linux 零拷贝技术 选例bu间的给loper 零拷贝的几种方法 原始数据拷贝操作 在介绍之前,先看看 Linux 原始的数据拷贝操作 是怎样的。如下图,假如一个应用需要从某个磁 盘文件中读取内容通过网络发出去,像这样: while((n = read(diskfd, buf, BUF -SIZE)) > 0)write(sockfd, buf , n); 那么整个过程就需要经历:1)read 将数据从磁

盘文件通过 DMA 等方式拷贝到内核开辟的缓冲

区; 2)数据从内核缓冲区复制到用户态缓冲区;

3)write 将数据从用户态缓冲区复制到内核协议

栈开辟的 socket 缓冲区; 4) 数据从 socket 缓

User 缓存

可见,整个过程发生了至少四次数据拷贝,其中

两次是 DMA 与硬件通讯来完成, CPU 不直接参

与,去掉这两次,仍然有两次 CPU 数据拷贝操

User 缓存

1) 这种方法只能适用于那些不需要内核缓冲区处

理的应用程序,这些应用程序通常在进程地址空

间有自己的数据缓存机制,称为自缓存应用程

2) 这种方法直接操作磁盘 I/O, 由于 CPU 和磁

盘 I/O 之间的执行时间差距, 会造成资源的浪

序,如数据库管理系统就是一个代表。

Сору

Сору

Socket 缓存

TactoudDeveloper

方法一: 用户态直接 I/O

CPU Copy

Socket 缓存

DMA Copy

网络

🤼 aCloudDeveloper

冲区通过 DMA 拷贝到网卡上发出去。

CPU Copy

Kernel 缓存

DMA Copy

磁盘文件

再谈Linux零拷贝技术

2019-08-04

篇类似的,看这里。

为什么需要零拷贝

数据包整个处理流程的57.1%。

什么是零拷贝

基于这个思想去做的优化。

释。

CloudDeveloper Linux云计算网络

PS: 其实这篇文章其实是我很早的原创(当时

公众号还叫aCloudDeveloper),但是之前的

文章删了, 但是网上却被流传开来, 有些还直

接标原创,有些把作者弄错,这里我重新发出

来,做一个说明。另外,我前两天还转载过一

本文借鉴并总结了几种比较常见的 Linux 下的零

拷贝技术,相关的引用链接见文后,大家如果觉

得本文总结得太抽象,可以转到链接看详细解

传统的 Linux 系统的标准 I/O 接口(read、

write) 是 基 于 数 据 拷 贝 的 , 也 就 是 数 据 都 是

copy_to_user 或者 copy_from_user, 这样做的

好处是,通过中间缓存的机制,减少磁盘 I/O 的

操作,但是坏处也很明显,大量数据的拷贝,用

户态和内核态的频繁切换,会消耗大量的 CPU 资

源,严重影响数据传输的性能,有数据表明,在

Linux内核协议栈中,这个拷贝的耗时甚至占到了

零拷贝就是这个问题的一个解决方案,通过尽量

避免拷贝操作来缓解 CPU 的压力。Linux 下常见

的零拷贝技术可以分为两大类:一是针对特定场

景,去掉不必要的拷贝;二是去优化整个拷贝的

过程。由此看来,零拷贝并没有真正做到"0"拷

贝,它更多是一种思想,很多的零拷贝技术都是

这种方法可以使应用程序或者运行在用户态下的 库函数直接访问硬件设备,数据直接跨过内核进 行传输,内核在整个数据传输过程除了会进行必 要的虚拟存储配置工作之外,不参与其他任何工 作,这种方式能够直接绕过内核,极大提高了性 能。

User

Kernel

Hardware

缺陷:

Kernel 缓存

磁盘文件

作。

User

Kernel

Hardware

费,解决这个问题需要和异步 I/O 结合使用。 方法二: mmap 这种方法,使用 mmap 来代替 read,可以减少 一次拷贝操作,如下: buf = mmap(diskfd, len); write(sockfd, buf, len); 应用程序调用 mmap, 磁盘文件中的数据通过 DMA 拷贝到内核缓冲区,接着操作系统会将这个

缓冲区与应用程序共享,这样就不用往用户空间

拷贝。应用程序调用write ,操作系统直接将数据

从内核缓冲区拷贝到 socket 缓冲区,最后再通

1) mmap 隐藏着一个陷阱, 当 mmap 一个文件

时,如果这个文件被另一个进程所截获,那么

write 系统调用会因为访问非法地址被 SIGBUS

信号终止,SIGBUS 默认会杀死进程并产生一个

coredump,如果服务器被这样终止了,那损失

破坏文件,这样 write 在被 SIGBUS 杀死之前,

会被中断,返回已经写入的字节数,并设置 er-

通常的做法是在 mmap 之前加锁,操作完之后解

从Linux 2.1版内核开始,Linux引入了sendfile,

—CPU Copy—→ Socket 缓存

DMA Copy

网络

🤼 aCloudDeveloper

过 DMA 拷贝到网卡发出去。

Kernel 缓存

DMA Copy

磁盘文件

MMAP Shared

User

Kernel

Hardware

缺陷:

解决这个问题通常使用文件的租借锁:首先为文 件申请一个租借锁, 当其他进程想要截断这个文 件时, 内核会发送一个实时的 RT_SIG-NAL_LEASE 信号,告诉当前进程有进程在试图

锁:

就可能不小了。

rno 为 success。

方法三: sendfile

#include<sys/sendfile.h>

ssize_t sendfile(int out-

也能减少一次拷贝。

Kernel

Hardware

缺陷:

序。

fd, int in fd, off t *offset, size t count); sendfile 是只发生在内核态的数据传输接口,没 有用户态的参与,自然避免了用户态数据拷贝。 它指定在 in_fd 和 out_fd 之间传输数据, 其中, 它规定 in_fd 指向的文件必须是可以 mmap 的, out_fd 必须指向一个套接字,也就是规定数据只 能从文件传输到套接字,反之则不行。sendfile 不存在像 mmap 时文件被截获的情况,它自带异 常处理机制。 User User 缓存

-CPU Copy-

1) 只能适用于那些不需要用户态处理的应用程

常规 sendfile 还有一次内核态的拷贝操作,能不

这种方法借助硬件的帮助, 在数据从内核缓冲区

到 socket 缓冲区这一步操作上,并不是拷贝数

据,而是拷贝缓冲区描述符,待完成后,DMA 引

擎直接将数据从内核缓冲区拷贝到协议引擎中

User 缓存

DMA Copy

1) 除了3.4 中的缺陷,还需要硬件以及驱动程序

splice 去掉 sendfile 的使用范围限制,可以用于

/* See fea-

2)只适用于将数据从文件拷贝到套接字上。

Socket 缓存

i aCtoudDeveloper

方法四: DMA 辅助的 sendfile

Kernel 缓存

DMA Copy

磁盘文件

能也把这次拷贝给去掉呢?

去,避免了最后一次拷贝。

Kernel 缓存

DMA Copy_

磁盘文件

方法五: splice

#define _GNU_-

SOURCE

任意两个文件描述符中传输数据。

ture_test_macros(7) */

ssize_t splice(int fd_in, lof-

f_t *off_in, int fd_out, lof-

*off out, size t len, un-

但是 splice 也有局限,它使用了 Linux 的管道缓

冲机制,所以,它的两个文件描述符参数中至少

splice 提供了一种流控制的机制,通过预先定义

的水印(watermark)来阻塞写请求,有实验表

明,利用这种方法将数据从一个磁盘传输到另外

一个磁盘会增加 30%-70% 的吞吐量,CPU负责

在某些情况下,内核缓冲区可能被多个进程所共

享,如果某个进程想要这个共享区进行 write 操

作,由于 write 不提供任何的锁操作,那么就会

对共享区中的数据造成破坏,写时复制就是 Lin-

写时复制,就是当多个进程共享同一块数据时,

如果其中一个进程需要对这份数据进行修改,那

么就需要将其拷贝到自己的进程地址空间中,这

样做并不影响其他进程对这块数据的操作,每个

这种方法完全改写 I/O 操作,因为传统 I/O 接口

都是基于数据拷贝的,要避免拷贝,就去掉原先

的那套接口,重新改写,所以这种方法是比较全

面的零拷贝技术,目前比较成熟的一个方案是最

先在 Solaris 上实现的 fbuf (Fast Buffer, 快速

Fbuf 的思想是每个进程都维护着一个缓冲区池,

这个缓冲区池能被同时映射到程序地址空间和内

核地址空间,内核和用户共享这个缓冲区池,这

fbufs 管理器

1/0子系统

1) 管理共享缓冲区池需要应用程序、网络软件、

Netmap 基于共享内存的思想,是一个高性能收

发原始数据包的框架,由Luigi Rizzo 等人开发完

成,其包含了内核模块以及用户态库函数。其目

标是,不修改现有操作系统软件以及不需要特殊

硬件支持, 实现用户态和网卡之间数据包的高性

缓冲区池A

缓冲区池B

程序B

地址空间

磁盘

文件

ুল: aGloudDeveloper

1) 同样只适用于不需要用户态处理的程序

2) 传输描述符至少有一个是管道设备。

方法六:写时复制

ux 引入来保护数据的。

#include <fcntl.h>

signed int flags);

有一个必须是管道设备。

也会减少一半。

缺陷:

User

Kernel

Hardware

缺陷:

支持。

答案就是这种 DMA 辅助的 sendfile。

Socket 缓存

Copy

TactoudDeveloper

DMA

进程要修改的时候才会进行拷贝,所以叫写时拷 贝。这种方法在某种程度上能够降低系统开销, 如果某个进程永远不会对所访问的数据进行更 改,那么也就永远不需要拷贝。 缺陷: 需要 MMU 的支持,MMU 需要知道进程地址空 间中哪些页面是只读的,当需要往这些页面写数 据时,发出一个异常给操作系统内核,内核会分 配新的存储空间来供写入的需求。 方法七:缓冲区共享

缓冲区)。

样就避免了拷贝。

应用程序

地址空间

操作系统内核

存储

程序A

以及设备驱动程序之间的紧密合作 2) 改写 API, 尚处于试验阶段。 高性能网络 I/O 框架——netmap

能传递。

传统系统中网络数据包映射关系

缺陷:

NETMAP数据结构 🔆 aCloudDeveloper netmap_ring,进行数据包的获取发送。 总结 1、零拷贝本质上体现了一种优化的思想 file, splice, 缓冲区共享, 写时复制......

在 Netmap 框架下,内核拥有数据包池,发送环 \接收环上的数据包不需要动态申请,有数据到达 网卡时, 当有数据到达后, 直接从数据包池中取 出一个数据包, 然后将数据放入此数据包中, 再 将数据包的描述符放入接收环中。内核中的数据 包池,通过 mmap 技术映射到用户空间。用户态 程序最终通过 netmap_if 获取接收发送环 2、直接 I/O, mmap, sendfile, DMA send-后台回复"<mark>加群</mark>",带你进入高手如云交流群 推荐阅读: <u>Linux 下实践 VxLAN</u> 10 个实战与面试【常用 Shell 脚本】编写 DockerFile 命令总结 <u>10个小技巧提高 Kubernetes 容器效率</u> 喜欢,就给我一个"在看", Linux云计算网络

云计算 | 网络 | Linux | 干货

获取学习大礼包后台 回复"1024"

加群交流后台回复"加群"

10T 技术资源大放送!包括但不限于:云计算、虚拟

化、微服务、大数据、网络、Linux、Docker、Ku-

bernetes、Python、Go、C/C++、Shell、PPT 等。

在公众号内回复「1024」,即可免费获取!!