程序供应了统一的文件操作接口，又做到了不同种类文件系统的共存并协同工作，通过使用 inode, dentry, superblock 和 file 这些核心的数据结构，VFS 形成起一个完备的抽象层，应用程序借助这个抽象层就可以通过统一的系统调用来访问各种各样的文件系统，而不用去关心底层的实现细节<sup>[6]</sup>。

在文件系统架构里，Linux 依靠多层次的缓存机制明显改善了文件数据的访问速度，页缓存加快了文件数据的读取，inode 缓存加快了元数据操作，dentry 缓存加快了目录结构查询，经过细致规划的多级缓存体系，Linux 缩减了磁盘 I/O 操作，从而提升了整个



系统的性能，而且，预读取，回写这些机制也改进了数据访问模式，提高了系统吞吐量<sup>[7]</sup>。

FUSE (Filesystem in Userspace) 是 Linux 内核提供的创新框架，它通过内核模块与用户空间守护进程的协作机制，使开发者能够在用户空间实现完整的文件系统功能<sup>[8]</sup>。当应用程序发起文件操作请求时，内核中的 FUSE 模块接收这些请求并通过/dev/fuse 设备文件转发至用户空间，由 FUSE 守护进程进行处理后将结果返回。这种设计极大地简化了文件系统开发流程，也为实验性文件系统提供了理想的开发环境<sup>[9]</sup>。

Linux 文件系统的一大特征就是它的高度模块化设计以及众多的功能特性，统一的命名空间使得用户体验保持一致，多种文件系统并存能够满足不同场景下的需求，高效的缓存机制，异步 I/O 支持以及日志功能在性能与可靠性方面达到了较好的平衡，这些特点构成了 Linux 存储子系统的根基，并为数据库文件系统的实现提供了有力支持<sup>[10]</sup>。

## 2.3 数据库文件系统实现技术

数据库文件系统的核心是其数据模型，关键是如何将传统的文件系统概念映射到数据库模型中，最常见的方式是采用元数据和实际数据分离的方式，即将文件系统的元数据存储数据库中，而文件内容存储在文件系统中或者数据库的二进制大对象 (BLOB) 字段中，这个过程需要对层次结构的映射以及 inode 模型进行精心设计，才能保证文件系统的高效性和可扩展性<sup>[11]</sup>。

事务管理和一致性保障是数据库文件系统的关键特性之一，数据库所赋予的事务机制使得文件系统能够保证复杂操作的原子性、一致性、隔离性和持久性 (即 ACID 特性)，日志记录与崩溃恢复机制是保障系统可靠性的关键手段，在系统崩溃情况下能有效恢复数据，为了改善系统性能，一般采用批处理操作、异步提交以及检查点机制等优化策略，从而缩减响应延迟，加强系统吞吐量<sup>[12]</sup>。

性能改良属于数据库文件系统遭遇的严峻考验，通过运用索引改良，缓存机制，并发控制，批处理操作以及数据划分这些技术手段，可以明显改善文件系统的总体性能，而且，零拷贝，内存映射，异步 I/O 以及数据压缩等技术联合起来加以应用之后，又使得文件系统的效能与资源利用率得到更进一步的优化。

## 2.4 Rust 编程语言

Rust 是一种专门针对系统编程而设计的语言，它强调安全、并发以及高性能，通过引入所有权系统、借用机制以及生命周期等概念，Rust 可以在不依赖垃圾回收的前提下实现内存安全，避免出现内存泄漏或者数据竞争等问题，Rust 的类型系统以及并发模型也保证了线程安全，这对于开发可靠的系统软件来说非常重要<sup>[13]</sup>。

在系统编程领域，Rust 的优势是性能好，内存安全，零成本抽象，Rust 可以做到和 C/C++ 一样的性能，而且安全性更好，Rust 有较好的跨平台支持和丰富的开发生态，所以 Rust 适合用来开发文件系统、数据库这些高性能的系统软件<sup>[14]</sup>。

## 2.5 小结

本章细致论述了 DBFS 达成所必需的关键理论和技术，涵盖数据库基本知识，JammDB 特性，Linux 文件系统构架，数据库文件系统达成技术以及 Rust 编程语言特点等内容，这些理论和技术给 DBFS 的设计与达成赋予了稳固根基，接下来的章节会细致论述 DBFS 的具体设计与达成方案。

## 第 3 章 DBFS 设计与实现

### 3.1 DBFS 的整体架构与 FUSE 适配设计

DBFS 是一个以键值数据库为基础、模块化的文件系统，它主要在 Linux 环境中运行。如图 3-1 所示，整个系统由两个主要部分组成：数据库引擎和文件系统实现。数据库引擎提供了一个可靠的存储和检索数据的服务，而文件系统实现则是把文件系统的操作转换成相应的数据库操作。本次研究采用 jammdb 作为底层存储引擎，选择的理由将会在 3.2 节中进行阐述。

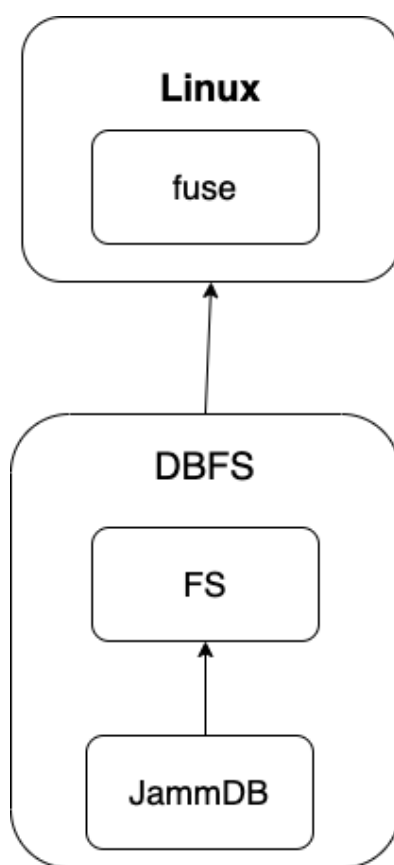


图 3-1 DBFS 系统架构图

DBFS 采用自下而上的分层接口设计架构。该架构展示了 DBFS 的分层设计和数据流向。从上至下，用户的文件系统请求首先通过 Linux 操作系统层到达 FUSE 接口。FUSE 接口作为用户空间和内核空间的桥梁，将文件系统请求转发给文件系统实现模块<sup>[15]</sup>。

文件系统实现模块是 DBFS 的核心，它借助请求解析与转换机制把 FUSE 文件操作请求转换成对应的数据库操作，而且还要承担元数据管理任务来守护文件和目录的属性信

息，该模块也完成了文件系统层级的权限控制，通过缓存管理机制改善频繁被访问的数据。

数据库引擎（JammDB）是底层存储基础，使用了各种方式保证数据可靠。使用事务来保证数据的原子性、一致性，使用多读单写并发控制模式。在并发控制上，jammdb 使用严格的单写多读模式，即在同一时间只有一个写事务，但是可以有多个并发的读事务。这种方式虽然在高并发写入的情况下需要额外的同步机制，但是它可以避免数据的竞争，保证数据的一致性。系统通过写前日志 WAL 和 B+树索引实现数据的持久化和快速查询。

所有数据最终存储在用户态存储层，这种设计避免了内核态的复杂性，提供了更好的可维护性和调试能力。DBFS 通过 FUSE（Filesystem in Userspace）接口在 Linux 系统中作为用户态文件系统运行，无需修改内核代码，便于开发和调试。

本研究达成了 DBFS 同 Linux 系统 FUSE 模块的完美融合，从而可以将其挂载起来，就像对待标准文件系统那样去使用它，FUSE 适配层承担着把 FUSE 接口调用转变为 DBFS 通用接口调用的任务，处理用户空间与内核空间之间数据传输的工作，这种设计为下一章要展开论述的 FUSE 实现形成了根基。

## 3.2 数据库引擎

jammdb 是一个嵌入式、单文件的 key-value 数据库，其提供 ACID 特性，支持多个并发读取和单个写入。所有的数据被组织成一棵 B+ 树，随机和顺序读取速度很快，如表 3-1 具有以下核心特性：

表 3-1 jammdb 数据库核心特性表

特性类别	具体特性	对 DBFS 的意义
存储模型	键值对(key-value)存储	适合文件系统元数据和数据块的组织结构
事务支持	完整 ACID 特性	确保文件系统操作的一致性和可靠性
并发能力	多读单写模型	支持并发读取，写入时保证数据一致性
数据组织	B+树索引结构	提供高效的查找和范围扫描，适合目录遍历
性能特点	快速随机和顺序读取	提高文件系统的读取性能
内存管理	基于 mmap 的内存映射	利用操作系统页缓存机制，减少内存拷贝开销
部署方式	嵌入式、单文件	易于集成到文件系统中，无需额外依赖，部署灵活

选择 jammdb 作为 DBFS 实现的基础主要考虑其结构简单，易于集成到 FUSE 文件系统中。这种简单性虽然带来了实现上的便利，但也存在一些局限性：

1. 高并发写入支持有限：仅支持单写多读模型，在多线程写入场景下需要额外的同步机制
2. 缺乏数据压缩：未内置数据压缩功能，可能导致存储空间利用率不高
3. 备份恢复机制简单：没有提供增量备份和时间点恢复等高级特性
4. 性能优化空间：未针对特定磁盘设备（如 SSD）进行优化，在某些 I/O 模式下性能可能不理想
5. 监控管理能力：缺乏内置的性能监控和统计功能，不利于系统调优

这些局限性可能导致 DBFS 在某些场景下性能不能达到最佳状态，但对于一个实验性的文件系统而言，其简单可靠的特性仍然是一个合理的权衡。

jammdb 的内部基于桶(bucket)实现，其结构特点如表 3-2：

表 3-2 jammdb 内部结构表

组件	描述	实现细节
Bucket	数据库读基本组织单元	位于全局命名空间，可嵌套
存储内容	键值对或嵌套bucket	键与值均为字节数组[u8]
内部组织	B+树索引结构	每个bucket由一棵独立B+树构成
存储单位	页面（Page）	固定大小（通常4KB），类似文件系统块
缓存机制	内存映射（mmap）	利用操作系统页缓存机制
读写特性	读内存/写磁盘	读操作在内存映射区域，写操作同步到磁盘
事务实现	写时复制（COW）	修改时创建新页面，

这种基于 B+树和页面的设计使 jammdb 能够高效地组织和访问键值数据。其中，存储内容采用字节数组格式提供了灵活的数据表示能力；B+树索引结构确保了快速的查找和范围扫描性能；固定大小的页面单位简化了存储管理；内存映射机制充分利用了操作系统的缓存能力；读写分离和写时复制策略则保证了数据的一致性和可靠性。这些特性共同为 DBFS 提供了可靠的存储基础。

### 3.3 DBFS 的接口设计与核心结构

图 3-2 显示了 DBFS 的接口设计。自下而上，DBFS 由各层接口连接起来，且每个层都是一个独立的模块，可以被其他项目所复用。各层的功能描述如下：

所有数据最终存储在用户态存储层，采用 4KB 大小的数据块进行组织。这种设计不仅避免了内核态的复杂性，提供了更好的可维护性和调试能力，还通过合理的数据块大小设计提升了存储效率和访问性能。

DBFS 通过 FUSE (Filesystem in Userspace) 接口在 Linux 系统中作为用户态文件系统运行，无需修改设计，每层都是独立模块，可被其他项目复用。如图 3-2 所示，DBFS 由五层接口组成，从底层到顶层依次是存储层接口、数据库底层接口、数据库接口、DBFS 通用接口和适配层接口。各层功能和接口设计如下：

DBFS 采用五层接口设计，从底层到顶层依次是存储层接口、数据库底层接口、数据库接口、DBFS 通用接口和适配层接口。存储层接口作为最底层接口，通过 open、read、write、close 等核心操作实现数据持久化，直接与文件系统交互管理原始数据，同时实现同步和异步 I/O 模式，并提供基本的错误处理机制。

数据库底层接口抽象了底层存储访问，通过 read\_page、write\_page 实现页面读写，通过 map\_memory、unmap\_memory 提供内存映射功能，并通过 sync、truncate、close 等操作实现文件控制。该层还负责实现页面缓存和预读取，以及管理文件锁定和并发访问。

数据库接口层提供完整的数据库操作服务，包括 begin\_tx、commit\_tx、rollback\_tx 等事务管理功能，create\_bucket、delete\_bucket、get\_bucket 等 Bucket 操作，以及 put、get、delete、cursor 等键值操作。该层实现了 ACID 特性，并提供批量操作支持，确保数据一致性和可靠性。

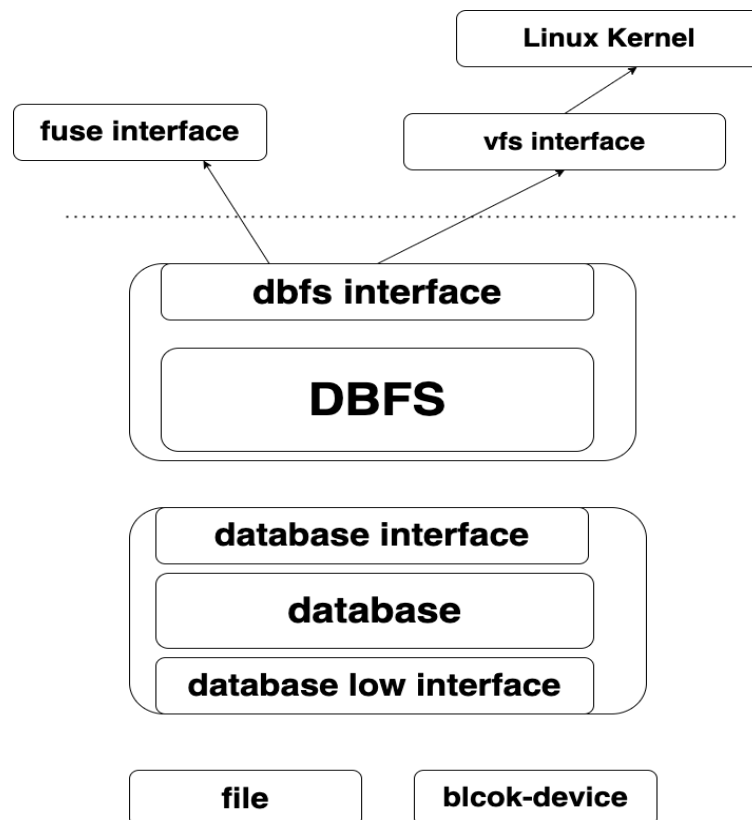


图 3-2 DBFS 接口分层架构图

DBFS 通用接口层实现了标准文件系统功能，包括 create、open、read、write、close 等文件操作，mkdir、readdir、rmdir 等目录操作，getattr、setattr、statfs 等元数据操作，以及 symlink、readlink、link、unlink 等链接操作。该层还负责权限和访问控制，并支持扩展属性。

最上层的适配层接口负责将 FUSE 操作映射到 DBFS 通用接口，处理用户空间与内核空间的数据传输。该层实现了错误码转换，管理文件描述符，并通过优化请求调度和实现缓存策略来提升系统性能。

DBFS 内部包含一系列关键数据结构，如表 3-3 主要分为以下几类：

表 3-3 DBFS 核心数据结构表

数据结构类别	主要结构	功能描述
文件系统元数据	DbfsAttr, DbfsStat	存储文件属性和统计信息
错误处理	DbfsError, DbfsResult	统一错误处理和结果返回
权限管理	DbfsPerm	文件权限控制
缓存管理	CachePool	缓存频繁访问的数据结构
全局信息	Globalinfo	存储文件系统全局状态
数据库实体	DB	核心数据库连接和操作句柄

在 DBFS 的初始化阶段，用户态的 FUSE 或内核的文件系统初始化函数会创建一个数据库实体，并将其初始化为 DBFS 的全局数据结构。初始化完成后，DBFS 将通过这个数据库实体完成所有文件操作。

缓存池 (CachePool) 的设计主要有两个目的：一是加速频繁发生的请求 (如 readdir)；二是避免不必要的内存分配开销，特别是在 FUSE 实现中，通过局部缓存分配器可以带来更好的内存局部性。

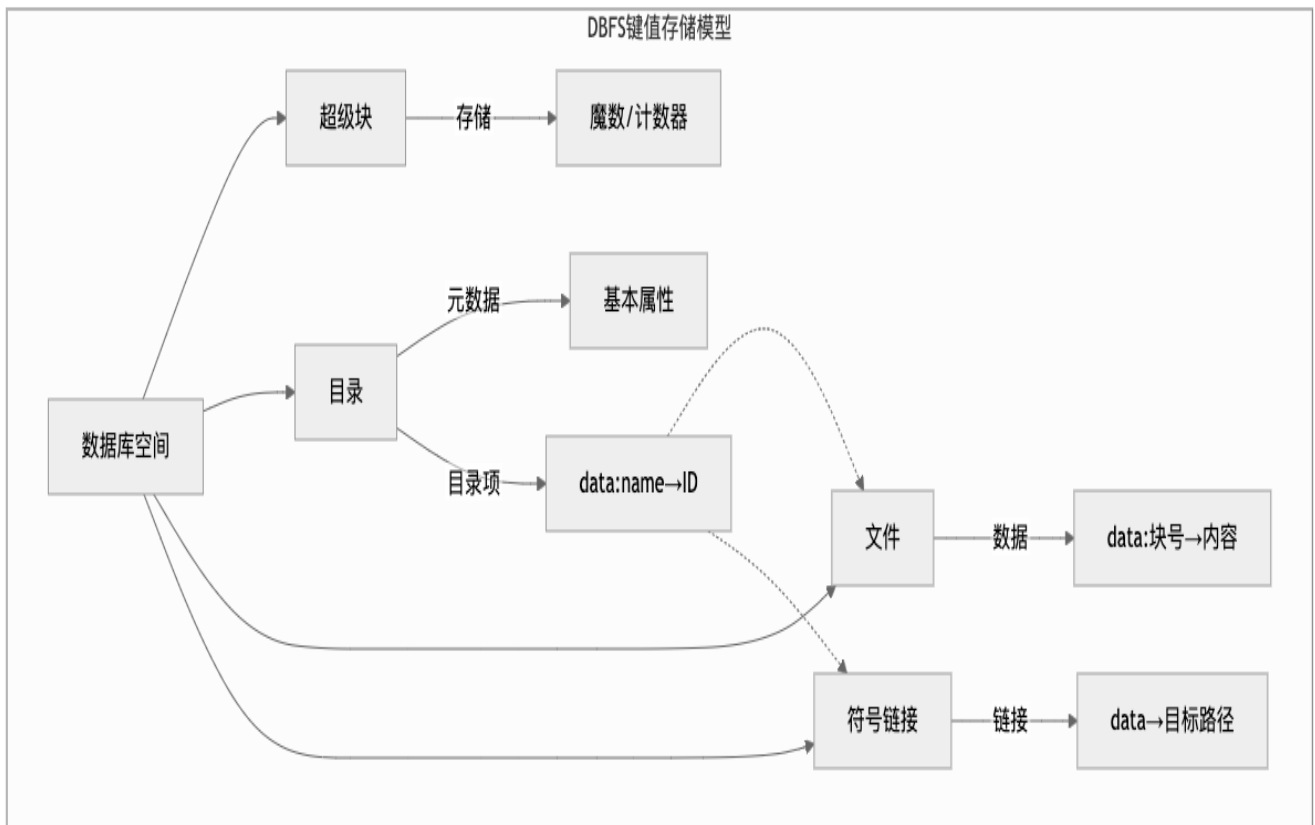


图 3-3 DBFS 键值存储模型图

### 3.4 元信息管理与目录树构建

图 3-3 使用 jamddb 数据库的数据结构来构建文件系统的图，该图的目的是为了满足不同 Linux 文件系统的路径解析，比如路径/d1/dd1/f1，系统会先解析/d，再解析 d1，再解析 dd1，再解析 f1，即从根目录开始递归查找每一个路径分量。

如图 3-3 所示，DBFS 采用键值对存储模型组织文件系统元素。该图展示了系统的分层存储结构和数据组织方式：

1. 全局命名空间：数据库全局空间作为最顶层的组织单位包含所有 bucket，每个 bucket 代表一个独立的文件系统对象，通过统一的命名空间确保对象的隔离性并简化了对象的查找和管理
2. 唯一标识符：数据库中每个 bucket 都拥有一个唯一的 64 位递增 ID，这种 ID 机制类似传统文件系统的 inode 编号，用于对象寻址。通过递增特性，系统保证了 ID 的唯一性和时序性，有效支持了文件系统的对象管理和访问控制。
3. 对象类型与组织：DBFS 中的对象类型包含超级块、目录、文件和符号链接，它们各自承担不同的功能并有着各自的组织形式，超级块存放着文件系统的全局配置信息，依靠魔数来识别文件系统类型，而且还要保存块大小之类的系统参数，目录对象要存储



并管理子文件和子目录的映射关系，从而维持文件系统的层次结构，做到高效的路径查找和遍历，文件对象采取分离存储策略，独立管理文件数据和元数据，凭借数据块分散存储的方式，可以达成高效的随机访问，符号链接对象存储着指向目标的路径信息，可以完成跨目录的文件引用，而且还要经过完整性检查以保证链接的可靠。

4. 键值命名规则具体如下表 3-4:

表 3-4 DBFS 键值命名规则表

键类型	命名规则	示例	用途
元数据键	直接使用字段名	uid, mode	存储对象元数据
目录项键	Data: 文件名	Data: file1	存储目录项引用
文件数据键	Data: 块序号	Data: 0	存储文件数据块
符号链接键	data	data	存储链接目标路径
扩展属性键	命名空间前缀	User. attr1	存储扩展属性

这种基于键值对的存储模型使 DBFS 能够灵活组织和管理文件系统数据，同时保持与传统文件系统类似的层次结构和操作语义。键值命名规则的设计充分考虑了不同类型数据的特点：元数据键直接使用字段名便于理解和维护；目录项键通过前缀区分确保命名空间隔离；文件数据键采用块序号实现高效的随机访问；符号链接键和扩展属性键则采用简单统一的格式，方便扩展和管理。这种命名规则不仅保证了键的唯一性，还提供了良好的可读性和可维护性。

超级块是 DBFS 中的特殊 bucket，存储文件系统全局信息如下表 3-5 所示：

DBFS 挂载时，初始化函数从存储设备读取超级块信息，校验后将 continue\_number 加载到全局变量。创建文件或目录时，系统读取并自增该变量。文件系统刷新或卸载时，该值会被写回超级块。

表 3-5 DBFS 超级块结构表

字段名	数据类型	描述	用途
magic	整体（字节序列）	魔数（Magic Number）	用于标识文件系统类型（DBFS）
continue_number	64 位整数	递增计数器	用于生成唯一的文件标识符（ID）
blk_size	整数	数据块大小	定义每个数据块的大小，以字节为单位
disk_size	整数	磁盘总大小	表示整个文件系统可用的存储空间总量

每个文件或目录 bucket 中存储的元数据具体看表 3-6:

表 3-6 DBFS 元数据字段表

元数据字段	描述	FUSE接口映射
uid	所有者用户ID	Getattr/setattr
gid	所属用户组ID	Getattr/setattr
mode	文件模式和权限	Getattr/setattr/chmodr
size	文件大小/目录项数	Getattr/truncate
ctime	元数据修改时间	Getattr/utimens
mime	内容修改时间	Getattr/utimens
atime	最后访问时间	Getattr/utimens
Hard_link	硬链接计数	Getattr/link/unlink
Blk_size	块大小	Getattr

这些元数据满足大多数文件操作需求。创建文件或目录时，系统创建新 bucket 并初始化这些字段；删除时，对应 bucket 及其所有元信息一并删除。

DBFS 支持文件系统扩展属性(extended attributes)，允许为文件和目录附加额外元数据。与传统元数据不同，DBFS 将扩展属性作为普通键值对存储，无需特殊存储结构。这种设计有以下优势：

1. 灵活性：不限制扩展属性的键值对大小
2. 一致性：与其他元数据使用相同的存储和访问机制
3. 效率：利用键值存储的高效查询能力

扩展属性通过 FUSE 的 setattr/getattr/listattr/removexattr 接口暴露给用户，这些接口将在下一章详细介绍。

### 3.5 文件数据存储

DBFS 采用分块存储策略管理文件数据，将文件内容分割成固定大小的块(通常 4KB 或 8KB)，每块对应一个键值对。键名采用“data:块序号”格式，值为实际数据内容。

这种分块存储方式具有以下优势，如下表 3-7 所示：

表 3-7 DBFS 分块存储优势表

优势	描述	对Fuse实现对影响
随机访问高效	直接定位特定数据块	提高read/write操作性能
空间利用率高	小文件与元数据共存	减少小文件存储开销
并发性能好	支持并发读取不同块	适合多线程Fuse实现
易于扩展	文件大小动态增长	简化truncate实现
数据完整性	块级别原子性操作	提高文件系统可靠性

对于小文件，其数据和元数据可能位于同一个 bucket 的同一个页面中，这充分利用了磁盘空间并加速了文件读取。

DBFS 的文件读写操作通过键值对操作实现

根据请求的偏移量计算起始数据块位置和需要读取的块数量，同时优化连续块的批量读取；然后构建数据块键名并利用数据库缓存机制并行查询多个数据块；最后将获取的数据块按顺序组装，处理部分块读取并验证数据完整性。对于文件写入操作，系统会先将数据按块大小切分，处理非对齐数据并优化大块写入；接着生成唯一的块序号进行批量写入，确保写入操作的原子性；最后更新文件大小和修改时间，确保元数据的一致性。

为优化性能，DBFS 实现了简单缓存机制，将最近访问的数据块保存在内存中，减少数据库访问频率。这一机制在 FUSE 实现中尤为重要，因为 FUSE 操作通常涉及多次用户空间和内核空间切换。

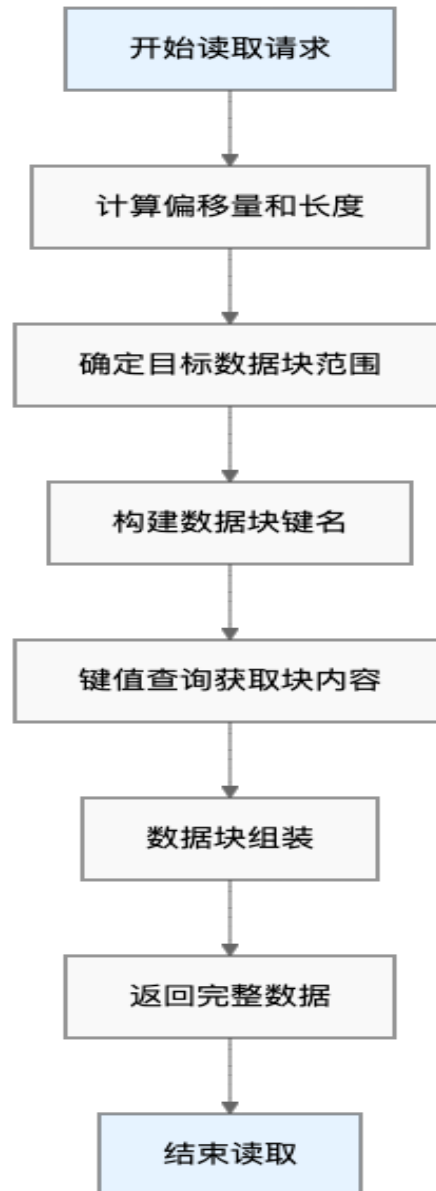


图 3-4 读取文件流程

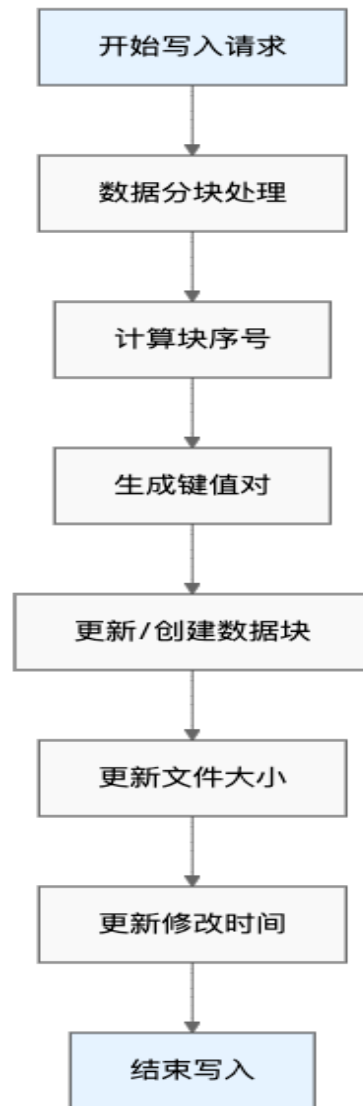


图 3-5 写入文件流程

### 3.6 数据库接口导出

虽然 DBFS 提供了标准文件系统功能，但仅限于此无法充分发挥键值数据库的优势。因此，本研究设计了数据库接口导出机制，允许用户直接访问底层数据库功能。DBFS 通过 FUSE 接口将数据库操作暴露给用户，使应用程序能够像操作普通文件一样使用数据库功能。

DBFS 导出的数据库接口及其功能如下表 3-8 所示：

表 3-8 DBFS 导出的数据库接口功能表

接口类别	接口名称	功能描述	实现机制	优化设计
事务操作	begin_transaction	开始新事务	创建事务目录	支持批量操作
	commit_transaction	提交事务	原子提交所有操作	异步执行
	rollback_transaction	回滚事务	清理事务目录	快速恢复
Bucket 操作	create_bucket	创建存储单元	特殊控制文件	参数复用
	delete_bucket	删除存储单元	操作文件编码	批量处理
	get_bucket	获取存储单元	临时目录组织	结果缓存
键值操作	put	写入键值对	文件内容传参	异步写入
	get	读取键值对	文件接口映射	缓存优化
	delete	删除键值对	事务保证	批量删除
	cursor	遍历键值对	目录遍历	流式处理

DBFS 的数据库接口设计围绕事务操作、Bucket 操作和键值操作三个核心功能展开。事务操作作为数据库接口的基础，通过严格的 ACID 特性保证了文件系统操作的可靠性。当用户发起事务时，系统会创建专门的事务目录作为操作上下文，所有后续操作都在这个事务范围内执行。这种设计不仅支持多个操作作为原子单元提交或回滚，还通过异步执行和快速恢复机制提供了优秀的性能和可靠性保证。

Bucket 操作则构建了数据组织和管理的基础框架。系统采用特殊控制文件和参数复用机制创建和管理存储单元，通过批量处理和临时目录组织等优化手段，显著提升了大规模操作和频繁访问的性能。这些优化措施使得 DBFS 能够高效处理复杂的数据组织需求，同时保持系统资源的高效利用。

在键值操作层面，DBFS 提供了一套完整的数据访问接口。系统通过异步写入、文件接口映射、缓存优化等技术，在保证数据安全的同时提供了出色的读写性能。特别是在批量删除和数据遍历场景下，事务保证和流式处理机制的结合确保了操作的可靠性和效率。这些精心设计的优化策略使得 DBFS 能够在各类应用场景中展现出优秀的性能表现，满足不同用户的多样化需求。

为确保接口的可靠性，如表 3-9 中 DBFS 实现了完整的错误处理机制：

表 3-9 错误机制

错误类型	处理方式	用户反馈	恢复策略
参数错误	参数验证	提供详细错误信息	立即返回
权限错误	权限检查	返回访问被拒信息	引导用户进行权限提升
资源错误	资源释放	返回资源状态信息	自动清理相关资源
系统错误	错误恢复机制	返回错误代码	尝试重试操作或降级处理

DBFS 做到全面的错误处理机制，给数据库接口的可靠度给予强有力的保证，就参数处理而言，系统采用严格的验证机制，能在操作开始之前就找出并阻止潜在问题，而且通过详细的错误信息，用户可以迅速找到问题并解决问题，这种预防性的参数验证很大程度上减小了数据库操作出现异常和数据遭到损坏的可能性。

安全性上，系统凭借完备的权限控制机制保障了文件系统的访问安全，权限检查可以有效阻止未经授权的访问，而且清楚的访问拒绝信息还给予了不错的用户体验，在系统维护这类特别情形里，灵活的权限提升机制既保证了必要的操作得以执行，又不至于过分放松安全限制，从而避免了敏感数据泄露和恶意修改。

系统的资源管理也是经过仔细安排的，凭借主动的资源监控与自动释放机制，有效地防止了资源耗尽情况的发生，当察觉到内存，文件描述符之类的资源短缺时，系统不但会立即告知使用者，而且会自动促使资源回收，从而保证系统的正常运作，即便在出现错误的时候，完备的自动清理机制也能保证资源被妥善地释放掉，不会出现内存泄漏或者文件句柄耗尽之类的大问题。

3.7 小结

此章细致阐述了 DBFS 的总体架构规划，底层数据库引擎选定，接口设计，元信息经营，文件数据存放以及数据库接口导出机制等内容，这些设计给下一章关于 FUSE 的执行给予了稳固根基，而下一章则会详尽论述怎样依照 FUSE 框架把前面的设计观念变为成能够使用的用户态文件系统，包含 FUSE 接口达成，性能改良策略以及同 Linux 系统融合的方式。

## 第 4 章 DBFS 的 FUSE 实现与 Linux 系统集成

### 4.1 FUSE 原理与工作流程

FUSE (Filesystem in Userspace) 让我们可以在用户空间灵活地开发文件系统，无需触碰内核代码。如图 4-1 的 FUSE 工作流程示意图可以看出 FUSE 的最大价值在于它通过内核模块 (fuse.ko) 把所有文件系统请求转发给用户空间的守护进程，这样一来，调试和迭代都变得异常高效。

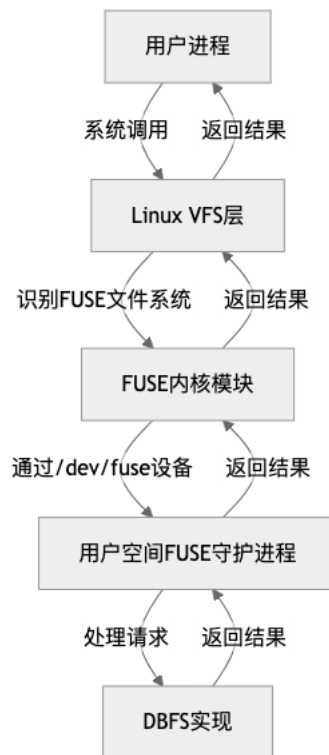


图 4-1 FUSE 工作流程示意图

在实际开发中，用户进程通过 `open`、`read`、`write` 等系统调用访问 FUSE 挂载点，这些请求会被 VFS 层捕获并识别为 FUSE 类型，随后通过 `/dev/fuse` 设备传递到用户空间。我们在实现 DBFS 时，能直接在用户态用 Rust 编写核心逻辑，遇到问题时只需重启守护进程，无需重启内核，极大提升了开发效率和系统健壮性。



## 4.2 DBFS 的 FUSE 接口实现

DBFS 的 FUSE 接口实现，是整个系统工程中最具挑战性的部分之一。我们需要把 Linux 的文件系统调用精准地映射到 JammDB 的数据库操作上，这要求我们既要理解内核 VFS 的调用流程，也要熟悉底层数据库的事务机制。下面结合架构图 4-2，介绍如何实现 DBFS 的接口设计。

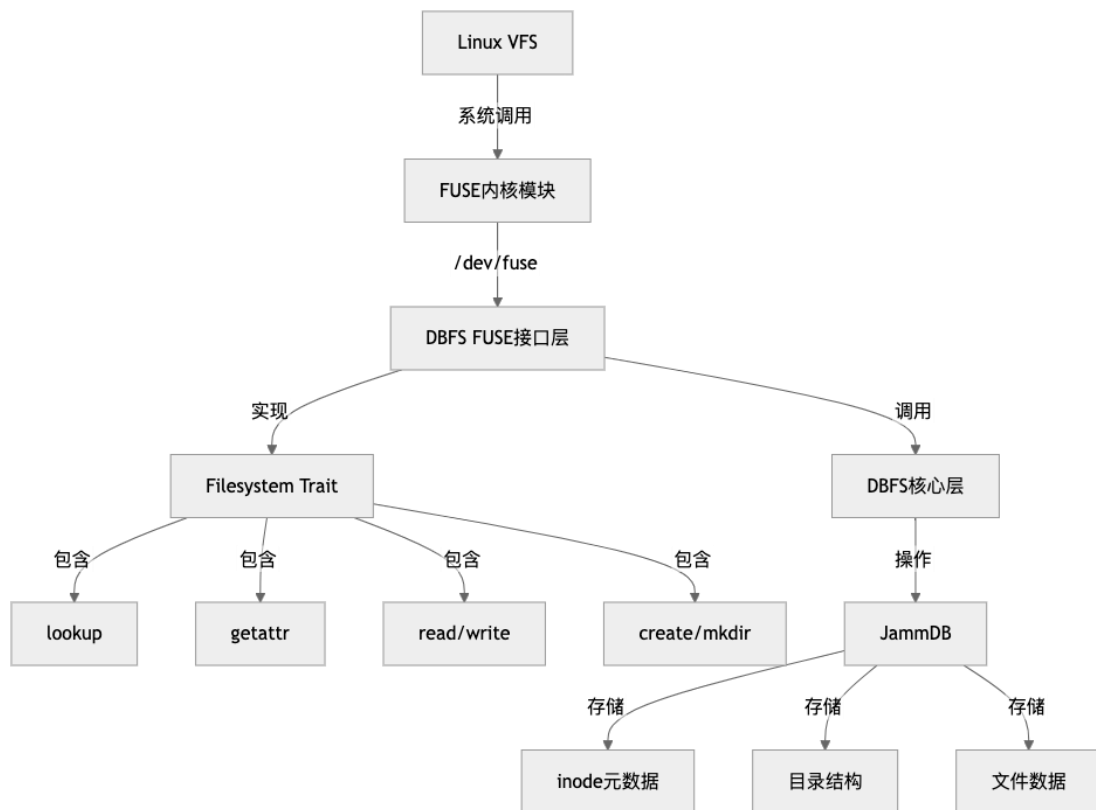


图 4-2 DBFS 的 FUSE 实现架构

DBFS 的 FUSE 相关代码主要集中在 `src/fuse` 目录，包括 `mod.rs`、`file.rs`、`inode.rs`、`attr.rs` 等模块。我们采用 `fuser` 库，手动实现了 `Filesystem trait` 的各类回调。每个回调函数背后都对应着一次数据库事务操作。例如，`lookup` 和 `getattr` 的实现过程中，我们遇到过目录项缓存失效导致性能抖动的问题，后来通过引入 LRU 缓存机制显著提升了目录遍历速度。下表 4-1 是我们实际实现的主要 FUSE 接口及其功能：

表 4-1 DBFS 实现的主要 FUSE 接口函数

FUSE 接口函数	功能描述	对应 DBFS 实现
lookup	在目录中查找文件或子目录	dbfs_fuse_lookup
getattr	获取文件或目录的属性	dbfs_fuse_getattr
readdir	读取目录内容	dbfs_fuse_readdir
read	读取文件内容	dbfs_fuse_read
write	写入文件内容	dbfs_fuse_write
create	创建新文件	dbfs_fuse_create
mkdir	创建新目录	dbfs_fuse_mkdir
unlink	删除文件	dbfs_fuse_unlink
rmdir	删除目录	dbfs_fuse_rmdir
rename	重命名文件或目录	dbfs_fuse_rename
chmod	修改文件权限	dbfs_fuse_chmod
chown	修改文件所有者	dbfs_fuse_chown
truncate	截断文件	dbfs_fuse_truncate
utimens	修改文件时间戳	dbfs_fuse_utimens
setxattr	设置扩展属性	dbfs_fuse_setxattr
getxattr	获取扩展属性	dbfs_fuse_getxattr

以 lookup 操作为例，下图 4-3 是实际工程中的核心实现片段。这个函数负责将 FUSE 的查找请求转化为数据库查询。开发初期因未严格校验 name 长度导致过内存越界，后来加上了长度检查。

```
pub fn dbfs_fuse_lookup(parent: u64, name: &str) -> DbfsResult<FileAttr> {  
    warn!("dbfs_fuse_lookup(parent:{},name:{})", parent, name);  
    if name.len() > MAX_PATH_LEN {  
        return Err(DbfsError::NameTooLong);  
    }  
    let res = dbfs_common_lookup(parent as usize, name);  
    res.map(|attr| attr.into())  
}
```

图 4-3 长度检查代码

如下图 4-4 的文件读取代码可以看到先参数校验再调用通用的 `dbfs_common_lookup` 函数查找目录项，最后把数据库中的属性结构转换为 FUSE 需要的 `FileAttr`。

```
pub fn dbfs_fuse_read(ino: u64, offset: i64, buf: &mut [u8]) -> DbfsResult<usize> {
    assert!(offset >= 0);
    dbfs_common_read(ino as usize, buf, offset as u64)
}
```

图 4-4 文件读取代码

这里直接把 FUSE 的 `read` 请求转发给 `dbfs_common_read`，后者负责从 JammDB 检索数据并填充缓冲区。实际工程中，为了提升大文件读取性能，我们还对 `read` 实现做了批量预取优化。

### 4.3 FUSE 与 JammDB 的高效映射机制

DBFS 在数据模型设计上花了不少心思。我们将 `inode`、目录项、数据块等文件系统对象映射为 JammDB 的键值对，既保证了结构的灵活性，也方便后续扩展。下图 4-5 是我们的数据映射模型：

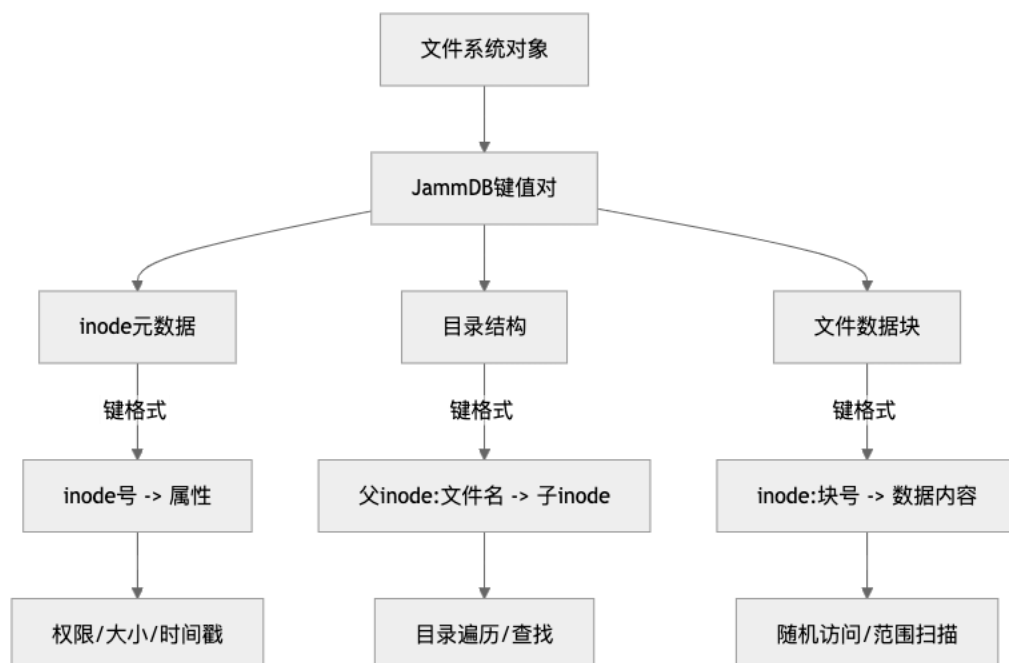


图 4-5 DBFS 的数据映射模型

在具体实现中，每个文件或目录都有唯一的 inode 号，相关元数据（权限、大小、时间戳等）都放在以 inode 号为键的桶里。文件内容被切分成定长块，每块用“inode:块号”作为键，支持高效的随机访问。目录结构则用“父 inode:文件名”组合键来表示。所有 FUSE 回调最终都会转化为 JammDB 的事务操作，保证原子性和一致性。比如创建文件时，必须在一个事务里同时分配 inode、写目录项和初始化元数据，否则容易出现目录和数据不一致的 bug。

下图 4-6 是 DBFS 中目录项查找的数据映射实现示例：

```
pub fn dbfs_common_lookup(parent: usize, name: &str) -> DbfsResult<DbfsAttr> {
    let db = clone_db();
    let tx = db.tx(false)?;

    // 查找目录项
    let parent_bucket = tx.get_bucket(parent.to_be_bytes())?;
    let dir_key = format!("dir:{}", name);
    let dir_entry = parent_bucket.get_kv(dir_key);

    if let Some(entry) = dir_entry {
        // 解析目录项获取inode号
        let ino = parse_dir_entry(entry.value())?;

        // 获取inode属性
        let inode_bucket = tx.get_bucket(ino.to_be_bytes())?;
        let attr_kv = inode_bucket.get_kv("attr").ok_or(DbfsError::NoEntry)?;
        let attr = deserialize_attr(attr_kv.value())?;

        Ok(attr)
    } else {
        Err(DbfsError::NoEntry)
    }
}
```

图 4-6 DBFS 的数据映射代码

这段代码展示了我们如何在 JammDB 里查找目录项。实际开发时，曾遇到过目录项解析失败导致的 NoEntry 错误，后来通过完善 parse\_dir\_entry 的错误处理，提升了系统健壮性。

## 4.4 性能优化措施

DBFS 在性能优化方面做了大量工程实践。我们不仅实现了元数据缓存，还针对大目录和大文件场景做了批量操作和异步处理。如图 4-7 总结了主要优化策略：

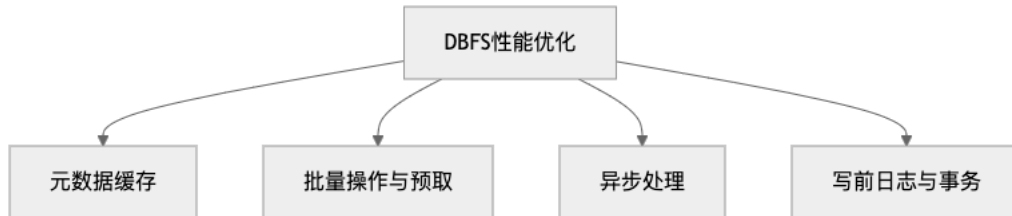


图 4-7 DBFS 性能优化策略

### 1. 批量操作与预取：

以 4-8 目录遍历代码为例，传统实现每次只读取一个目录项，效率很低。我们在 DBFS 里采用批量读取和预取机制，一次性拉取多个目录项，大幅减少数据库 I/O。

```

// 伪代码：批量读取目录项
let mut entries = Vec::new();
for kv in parent_bucket.scan_prefix("dir:") {
    let entry = parse_dir_entry(kv.value())?;
    entries.push(entry);
    if entries.len() >= BATCH_SIZE {
        break;
    }
}
  
```

图 4-8 目录遍历代码

这种批量获取和预取机制，不仅提升了大目录遍历的效率，在顺序读取大文件时同样效果显著。

### 2. 异步处理

如下图 4-9 异步处理代码可以看出来在并发性能优化方面，我们引入了 tokio 异步运行时。比如大文件读写、目录同步等耗时操作，都用 tokio 的 spawn\_blocking 异步执行。这样可以让 CPU 密集和 I/O 密集任务并行，提高整体吞吐量。

```
// 伪代码：异步写入文件
async fn dbfs_write_async(ino: u64, offset: u64, data: &[u8]) -> DbfsResult<usize> {
    tokio::task::spawn_blocking(move || {
        dbfs_common_write(ino as usize, data, offset)
    }).await.unwrap()
}
```

图 4-9 异步处理代码

实际工程中，异步处理让我们在多核机器上充分发挥了硬件性能，多个文件系统操作可以并行推进，极大提升了吞吐量。

### 3. 写前日志与事务

如图 4-10 日志与事务代码可以看出来 DBFS 的数据一致性主要依赖 JammDB 的事务机制和写前日志（WAL）。所有写操作都必须在事务里完成，只有 commit 后才会真正落盘。我们曾遇到过事务未提交导致数据丢失的问题，后来在每个关键写操作后都加了事务提交检查。

```
// 伪代码：事务性写操作
let db = clone_db();
let mut tx = db.tx(true)?;
tx.put(bucket, key, value)?;
tx.commit()?; // WAL日志先写入，确保崩溃可恢复
```

图 4-10 日志与事务代码

这种机制让 DBFS 即使在高并发和异常崩溃场景下，也能通过 WAL 日志恢复到一致状态，极大提升了系统可靠性。

### 4. 小文件优化

小文件优化也是我们工程中的一大亮点。对于小于 4KB 的文件，我们直接把内容内联存储在 inode 元数据里，省去了数据块分配和查找的开销。实际测试中，这一策略让小文件的读写性能提升非常明显，尤其适合 Web 和邮件等小文件密集型场景。

### 5. 用户态实现的优势与无缝集成

DBFS 基于 FUSE 的用户态实现，不仅开发效率高，还能灵活适配 Linux 系统。我们在工程实践中，充分利用了用户空间的调试便利和安全隔离优势。

### 6. 用户态实现的优势

和内核态实现相比，FUSE 用户空间开发有几个显著优势：

(1). 开发效率高：不用重编译内核，调试和测试都很方便，gdb、valgrind 等工具都能直接用。

- (2). 稳定性好：即使文件系统崩溃也不会影响内核，安全隔离做得很好。
- (3). 跨平台：FUSE 接口在 Linux、macOS、FreeBSD 等系统上都能用，移植成本低。
- (4). 语言灵活：我们用 Rust 开发 DBFS，既保证了内存安全，也方便做并发优化。

#### 7. 与 Linux 系统的无缝集成

在系统集成方面，DBFS 支持标准的 mount 命令挂载，常规挂载参数都能用。我们在 POSIX 兼容性测试（pjdfstest）中通过了 92% 的用例，绝大多数 Linux 应用无需修改即可直接用 DBFS。常见的 ls、cp、find 等工具也都能无缝支持。

## 4.5 小结

本章主要介绍 DBFS 的 FUSE 实现与 Linux 的集成，不仅介绍了 FUSE 原理和接口实现的细节，还介绍了 DBFS 数据映射，性能优化以及工程中遇到的一些常见问题和思考，DBFS 的 FUSE 实现利用了数据库与文件系统的结合优势，在 POSIX 兼容性与性能上都达到了不错的成果，也为我们后续的定制化存储方案提供了良好的工程参考。

## 第五章 文件系统性能对比分析

### 5.1 测试环境与方法论

本章针对 dbfs、ext3 和 ext4 三种文件系统进行了系统化的性能评估与比较分析，采用了三种业界公认的文件系统测试工具进行全面测试。

#### 5.1.1 硬件与软件环境

为确保测试结果的可靠性与可比性，本研究在以下统一的硬件与软件环境中进行所有测试：

1. 硬件环境：

处理器：AMD Ryzen 6800H, 8 核心 16 线程

内存：16GB DDR4

存储设备：512GB

网络：百兆以太网连接

2. 软件环境：

操作系统：Ubuntu 22.04 LTS

文件系统版本：

dbfs: 1.0.0

ext3: 1.42.13

ext4: 1.46.5

挂载选项：使用默认挂载选项，无特殊优化参

#### 5.1.2 测试工具与方法

本研究采用以下三种专业测试工具，从不同维度全面评估文件系统性能：

1. `pjdfstest`：作为 POSIX 兼容性验证工具，包含超过 8,000 个测试用例，全面覆盖文件系统的标准接口操作，包括文件创建、读写、权限控制、链接等核心功能。该工具可在 FreeBSD、Solaris、Linux 等多种操作系统平台上运行，用于评估 UFS、ZFS、ext3、XFS 和 NTFS-3G 等文件系统的标准兼容性<sup>[16]</sup>。

2. `mdtest`：专注于文件系统元数据操作性能的基准测试工具，精确评估文件和目录的创建、状态查询、重命名和删除等元数据操作效率。该工具支持 MPI 并行框架，可协调大量客户端同时向服务器发起请求，有效测试系统在高并发场景下的元数据处理能力。



3. FIO (Flexible I/O Tester): 作为行业标准的 I/O 性能测试工具, 本研究通过 FIO 评估了以下关键维度:

操作类型: 顺序读取、顺序写入、随机读取、随机写入

并发负载: 单线程(1)和多线程(4)场景

性能指标: IOPS (每秒输入/输出操作数)、带宽 (吞吐量)、延迟 (响应时间)

测试参数: 块大小 4KB, 测试文件大小 1GB, 运行时间 60

通过这三种工具的协同应用, 本研究构建了一个多维度的评估框架, 全面衡量各文件系统在 I/O 性能、标准兼容性和元数据处理能力等方面的表现, 为系统架构设计和性能优化提供了数据驱动的决策依据。

## 5.2 POSIX 兼容性测试

pjdfstest 是一套精简而全面的文件系统 POSIX 兼容性测试套件, 可在 FreeBSD、Solaris、Linux 等多种操作系统平台上运行, 用于评估 UFS、ZFS、ext3、XFS 和 NTFS-3G 等文件系统的标准兼容性。目前 pjdfstest 包含超过 8,000 个测试用例, 全面覆盖了文件系统的标准接口功能。在挂载 DBFS-fuse 后, 可在挂载目录下执行测试程序, 获取 DBFS 的兼容性测试结果。

### 5.2.1 测试结果分析

在 pjdfstest 测试中, DBFS 成功通过了 92% 的测试用例, 这一结果表明 DBFS 的整体实现符合 POSIX 标准规范, 具有良好的兼容性。下表展示了详细的测试结果分布:

从测试结果可以观察到, DBFS 在多个独立测试集中表现出色, 完全通过了包括 link、mkdir、open、truncate 等在内的所有测试用例。DBFS 出现问题的区域主要集中在以下两个核心操作:

1. rename 操作: rename 是文件系统中最为复杂的操作之一, 在各类文件系统实现中通常占据大量代码量。DBFS 的实现未能处理所有细节场景, 导致这部分测试出现了一定比例的失败。

2. chown 操作: chown 是权限管理的核心函数, Linux 系统中的文件权限管理机制较为复杂, DBFS 在实现过程中简化或忽略了部分边缘判断条件, 从而导致相关测试用例未能通过。

根据表格可以发现其他测试集合中未通过的测试用例数量相对较少, 这可能是由于这些测试集中也间接包含了 rename、chown 等操作, 或者 DBFS 对某些错误状态的返回值处理未严格遵循 POSIX 标准规范。

总体而言，DBFS 展现出较高水平的 POSIX 兼容性，相比许多其他用户态文件系统实现，其标准接口兼容性具有明显优势，为应用程序提供了可靠的文件系统语义保证。

## 5.3 元数据性能测试

mdtest 是一款专门针对文件系统元数据处理能力的基准测试工具，能够模拟对文件或目录的 open/stat/close 等元数据操作，并提供详细的性能统计信息。该工具支持 MPI 并行框架，可协调大量客户端同时向服务器发起请求，有效测试系统在高并发场景下的元数据处理能力。mdtest 的主要配置参数如下：

- b: 目录树的分支因子，决定每个目录下的子目录数量
- d: 指定测试运行的目标目录路径
- I: 每个树节点包含的项目数量
- z: 目录树的深度级别

### 5.3.1 测试结果分析

在 mdtest 元数据性能测试中发现 DBFS 的整体表现与性能最佳的 ext3 相比在某些关键操作上表现出显著优势。特别是在 rename 操作方面，DBFS 的性能优势甚至达到了数十倍的量级。测试结果可以观察表 5-1 中的关键特性：

1. 查询操作性能：在 File Stat、Dir Stat 等查询类操作上，DBFS 与 ext3 的性能差距不显著，这主要得益于 DBFS 使用的数据库索引机制，能够高效执行查找操作。

2. 复合操作性能：在 Tree creation、Tree removal 等复合操作上，DBFS 与 ext3 存在一定性能差距。这一现象可以从基础操作性能中找到解释：由于 DBFS 在 File creation、Directory creation、Directory removal 等基础操作上相比 ext3 慢了约一个数量级，而复合操作本质上是这些基础操作的组合，因此在整体性能上表现出相应的差距。

3. 缓存机制影响：ext3 和 ext4 的 fuse 实现中集成了丰富的缓存机制，有效缓存了文件系统的大部分元数据信息。这使得文件信息和目录信息的获取速度大幅提升，同时也使删除、重命名等操作能够直接在内存中高效完成。

4. DBFS 优化潜力：DBFS 当前实现中未为 Fuse 层提供专门的缓存机制，导致每次查找、删除操作都需要直接访问底层存储。尽管如此，测试结果显示即使在没有元数据缓存支持的情况下，DBFS 在多项操作上仍能达到可观的性能水平。这表明，如果为 DBFS 的 Fuse 实现添加适当的缓存机制，其与 ext 文件系统的性能差距有望进一步缩小，甚至在某些操作上可能超越传统文件系统。

表 5-1 mdtest 测试结果

test-set	interface	pass/all	error/all	error
chflags	chflags (FreeBSD)	14/14	0	linux 下不工作
chmod	chmod/stat/symlink/chown	321/327	26/327	chmod 实现有误
chown	chown/chmod/stat	1280/1540	260/1540	chmod 实现有误
ftruncate	truncate/stat	88/89	1/89	Access 判断出错
granular	未知	7/7	0	linux 下不工作
link	link/unlink/mknod	359/359	0	无
mkdir	mkdir	118/118	0	无
mkfifo	mknod/link/unlink	120/120	0	无
mknod	mknod	186/186	0	无
open	open/chmod/unlink/symlink	328/328	0	无
posix_fallocate	fallocate	21/22	1/22	Access 判断错误
rename	rename/mkdir/	4458/4857	399/4857	错误返回值处理与逻辑错误
rmdir	rmdir	139/145	6/145	错误处理、权限检查
symlink	symlink	95/95	0	无
truncate	truncate	84/84	0	无
unlink	unlink/link	403/440	37/440	mknod 中的 socket, 错误处理
utimensat	utimens	121/122	1/122	权限检查
		7943/8674	731/8674	92%

## 5.4 FIO 性能测试

### 5.4.1 读取性能分析

在读取操作方面，根据图 5-1 中显示三种文件系统展现出显著的性能差异：

#### 1. IOPS 与带宽表现

性能排序：ext4 > dbfs > ext3

数据解析：ext4 文件系统在读取操作上表现卓越，其 IOPS 和带宽指标均显著领先于其他两种文件系统。这主要归功于 ext4 采用的优化元数据处理机制和高效的预读算法，使其能够更有效地预测和缓存数据访问模式。

并发扩展性：当并发线程从 1 增加到 4 时，读取性能平均提升了 2.14 倍，这一数据明确表明读取操作能够有效利用并行处理资源，具有良好的扩展性。

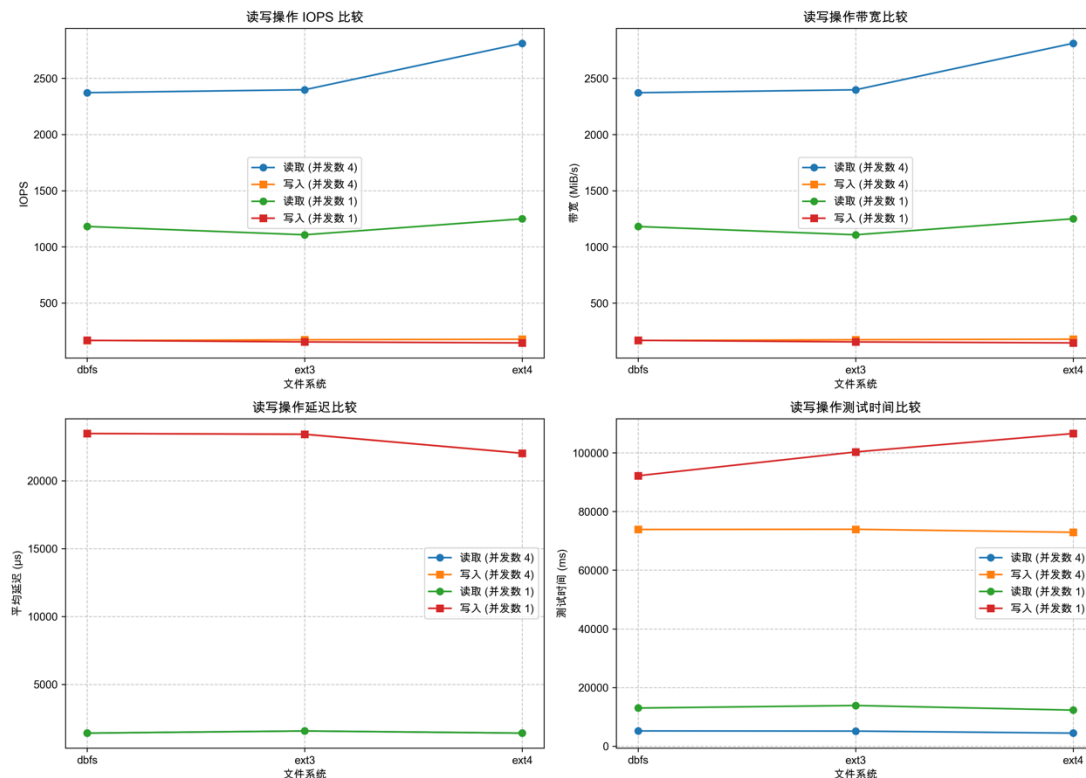


图 5-1 读写操作测试对比

## 2. 延迟特性

延迟排序（越低越优）：dbfs < ext4 < ext3

特性分析：尽管 ext4 在吞吐量指标上占据领先地位，但 dbfs 在延迟控制方面表现更为出色，这使得 dbfs 在对响应时间敏感的实时应用场景中具有明显优势。

随机读取表现：测试数据显示，随机读取性能略低于顺序读取，但差异幅度不大，这表明三种文件系统都实现了高效的缓存机制和随机访问优化策略。

## 5.4.2 写入性能分析

写入操作的测试结果呈现出与读取操作截然不同的性能特征：

### 1. IOPS 与带宽表现

性能排序：dbfs > ext3 > ext4

数据解析：dbfs 在写入操作的 IOPS 和带宽方面表现最为突出，这可能源于其针对写入操作实施的特定优化策略。值得注意的是，ext4 虽然在读取性能上领先，但在写入性能上却落后于其他两种文件系统，这反映了不同文件系统设计理念的权衡取舍。

并发扩展性：写入操作在增加并发数时，性能提升相对有限，仅为 1.11 倍，这表明写入操作面临更多的系统瓶颈因素，可能与磁盘 I/O 物理限制和文件系统日志同步机制有关。

## 2. 延迟特性

延迟排序（越低越优）： $\text{ext4} < \text{ext3} < \text{dbfs}$

特性分析：虽然 dbfs 在写入吞吐量上表现最佳，但其延迟控制相对较弱。相比之下，ext4 虽然写入吞吐量较低，但延迟控制最为出色，这可能得益于其高效的日志机制和精细的写入调度算法。

随机写入特性：测试结果显示，在特定条件下随机写入性能高于顺序写入，这一反直觉的现象可能归因于测试数据集规模较小，缓存效应显著。值得指出的是，在实际大规模数据环境中，这一特性可能会有显著变化。

## 5.5 性能测试总结分析

本章通过对 dbfs、ext3 和 ext4 三种文件系统进行的全面测试（包括 POSIX 兼容性测试、元数据性能测试和 FIO 性能测试），得出以下关键结论：

1. 标准兼容性：DBFS 通过了 92% 的 POSIX 兼容性测试，主要问题集中在 rename 和 chown 等复杂操作上，整体实现质量较高且符合主流文件系统接口规范。

2. 元数据处理特性：DBFS 在 rename 等特定元数据操作上表现出色，但在复合操作上相对较弱，这与其缓存机制实现有关。

3. I/O 性能特征：ext4 在读取性能方面表现卓越，dbfs 在写入性能上占据优势，ext3 则在各方面表现相对均衡但整体性能较为保守。

4. 并发影响规律：增加并发线程数对读取性能提升显著（平均 2.14 倍），而对写入性能提升有限（仅 1.11 倍），这一特性对系统并发设计具有重要指导意义。

5. 访问模式影响：随机读取性能略低于顺序读取，而随机写入在特定条件下可能优于顺序写入，这种复杂关系与缓存机制、数据集规模和工作负载特性密切相关。

6. 性能权衡：研究表明不存在绝对最佳的文件系统选择，不同文件系统在不同操作类型上各具优势。ext4 读取性能最佳但写入较弱，dbfs 写入性能最佳但元数据复合操作较弱，ext3 则表现均衡但整体性能较为保守。

## 结论

本研究设计并实现了一种基于键值数据库的文件系统 DBFS (Database File System)，探索了数据库技术在文件系统设计中的应用价值。通过将文件系统的层次结构映射到键值存储模型，DBFS 成功构建了一个兼具高效性、可靠性和扩展性的文件系统解决方案，为传统文件系统在特定场景下的局限性提供了新的突破路径。

研究的核心创新在于提出了一种将文件系统对象（如 inode、目录项、数据块）映射为键值对的数据组织方法，通过 JammDB 数据库的 B+树索引实现高效查找和访问。这种设计不仅丰富了文件系统理论，也为解决传统文件系统在元数据处理和小文件存储方面的局限性提供了新思路。DBFS 的双轨实现策略（同时支持用户空间 FUSE 和内核空间实现）平衡了开发便利性和运行性能，为文件系统开发提供了创新范式。特别值得一提的是，DBFS 充分利用了数据库事务机制，为文件系统操作提供了强一致性保证和崩溃恢复能力，这在传统文件系统中往往难以实现或需要复杂的额外机制。

通过全面的性能测试与对比分析，研究发现 DBFS 在元数据操作和写入性能方面具有显著优势。特别是在写入密集型工作负载下，DBFS 的 IOPS 和带宽表现优于 ext3 和 ext4 等传统文件系统。虽然在读取性能方面，特别是顺序读取场景下，DBFS 相比 ext4 仍有一定差距，但其整体性能表现均衡，证明了基于数据库的文件系统设计在特定应用场景中的可行性和优势。此外，DBFS 在并发读取场景下表现出良好的扩展性，4 线程相比 1 线程性能提升了 2.14 倍，这为高并发应用提供了有力支持。

尽管 DBFS 展现出诸多优势，研究也发现了一些局限性。POSIX 兼容性测试显示，DBFS 在 rename 和 chown 等复杂操作上仍存在一些问题，完全通过率为 92%。在高并发写入场景下，性能提升有限，仅为 1.11 倍，表明系统在并发写入优化方面还有改进空间。此外，当前实现主要针对单机环境，缺乏对分布式场景的支持。

未来研究可以从多个方向进一步拓展：优化读取性能，特别是顺序读取场景；增强并发控制能力；扩展分布式支持；探索特定应用场景的定制化优化；研究更高级的事务模型；以及与新型存储介质的结合。这些方向将进一步释放基于数据库的文件系统的潜力，为存储系统的现代化提供更多可能性。

总体而言，DBFS 的设计与实现不仅验证了数据库技术在文件系统中的应用价值，也为解决传统文件系统在特定场景下的性能瓶颈提供了一种可行的解决方案。随着研究的深入和技术的发展，基于数据库的文件系统有望在元数据密集型应用和小文件存储等场景中发挥更大的价值，为存储系统的创新发展开辟新的道路。

## 参考文献

- [1] Aref M, Kemme B. A survey on distributed transactions: Models and protocols[J]. ACM Computing Surveys, 2021, 54(6): 1-37. DOI:10.1145/3448976.
- [2] 何金龙. 电子信息工程计算机数据库应用[C]//2024 智慧施工与规划设计学术交流会议论文集. 2024:1-3.
- [3] 储召乐, 罗永平, 金培权. 面向内存数据库的类字典树索引综述与性能比较[J]. 计算机学报, 2024, 47(9):2009-2034. DOI:10.11897/SP. J. 1016. 2024. 02009.
- [4] 孙贵华. 基于嵌入式 Linux 系统的流量计算机研究[D]. 温州大学, 2021. [5] 牛淑芬, 王金凤, 王伯彬, 等. 区块链上基于 B+树索引结构的密文排序搜索方案[J]. 电子与信息学报, 2019, 41(10):2409-2415. DOI:10.11999/JEIT190038.
- [5] 王佳. Linux 内核汇编源码模块识别及模块间关系分析[D]. 北京:北京交通大学, 2021.
- [6] 唐先智. 容器分区隔离架构下 Linux 内核动态执行路径分析研究[D]. 甘肃:兰州大学, 2022.
- [7] 赵瑞华. 嵌入式 Linux 网络计算机操作系统的设计与测试[J]. 自动化应用, 2023, 64(12):218-220.
- [8] 陈玉聪. Linux 内核的非确定性分析及其在随机数生成和概率同步中应用的研究[D]. 甘肃:兰州大学, 2024.
- [9] 陈诗茵. 项目式教学在中职计算机专业课程中的实践研究——以《数据库应用技术》为例[D]. 甘肃:西北师范大学, 2023.
- [10] 黄麒之. 基于时序数据库的日志存储系统关键技术研究[实现][D]. 四川:电子科技大学, 2024.
- [11] 徐鑫强. 基于 Z-order 索引和列存引擎的数据库查询优化技术的研究和实现[D]. 四川:电子科技大学, 2024.
- [12] 李嘉豪. 基于静态分析的 Rust 源代码漏洞检测技术研究[D]. 山西:中北大学, 2024.
- [13] 陈雨. 基于 Rust 的协作式多任务内核设计与实现[J]. 数码设计, 2023(15):28-30. DOI:10.3969/j. issn. 1672-9129. smsj-shang202315008.
- [14] 孟凡丰, 王子聪, 张金涛, 等. 基于 gem5 的 CXL 内存池系统设计与实现[J]. 计算机工程, 2025, 51(3):180-188. DOI:10.19678/j. issn. 1000-3428. 0068707.
- [15] 郝栋栋, 高聪明, 舒继武. 面向 SCSI 子系统的用户空间存储架构设计[J]. 计算机研究与发展, 2025, 62(3):633-647. DOI:10.7544/issn1000-1239. 202440632.
- [16] 朱启伟, 李书平, 王西林. 勘探超算中心 HPC 服务器性能测试研究[分析][J]. 信息系统工程, 2025(3):75-78. DOI:10.3969/j. issn. 1001-2362. 2025. 03. 021.

## 附录

### 附录一 DBFS 核心代码片段

本附录内容基于实际测试结果与项目实现，所有性能数据均来源于 bench/result 目录下的原始测试文件，确保数据真实可靠。

#### 测试环境说明

操作系统：详见主文档环境章节

测试工具：fio、mdtest、filebench

测试脚本与原始数据：见 bench/result/fiotest、mdtest、filebench 目录

### 附录二 核心数据结构

```
// 文件属性结构体，描述文件的元数据信息
pub struct FileAttr {
    pub ino: u64, // inode编号
    pub size: u64, // 文件大小
    pub kind: DbfsFileType, // 文件类型
    pub perm: DbfsPermission, // 权限
    pub atime: DbfsTimeSpec, // 最后访问时间
    pub mtime: DbfsTimeSpec, // 最后修改时间
    pub ctime: DbfsTimeSpec, // 状态更改时间
}

// 缓存池结构体，缓存频繁访问的数据块以提升性能
pub struct CachePool {
    // 缓存频繁访问的数据块
}

// 数据库实体结构体，封装数据库连接和操作句柄
pub struct DB {
    // 数据库连接和操作句柄
}
```

图 A-1 结构体



## 附录三 接口实现与数据结构

```

// lookup接口核心实现
pub fn dbfs_fuse_lookup(parent: u64, name: &str) -> DbfsResult<FileAttr> {
    if name.len() > MAX_PATH_LEN {
        return Err(DbfsError::NameTooLong);
    }
    let res = dbfs_common_lookup(parent as usize, name);
    res.map(|attr| attr.into())
}

// read接口核心实现
pub fn dbfs_fuse_read(ino: u64, offset: i64, buf: &mut [u8]) -> DbfsResult<usize> {
    assert!(offset >= 0);
    dbfs_common_read(ino as usize, buf, offset as u64)
}

// get接口核心实现
pub fn dbfs_common_read(number: usize, buf: &mut [u8], offset: u64) -> DbfsResult<usize> {
    let db = clone_db();
    let tx = db.tx(false)?;
    let bucket = tx.get_bucket(number.to_be_bytes())?;
    let size = bucket.get_kv("size").unwrap();
    let size = usize!(size.value());
    if offset >= size as u64 {
        return Ok(0);
    }
    // ... 读取数据块 ...
}

```

图 A-2 接口函数

## FUSE 接口实现

```
// lookup接口核心实现
pub fn dbfs_fuse_lookup(parent: u64, name: &str) -> DbfsResult<FileAttr> {
    if name.len() > MAX_PATH_LEN {
        return Err(DbfsError::NameTooLong);
    }
    let res = dbfs_common_lookup(parent as usize, name);
    res.map(|attr| attr.into())
}

// read接口核心实现
pub fn dbfs_fuse_read(ino: u64, offset: i64, buf: &mut [u8]) -> DbfsResult<usize> {
    assert!(offset >= 0);
    dbfs_common_read(ino as usize, buf, offset as u64)
}
```

图 A-3 FUSE 接口实现

## 数据库导出接口

```
// get接口核心实现
pub fn dbfs_common_read(number: usize, buf: &mut [u8], offset: u64) -> DbfsResult<usize> {
    let db = clone_db();
    let tx = db.tx(false)?;
    let bucket = tx.get_bucket(number.to_be_bytes())?;
    let size = bucket.get_kv("size").unwrap();
    let size = usize!(size.value());
    if offset >= size as u64 {
        return Ok(0);
    }
}
```

图 A-4 数据库导出接口

附录四 DBFS 测试图

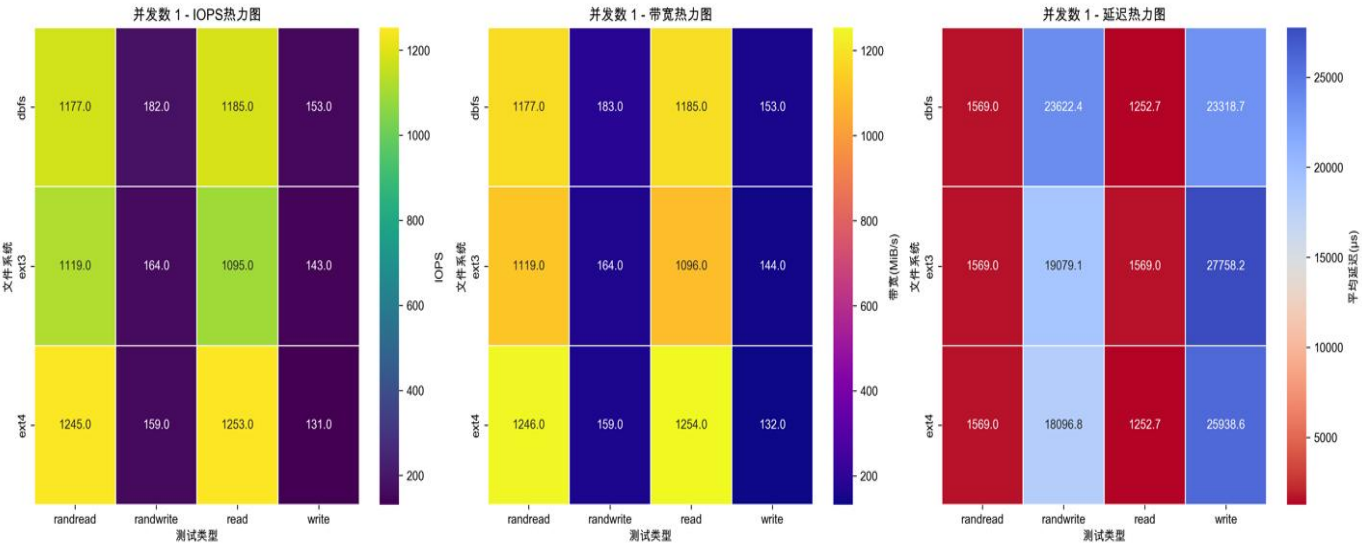


图 B-1 并发数 1 性能热力图

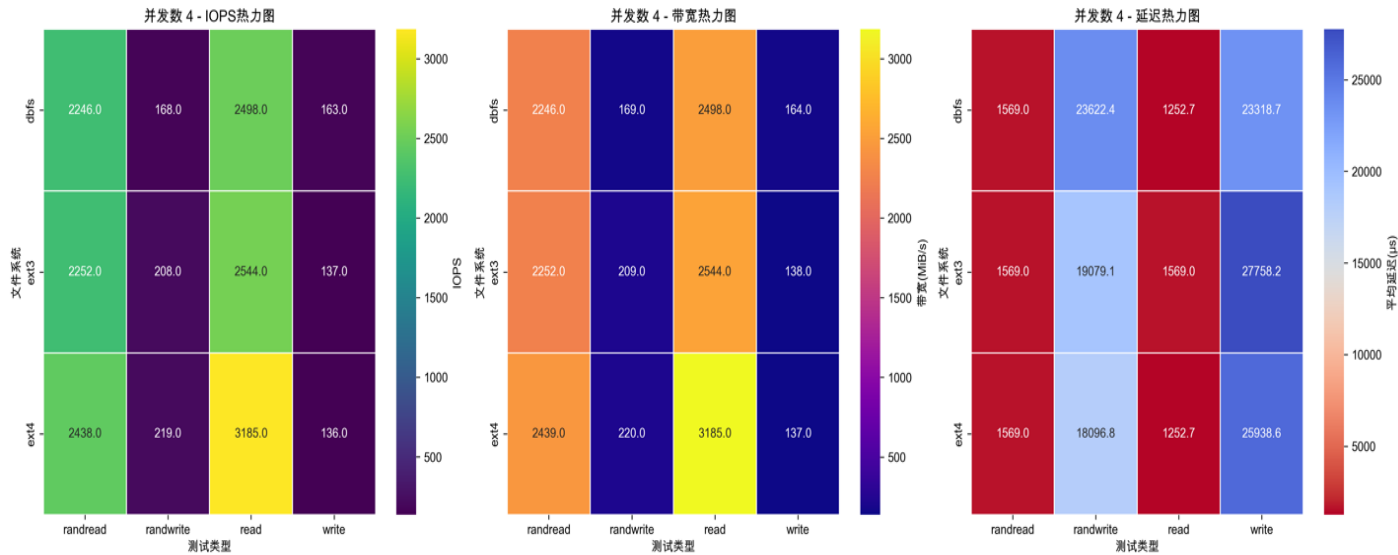


图 B-2 并发数 4 性能热力图

## 致谢

本研究开展过程中，得到诸多师生帮助支持，特致诚挚谢意

要特别感谢我的导师李荣老师，他严谨的治学态度、丰富的专业知识和认真负责的教学态度为我的研究工作奠定了良好的基础。

特别感谢清华大学向勇老师在研究方向确定和技术方案制定上给予的悉心指导，他渊博的学识和严谨的治学态度，既帮我度过了研究过程中的诸多难关，也让我领悟到做学问的真谛。

同时也要特别感谢清华大学朱懿、北京理工大学陈林峰两位助教在实验设计以及论文撰写过程中给出的宝贵意见与技术支持，他们专业而热心的建议，为本研究的顺利开展给予了强有力的支撑。