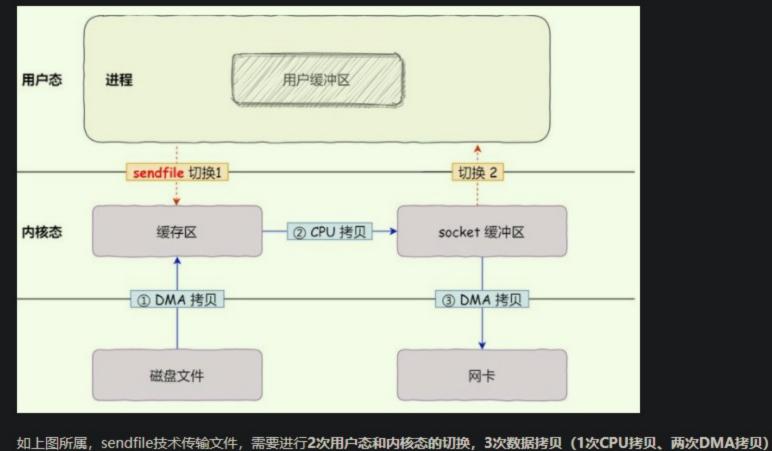


相对于传统数据传输,mmap减少了一次CPU拷贝,其具体过程如下: 1. 应用进程调用 mmap() , DMA 会把磁盘的数据拷贝到内核的缓冲区里,应用进程跟操作系统内核「共享」这个缓冲区

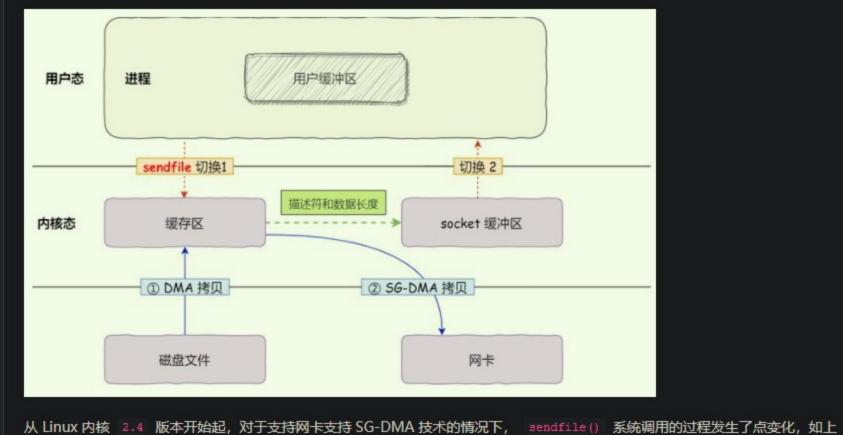
- 2. 应用进程再调用 write(),操作系统直接将内核缓冲区的数据拷贝到 socket 缓冲区中,这一切都发生在内核态,由 CPU 来搬 3. 最后,把内核的 socket 缓冲区里的数据,拷贝到网卡的缓冲区里,这个过程是由 DMA 搬运的
- 显然仅仅减少一次数据拷贝,依然难以满足要求

sendfile



相对于mmap, 其又减少了两次上下文的切换, 具体过程如下:

- 1. 应用调用sendfile接口,传入文件描述符,应用程序切换至内核态,并通过 DMA 将磁盘上的数据拷贝到内核缓冲区中 2. CPU将缓冲区数据拷贝至Socket缓冲区
- 3. DMA将数据拷贝到网卡的缓冲区里,应用程序切换至用户态 sendfile其实是将原来的两步读写操作进行了合并,从而减少了2次上下文的切换,但其仍然不是真正意义上的"零"拷贝
- sendfile + SG-DMA



贝) 具体过程如下:

图所示, sendfile + SG-DMA技术传输文件, 需要进行2次用户态和内核态的切换, 2次数据拷贝(1次DMA拷贝, 1次SG-DMA拷

- 1. 通过 DMA 将磁盘上的数据拷贝到内核缓冲区里; 2. 缓冲区描述符和数据长度传到 socket 缓冲区,这样网卡的 SG-DMA 控制器就可以直接将内核缓存中的数据拷贝到网卡的缓冲
- 区里,此过程不需要将数据从操作系统内核缓冲区拷贝到 socket 缓冲区中,这样就减少了一次数据拷贝; 此种方式对比之前的,真正意义上去除了CPU拷贝,CPU 的高速缓存再不会被污染了,CPU 可以去执行其他的业务计算任务,同时

和 DMA 的 I/O 任务并行,极大地提升系统性能。 但他的劣势也很明显, 强依赖于硬件的支持

splice Linux 在 2.6.17 版本引入 splice 系统调用,不再需要硬件支持,同时还实现了两个文件描述符之间的数据零拷贝。

splice 系统调用可以在内核空间的读缓冲区(read buffer)和网络缓冲区(socket buffer)之间建立管道(pipeline),从而避免了 用户缓冲区和Socket缓冲区的 CPU 拷贝操作。

基于 splice 系统调用的零拷贝方式,整个拷贝过程会发生 2次用户态和内核态的切换,2次数据拷贝(2次DMA拷贝),具体过程如

- 1. 用户进程通过 splice() 函数向内核(kernel)发起系统调用,上下文从用户态(user space)切换为内核态(kernel space) . 2. CPU 利用 DMA 控制器将数据从主存或硬盘拷贝到内核空间(kernel space)的读缓冲区(read buffer)。
- 3. CPU 在内核空间的读缓冲区 (read buffer) 和网络缓冲区 (socket buffer) 之间建立管道 (pipeline)。 4. CPU 利用 DMA 控制器将数据从网络缓冲区 (socket buffer) 拷贝到网卡进行数据传输。
- 5. 上下文从内核态(kernel space)切换回用户态(user space),splice 系统调用执行返回。
- 件描述符中传输数据,但是它的两个文件描述符参数中有一个必须是管道设备

splice 拷贝方式也同样存在用户程序不能对数据进行修改的问题。除此之外,它使用了 Linux 的管道缓冲机制,可以用于任意两个文

总结 本文简单介绍了 Linux 中的几种 Zero-copy 技术,随着技术的不断发展,又出现了诸如:写时复制、共享缓冲等技术,本文就不再

赘述。 广义的来讲,Linux 的 Zero-copy 技术可以归纳成以下三大类:

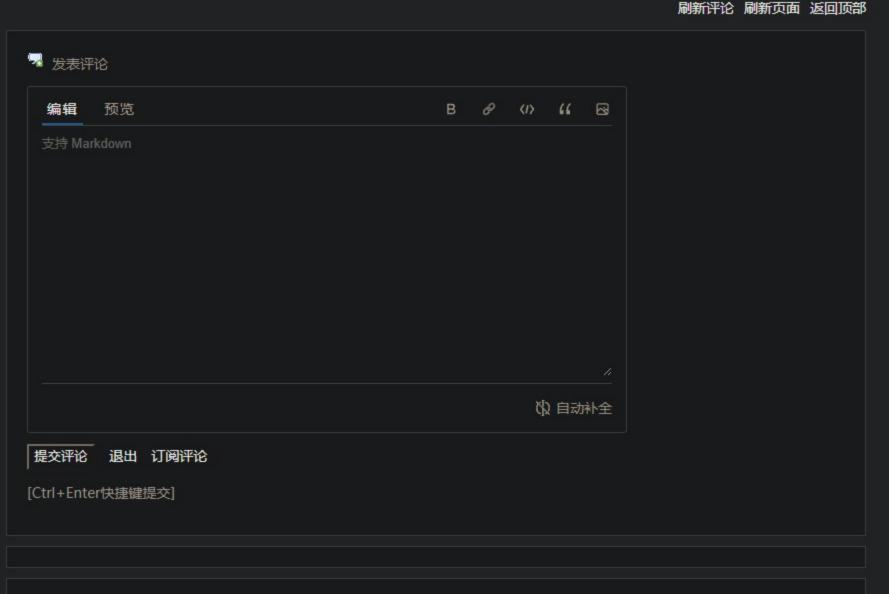
• 减少甚至避免用户空间和内核空间之间的数据拷贝:在一些场景下,用户进程在数据传输过程中并不需要对数据进行访问和处 理,那么数据在 Linux 的 Page Cache 和用户进程的缓冲区之间的传输就完全可以避免,让数据拷贝完全在内核里进行,甚

- 至可以通过更巧妙的方式避免在内核里的数据拷贝。这一类实现一般是是通过增加新的系统调用来完成的,比如 Linux 中的 mmap(), sendfile()以及 splice()等。 • 绕过内核的直接 I/O: 允许在用户态进程绕过内核直接和硬件进行数据传输,内核在传输过程中只负责—些管理和辅助的工 作。这种方式其实和第一种有点类似,也是试图避免用户空间和内核空间之间的数据传输,只是第一种方式是把数据传输过程
- 放在内核态完成,而这种方式则是直接绕过内核和硬件通信,效果类似但原理完全不同。 • 内核缓冲区和用户缓冲区之间的传输优化:这种方式侧重于在用户进程的缓冲区和操作系统的页缓存之间的 CPU 拷贝的优 化。这种方法延续了以往那种传统的通信方式,但更灵活。

分类: IO, Linux 标签: Java, Linux

本文来自博客园,作者:十三,转载请注明原文链接: https://www.cnblogs.com/hystrix/p/15213700.html





编辑推荐: · 五个维度打造研发管理体系 ·不会SQL也能做数据分析?浅谈语义解析领域的机会与挑战 ·Spring IoC Container 原理解析 前端实现的浏览器端扫码功能 · ASP.NET Core Filter 与 IOC 的羁绊 最新新闻: ·CV进入三维时代! Facebook在ICCV 2021 发布两个3D模型,自监督才是终极答案? (2021-10-20 14:22)

· 「老司机」盖茨再上热搜!被曝婚外情史将近20年,公然给女下属发邮件调情(2021-10-20 14:11) ·苹果M1「徒有其表」?「地表最强」芯只能剪视频引知乎热议(2021-10-20 14:00) · 物理学家获得至今对中子寿命的最精确测量 (2021-10-20 13:45) ·NeurIPS 2021论文公开评审被删除,Reddit网友猜测:观众太蠢,我决定隐藏讨论(2021-10-20 13:33) » 更多新闻...

1. Nacos 自动更新配置不生效问题(689) 2. 深入理解零拷贝技术(358) 3. Netty入门 (—): ByteBuf(289) 4. Netty入门(三): EventLoop(190) 5. Java基础 (一): I/O多路复用模型及Lin ux中的应用(89)

1. Netty入门(三): EventLoop(1) 推荐排行榜 1. 深入理解零拷贝技术(1)

评论排行榜

观,我也想学一下

2. Netty入门(三): EventLoop(1) 最新评论 1. Re:Netty入门(三): EventLoop ###请问博主这个排版是markdown, 很美

--懂得了才能做一些改变