

Lustre 文件系统性能评测与分析

陈四建¹ 吴庆波¹ 周恩强¹

¹(国防科学技术大学计算机学院 长沙 410073)

(hawk10242003@hotmail.com)

【摘要】 传统网络文件系统难以满足高性能计算集群系统的 I/O 需求, 基于对象存储的全局并行文件系统可以有效缓解传统文件系统在可扩展性、可用性和性能上的问题。本文针对 Lustre 文件系统, 与 NFS 文件系统进行对比测试。首先介绍了 Lustre 的结构及其优势; 然后从带宽、元数据处理性能和文件系统可扩展性对其进行测试, 并与 NFS 的测试结果进行比较分析, 说明 Lustre 具有带宽高, 吞吐率稳定, 可扩展性好的优点。最后, 针对 Lustre 的不足, 提出了改进意见。

【关键词】 Lustre; 文件系统; NFS; 性能测试

【中图法分类号】 TP302.7

Performance Evaluation and Analysis of Lustre File System

CHEN Si-Jian¹, WU Qing-Bo¹, ZHOU En-Qiang¹

¹(School of Computer, National University of Defense Technology, Changsha 410073)

Abstract Traditional network file systems can't reach the achievement of high-performance computing Cluster systems, but object-based global parallel file system—Lustre—avoids the problems of traditional file systems about scalability, availability, performance. First, this paper introduces the architecture and advantages of Lustre. Second, we test the bandwidth, performance of handling metadata, scalability of Lustre, and compare the test results with NFS in the same environment. Both theoretic analysis and the test results show that Lustre has more bandwidth, steadier throughout and better scalability. Finally, we give some improvements according to the pitfall of Lustre.

Key words Lustre; file system; NFS; performance evaluation

1 引言

随着高性能计算由传统的主机方式向网络化集群演变, 传统的基于主机的存储架构已成为新的瓶颈, 网络化存储的研究已成为近来的研究热点。SAN 系统由于具有高带宽、低延迟的优势, 在高性能计算中占有一席之地, 如 SGI 的 CXFS^[1] 文件系统就是基于 SAN 实现高性能文件存储的,

但是由于 SAN 系统的价格较高, 可扩展性较差, 已不能满足成千上万个 CPU 规模的系统。NAS 技术通过 TCP/IP 实现网络化存储, 可扩展性好、价格便宜、用户易管理, 如目前在集群计算中应用的较多的 NFS 文件系统, 但由于 NAS 的协议开销高、带宽低、延迟长, 不利于在高性能集群中的应用。目前国外已开始了新型文件系统的研究, 希望它能有效结合 SAN 和 NAS 系统的优点, 如

File System 公司的 Lustre、Panasas 公司的 ActiveScale^[2] 文件系统等。

Lustre^[3,4] 在美国能源部 (U. S. Department of Energy: DOE)、Lawrence Livermore 国家实验室、Los Alamos 国家实验室、Sandia 国家实验室、Pacific Northwest 国家实验室等处的 HP 系统中得到了初步的应用^[5]。HP 多年来一直主持由 Sandia、Los Alamos 与 Lawrence Livermore 三个国家实验室共同研发的 PathForward 计划。该计划的目的是让 Lustre 能够与集群数据系统与 Intel 结合。根据此计划, 未来 HP 也将在高性能计算集群市场上推出 Lustre 相关产品。

本文基于相同的硬件环境, 使用 Bonnie++^[6] 和 PostMark^[7], 对 NFS 文件系统和 Lustre 1.2.1 版本文件系统的带宽和吞吐率进行了测试; 同时, 测试 Lustre 使用多个存储服务器的可扩展性, 并对测试结果进行了分析。

2 Lustre 文件系统介绍

Lustre 文件系统是由 Cluster File Systems^[4] 公司开发的一个开源、高性能的分布式并行全局文件系统。它来源于卡耐基梅隆大学的 Coda 项目研究工作。Cluster File Systems 公司在 2003 年 12 月发布了 Lustre 1.0 版, 预计在 2005 年发布 2.0 版。目前 Lustre 1.0.4 版 for Linux 2.4.x 可以免费从 Cluster File Systems 公司主页下载, 但是需要向 Cluster File Systems 公司支付一定的费用才可以得到 1.0.4 的后续版本。

Lustre 针对大文件的读写作了优化, 可以为集群系统提供高性能的 I/O 吞吐率、全局数据共享环境、数据存储位置独立性和对节点失效提供冗余机制, 以及当集群重配置或者服务器和网络失效时的快速恢复服务, 较好地满足了高性能计算集群系统的需要。

Lustre 使用了基于对象的存储技术, 基于意图的分布式锁管理机制, 元数据和存储数据相分离的解决方案, 提供了一个全局命名空间, 并融合了传统分布式文件系统 (比如: AFS 和 Locus CFS) 的特色和传统共享存储集群文件系统 (比如: Zebra、Berkeley XFS、GPFS、Calypso、InfiniFile

和 GFS) 的设计思想, 消除了传统文件系统在可扩展性、可用性和性能上的问题^[8]。

Lustre 由客户端 (client)、存储服务器 (OST) 和元数据服务器 (MDS) 三个主要部分组成。Lustre 的客户端运行 Lustre 文件系统, 它和 OST 进行文件数据 I/O 的交互, 和 MDS 进行命名空间操作的交互。当 Client、OST 和 MDS 是分离的时候, Lustre 表现为一个有文件管理器的集群文件系统, 但是由于其设计的均匀性, 这些子系统可以运行在同一个系统中^[9]。其三个主要部分如图 1 所示。

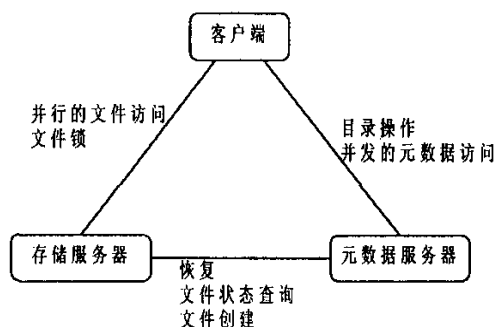


图1

Lustre 是一个透明的全局文件系统, 客户端可以透明地访问集群文件系统中的数据, 而无需知道这些数据的实际存储位置。客户端通过网络读取服务器上的数据, 存储服务器负责实际文件系统的读写操作以及存储设备的连接, 元数据服务器负责文件系统目录结构、文件权限和文件的扩展属性以及维护整个文件系统的数据一致性和响应客户端的请求。

Lustre 把文件当作由元数据服务器定位的对象, 元数据服务器指导实际的文件 I/O 请求到存储服务器, 存储服务器管理在基于对象的磁盘组上的物理存储^[6]。高速网络和硬盘技术的发展为集群文件系统的扩展提供了技术保证。Lustre 支持多种网络协议, 其逻辑结构如图 2 所示。

由于采用元数据和存储数据相分离的技术, 可以充分分离计算和存储资源。它使得客户端计算机可以专注于用户和应用程序的请求; 存储服务器和元数据服务器专注于读、传输和写数据。服务器端的数据备份和存储配置以及存储服务器扩充等操作不会影响到客户端, 存储服务器和元数据服务器均不会成为性能瓶颈^[9]。

Lustre 的全局命名空间为文件系统的所有客户端提供了一个有效的全局唯一的目录树,并将数据条块化,再把数据分配到各个存储服务器上面,提供了比传统 SAN 的“块共享”更为灵活的共享访问方式.全局目录树消除了在客户端的配置信息,并且在配置信息更新时仍然保持有效.同时,客户端请求自包含也减轻了 MDS 的负担.

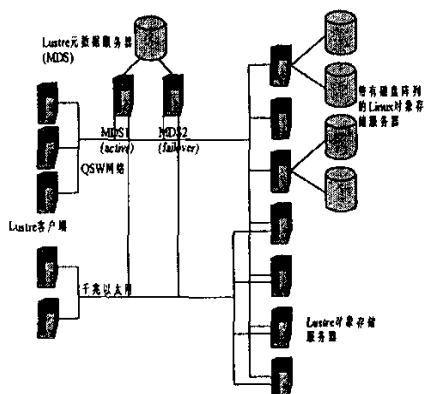


图 2

3 Lustre 和 NFS 的性能比较

3.1 测试环境

测试共使用了 9 台服务器,配置均为双 CPU 安腾 1.4G,内存 8G,使用的网络设备是千兆网卡和千兆交换机,操作系统安装的是 Red Hat Advanced Server 1A64 2.1 版^[10],使用的 Lustre 文件系统是 1.2.1 版.

使用 Bonnie++ 作为带宽测试工具.bonnie++ 是针对文件并发访问数量较少,但单个文件相对较大的应用而设计的.

为了消除文件系统缓存策略对系统性能的影响,每组块读/写测试中每个客户端读写的文件长度都是服务器内存大小的两倍,即 16G,带宽测试结果为聚集带宽.在硬盘的基准性能测试中使用 Ext3 本地文件系统,测试 16G 的文件,得到的结果为:块写带宽为 67476KB/s,块读带宽为 57831KB/s.

由于千兆交换机有 100MB/s 以上的通信带宽,而本地硬盘仅能提供 56~66MB/s 的数据传输率,

因此可以认为,所测得数据的差异能真实反映分布式文件系统的性能^[11].

3.2 传输带宽

使用同一台服务器分别作为 NFS 服务器和 Lustre 存储服务器,客户端数量从 1 增加到 5,测得的 NFS 和 Lustre 块读/写带宽结果如图 3、图 4 所示.横轴坐标表示客户端的数量,纵轴坐标表示读写带宽,单位是 KB/s.

NFS 的块读带宽随着客户端的数量的增长,聚集带宽下降明显;块写带宽在 3 个客户端时达到最大值,然后随着客户端的增多,聚集带宽下降.预计随着客户端的数量的增长,NFS 的块读和块写带宽将持续下降.

Lustre 的块读带宽和块写带宽随着客户端的数量的增长,聚集带宽一直在上升.可以预计随着客户端的数量的增长,块读和块写带宽将持续上升,达到一个最大值后开始持续的下降.

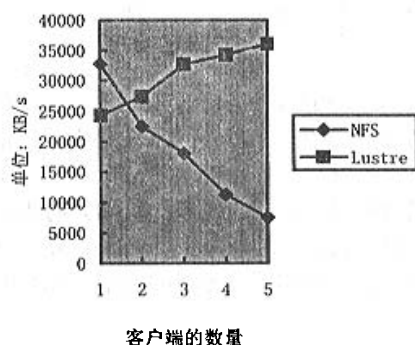


图 3 NFS 和 Lustre 的块读比较

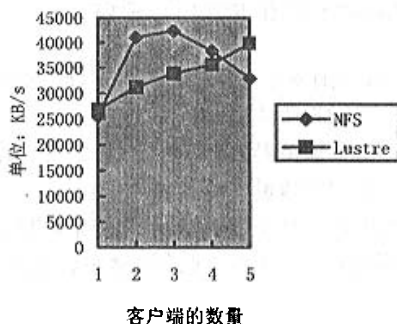


图 4 NFS 和 Lustre 的块写比较

3.3 小文件吞吐率

使用 PostMark 对 NFS 和 Lustre 的元数据处理能力进行测试。PostMark 模拟的是需要频繁、大量地存取小文件的应用。

Postmark 的测试原理是创建一个测试文件池, 文件的数量和最大、最小长度以及处理文件的事务可以设定。创建完成后, Postmark 对文件池进行一系列的事务操作, 根据从实际应用中统计的结果, 设定每一个事务包括一次创建或删除操作和一次读或添加操作, 文件系统的缓存策略可能对

性能造成影响, Postmark 可以通过对创建/删除以及读/添加操作的比例进行修改来抵消这种影响。事务操作进行完毕后, PostMark 对文件池进行删除操作, 并结束测试, 输出结果。Postmark 是用随机数来产生所操作文件的序号, 从而使测试更加贴近于现实应用^[7]。

输出结果中比较重要的输出数据包括测试总时间、每秒钟平均完成的事务数、平均每秒创建和删除的文件数, 以及读和写的平均传输速度。测试结果如表 1 所示。

表 1 NFS 和 Lustre 事务处理能力

文件数量(个)	事务数量(个)	事务处理 (个/秒)		文件创建 (个/秒)		文件删除 (个/秒)		读(KB/s)		写(KB/s)	
		NFS	Lustre	NFS	Lustre	NFS	Lustre	NFS	Lustre	NFS	Lustre
10,000	50,000	276	100	588	192	785	392	735.86	268.63	1055.79	385.41
20,000	50,000	256	100	434	192	747	388	550.42	224.87	1038.79	424.39
40,000	100,000	213	98	400	191	443	388	450.04	224.00	844.59	420.39

因为目前的 Lustre 元数据服务器并没有针对小文件进行优化, 所以 Lustre 的小文件处理能力比 NFS 差。但是由于 Lustre 采用了元数据和存储数据相分离的技术, 基于意图的分布式的锁机制, 从表中可以看出, 随着文件数量的增加, Lustre 的小文件处理能力并没有象 NFS 一样有明显下降, 可以预计 Lustre 能满足大量客户端同时读写大量文件的需求。

客户端可以并行地从多个存储服务器上读写。因此, 多个存储服务器的性能明显优于单个服务器的性能, 整体 I/O 性能随着存储服务器数量的增加而增强。

4 Lustre 的扩展性能测试

本测试针对 Lustre 的扩展性进行, 环境如下:

- 元数据服务器 1 台;
- 存储服务器数量分别为 1、2、3、4;
- 客户端数量分别为 1、2、3、4;

共进行了十六组测试, 测试结果如图 5、图 6 所示。横轴坐标表示存储服务器数量, 纵轴坐标表示块读/写带宽, 单位是 KB/s。

Lustre 使用了“stripe”^[12]技术, 把大文件分成“条块”分散到多个存储服务器上存储, 这样

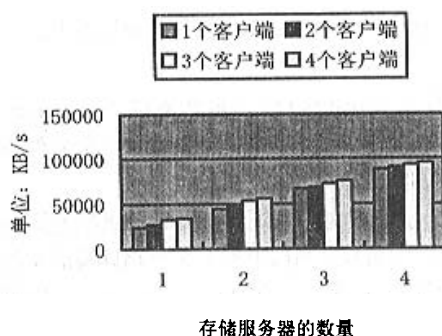


图 5 块读性能

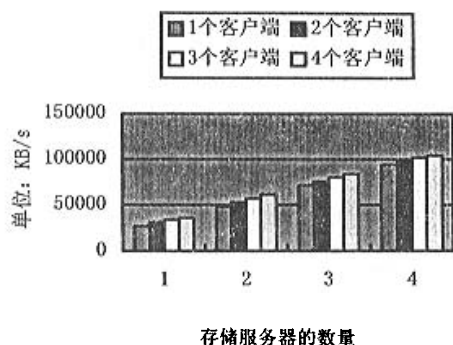


图6 块写性能

Lustre 的块读写性能随着客户端和存储服务器的增加基本上呈线性增长，表现出良好的可扩展性。Lustre 的块读性能一直低于块写性能，是由于在 8G 内存的机器上测试 16G 的文件时，写操作能够部分利用系统缓存，而读操作则无法利用。

由于每一个客户端测试都需要 16G 硬盘空间，本次测试受限于硬盘空间和服务器数量无法做进一步测试。在 GUPS 项目开发组针对 Lustre1.04 版的测试中，可以得知：在只有一台存储服务器的情况下，客户端数量为 9 时，块读和块写的聚集带宽达到最大值^[13]。因此可以预测随着存储服务器和客户端数量的增加，读写带宽会越来越大，在存储服务器和客户端的数量为大约为 1:9 时，系统的聚集带宽达到最大。

5 改进建议

Lustre 的出色设计使得它对大文件的读写性能优于 NFS 文件系统。Lustre 的元数据处理能力比 NFS 文件系统差，但是由于它采用了元数据和存储数据相分离的技术，基于意图的锁这一新颖的机制，随着吞吐率的增加，其元数据处理能力并没有明显的下降。

Lustre 的所有元数据由一个元数据服务器集中控制，这使得整体的可扩展性受到了限制，可以考虑通过实现元数据服务器集群、元数据服务器写回 cache 和设计更好的预取 cache 以及更为精巧的分布式的锁的方法，增强系统对元数据的处理能力；同时，如果存储服务器采用协作式 cache 将会提高系统整体的性能。

Lustre 文件系统在卸载时经常出现问题，设计更为方便的配置工具，进而开发图形化的配置工具，使得 Lustre 可以对网络的暂时失效很好的处理，可以象 InterMezzo^[14]文件系统一样具有断开连接下文件操作的能力，可以动态的增加新的存储服务器和删除需要淘汰的存储服务器，从而能够更方便的扩展系统以满足不断增长的应用的需求，这些都是需要进一步改进和研究的内容。

参考文献

- [1] <http://www.sgi.com/products/storage/cxfs.html>
- [2] <http://www.panasas.com/>
- [3] <http://www.lustre.org/>
- [4] <http://www.clusterfs.com/>
- [5] HP 为美国政府打造高质量 Linux 群集系统。
<http://www.linuxaid.com.cn/>
- [6] <http://www.coker.com.au/bonnie++/>
- [7] NAS: 文件服务器的替代者。
<http://www.pcworld.com.cn/>
- [8] Peter J. Braam. The Lustre Storage Architecture, Cluster File Systems, Inc. <http://www.lustre.org/docs/lustre.pdf> June 2, 2004
- [9] <http://www.lustre.org/docs/whitepaper.pdf>
- [10] <http://www.redhat.com/>
- [11] 曹立强, 熊劲. Global File System 性能评测与分析. 高性能计算技术. 12, 2003(165):2~3.
- [12] Philip Schwan. Lustre: Building a File System for 1,000-node Clusters. Cluster File Systems, Inc <http://www.lustre.org/docs/ols2003.pdf>
- [13] <http://hpcf.nersc.gov/projects/GUPFS/index.php>
- [14] <http://www.inter-mezzo.org/>

作者：[陈四建](#)，[吴庆波](#)，[周恩强](#)
作者单位：[国防科学技术大学计算机学院, 长沙, 410073](#)

相似文献(10条)

1. 期刊论文 [聂刚. 卿秀华. NIE Gang. QING Xiu-hua 基于对象存储的Lustre文件系统的研究 -信息技术2007, "" \(9\)](#)

随着高性能计算网络集群系统的高速发展, 传统的网络存储架构—网络附加存储NAS和存储区域网SAN已越来越不能够满足系统对存储性能的要求. 针对SAN和NAS的不足, 新一代的网络存储技术—基于对象存储OBS成为了研究的热点. 重点论述了基于对象存储的架构和特点, 并针对基于对象存储的Lustre文件系统进行了初步测试, 通过和传统的NFS文件系统对比, 分析了对象存储文件系统在可扩展性、性能、易用性和安全性等方面的优越性能.

2. 学位论文 [林松涛 基于Lustre文件系统的并行I/O技术研究 2004](#)

典型的并行程序都要进行数据量非常大的I/O访问, 但是这些I/O访问序列很多是由大量的小块数据的I/O访问组成. 这样, 虽然并行文件系统在大块数据传输上能够达到很大的I/O带宽, 但是由于并行程序传递给并行文件系统层次的调用是频繁的I/O启动操作和小块数据的I/O数据传输操作, 这样使得在实际中并行程序能够用到的I/O带宽远远低于并行文件系统所能提供的最大带宽.

一个重要的解决方法是在从上层到下层传递I/O请求时尽量的保持住应用程序的I/O访问模式的信息, 并且在I/O操作时尽量利用这些信息. 常见的实现方式是在并行文件系统上层再加入一个可以获取并行程序数据访问模式的层次—并行I/O库. I/O库获得并行应用程序的独立I/O的数据访问分布, 以及程序总体I/O行为的信息, 利用这些信息把很多多个数据块的I/O请求转化为对并行文件系统的少量的大块数据的实际I/O操作, 从而可以最大程度的利用并行I/O/文件系统所提供的带宽, 但是也将带来更多的I/O开销. 如果底层并行文件系统能直接提供对于这些小块数据访问的高性能接口和实现, 可以直接高效地读写小块数据, 那么即简化了并行I/O库的实现, 又将提高并行I/O的性能. 本文首先详细地描述了Lustre文件系统的I/O结构及其实现, 在研究基于Lustre文件系统的并行、I/O技术的基础上, 实现了细粒度的“DirectI/O”读写方式及接口, 并在MPI-IO中实现了基于这些接口的并行I/O方式.

本文还对改进后的Lustre文件系统的并行I/O模式的性能进行了读写带宽测试. 评测结果表明基于新的I/O接口的I/O模式将大大的提高并行I/O性能, 特别是访问结构化数据集的时候.

3. 期刊论文 [谢旻. 卢宇彤. 周恩强. 曹宏嘉. 杨学军. Xie Min. Lu Yutong. Zhou Enqiang. Cao Hongjia. Yang Xuejun](#)

[基于Lustre文件系统的MPI检查点系统实现技术与性能测试 -计算机研究与发展2007, 44 \(10\)](#)

基于协同式检查点的回滚恢复是在大规模并行计算机系统中得到采用的一项重要容错技术, 其性能开销主要为协同协议和检查点映像存储所决定. 描述了一个在EMPICH2中实现的应用透明的并行检查点系统, 相比已有的技术, 该系统有以下特点: 1) 协同协议操作利用了并行应用的近邻通信特性, 通过虚连接方法减少协议的处理开销; 2) 采用Lustre文件系统简化检查点映像文件管理的复杂性; 3) 通过并行I/O操作提高性能, 优化检查点映像的存储过程. 实际应用的测试表明, 该检查点系统具有较小的运行时间开销和良好的可扩展性.

4. 学位论文 [邹丹 基于对象存储的固态硬盘存储加速技术研究 2008](#)

随着用户对数据访问需求的增长, 传统的外部存储系统的结构以及磁盘的I/O延时限制了存储系统的性能. 为了缓解系统的I/O瓶颈, 一方面需要新的存储设备, 另一方面需要新的存储结构.

在存储设备方面, 高性能存储设备固态硬盘(Solid State Disk, SSD)逐渐成为关注的热点. 存储结构方面, 对象存储具有传统存储结构难以比拟的优势. 在这种背景下, 本文对固态硬盘设备和对象存储进行了系统研究, 设计了基于对象存储的固态硬盘存储加速系统.

本文的研究工作主要包括以下几个方面:

(1) 研究了闪存型固态硬盘和DRAM型固态硬盘的基本原理、组成结构、存储特性和应用现状, 分析对比了传统的磁盘和固态硬盘系统的性能;

(2) 研究了存储结构的演变, 重点研究了对象存储的基本理论和基于对象存储的文件系统;

(3) 以Lustre文件系统为基础, 针对不同存储设备的性能特性, 利用对象粒度的灵活性以及对象存储接口丰富的功能特性, 设计了可应用于基于对象存储的Lustre文件系统的固态硬盘存储加速系统, 将对象迁移到不同的OST上, 以提高系统性能;

(4) 研究了I/O访问模式特点、多专家系统原理及决策算法, 分析了传统文件Cache的替换及预取策略. 针对实际应用中I/O访问模式的多样性和变化性, 设计了基于多专家决策的对象调度算法;

(5) 实现了基于多专家对象调度策略的对象Cache原型, 集成了基于FIFO、LRU、LFU、MRU替换算法的替换专家模块, 在不同类型的I/O负载下进行了测试, 证明了多专家对象调度算法在各种负载下的性能均能接近或超过最优算法, 具有较强的自适应性;

(6) 实现了混合型OST应用模型的原型系统, 在随机I/O负载下分别对各种混合型OST的性能进行了分析比较.

本文研究工作中一些设计思想和关键技术, 对其它新型外部存储系统的研究具有参考价值.

5. 期刊论文 [张媛. 卢泽新. 刘亚萍. ZHANG Yuan. LU Zexin. LIU Yaping NFS over Lustre性能评测与分析 -计算机工程2007, 33 \(10\)](#)

传统的网络文件系统难以满足高性能计算系统的I/O需求. 基于对象存储的全局并行文件系统Lustre可以有效地解决传统文件系统在可扩展性、可用性和性能上存在的问题. 该文介绍了Lustre文件系统的结构及其优势, 对NFS over Lustre进行了性能测试, 并将测试结果与Lustre文件系统、NFS网络文件系统及本地磁盘Ext3文件系统的性能进行了比较分析, 给出了性能差异的原因, 提出了一种可行的解决方法.

6. 会议论文 [林松涛. 周恩强. 廖湘科 Lustre文件系统I/O性能的分析 and 改进](#)

Lustre是一种高性能的分布式文件系统, 缓解了传统文件系统在可扩展性、可用性和性能上的问题. 文章介绍了Lustre的结构和采用的关键技术, 并对其I/O性能进行详细的测试, 分析了影响系统I/O性能的关键因素; 最后根据Lustre在I/O上性能的问题, 提出了改进方案, 并给出了测试数据.

7. 学位论文 [李柱 分布式文件系统小文件性能优化技术研究 with 实现 2008](#)

分布式文件系统以其高可靠性、高可扩展性以及高性能和高性价比成为高性能计算平台存储系统的首选, 已经在军事技术、天气预报等环境中得到广泛应用. 相比其它文件系统, 它具有两个特点: 一是通过数据的分布存储, 来提供更大的存储空间, 并利用并行的I/O服务模式提供更高的I/O带宽; 二是通过使用各种新颖的分布式存储体系结构, 来为应用程序提供更丰富的I/O模式. 比如通过使用对象存储技术, 为应用程序提供面向对象的数据存储格式, 并提供Peta级大小的存储空间.

Lustre是典型的基于对象存储体系结构的并行文件系统, 它起源于卡耐基梅隆大学的Coda项目研究工作, 已经成为当前高性能计算领域使用最广泛的并行文件系统之一. Lustre具有良好的大文件I/O性能, 但是由于Lustre使用分布式的存储体系结构, 文件元数据和数据分开存储, 它的小文件I/O性能低下, 甚至不如本地文件系统. 本文以Lustre为具体研究对象, 通过研究Lustre的存储体系结构和实现原理, 在Lustre的OST组件中设计并实现了一种分布独立式的小文件Cache结构: Filter Cache. 该方法通过扩展Lustre的OST端的数据通路, 在原有数据通路的基础上, 增加对小对象I/O的缓存措施, 以此来改善Lustre的小文件性能. 测试表明: 使用Filter Cache方法之后, Lustre的小文件I/O性能得到了很好的改善, 在Cache资源全命中时, 读性能最大能够提高65%.

命中率 and 访问延迟是Cache系统中最重要的两个指标. 本文研究了Cache技术的设计思想和实现技术, 设计了对Filter Cache方法的优化方案. 优化方案主要针对方法使用的资源结构、Cache置换算法和Cache读写流程. 本文下一步工作将进一步完善这些优化措施的设计, 并进行实现.

最后, 本文对分布式文件系统中的另一种Cache结构: 协作—对象Cache进行了研究, 详细介绍了其特点和实现, 对比了该Cache结构和Filter Cache方法的不同点, 根据它的优点提出了两点对Filter Cache方法的改进思想.

8. 期刊论文 [李柱, 周恩强, 廖湘科, Li Zhu, Zhou Enqiang, Liao Xiangke Filter Cache:一种提高Lustre I/O性能的方法](#)

[-计算机研究与发展](#)2009, 46(z2)

高性能计算系统需要一个可靠高效的并行文件系统. Lustre集群文件系统是典型的基于对象存储的集群文件系统, 它适合大数据量聚合I/O操作. 大文件I/O操作能够达到很高的带宽, 但是小文件I/O性能低下. 针对导致Lustre的设计中不利于小文件I/O操作的两个方面, 提出了Filter Cache方法. 在Lustre的OST组件中设计一个存放小文件I/O数据的Cache, 让OST端的小文件I/O操作异步进行, 以此来减少用户感知的小文件I/O操作完成的时间, 提高小文件I/O操作的性能.

9. 会议论文 [谢旻, 万国伟, 卢宇彤, 周恩强, 曹宏嘉 基于Lustre文件系统的MPI检查点系统实现技术与性能测试](#) 2007

基于协同式检查点的回卷恢复是在大规模并行计算机系统中得到采用的一项重要容错技术, 其性能开销主要为协同协议和检查点映像存储所决定. 描述了一个在MPICH2中实现的应用透明的并行检查点系统, 相比已有的技术, 该系统有以下特点: (1) 协同协议操作利用了并行应用的近邻通讯特性, 通过虚连接方法减少协议的处理开销; (2) 采用Lustre文件系统简化检查点映像文件管理的复杂性; (3) 通过并行I/O操作提高性能, 优化检查点映像的存储过程. 实际应用的测试表明, 该检查点系统具有较小的运行时间开销和良好的可扩展性.

10. 会议论文 [杨昕, 沈文海 Lustre并行文件系统的发展及在气象领域的应用前景](#) 2006

Lustre是一个在Linux集群环境中广为应用的并行文件系统. 本文介绍Lustre相关技术, 并针对气象应用的特点和国家气象信息中心的现状和发展, 设计一个基于Lustre的多集群全局文件系统整合的模型, 通过分析, 我们认为这一模型可以为气象应用提供高效、灵活和统一的在线存储资源, 也能很好地满足总体计算环境扩展的需求.

本文链接: http://d.g.wanfangdata.com.cn/Conference_6036255.aspx

授权使用: 中科院计算所(zkyjsc), 授权号: e869773d-42a5-4dbf-b0e3-9e400107c147

下载时间: 2010年12月2日