

跨语言用户态文件系统框架读写性能优化

顾荣^{1,2}, 罗义力², 仇伶玮², 王肇康^{1,3}, 戴海鹏^{1,2}, 黄宜华^{1,2}

(1. 计算软件新技术国家重点实验室(南京大学), 江苏南京 210023; 2. 南京大学计算机科学与技术系, 江苏南京 210023;
3. 南京航空航天大学计算机科学与技术学院, 江苏南京 211106)

摘要: 以深度学习为代表的数据分析应用越来越多依赖分布式文件系统存储管理大规模数据集。为了增强数据访问的兼容性, 现有分布式文件存储系统通常需提供标准 POSIX 接口, 以支持深度学习等应用的无缝对接。然而, 以内核模块形态开发提供 POSIX 接口的文件系统非常复杂耗时。近年来, 用户态文件系统(Filesystem in Userspace, FUSE)框架大幅简化了文件系统的开发工作, 已被 Alluxio 和 Ceph 等诸多知名分布式文件系统使用。目前常用的用户态 FUSE 库 libfuse 仅提供 C 语言编程接口, 但现有大数据分布式文件系统基本都是基于 Java 语言开发的(例如 HDFS 和 Alluxio 等), 为了使基于 Java 语言开发的分布式文件系统可以对接 C 语言开发的 FUSE 库, 需采用跨语言 FUSE 框架作为中介。跨语言 FUSE 框架利用跨编程语言的函数回调机制, 使操作系统 FUSE 库的 C 语言函数可以跨语言的调用分布式文件系统提供的 Java 语言编程接口, 从而为大数据分布式文件系统提供标准 POSIX 接口的访问能力。但在数据密集型应用中, 现有跨语言 FUSE 框架的执行效率低, 导致数据密集型作业(深度学习、大数据分析等)中数据 I/O 耗时占据了显著的性能开销, 成为新的潜在性能瓶颈。针对此问题, 本文首先评估分析了重要且广为使用的跨语言 FUSE 框架 JNR-FUSE 的性能, 发现并定位其在高并发和小文件场景下存在的性能瓶颈; 接着从多方面剖析性能瓶颈根因, 进而总结出高效跨语言 FUSE 框架的性能优化方向, 并面向 Java 语言设计实现了跨语言 FUSE 框架 JNI-FUSE。JNI-FUSE 利用延迟分离和元信息缓存等优化技术降低跨语言函数回调开销, 从而提升跨语言 FUSE 框架的性能。实验结果表明, 对比当前性能最好的 Java FUSE 框架 JNR-FUSE, 本文提出的 JNI-FUSE 带来了 1.15~6.04 倍的 FUSE 框架性能提升和 1.90~2.71 倍的文件系统端到端性能提升, 并为上层深度学习训练任务带来了 1.06~1.73 倍的训练加速。本文设计提出的 JNI-FUSE (Java Native Interface-Filesystem in User Space) 因性能优势, 已被知名开源分布式文件系统 Alluxio 官方接受集成。

关键词: POSIX; 用户态文件系统; 跨语言; 性能优化; Java 原生接口

基金项目: 国家自然科学基金面上项目(No.62072230, No.61872178); 江苏省科技厅重点项目(No.BE2021729); 软件新技术国家重点实验室开放课题(No.KFKT2021B33); 软件新技术与产业化协同创新中心项目

中图分类号: TP31 **文献标识码:** A **文章编号:** 0372-2112(2023)06-1590-17

电子学报 URL: <http://www.ejournal.org.cn>

DOI: 10.12263/DZXB.20220372

Reading and Writing Performance Optimization of Cross-Language FUSE Framework

GU Rong^{1,2}, LUO Yi-li², QIU Ling-wei², WANG Zhao-kang^{1,3}, DAI Hai-peng^{1,2}, HUANG Yi-hua^{1,2}

(1. State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing, Jiangsu 210023, China;

2. Department of Computer Science and Technology, Nanjing University, Nanjing, Jiangsu 210023, China;

3. College of Computer Science and Technology, Nanjing University of Aeronautics and Astronautics, Nanjing, Jiangsu 211106, China)

Abstract: Big data analytical applications like deep learning-based AI applications more and more rely on distributed file system to store and manage large scale data sets. File systems often need to provide standard POSIX interfaces to enhance their access compatibility with upper-layer applications. However, it is complicated to develop POSIX-compatible file systems in kernel space. In recent years, FUSE (File System in User Space) has been used by many well-known file systems, including Alluxio, Ceph, etc., because it significantly simplifies the file system development. The popular FUSE library libfuse is developed in the C language. However, the popular distributed file systems (like HDFS and Alluxio) for big data applications are developed with the Java language. To make the Java-based distributed file systems use the FUSE

mechanism, the cross-language FUSE frameworks are needed to bridge the gap, which becomes a potential performance bottleneck. The cross-language FUSE framework uses a cross-programming language function callback mechanism to enable the C functions of the FUSE library to call the programming interface provided by the distributed file system in Java. In this way, we can provide the access to the standard POSIX interface for the Java-based distributed file systems. However, the existing cross-language FUSE frameworks are inefficient in performance. It makes the data I/O in data-intensive applications (like deep learning and big data analysis) occupy noticeable proportion of their execution costs. To address the problem, we first systematically evaluate the performance of the widely-used cross-language FUSE framework, and find the bottlenecks of throughput performance in high concurrency and small file scenarios. We then analyze the bottlenecks of the cross-language FUSE framework from multiple perspectives, and propose several directions for optimizing the cross-language FUSE framework. According to the optimization directions, we design and implement JNI-FUSE (Java Native Interface-Filesystem in User Space), an efficient cross-language FUSE framework. In JNI-FUSE, we propose the defer detach and meta cache techniques to reduce execution costs of cross-programming language function callbacks. Experiment results show that JNI-FUSE improves the average framework performance from 1.15 times to 6.04 times compared to the cutting-edge cross-language FUSE framework JNR-FUSE. JNI-FUSE improves the end-to-end performance by 1.90 time to 2.71 times, and accelerates the deep learning training by 1.06 times to 1.73 times compared to JNR-FUSE. JNI-FUSE has been accepted and integrated by the well-known open-source distributed file system Alluxio due to its good performance.

Key words: portable operating system interface; user-space file system; cross-language; performance optimization; java native interface

Foundation Item(s): National Natural Science Foundation of China (No.62072230, No.61872178); Key Project of Jiangsu Provincial Science and Technology Department (No.BE2021729); Open Project of State Key Laboratory for Novel Software Technology (No.KFKT2021B33); The project of Collaborative Innovation Center of Novel Software Technology and Industrialization

1 引言

人工智能和大数据分析应用通常依赖大规模数据集训练和分析。基于磁盘的传统单机本地文件系统难以满足此类应用对于性能和可用性的要求。面向数据密集型应用而设计的分布式数据存储系统(例如 HDFS, Alluxio^[1], M-Cloud^[2])通过可扩展的系统架构,可以较好地解决此问题。但许多现有人工智能和数据分析应用(例如基于 PyTorch 的深度学习应用)只支持从提供 POSIX^[3]接口的传统单机本地文件系统读写数据,面临着难以直接读写分布式数据存储系统数据的问题。

为了支持现有人工智能和数据分析应用快速简单地访问数据,众多存储系统厂商引入 POSIX 对接层,将分布式文件系统模拟为提供 POSIX 接口的本地文件系统,以帮助现有人工智能和大数据应用在不修改现有文件 I/O 接口的情况下透明读取数据。因此,POSIX 对接层的数据访问效率直接影响了上层人工智能和数据分析应用的总体执行效率。

POSIX 对接层需要借助 Linux 操作系统提供的 POSIX 兼容文件系统框架,其可进一步由内核态或用户态两种方式之一实现。以内核态形式提供 POSIX 兼容的文件系统时,通常需要修改操作系统内核源码,开发和维护难度较高,而用户态文件系统(Filesystem in Userspace, FUSE)^[4]允许在不更改操作系统内核的情况

下对接自定义的文件系统,大幅降低了文件系统的设计开发门槛,因此成为很多文件系统 POSIX 对接层的实现方式。各种用户态文件系统框架^[5,6]被广泛使用在 Alluxio^[1], Ceph^[7] 和 Gluster^[8]等主流分布式文件系统中,它们皆借助 FUSE 提供 POSIX 访问接口。与内核态文件系统相比,FUSE 虽然简化了文件系统开发过程,但也牺牲了一定的数据访问性能。

目前常用的用户态 FUSE 库 libfuse 仅提供 C 语言编程接口,但分布式文件系统可能采用非 C 语言开发(例如 Java)。为了使基于其他语言开发的分布式文件系统可以对接基于 C 语言开发的 FUSE 库,需要额外引入跨语言 FUSE 框架,以实现文件系统编程接口与基于 C 语言的 FUSE 库之间的互相调用。鉴于现有大数据分布式文件系统基本都是基于 Java 语言开发(例如 HDFS, Alluxio 等),因此跨语言 FUSE 当前主要关注的问题就是面向 Java 的跨语言 FUSE 实现。

跨语言 FUSE 框架的引入进一步放大了用户态文件系统的性能问题。在深度学习训练、大数据分析等数据密集型作业中,数据 I/O 的速度会显著影响作业执行时间,进而影响用户的使用成本。然而,跨语言 FUSE 框架从系统内核层到语言实现层横跨多个系统层次,比单一语言 FUSE 框架实现更为复杂^[9],导致其性能瓶颈尚未得到充分的分析,现有跨语言 FUSE 框架存在严重的性能问题。

针对现有面向 Java 的跨语言 FUSE 框架存在的性

能问题,本文首先以实验实证的方式分析了影响跨语言 FUSE 框架性能的关键因素,从而总结了跨语言 FUSE 框架的性能优化方向.在优化方向的指导下,本文面向 Java 语言提出并实现了基于高效跨语言函数回调的跨语言 FUSE 框架 JNI-FUSE.具体地,本文贡献可以归纳为如下三个方面:

(1)通过设计实验分析,深度评估了目前广为使用的跨语言 FUSE 框架 JNR-FUSE 在数据密集型作业下的性能表现,发现其在高并发和小文件场景下存在明显的性能瓶颈.进一步地,本文从原理与实验两方面分析了性能瓶颈的产生原因,并提出高效跨语言 FUSE 框架的性能优化方向.

(2)根据上述总结的性能优化方向,本文面向 Java 提出并实现了高效跨语言 FUSE 框架 JNI-FUSE,该框架通过延迟分离和元信息缓存等优化技术,降低了跨语言 FUSE 过程中函数回调开销.本文优化方法可进一步应用于其他基于 JVM 编程语言开发的文件系统的跨语言 FUSE 场景.

(3)通过大量实验评估发现,对比当前性能最好的 Java FUSE 框架 JNR-FUSE,本文提出的 JNI-FUSE 读写性能提升 1.15~6.04 倍、文件系统端到端性能提升 1.90~2.71 倍,并可为深度学习训练任务带来 1.06~1.73 倍的训练加速.本文设计提出的 JNI-FUSE 因性能优势,已被知名开源分布式文件系统 Alluxio 官方接受集成,且作为默认配置使用.

本文的组织结构如下.第 2 节介绍 FUSE 框架和跨语言 FUSE 框架的背景知识及相关工作,第 3 节分析跨语言 FUSE 框架存在的性能瓶颈并总结高效跨语言 FUSE 框架的优化方向,第 4 节介绍 JNI-FUSE 框架的设计与实现,在第 5 节通过实验评估了 JNI-FUSE 的性能,第 6 节总结全文工作.

2 相关工作与背景知识

2.1 相关工作

跨语言 FUSE 框架相关工作可分为 FUSE 性能分析、FUSE 性能优化与跨语言 FUSE 框架设计等三类.在 FUSE 性能分析方面,用户态文件系统 FUSE 相比内核态文件系统提供了更高的安全性和可靠性^[10],但它也带来了一定的性能损耗^[9],其中最为突出的瓶颈是内核态和用户态的多次通信.Vangoor 等人^[9]以本地文件系统 Ext4 为基础,详尽分析了 FUSE 在不同负载下的性能损耗情况.实验结果表明 FUSE 最高可带来超过 83% 的性能损耗,尽管通过优化 FUSE 配置,大部分工作负载的性能损耗可以降低到 5% 以内,但仍存在部分负载(尤其是元数据操作频繁的工作负载)性能大幅弱于内核态文件系统.

在 FUSE 性能优化方面,FUSE 因为横跨内核态与用户态,其性能优化更加困难.Ishiguro 等人^[11]和 Narayan 等人^[12]曾尝试在特定情况下减少 FUSE 内核态和用户态的上下文切换.为进一步优化 FUSE 性能,Bijlani^[10]等提出允许 FUSE 文件系统针对某些请求自定义一些特殊优化的处理逻辑,从而在不改变原有复杂逻辑的基础上提升性能.以上 FUSE 框架优化工作均只面向文件系统与 FUSE 库采用相同语言实现的情况,重点关注减少内核态和用户态间的消息通信开销,但未考虑跨语言 FUSE 中跨语言函数调用带来的性能问题,本文在这方面将开展研究.

在跨语言 FUSE 框架设计方面,目前相关研究较少,主要有 JNR-FUSE^[13]和 FUSE-J^[14],两者都提供了 FUSE 编程接口的 Java 语言绑定.相比 FUSE-J,JNR-FUSE 的应用范围更广且性能更好,并仍处于活跃的维护状态.

2.2 用户态文件系统 FUSE

图 1 展示了用户态文件系统的设计架构.FUSE 通过 Linux 的虚拟文件系统(Virtual File System, VFS)机制提供了统一的 POSIX 访问接口.它包含一个内核模块(FUSE Driver)和一个用户态函数库 libfuse^[15],两个模块通过读写一个块设备文件/dev/fuse 交互.FUSE Driver 负责将应用从 VFS 传过来的文件系统请求转换为 FUSE 请求,并加入 FUSE 处理队列;libfuse 通过读取/dev/fuse,从 FUSE 队列拉取 FUSE 请求,并交给用户态的工作线程处理,工作线程会调用用户定义的回调函数处理文件系统的访问请求,最后将处理结果通过/dev/fuse 写回内核.

为了处理应用程序对用户态文件系统的并发访问请求,libfuse 会创建多个工作线程(Daemon 线程)并发处理 FUSE 请求.FUSE 的工作线程管理机制的优化原则是尽可能利用高并发度减少 FUSE 请求的总等待和处理时间.具体而言,文件系统通常会读写硬盘等较慢速的设备文件,涉及大量 I/O 请求,导致性能瓶颈通常在于 I/O.为了减少 I/O 等待带来的计算资源浪费、降低所有请求的总等待时间,FUSE 的工作线程管理机制通常会创建足够多的工作线程,以保证所有文件系统调用请求都尽可能被及时处理.FUSE 在初始化时仅会创建一个工作线程,当应用程序对文件系统发起并发访问时,应用程序发出的 FUSE 请求将会在 libfuse 库内排队等待工作线程处理.当 libfuse 库发现排队等待的 FUSE 请求数量过多时,表明当前工作线程数量不足,libfuse 库会创建额外的工作线程以加速处理 FUSE 请求.在典型配置下,如果存在超过 10 个工作线程处于空闲状态,FUSE 会回收空闲工作线程以释放资源.在高并发场景下,受跨语言 FUSE 高执行开销的影响,

FUSE 工作线程的实际处理速度低于 I/O 请求的到达速度,导致 I/O 请求堆积、大量工作线程被创建,这些工作线程在对应的 I/O 请求完成后又被销毁,从而带来大量线程管理开销。

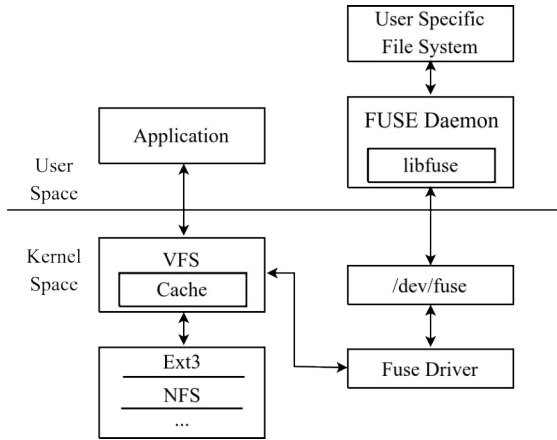


图1 用户态文件系统架构

2.3 跨语言 FUSE 框架 JNR-FUSE

目前大数据生态中的常用分布式文件系统(例如 HDFS, Alluxio 等)多以 Java 语言实现,而 FUSE 的用户态函数库 libfuse 采用 C 语言实现。这些文件系统为了能够基于 FUSE 机制提供 POSIX 接口,必须依赖跨语言 FUSE 框架。

跨语言 FUSE 框架的核心在于实现从 libfuse 库中的 C 语言函数,回调文件系统以 Java 语言提供的应用程序接口。C 语言函数通过 JVM 提供的 Java 原生接口 (Java Native Interface, JNI)^[16] 机制与 JVM 虚拟机交互,实现对象创建、方法调用等。但 JNI 使用方式较为烦琐,为了降低 JNI 编程难度, JNA (Java Native Access)^[17] 和 JNR (Java Native Runtime)^[18] 等技术基于 JNI 提供了更为简洁的跨语言函数调用方式。

Tselovalnikov 等人^[13]基于 JNR 技术提出了 Java 平台上最为常用的跨语言 FUSE 框架 JNR-FUSE。JNR-FUSE 基于 JNR 实现了 FUSE 接口的 Java 绑定,继承了 JNR 跨平台、易于实现的优点。图 2 是 JNR-FUSE 架构图。JNR 为了实现跨平台的通用性、减少用户的编程负担,引入了多层封装结构,这使得 JNR-FUSE 中跨语言函数调用链更长,性能瓶颈更难以分析。在第 3 节,本文将 JNR-FUSE 为例,分析 FUSE 在跨语言使用时存在的性能瓶颈。

3 跨语言 FUSE 框架 JNR-FUSE 性能分析

3.1 性能分析设置

跨语言 FUSE 过程的端到端总体性能一般由跨语言 FUSE 框架、FUSE 机制本身、文件系统与底层存储介

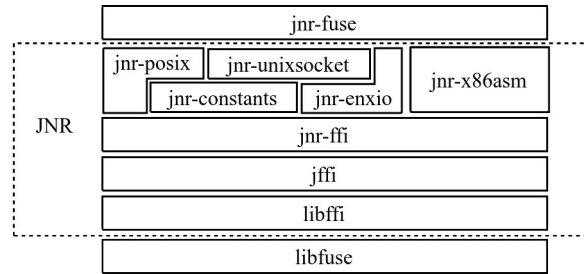


图2 JNR-FUSE 框架架构

质等因素共同决定。本文采用 Java 平台上最为成熟且广为使用的跨语言 FUSE 框架 JNR-FUSE^[13]作为评估对象,进行性能评估与现象分析。

为屏蔽文件系统与底层存储介质的干扰影响,本文设计并实现了具备低管理开销的伪内存文件系统 MemFS 作为 FUSE 的目标文件系统。为了降低 MemFS 本身的管理开销,MemFS 存储的所有文件系统元信息皆保留在内存中,并且不实际存储文件数据,而仅记录文件长度。为降低 MemFS 的内存空间开销,MemFS 的写入操作不会触发任何额外的磁盘或内存数据拷贝,MemFS 的读取操作会返回预先填充的虚拟数据。因此采用 MemFS 作为目标文件系统能更准确地反映 JNR-FUSE 框架自身的性能开销情况。

在深度学习训练、大数据分析等数据密集型作业中,数据 I/O 的速度会极大地影响作业的执行时间。高并发和小文件是数据密集型作业的两个重要特征。为了提高深度学习等模型训练速度,往往会采用多个线程进行高并发读写;当数据处理对象是语音片段、图片等,往往需读写由大量小文件组成的数据集^[19]。本文分别测试并收集了 JNR-FUSE 框架在高并发和小文件两种场景下访问 MemFS 的性能数据,实验中采用的服务器软硬件配置见第 5 节性能评估部分的表 3 (如无特殊说明,下同)。

高并发场景使用不同数量的线程并发读取 81 920 个文件,每个文件规模固定为 128 KB。JNR-FUSE 的吞吐量实验结果如图 3(a)所示。在较低线程数(不超过 8 线程)的情况下,吞吐量随着线程数量增多而增长。但当线程数量从 8 扩展到 16 时,JNR-FUSE 出现了严重的性能下降。16 线程的吞吐量只有 8 线程的 40%,大于 16 线程时,继续增加线程数量,吞吐量有一定上升,但始终达不到 8 线程时的吞吐量。

小文件场景使用 32 个线程并发读取不同规模的大量小文件,性能实验结果如图 3(b)所示。JNR-FUSE 的吞吐量随着文件规模的增大而增长。当单个文件大小为 1 MB 时,吞吐量增长为 1 598 MB/s;当单个文件大小减小到 4 KB 时,吞吐量降低为 19 MB/s,仅为 1 MB 情况的 1.2%。

以上实验结果表明,JNR-FUSE 在高并发和小文件

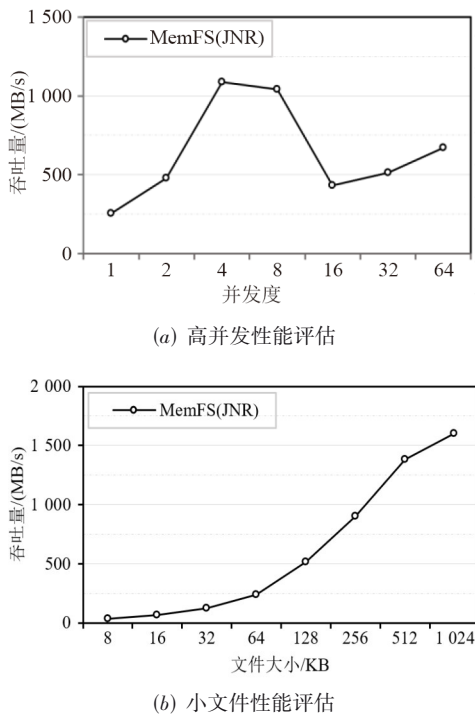


图3 JNR-FUSE框架访问MemFS的性能评估

场景下存在严重的性能问题。JNR-FUSE 在高并发场景下性能会异常下降；而在读取大量小文件时，JNR-FUSE 的性能远低于大文件读取。在深度学习训练、大数据分析等数据密集型作业中，经常会有高并发和小文件的工作负载，JNR-FUSE 将面临严峻的性能挑战。

3.2 性能瓶颈分析

基于上节发现的 JNR-FUSE 性能实验结果，本节重点分析其性能瓶颈产生原因。跨语言 FUSE 框架的时间开销主要可分为三部分：跨语言函数回调开销、跨语言数据拷贝开销、FUSE 工作线程管理开销。本节从以上三个维度对以 JNR-FUSE 为代表的跨语言 FUSE 框架的性能瓶颈进行分析。

3.2.1 跨语言函数回调

FUSE 依靠文件系统开发者定义的回调函数完成文件系统功能。FUSE 库 libfuse 使用 C 语言实现，当文件系统与 libfuse 使用不同语言开发时，则存在跨语言函数回调的过程。以 JNR-FUSE 为例，文件系统开发者用 Java 语言开发文件系统的各种操作，各种操作以回调函数的形式被 FUSE 的工作线程 (Daemon 线程) 调用，中间依次经过 FUSE 层、跨语言调用层 (例如 JNR Native Libraries)、JVM 层等，如图 4 所示。

从 libfuse 库本地工作线程 (Native Thread) 回调文件系统提供的 JVM 方法的流程包含以下步骤。

(1) 线程链接 (Attach)^[20]：将本地线程链接到一个

JVM 线程，本地线程成为能被 JVM 统一调度管理的 JVM 线程，并获取能和 JVM 交互的 JNIEnv 句柄。

(2) 元信息初始化：使用 JNIEnv 句柄获取待调用方法的 MethodID 等调用元信息，这些元信息在之后的调用过程中将被用于定位待调用的 Java 方法。

(3) 函数参数初始化：使用 JNIEnv 句柄将 C 语言的本地数据类型转化为 JVM 数据类型，对于本地的结构体等复合类型，还需额外通过 JNIEnv 句柄调用 JVM 中的对应构造函数，以将结构体转换为 JVM 中对应的临时对象。

(4) 跨语言调用：使用 JNIEnv 句柄、调用方法元信息和转换后的函数参数，通过 JNI 机制调用 JVM 中的对应方法。

(5) 环境清理：清理跨语言调用使用的临时对象，释放无用的元信息、临时对象等。

(6) 线程分离 (Detach)：将本地线程从 JVM 中分离，释放 JNIEnv 句柄，使当前线程再次成为独立调度管理的本地线程。

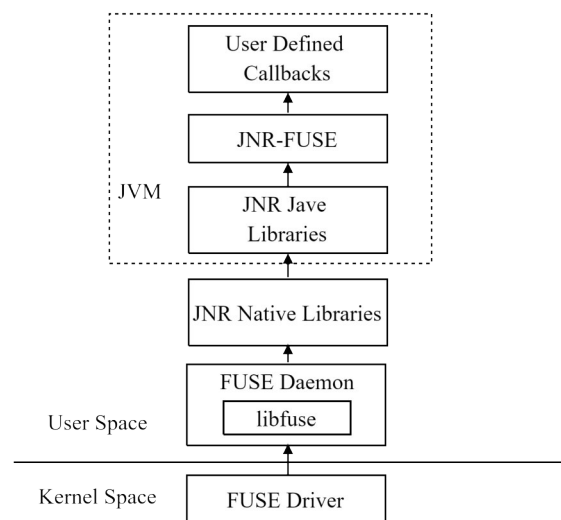


图4 基于 JNR-FUSE 的回调函数调用链

相比同语言内的函数调用，跨语言函数调用包含了复杂的准备步骤和清理步骤，会大幅影响系统整体性能。当使用跨语言函数回调的方式实现 FUSE 时，数据密集型工作负载会频繁触发跨语言函数回调。在大量小文件读取场景中，跨语言函数调用次数将大幅增加。当读取总大小为 10 GB 的大量小文件时，触发的跨语言调用次数与文件大小的关系如表 1 所示。大量跨语言函数回调显著增加了文件系统的额外开销，从而造成了 JNR-FUSE 在小文件场景下出现性能下降的问题。

3.2.2 跨语言数据拷贝

数据拷贝也是文件系统设计中不可忽视的开销因

表 1 跨语言调用次数与文件大小关系

文件大小	跨语言调用次数	相对比率
128 KB	573 441	1.0
32 KB	1 966 081	3.43
8 KB	7 864 321	13.71

素. 当 FUSE 和用户自定义文件系统以动态链接的形式使用同一个内存地址空间时, FUSE 所产生的数据拷贝是简单的内存间数据移动. 比如在 read 和 write 等系统调用中, 文件系统仅需通过 memcpy 函数把文件数据移动到 FUSE 的缓冲区即可. 但对于采用 Runtime 管理内存的 Java 语言而言, JNR-FUSE 等跨语言 FUSE 框架存在着额外的跨语言数据拷贝问题, 如图 5 所示. 例如 read 系统调用要求将存储在 JVM 堆上空间的用户态文件的文件数据拷贝到由 FUSE 管理的堆外原生的内存缓冲区中; 对于 write 系统调用而言, 则需要相反的过程.

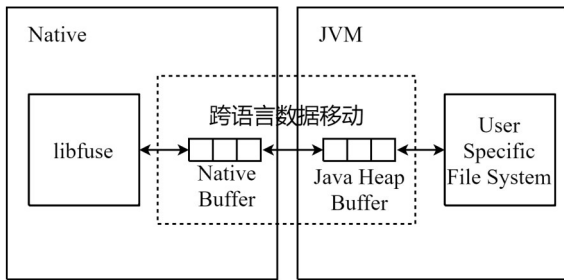


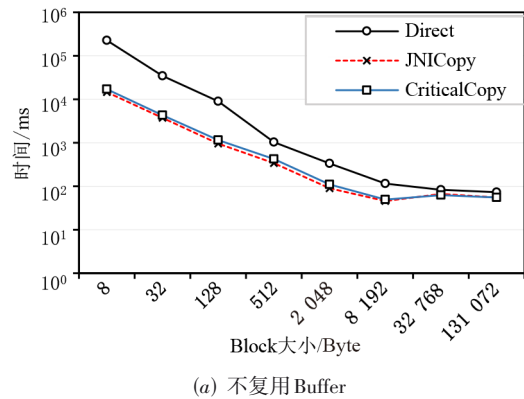
图 5 跨语言数据拷贝

JNR-FUSE 采用 JNI 拷贝 (JNICopy)^[21] 的方式实现跨语言数据拷贝, 其通过调用 JNI 方法实现堆外数据缓冲区和堆上缓冲区间的数据拷贝. 每次数据拷贝均涉及一次 JNI 调用, 采用单独的缓冲区完成, 带来额外开销. 因此当需要拷贝的文件越小、数量越多时, JNI 调用越频繁, 数据拷贝效率越差.

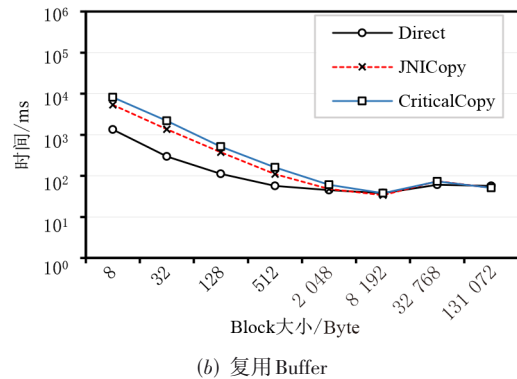
除了 JNICopy, JVM 还提供了 DirectByteBuffer (Direct)^[22] 与临界区拷贝 (CriticalCopy)^[23] 两种方式实现堆上空间和堆外空间的数据拷贝. Direct 方式通过 Java 内置的直接缓冲区编程接口, 在 Java 代码中直接读写堆外内存数据. Direct 方式是 Java 原生支持的堆外数据访问方式, 由 Java 管理缓冲区空间, 因此其安全性最高. 虽然该方式可避免 JNI 调用, 但创建 DirectByteBuffer 对象会带来额外的执行开销. CriticalCopy 方式通过调用 JNI 获取堆上缓存区的内存地址, 然后使用 memcpy 函数直接将堆上的数据复制到堆外缓冲区中. 然而此过程会阻塞 JVM 垃圾回收过程, 并且直接暴露堆上对象的内存地址, 安全性最低.

因为系统实现机制不同, 上述三种数据拷贝方式的效率不同. 本文通过实验测试了这三种方式通过不

同容量的缓冲区拷贝 10 GB 数据的总耗时. 由于缓冲区的创建和销毁会产生额外开销, 本文分别测试了不复用缓冲区和复用缓冲区时的性能, 实验结果如图 6 所示. 实验结果表明在复用缓冲区的情况下 DirectByteBuffer 的数据拷贝速度最快, 但 DirectByteBuffer 对象的创建和销毁会产生极大的额外开销, 导致其在不复用缓冲区时显著慢于其他两种方式. 同时, 对所有的数据拷贝方式而言, 增大缓冲区容量均能显著优化性能. 当缓冲区容量增大到 128 KB 时, 三者几乎不存在性能差距. 相反, 减小缓冲区容量会导致数据拷贝次数增多, 使数据拷贝总用时增长.



(a) 不复用 Buffer



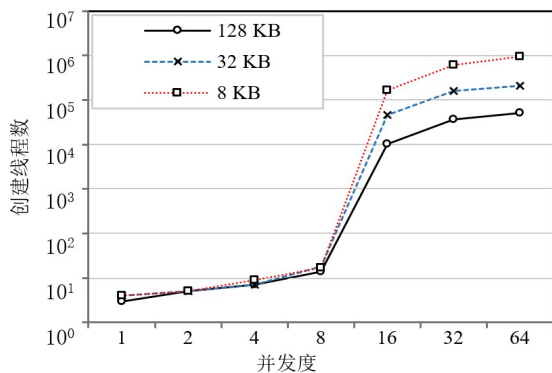
(b) 复用 Buffer

图 6 跨语言数据移动性能评估(纵轴采用对数坐标)

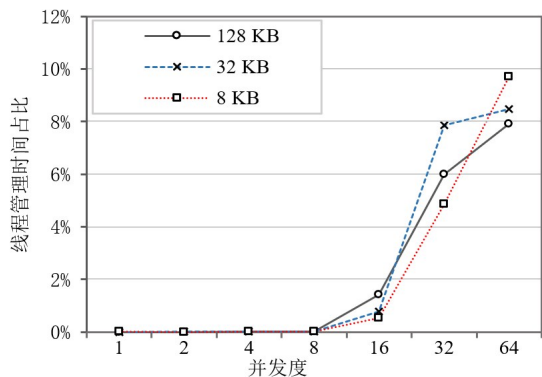
实验结果表明当缓冲区容量较小(小于 128 Byte)时, 三种跨语言数据移动方式存在较为明显的性能差异. 但在实际应用中, 这部分开销难以构成主要性能瓶颈. 首先, 实际应用中极少出现小于 512 Byte 的文件, 并且数据拷贝基本都优先以大块方式进行, 因此三种方式性能差距较小. 其次, 三种跨语言数据拷贝执行时间较短, 其开销在整个跨语言 FUSE 过程中占比不高. 如图 6 所示, 当缓冲区容量大于 128 Byte 时, 三种数据拷贝方式均能于 10s 内完成 10 GB 数据拷贝, 数据拷贝带宽超过 1 GB/s, 远大于文件系统自身数据拷贝所能达到的实际吞吐量.

3.2.3 FUSE 工作线程管理

高并发场景下的实验结果表明 FUSE 现有的工作线程管理机制难以适应大量线程并发读写的场景. 如图 7(a)所示,随着读取线程数量增加,在读取并发度不超过 8 时,FUSE 创建的线程数量很少,而当读取并发度达到 16 及其以上时,FUSE 创建的线程数量呈指数增长. 小文件场景会进一步加重此问题. 在读取 32 KB 和 8 KB 的文件时,创建线程数量分别达到了读取 128 KB 文件时的 4.1 倍和 18.6 倍. 图 7(b)进一步展示了 FUSE 工作线程管理开销对系统性能的影响. 当读取并发度为 1~8 时,因为只创建了很少的线程,线程管理时间开销基本占比可忽略不计;而当读取并发度从 16 增加到 64 时,线程管理开销占比不断增加,在 8 KB 文件规模和 64 读取并发度设置下,约有 10% 的执行时间被浪费在线程管理上.



(a) 创建线程数量(纵轴采用对数坐标)



(b) 线程管理开销占总时间比例

图 7 FUSE 工作线程管理开销分析

造成上述问题的原因在于 FUSE 的工作线程管理机制. 文件系统通常会读写硬盘等较慢速的设备文件, 涉及大量 I/O 请求, 导致 FUSE 的工作线程通常被阻塞在 I/O 操作上. FUSE 的工作线程管理机制通过创建足够多的工作线程(可超过 CPU 核数的限制), 以提升文件系统的总 I/O 吞吐量, 使尽可能多的文件系统调用可

以被及时处理. FUSE 在初始化时仅会创建一个工作线程, 当 FUSE 发现排队等待的 FUSE 请求数量过多、当前工作线程数量不足时, FUSE 会创建额外的工作线程以并发处理 FUSE 请求. 然而在高并发场景下, 当 FUSE 工作线程的实际处理速度低于 I/O 请求的到达速度时, 会造成 I/O 请求堆积, FUSE 会创建大量工作线程, 这些工作线程在对应的 I/O 请求完成后又被销毁, 从而带来不可忽视的线程管理开销.

实际上, 上述问题不只限于跨语言 FUSE 场景, 在 C 语言原生的 FUSE 调用场景也可能发生. 该问题产生的根源在于文件系统的回调函数处理速度慢于 I/O 请求的到达速度, 造成大量 I/O 请求堆积. 该问题与文件系统采用的程序设计语言无关. 但因为跨语言 FUSE 相比原生 C 语言函数回调额外增加了跨语言调用开销, 该问题在实际应用中变得更加显著.

从上述分析可知, 降低 FUSE 工作线程管理开销的方法包括直接方式和间接方式两类. 在直接方式方面, 可优化 FUSE 的工作线程管理机制, 例如限制最大工作线程数量、推迟工作线程销毁时机, 从而降低 FUSE 的工作线程管理开销. 但这类方法需要修改用户态 FUSE 库 libfuse 的源代码, 需重新编译和部署操作系统环境, 维护开销较高, 还存在不兼容现有上层应用的风险, 难以快速应用于企业的生产环境. 与上述直接方式不同, 间接方式则是通过降低跨语言 FUSE 的回调函数开销, 提升跨语言 FUSE 请求的处理速度, 降低 I/O 请求堆积的风险, 从而间接减少 FUSE 工作线程管理机制对整体性能的负面影响.

3.3 性能分析总结

上述的性能分析剖析了现有跨语言 FUSE 框架 JNR-FUSE 在高并发和小文件场景下存在的性能瓶颈及其产生的原因: 第一, JNR-FUSE 会频繁触发跨语言函数回调, 在小文件场景下会大幅增加函数回调次数; 第二, 堆内内存和堆外内存间频繁的数据拷贝开销在极端小文件场景(文件小于 128 Byte)下也不容忽略; 第三, FUSE 本身的工作线程管理机制会导致在高并发场景下易造成工作线程的频繁创建和销毁.

基于上述分析, 本文提出优化跨语言 FUSE 框架的性能可遵循的三个优化方向: (1) 设计低开销的跨语言函数回调机制, 从而实现对大量小文件的高效访问; (2) 选择低开销且高安全性的跨语言数据拷贝方式, 从而实现对大量文件的高吞吐、高安全访问; (3) 设计轻量级的工作线程管理机制, 从而降低高并发场景下工作线程创建与销毁所带来的额外开销. 在企业生产环境下, 因为 FUSE 库 libfuse 通常以操作系统内置库的形式提供, 在不修改 libfuse 库源代码的情况下, FUSE 库的工作线程管理机制难以修改, 因此在设计系统兼容性

强的跨语言 FUSE 框架时应更关注前两个优化方向。

4 JNI-FUSE 框架设计实现

根据跨语言 FUSE 框架设计的优化方向,本文面向以 Java 语言开发的分布式文件系统,设计了跨语言函数回调开销小、跨语言数据拷贝安全性高的跨语言 FUSE 框架 JNI-FUSE,以加速数据密集型作业在分布式文件系统上的运行速度。

本文设计实现的 JNI-FUSE 框架原型总共包含约 3 700 行代码,其中包括 C 语言实现的本地代码约 1 400 行和 Java 语言实现的 JVM 接口代码约 2 300 行。JNI-FUSE 的源代码^①已被知名分布式文件系统 Alluxio 接受并集成作为 Alluxio 默认的 FUSE 实现方式。

4.1 框架架构

作为跨语言 FUSE 框架,JNI-FUSE 连通了操作系统 FUSE 库 libfuse 的本地 C 语言环境与目标分布式文件系统的 Java 语言环境,使 FUSE 库 libfuse 可以通过 JNI-FUSE 框架的跨语言函数回调机制调用分布式文件系统的相关功能,而分布式文件系统则无须感知 libfuse 库与内核 FUSE 驱动的技术细节即可通过 JNI-FUSE 提供 POSIX 文件访问接口。

4.1.1 架构设计

JNI-FUSE 由本地代码层和 JVM 层两个模块组成,其架构如图 8 所示。本地代码层向下和操作系统提供的 libfuse 库对接,向上则调用 JVM 层文件系统提供的 Java 语言编程接口,以实现相关的文件操作。JVM 层则向下对接本地代码层,向上为用户提供统一便捷的文件系统接口 AbstractFuseFileSystem。该接口将文件系统的典型操作抽象为 read, write, release 等方法,与 POSIX 文件访问接口保持一致。分布式文件系统的开发者需要继承该接口,并在该接口的各方法中调用分布式文件系统的对应编程接口,从而对接各种用户态分布式文件系统。JNI-FUSE 实现了 FUSE 库与分布式文件系统的编程接口之间的解耦,分布式文件系统的开发者无须了解 FUSE 编程的具体细节即可通过 JNI-FUSE 框架对接操作系统的 FUSE 库,从而提供 POSIX 文件接口的访问能力。

根据第 3.3 节的性能瓶颈分析,跨语言 FUSE 的可能性能优化方向包括降低跨语言函数回调开销、降低跨语言数据拷贝开销与改进工作线程管理机制。其中前两个优化方向只与跨语言 FUSE 框架本身相关,通过改进跨语言 FUSE 框架自身的架构与实现方法即可达成。而第三个优化方向与操作系统提供的 libfuse 库实现相关,需要修改与重新编译相关系统库,使其难以部

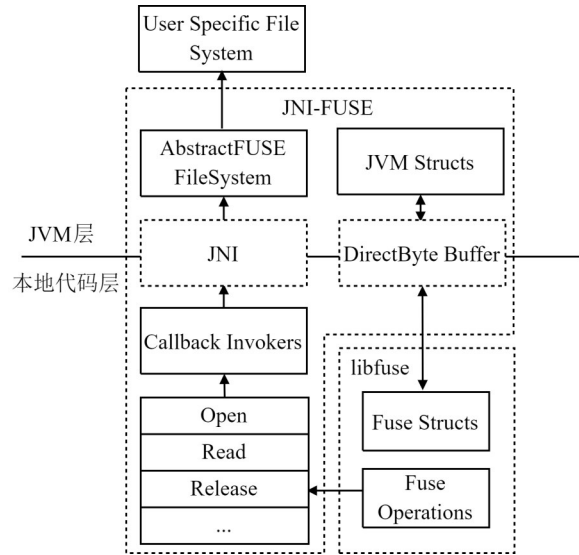


图 8 JNI-FUSE 框架

署于注重稳定的企业生产环境。

因此为了避免修改操作系统提供的 libfuse 库、提升 JNI-FUSE 框架对不同 libfuse 库版本的适应性,本文将 JNI-FUSE 框架设计为独立于操作系统 libfuse 库的第三方软件模块,通过延迟分离、元信息缓存等技术降低跨语言函数回调的执行开销,通过降低 I/O 请求堆积的风险,从而间接减少 FUSE 工作线程管理机制的负面影响。对跨语言 FUSE 的性能优化与对 FUSE 工作线程管理机制的优化之间是互相独立的,可以互相叠加。本文的研究主题更关注于跨语言 FUSE 特有的性能优化问题。

JNI-FUSE 框架的独立性使得本文将跨语言 FUSE 性能优化的重点聚焦于降低跨语言函数回调的执行开销,通过减少 I/O 请求堆积的发生来间接降低 FUSE 工作线程管理机制的负面影响。

4.1.2 工作流程

JNI-FUSE 处理一次来自应用程序的 POSIX 文件系统调用的工作流程如图 9 所示。运行在用户空间的应用程序通过 POSIX 文件系统调用(例如 read, write 等)的方式启动一次基于 FUSE 的 POSIX I/O 请求。应用程序的系统调用由运行在内核中的 FUSE 驱动处理, FUSE 驱动将 I/O 请求发送至运行在用户态的 libfuse 库的 FUSE 工作线程处理。libfuse 库的 FUSE 工作线程调用 JNI-FUSE 框架本地代码层中的对应函数,将文件系统的 I/O 请求转发至 JNI-FUSE 框架处理。

具体地,libfuse 库的工作线程首先通过线程链接步骤(步骤 1)将 FUSE 工作线程链接(Attach)为 JVM 线程^[20]。之后,元信息初始化步骤(步骤 2)会获取跨语言函数调用所需的调用方法元信息(包括 JVM 对应方法

① <https://github.com/Alluxio/alluxio/tree/release-2.6.0/integration/jnifuse>

的 MethodID 等)。跨语言数据拷贝步骤(步骤 3)将 C 语言缓冲区内的数据拷贝至 Java 语言的内存缓冲区中。基于跨语言函数调用元信息与 Java 内存缓冲区, JNI-FUSE 框架通过 JNI 机制调用 JVM 层由分布式文件系统开发者实现的 AbstractFuseFileSystem 接口的相应方法(步骤 4), 执行流转入 JVM 层。JVM 层在 AbstractFuseFileSystem 接口方法中具体地调用目标分布式文件系统的相关函数, 从而完成对分布式文件系统的操作。AbstractFuseFileSystem 接口调用完成之后, JNI-FUSE 框架清理跨语言函数调用产生的临时数据与元信息(步骤 5), 最后将本地工作线程从 JVM 分离(步骤 6)。

一次完整的跨语言 FUSE 函数回调链路很长, 这个过程存在大量冗余的准备和清理工作。本文通过引入延迟分离(Defer Detach)和元信息缓存(Meta Cache)技术, 跳过跨语言 FUSE 函数回调中回调准备和清理的相关步骤(包含链接线程、元信息初始化、环境清理与线程分离等四个步骤), 从而降低跨语言 FUSE 的函数

回调开销。

4.2 延迟分离优化

延迟分离优化主要针对跨语言 FUSE 过程中执行开销较高的线程链接与线程分离步骤, 如图 9 所示。这两个步骤的执行开销主要来自 JVM 的线程管理开销。当一个本地线程需要从 C 语言函数中调用 Java 的方法时, 首先需要通过 JNI 提供的线程链接机制将本地线程链接(Attach)为一个 Java 线程, 从而获得 JNIEnv 句柄并通过该句柄调用 Java 方法。在线程链接过程中, JVM 需要执行一系列 JVM 线程初始化操作(例如关闭浮点数异常、调整线程优先级等), 以将本地线程包装为一个 JVM 线程, 由 JVM 运行时统一管理; 线程分离过程则需要销毁对应的 JVM 线程, 从而实现本地工作线程与 JVM 线程的分离(Detach)。JVM 线程的初始化与销毁给两个步骤带来了显著的额外开销, 当跨语言文件系统函数调用过程(步骤 4)执行开销较小时, JVM 线程链接与分离带来的额外开销将变得更加显著。

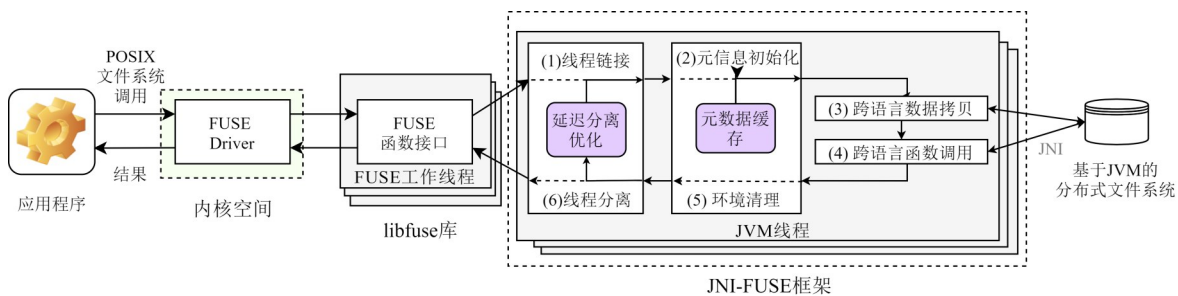


图9 JNI-FUSE框架工作流程

延迟分离技术的总体思路是在跨语言调用时推迟本地线程从 JVM 分离的时间, 以便在后续的跨语言函数回调中复用已链接的线程, 跳过程序链接(步骤 1)与线程分离(步骤 6)步骤, 如图 9 所示。具体地, 在从 libfuse 库本地线程首次进行跨语言函数回调的过程中, 在线程分离步骤暂不释放当前的 JVM 线程, 而是保持当前线程与 JVM 的链接关系, 将线程链接产生的 JNIEnv 句柄^[20]缓存在线程局部变量中。在之后的跨语言函数回调过程中, 线程链接步骤直接读取缓存的 JNIEnv 句柄, 跳过 JVM 线程链接过程。跨语言函数调用之后的线程分离步骤则同样被跳过。采用延迟分离技术后, JNI-FUSE 只会第一次跨语言函数回调时将本地线程链接到 JVM, 直到本地线程被销毁时才从 JVM 线程分离出来, 从而大幅降低 JVM 线程链接与分离带来的额外开销。

延迟分离技术要求各线程缓存的 JNIEnv 句柄应在 FUSE 工作线程销毁之前执行释放句柄的操作, 以保证 JVM 线程管理机制正常工作。但 FUSE 工作线程的销毁由 libfuse 库触发执行, JNI-FUSE 框架难以控制与掌

握。因此如果采用全局变量缓存各线程的 JNIEnv 句柄, 无法保证 JNIEnv 句柄随工作线程的结束而正确释放。为了解决延迟分离过程中 JNIEnv 句柄的释放问题, JNI-FUSE 框架采用了 pthread 多线程编程技术提供的线程特有数据(thread-specific data)机制来管理每个线程缓存的 JNIEnv 句柄。JNI-FUSE 利用该机制将每个工作线程的 JNIEnv 句柄声明为该线程特有数据对象, 并为该数据对象注册特定析构函数。在工作线程被操作系统销毁之前, pthread 会自动调用线程特有数据对象的注册析构函数, 从而正确释放 JNIEnv 句柄。

延迟分离技术因为缓存 JNIEnv 句柄、保持 JVM 线程链接, 而额外增加了与工作线程数成正比的空间开销, 但该部分额外空间开销的总量较小。该空间开销主要由两部分组成: (1)工作线程缓存的 JNIEnv 句柄所占空间, (2)工作线程对应的 JVM 内部线程对象所占空间。具体而言, 每个 JNIEnv 句柄包含一系列 JVM 函数指针与数据指针信息。此外, 工作线程链接到 JVM 时, JVM 会为工作线程创建对应的 JVM 线程对象, 并在 JVM 栈区与堆区申请与初始化相应的内存空间。延迟

分离优化会一直保持工作线程与JVM的链接状态. 上述JVM内部线程的存储开销在工作线程销毁后即被释放. 在实际应用中,JNI-FUSE在高并发场景下通常产生20个左右的工作线程(详见第5.2.1节的实验),单个线程延迟分离优化所增加的平均物理内存开销为636.52 KB,其产生的总物理内存开销约12.43 MB,因此额外存储空间开销不大.

4.3 元信息缓存

元信息缓存技术主要针对跨语言函数回调过程中的元信息初始化(步骤2)与环境清理(步骤5)步骤. 元信息初始化步骤需要为跨语言函数调用准备调用元信息,包括获得JVM中对应方法的MethodID等. 因为同一个Java方法的JNI调用元信息是静态的、不随线程而改变,所以当同一个文件系统操作被应用程序重复调用多次时,重复的元信息初始化与环境清理将带来冗余的执行开销.

元信息缓存技术的总体思路是在框架初始化时缓存各种文件系统操作(例如open, read, release等)所需的元信息,在后续的跨语言函数回调过程中直接复用相关信息,从而跳过元信息初始化与环境清理的步骤. 因为各方法的调用元信息是用户不可更改的只读信息,且元信息只与回调函数相关,与工作线程或JVM线程无关,因此跨语言函数调用元信息只须初始化一次,即可在多个工作线程中复用,而无须考虑多线程访问一致性的问题.

具体地,JNI-FUSE框架在框架初始化时,会为用户实现的AbstractFuseFileSystem接口内每种支持的POSIX文件系统调用预取对应的跨语言函数回调相关元信息,并缓存在内存中(图10). 缓存的调用元信息在所有线程之间共享,从而避免多个线程之间冗余的元信息获取开销. 在应用程序的一次跨语言FUSE函数调用的过程中,JNI-FUSE框架在元信息初始化步骤中从内存中提取对应的跨语言函数调用元信息. 基于提取出的元信息,JNI-FUSE即可利用JNI机制发起跨语言函数调用.

值得注意的是,本文提出的元信息缓存技术缓存的并不是文件系统的元信息,而是跨语言函数调用的元信息(例如MethodID、调用句柄等). 文件系统自身的元信息缓存通常由文件系统客户端管理,其与JNI-FUSE框架的元信息缓存技术是独立的,两者可以联合使用以提升文件系统的跨语言FUSE性能.

元信息缓存技术同样引入了额外的内存空间开销. 元信息缓存需存储各类文件系统调用对应JVM方法的MethodID信息,其存储开销与调用的JVM方法数量成正比. 对于单个JVM方法而言,本文实验所采用的64位HotSpot JVM中,每个MethodID包含了一个指向JVM内部方法对象的指针,占用8个Byte. JNI-FUSE调用的JVM方法数量的上限由AbstractFileSystem接口包含的

POSIX文件系统调用数量决定,而实验所采用的v2.9.5版本的libfuse库最多支持44类文件系统调用,因此元信息缓存所需的额外空间至多为352 Byte,开销很小.

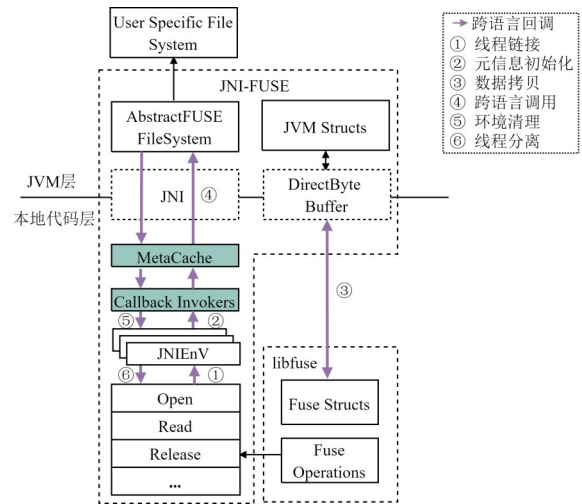


图10 元信息缓存优化

4.4 跨语言数据拷贝机制选择

目前在Java语言下可实现跨语言数据拷贝的机制有JNICopy, CriticalCopy和DirectByteBuffer三种. 表2从安全性、通用性与性能三个方面对比了三种跨语言数据拷贝方式. 在三种方式之中,DirectByteBuffer的安全性最高,DirectByteBuffer是Java内置的数据类型,无须用户直接管理缓冲区空间,可在规避内存泄漏等难以察觉和调试的问题的同时,避免直接访问指针可能带来的访问错误. 此外,DirectByteBuffer的通用性最好,其屏蔽了底层硬件和软件环境的细节,尽可能地避免了因硬件和软件环境变化所带来的适配问题,保持了JVM的“Write once, run anywhere”跨平台优势.

作为对比,JNR-FUSE等框架采用了JNICopy等通用性较弱的拷贝机制. 为保持跨语言FUSE框架的跨平台通用性,JNR-FUSE等框架额外封装了多层胶水代码,最终带来较大性能损失,得不偿失. 从性能角度比较,第3.2.2节的实验结果表明:DirectByteBuffer的数据拷贝性能在数据块较小时稍弱于JNICopy和CriticalCopy两种方式;但在FUSE的应用场景中数据多采取较大的数据缓冲区(128 KB),DirectByteBuffer与其他跨语言数据移动方式相比并没有明显性能差距. 综合考虑性能、安全性和通用性,JNI-FUSE选择采用DirectByteBuffer实现跨语言数据拷贝.

表2 不同跨语言数据拷贝机制比较

跨语言数据拷贝机制	安全性	通用性	性能
JNICopy	中	低	高
CriticalCopy	低	中	高
DirectByteBuffer	高	高	中

4.5 跨语言适用性分析

本文所提出的延迟分离、元信息缓存优化方法以及选择的跨数据拷贝机制适用于所有基于 JVM 的编程语言(例如 Java, Scala, Kotlin 等)开发的分布式文件系统,包括大数据软件栈中广为使用的 HDFS, Alluxio 等. 其原因是基于 JVM 的文件系统均需要借助 JNI 机制,实现 libfuse 库回调 JVM 中分布式文件系统提供的相应方法的功能. 本文提出的优化方法针对 JNI 调用中相关初始化与清理步骤设计,适用于所有基于 JVM 的编程语言,具备较好的跨语言适用性.

5 性能评估

5.1 实验设置

本节从 JNI-FUSE 框架(本文方法)自身性能、文件系统端到端读取性能、应用层深度学习训练性能三个方面评估 JNI-FUSE 的性能. 本节采用目前最为成熟稳定、性能最佳、使用最广的 Java 端 FUSE 框架 JNR-FUSE^[13]作为对比框架. 为保证测试的准确性和公平性,所有测试均在 FUSE 的 `direct_io` 模式(不使用 kernel cache)下进行. 如无特别说明,保留其余参数为默认配置.

考虑到很多大数据深度学习作业都基于容器云平台进行,本节实验在 Kubernetes^[24]集群中进行,测试程序运行在 Docker 容器中,通过 CSI^[25]对 FUSE. 第 5.2 节和第 5.3 节的实验在无 GPU 的普通计算服务器上进行测试,第 5.4 节的深度学习训练实验在一台具有两张 Tesla V100 GPU 的深度学习服务器上进行测试,两台计算服务器的配置如表 3 和表 4 所示.

表 3 普通计算服务器软硬件配置信息

参数	配置信息
CPU	Intel(R) Xeon(R) Silver 4114 CPU @ 2.20 GHz (40 cores)
Memory	96 GB DDR4
Disk	240 GB SSD
Network Bandwidth	10 Gbps
OS	CentOS Linux release 7.6.1810
JVM Version	Java 1.8.0
FUSE Version	2.9.5
Alluxio Version	2.5.0

5.2 框架性能评估

为评估 JNI-FUSE 框架(本文方法)自身带来的额外开销,本文以第 3.1 节所述的伪内存文件系统 MemFS 作为 FUSE 的目标文件系统,分别基于 JNR-FUSE^[13]和 JNI-FUSE 框架实现了该文件系统的 POSIX 接口,分别记为 MemFS(JNR-FUSE)和 MemFS(JNI-FUSE).

表 4 GPU 深度学习服务器软硬件配置信息

参数	配置信息
CPU	Intel(R) Xeon(R) Silver 4214 CPU @ 2.20 GHz (48 cores)
GPU	Tesla V100x2
Memory	256 GB DDR4
Disk	4.7 TB HDD
Network Bandwidth	10 Gbps
OS	CentOS Linux release 7.6.1810
JVM Version	Java 1.8.0
FUSE Version	2.9.5
Alluxio Version	2.5.0
CUDA Version	11.0
PyTorch Version	1.7.1

5.2.1 读取性能评估

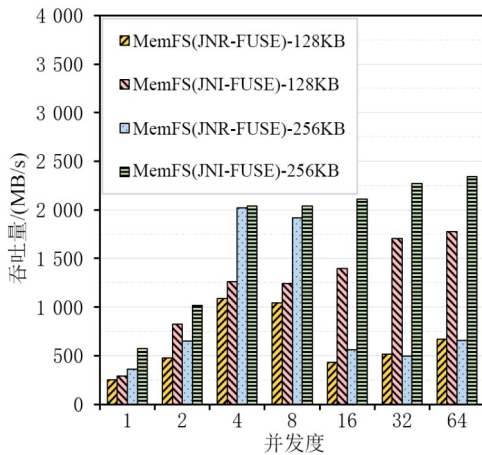
本文在普通计算服务器上分别对比了 MemFS (JNR-FUSE)和 MemFS(JNI-FUSE)在高并发和小文件场景下的文件读取性能.

高并发场景中,使用不同数量的线程并发从 MemFS 读取 81 920 个文件,每个文件大小先后设置为 128 KB 和 256 KB,总大小为 10 GB 和 20 GB. 高并发场景的读取性能测试结果如图 11(a)所示. 随着读取线程数的增长,MemFS(JNR-FUSE)的吞吐量先快速增长,但当并发度超过 8 时,MemFS(JNR-FUSE)的吞吐量出现骤降,尤其是在并发度从 8 增加到 16 时,吞吐量降低了 50% 以上,其主要受大量工作线程管理开销影响. 而 MemFS(JNI-FUSE)的吞吐量则随并发线程数的增长保持增长趋势. 当每个文件为 128 KB 时,MemFS (JNI-FUSE)相比 MemFS (JNR-FUSE)带来了平均约 2.06 倍、最高约 3.32 倍的性能提升.

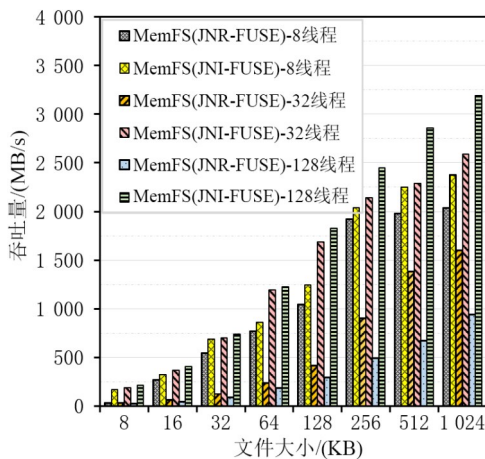
在 4 线程当每个文件为 256 KB 时,MemFS (JNI-FUSE)相比 MemFS(JNR-FUSE)带来了平均约 3.71 倍、最高约 5.78 倍的性能提升. MemFS (JNR-FUSE)和 MemFS(JNI-FUSE)分别在 4 线程与 64 线程并发度下取得最大吞吐量,MemFS(JNI-FUSE)的最大吞吐量相比 MemFS(JNR-FUSE)在 128 KB 和 256 KB 文件大小下分别提高 63.5% 和 16%. 在高并发场景下,JNI-FUSE 框架显著降低了跨语言函数回调开销,使 MemFS (JNI-FUSE)可以以较少的 FUSE 工作线程数处理高并发读取,大幅削弱了 FUSE 工作线程管理所带来的影响,从而获得比 JNR-FUSE 更好的性能.

本文进一步通过 8 线程和 64 线程并发度读取大量 32 KB 小文件的场景,对比了 JNI-FUSE 和 JNR-FUSE 所使用的 FUSE 工作线程数量的差异,实验结果如图 12 所示. 在 8 线程和 64 线程并发度下,JNI-FUSE 使用的工作线程数全程分别稳定在 10 和 20 左右. 而 JNR-FUSE

使用的工作线程数在 8 线程并发度下同样稳定在 10 个工作线程,但在 64 线程并发度下则呈现大幅波动,每次波动意味着大量工作线程的创建和销毁,其最高时甚至同时存在 120 个工作线程. 实验结果表明 JNI-FUSE 比 JNR-FUSE 的读取速度更快,在同样读取总规模为 10 GB 文件的负载下,在 8 线程和 64 线程并发度下 JNI-FUSE 均在第 50 s 时提前完成测试,而 JNR-FUSE 则分别需要 60 s 与 256 s 才能完成,JNI-FUSE 的性能显著优于 JNR-FUSE,尤其是在高并发场景下.



(a) 高并发读取性能评估



(b) 小文件读取性能评估

图 11 不同跨语言 FUSE 框架下 MemFS 文件系统的读取性能评估

在小文件场景中,本文在保持 MemFS 数据总量为 10 GB 的情况下生成多组不同规模的小文件,并分别使用 8 个、32 个和 128 个线程并发读取. 图 10(b)展示了小文件读取场景的性能测试结果. 随着单个文件规模的增长,MemFS(JNI-FUSE)和 MemFS(JNR-FUSE)的吞吐量都有所增长,且在较大文件规模下 MemFS(JNI-FUSE)的吞吐量大幅领先于 MemFS(JNR-FUSE). 当文件越小时,MemFS(JNI-FUSE)的性能优势更加明显.

使用 8 个线程并发读取时,MemFS(JNI-FUSE)取得了平均 1.16 倍、最高约 5.44 倍的性能提升;使用 32 个线程并发读取时,MemFS(JNI-FUSE)取得了平均约 3.98 倍、最高约 5.77 倍的性能提升;使用 128 个线程并发读取时,MemFS(JNI-FUSE)取得了平均约 6.04 倍、最高约 9.03 倍的性能提升. 读取的文件越小,文件系统元数据操作在总耗时中的比例就越高,跨语言 FUSE 框架本身性能开销所占比重越大,JNI-FUSE 相比 JNR-FUSE 的优势更加明显.

5.2.2 写入性能评估

本文在普通服务器上分别对比了 MemFS(JNI-FUSE)与 MemFS(JNR-FUSE)在高并发和小文件场景下的写入性能. MemFS 初始时不存储任何数据,文件数量为 0. 在随后的测试中,测试用写入线程将会持续向 MemFS 中写入文件. 在高并发场景中,分别限制单个文件大小为 128 KB 和 256 KB,总量对应为 10 GB 和 20 GB;在小文件场景中,固定使用 32 个和 128 个写入线程,每次写入不同大小的文件,保持 MemFS 数据总量为 10 GB.

图 13(a)展示了高并发写入场景的性能测试结果,其整体增长趋势类似于高并发读取场景;图 13(b)展示了小文件写入场景性能测试结果. 小文件写入时,文件越小两者性能差距反而越小,这与小文件读取相反. 这是因为写入操作是元数据高度密集型操作^[9],性能瓶颈主要在 FUSE 层,无法体现 JNR-FUSE 和 JNI-FUSE 框架本身的性能差距. 在高并发写入场景中,当每个文件为 128 KB 时,MemFS(JNI-FUSE)相比 MemFS(JNR-FUSE)带来了平均约 1.15 倍、最高约 1.36 倍的性能提升;当每个文件为 256 KB 时,MemFS(JNI-FUSE)相比 MemFS(JNR-FUSE)带来了平均约 3.26 倍、最高约 5.07 倍的性能提升. 从两个框架在实验中所取得的最高写入吞吐量来看,MemFS(JNI-FUSE)在 128 KB 和 256 KB 文件大小下的最高吞吐量相比 MemFS(JNR-FUSE)分别提高了 5.0% 和 256.8%.

在小文件写入场景中,使用 8 个线程并发写入时,MemFS(JNI-FUSE)相比 MemFS(JNR-FUSE)带来了平均约 2.41 倍、最高约 4.86 倍的性能提升;使用 32 个线程并发写入时,MemFS(JNI-FUSE)相比 MemFS(JNR-FUSE)带来了平均约 1.49 倍、最高约 3.08 倍的性能提升;使用 128 个线程并发写入时,MemFS(JNI-FUSE)相比 MemFS(JNR-FUSE)带来了平均约 5.02 倍、最高约 6.59 倍的性能提升.

5.2.3 优化技术效果评估

JNI-FUSE 框架采用了延迟分离优化与元信息缓存优化等两种优化技术来提升性能. 本节通过比较优化前后跨语言 FUSE 过程的性能,以评估本文所提优化技

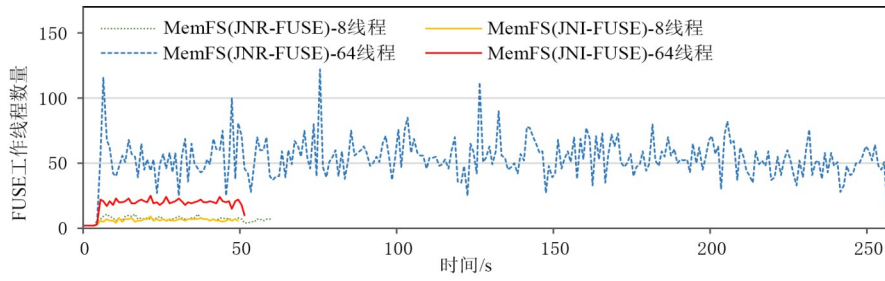
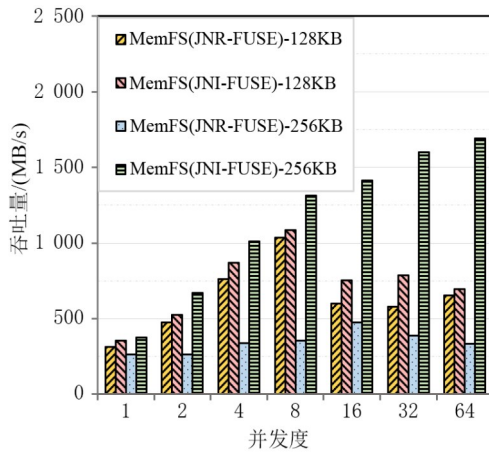
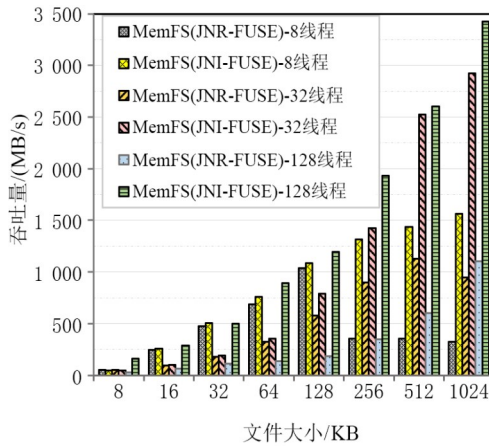


图 12 不同跨语言 FUSE 框架下 FUSE 工作线程数量对比



(a) 高并发写入测试



(b) 小文件写入测试

图 13 不同跨语言 FUSE 框架下 MemFS 文件系统的写入性能评估

术的有效性。以是否采用延迟分离技术与元信息缓存技术为标准,跨语言函数回调方式可以划分为 Baseline, Cache, Defer 和 Cache+Defer 四种。Baseline 方式采用 JNI 提供的基本跨语言调用过程,包含工作流程中步骤 1 至步骤 6 的所有过程。Cache 方式在 Baseline 的基础上使用元信息缓存优化。Defer 方式在 Baseline 的基础上,额外加入延迟线程分离优化。Cache+Defer 方式在 Baseline 的基础上,同时使用元信息缓存优化和延迟线程分离优化。

图 14 对比了在不同并发度下,分别采用以上四种方式执行 10 万次跨语言调用的总体耗时。图 14 实验结果表明不同方式的总体耗时差异较大:元信息缓存(Cache)和延迟分离优化(Defer)都能提升跨语言调用的性能,延迟分离优化效果最为明显,可带来两个数量级的性能提升。对比 Baseline,元信息缓存和延迟分离优化联合使用能够带来平均约 150 倍的性能提升。现有 JNR-FUSE 采用的跨语言函数回调是未充分优化的,它并未缓存元信息和延迟线程分离,因此只能达到 Baseline 的性能。JNI-FUSE 通过优化跨语言函数回调机制,可获得显著的性能提升。

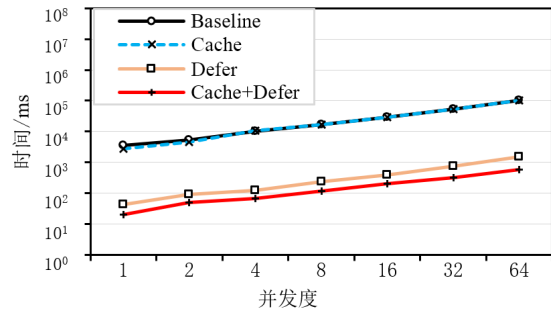


图 14 跨语言函数回调优化技术性能对比(纵轴采用对数坐标)

5.3 分布式文件系统端到端性能评估

本文在普通服务器上测试了广为使用的开源分布式文件系统 Alluxio 基于 JNI-FUSE 提供的 POSIX 接口的端到端性能。本文将 JNI-FUSE 集成到 Alluxio 中(记为 Alluxio(JNI-FUSE)),并采用 Alluxio 原本基于 JNR-FUSE 的实现^[26](记为 Alluxio(JNR-FUSE))作为对比。实验中采用大型真实图片数据集 SVHN^[27]的 extra 部分和大型真实图片数据集 COCO^[28]的测试集部分作为工作负载。SVHN 的 extra 部分包含 202 355 张图片文件,总大小约 3.6 GB,平均每张图片 19 KB,其中最小的图片不到 1 KB。COCO 的测试集部分包含 40 504 张图片文件,总大小约 6.2 GB,平均每张图片大小 160 KB,其中最小的图片不到 8 KB,最大图片接近 900 KB。Alluxio Worker 缓存大小设置为 30 GB,因此这两个数据集都可以完全缓存在内存中。为了加快 Alluxio 读取数据的速度,本文开启了 Alluxio 的元数据缓存优化,以减

少元数据操作 RPC 请求。

在测试时,本文预先将两种数据集分别加载进 Alluxio 中,接着分别用 Alluxio(JNR-FUSE)和 Alluxio(JNI-FUSE)把数据集挂载到本地文件系统以测试读取性能,测试结果如图 15 所示。在并发度不大于 8 时,两者的吞吐量基本与并发度呈线性增长关系,两个框架的性能基本持平;当并发度大于 8 时,Alluxio(JNR-FUSE)出现了严重的性能下降,而 Alluxio(JNI-FUSE)仅是降低了吞吐量随并发度的增长速度。因为 JNI-FUSE 比 JNR-

FUSE 更为高效稳定,最终在读取 SVHN 的 extra 部分时,Alluxio(JNI-FUSE)相比 Alluxio(JNR-FUSE)取得了平均 2.71 倍、最高 7.94 倍的性能提升;读取 COCO 的测试集部分时,对比 Alluxio(JNR-FUSE),Alluxio(JNI-FUSE)取得了平均 1.90 倍、最高 3.90 倍的性能提升。在两个数据集上,Alluxio(JNR-FUSE)均在并发度为 8 时取得最大吞吐量,而 Alluxio(JNI-FUSE)则均在 64 线程并发时取得最大吞吐量,Alluxio(JNI-FUSE)的最大吞吐量相比 Alluxio(JNR-FUSE)分别提升 19.4% 和 64.3%。

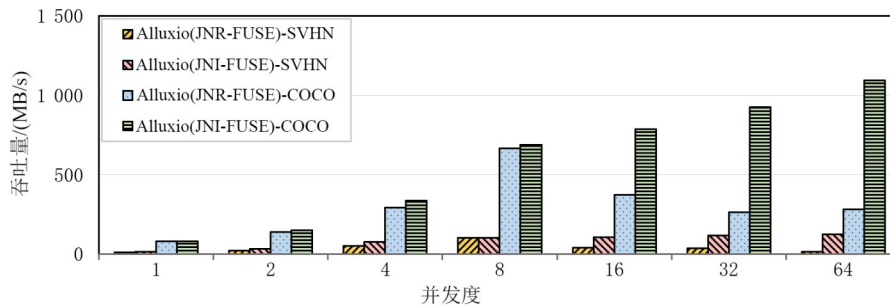


图 15 Alluxio 端到端性能测试

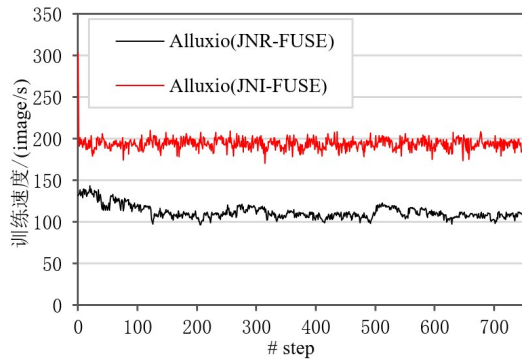
5.4 深度学习训练应用性能评估

最后本文在 Kubernetes 环境下的 GPU 计算服务器上利用大规模深度学习训练应用评估了本文优化后的跨语言 FUSE 框架对整个训练应用的性能影响。本文使用深度学习训练框架 PyTorch^[29]官方提供的 PyTorch 容器,运行其自带测试 Demo 程序^[30]对 ImageNet 数据集^[31] ILSVRC 子集的 CLS-LOC 部分进行训练;使用深度学习框架 Horovod^[32]官方提供的 Horovod 容器,运行其自带测试 Demo 程序^[33]对 CIFAR-10 数据集^[34]进行训练。ImageNet 和 CIFAR-10 都是大规模图像分类数据集。ImageNet 数据集 ILSVRC 子集的 CLS-LOC 部分总大小约 150 GB,包含了约 120 万张图片。CIFAR-10 数据集处理为 Horovod 自带测试 Demo 所需的格式后,总大小为 7.2 GB,包含了 60 000 张图片。本文将两种数据集分别存入 Alluxio 分布式文件系统中,并分别利用 Alluxio(JNR-FUSE)和 Alluxio(JNI-FUSE)挂载 Alluxio 中的数据集到本地文件系统,然后指定挂载点为数据源,并进行深度学习训练。深度学习框架在训练过程中采用的并发读取线程数通常由硬件配置(例如 CPU 数量)、用户给出的训练参数(例如 batch size 等)、深度学习框架配置(例如最大并发线程数等)共同确定。

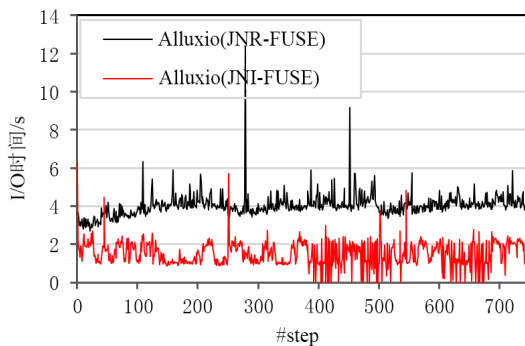
在 PyTorch 的实验中,训练过程中采用的数据并发读取线程数由 PyTorch 框架根据硬件配置、训练参数等计算得到,在本实验中为 8 个并发线程。因为完整的训练过程耗时很长,但前 3 个 Epoch 训练过程已经足够稳

定,之后的 Epoch 不会产生大的性能差异,因此在评估中取前 3 个 Epoch 的性能进行对比。图 16(a)比较了每 10 步迭代的平均训练速度,图 16(b)对比了每 10 步迭代中花费在数据读取上的平均耗时。在整个训练过程中,JNR-FUSE 平均每步训练有约 88.1% 的时间花费在了数据 I/O 上,而 JNI-FUSE 将这部分数据 I/O 的时间降低到了 55.6%,JNI-FUSE 的每步数据 I/O 时间比 JNR-FUSE 更短,因此使用 JNI-FUSE 时 GPU 资源利用率更高,训练速度更快。在深度学习训练任务中,JNI-FUSE 相比 JNR-FUSE 平均减少了 62.7% 的数据 I/O 时间,最终取得平均 1.73 倍的训练速度提升。

在 Horovod 测试中,选取前 100 次迭代的性能进行对比,控制数据并发读取线程数分别为 8 线程与 32 线程。图 17(a)比较了每步迭代的平均训练速度,图 17(b)对比了每步迭代中花费在数据读取上的平均耗时。在整个训练过程中,当并发预取线程数分别为 8 和 32 线程时,JNR-FUSE 平均每步训练有分别约 14.1% 和 59.0% 的时间花费在了数据 I/O 上,而 JNI-FUSE 将这部分数据 I/O 的时间降低到了 7.0% 和 11.1%,JNI-FUSE 的每步数据 I/O 时间比 JNR-FUSE 更短,因此使用 JNI-FUSE 时 GPU 资源利用率更高、训练速度更快。在深度学习训练任务中,8 线程和 32 线程并发读取情况下,JNI-FUSE 相比 JNR-FUSE 平均分别减少了 54.1% 和 86.7% 的数据 I/O 时间,最终取得了平均 1.06 倍和 1.50 倍的训练速度提升。

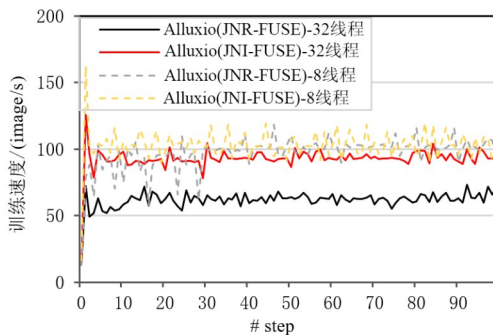


(a) 训练速度

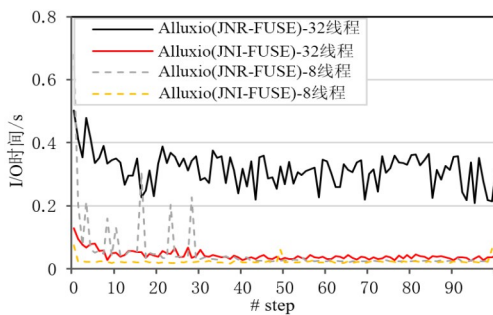


(b) I/O耗时

图16 PyTorch 自带 Demo 程序在 ImageNet 数据集 ILSVRC 子集的 CLS-LOC 部分上的训练性能对比



(a) 训练速度



(b) I/O耗时

图17 Horovod 自带 Demo 程序在 CIFAR-10 数据集上的训练性能对比

6 总结

现在很多大数据文件存储系统依赖 FUSE 提供 POSIX 兼容性,但跨语言 FUSE 框架的性能瓶颈缺乏系统性的研究与分析. 本文系统性地评估分析了现有分布式文件系统采用的跨语言 FUSE 框架 JNR-FUSE 的性能问题,进而总结出提升跨语言 FUSE 框架性能的可能优化方向,包括降低跨语言函数回调开销、跨语言数据拷贝开销与优化 FUSE 工作线程管理机制等.

基于上述优化方向,本文面向 Java 语言设计并实现了一个轻量化的跨语言 FUSE 框架 JNI-FUSE. 相比现有跨语言 FUSE 框架 JNR-FUSE, JNI-FUSE 能够显著提升小文件和高并发场景下的吞吐量. 本文将 JNI-FUSE 集成到了分布式文件系统 Alluxio,并在真实 Kubernetes 环境中测试了端到端性能及其对深度学习训练负载的影响. 实验结果显示本文提出的方法显著提升了文件读写吞吐量,降低了深度学习训练耗时. 本文设计实现的 FUSE 框架已开源,并被 Alluxio 官方社区接受集成并作为默认配置使用.

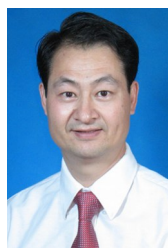
在未来的工作中,将进一步面向更多应用场景优化 FUSE 库的运行性能. 此外,还将继续深入研究其他编程语言的跨语言 FUSE 性能优化.

参考文献

- [1] LI H, GHODOSI A, ZAHARIA M, et al. Tachyon: reliable, memory speed storage for cluster computing frameworks[C]//Proceedings of the 2014 ACM Symposium on Cloud Computing (SoCC). New York: ACM, 2014: 6.1-6.15.
- [2] 孙勇, 林菲, 等. 面向云计算的键值型分布式存储系统研究[J]. 电子学报, 2013, 41(7): 1406-1411.
SUN Y, LIN F, WANG B J. Study on the key-value distributed storage system for cloud computing[J]. Acta Electronica Sinica, 2013, 41(7): 1406-1411. (in Chinese)
- [3] GNU. POSIX (The Portable Operating System Interface) [EB/OL]. (2016-04-11)[2021-09-28]. https://www.gnu.org/software/libc/manual/html_node/POSIX.html.
- [4] Linux Kernel Development Community. FUSE - the linux kernel documentation[EB/OL]. (2021-08-23)[2021-09-28]. <https://www.kernel.org/doc/html/latest/filesystems/fuse.html>.
- [5] ALEXANDROV A D, IBEL M, SCHAUSER K E, et al. Extending the operating system at the user level: the UFO global file system[C]//Proceedings of the 1997 USENIX Annual Technical Conference (USENIX ATC). California: USENIX, 1997: 77-90.
- [6] MAZIERES D. A toolkit for user-level file systems[C]//

- Proceedings of the 2001 USENIX Annual Technical Conference (USENIX ATC). California: USENIX, 2001: 261-274.
- [7] WEIL S A, BRANDT S A, MILLER E L, et al. Ceph: a scalable, high-performance distributed file system[C]//Proceedings of the 7th Symposium on Operating Systems Design and Implementation (OSDI). California: USENIX, 2006: 307-320.
- [8] Red Hat, Inc. Gluster[EB/OL]. (2021-09-28)[2022-04-08]. <https://www.gluster.org/>.
- [9] VANGOOR B K R, TARASOV V, ZADOK E. To FUSE or not to FUSE: performance of user-space file systems [C]//Proceedings of the 15th USENIX Conference on File and Storage Technologies (FAST). California: USENIX, 2017: 59-72.
- [10] BIJLANI A, RAMACHANDRAN U. Extension framework for file systems in user space[C]//Proceedings of the 2019 USENIX Conference on USENIX Annual Technical Conference (USENIX ATC). California: USENIX, 2019: 121-134.
- [11] ISHIGURO S, MURAKAMI J, OYAMA Y, et al. Optimizing local file accesses for fuse-based distributed storage[C]//Proceedings of the 2012 SC Companion: High Performance Computing, Networking Storage and Analysis (SC). New York: IEEE, 2012: 760-765.
- [12] NARAYAN S, MEHTA R K, CHANDY J A. User space storage system stack modules with file level control[C]//Proceedings of the 12th Annual Linux Symposium in Ottawa. Ottawa: OLS, 2010: 189-196.
- [13] TSELOVALNIKOV S. JNR-FUSE[EB/OL]. (2021-04-30)[2021-09-28]. <https://github.com/SerCeMan/jnr-fuse>.
- [14] LEVART P. FUSE-J: java bindings for FUSE[EB/OL]. (2015-08-06) [2021-09-28]. <https://sourceforge.net/projects/fuse-j/>.
- [15] libfuse: The reference implementation of the linux FUSE interface[EB/OL]. (2021-09-06) [2021-09-28]. <https://github.com/libfuse/libfuse/>.
- [16] Oracle, Inc. Java native interface[EB/OL]. (2021-09-28) [2022-04-08]. <https://docs.oracle.com/javase/8/docs/technotes/guides/jni/>.
- [17] DOUBROVKINE D, BLÄSING M. Java native access (JNA)[EB/OL]. (2021-08-23)[2021-09-28]. <https://github.com/java-native-access/jna>.
- [18] SIEGER N, NUTTER C O, JENVEY P, et al. The Java native runtime project[EB/OL]. (2021-03-09) [2021-09-28] <https://github.com/jnr>.
- [19] WANG L, YE S, YANG B, et al. DIESEL: A dataset-based distributed storage and caching system for large-scale deep learning training[C]//Proceedings of the 49th International Conference on Parallel Processing (ICPP). New York: ACM, 2020: 20.1-20.11.
- [20] Oracle, Inc. The invocation API of JNI[EB/OL]. (2021-09-26)[2021-09-28]. <https://docs.oracle.com/javase/8/docs/technotes/guides/jni/spec/invocation.html>.
- [21] Oracle, Inc. JNI copy array[EB/OL]. (2021-09-06)[2021-09-28]. https://docs.oracle.com/javase/8/docs/technotes/guides/jni/spec/functions.html#Get_PrimitiveType_ArrayRegion_routines.
- [22] Oracle, Inc. DirectByteBuffer of Java[EB/OL]. (2021-09-21) [2021-09-28]. <https://docs.oracle.com/javase/8/docs/api/java/nio/ByteBuffer.html>.
- [23] Oracle, Inc. JNI critical copy array[EB/OL]. (2018-02-21) [2021-09-28]. https://docs.oracle.com/javase/8/docs/technotes/guides/jni/spec/functions.html#GetPrimitiveArrayCritical_ReleasePrimitiveArrayCritical.
- [24] Kubernetes: Production-grade container orchestration[EB/OL]. (2020-09-06)[2021-09-28]. <https://kubernetes.io/>.
- [25] Kubernetes container storage interface (CSI) documentation [EB/OL]. (2018-03-03) [2021-09-28]. <https://kubernetes-csi.github.io/docs/introduction.html>.
- [26] Alluxio, Inc. FUSE-based POSIX API[EB/OL]. (2021-03-09)[2021-09-28]. <https://docs.alluxio.io/os/user/stable/en/api/POSIX-API.html>.
- [27] NETZER Y, WANG T, COATES A, et al. Reading digits in natural images with unsupervised feature learning[EB/OL]. (2020-09-20)[2021-09-28]. <http://ufldl.stanford.edu/housenumbers>.
- [28] COCO: Common objects in context[EB/OL]. (2020-09-26)[2021-09-28]. <https://cocodataset.org/>.
- [29] PyTorch[EB/OL]. (2020-09-02) [2021-09-28]. <https://pytorch.org/>.
- [30] PyTorch. ImageNet training in PyTorch[EB/OL]. (2020-09-02)[2021-09-28]. <https://github.com/pytorch/examples/tree/master/imagenet>.
- [31] DENG J, DONG W, SOCHER R, et al. ImageNet: a large-scale hierarchical image database[C]//Proceedings of 2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. New York: IEEE, 2009: 248-255.
- [32] SERGEEV A, BALSO M D. Horovod: fast and easy distributed deep learning in TensorFlow[J/OL]. (2018-02-01) [2021-09-28]. <https://arxiv.org/abs/1802.05799>.

- [33] Tensorflow Horovod2 benchmark[EB/OL]. (2020-09-26) [2021-09-28]. https://github.com/horovod/horovod/blob/master/examples/tensorflow2/tensorflow2_synthetic_benchmark.py.
- [34] KRIZHEVSKY A. The CIFAR-10 dataset[EB/OL]. (2021-09-28) [2021-09-28]. <http://www.cs.toronto.edu/kriz/cifar.html>.



黄宜华 男,1962年出生,江苏泰州人. 博士、教授、博士生导师. 主要研究方向为大数据系统、自动化机器学习、文本分析处理技术.
E-mail: yhuang@nju.edu.cn

作者简介



顾荣 男,1988年出生,江苏泰州人. 博士,南京大学特聘研究员. 主要研究方向为大数据与云计算系统,分布式缓存与高效索引系统.
E-mail: gurong@nju.edu.cn



罗义力 男,1998年出生. 南京大学计算机科学与技术系硕士研究生. 主要研究方向为分布式存储系统.
E-mail: luoyl@smail.nju.edu.cn



仇伶俐 男,1994年出生. 南京大学计算机科学与技术系硕士研究生. 主要研究方向为深度学习系统.
E-mail: mp1933003@smail.nju.edu.cn



王肇康(通讯作者) 男,1990年出生,河南郑州人. 博士、讲师、硕士生导师. 主要研究方向为分布式图计算、分布式数据处理、云计算技术.
E-mail: wangzhaokang@nuaa.edu.cn



戴海鹏 男,1985年出生,湖南娄底人. 博士、副教授、博士生导师. 主要研究方向为高效数据索引、物联网.
E-mail: haipengdai@nju.edu.cn