

Feeling odd

Using new OS and CPU features to speed up large file copying

Sebastian Kraemer *stealth [at] openwall net*

June 30, 2012

Abstract

Recent multi-core CPU setups and operating system features such as pipe data splicing offer interesting new techniques to speed up the copy process of large files. These techniques are currently unused in common file copy utilities, such as `dd` or `cp`.

Search-Engine-Tag: SET-feeling-odd-2012.

1 Introduction

The backup process of large data or device files can take serious amount of time. Commonly the `dd` tool [1] is used to copy large files or block devices, since it offers some flexibility such as the specification of the block size. In this paper I will describe a optimized version `odd` for Linux, which uses the following features that appeared during the last years:

- the `splice()` system call
- a variant of using `splice()` with two CPU cores
- the `mmap()` system call
- the `sendfile()` system call

These optimizations must be compared against a common `read()/write()` loop, assuming the best block size that can be used to copy files. The constant overhead, such as the time to `open()` or `close()` the files is neglectable if files of Gigabyte or Terabyte size are copied. To avoid page caching effects, the amount of copied data is large compared to the size of available RAM. Therefore, cached blocks will be discarded if a LRU algorithm is assumed. Nevertheless, the page cache has been dropped between each copy operation by issuing the

```
echo 1 > /proc/sys/vm/drop_caches
```

command. The effect of inode caches is neglectable since just one file is copied and can be even more reduced by calling `stat()` before any benchmark.

Speedup results against `dd` are shown in the last section.

2 The *splice()* system call

The *splice* system call "allows to move data between two file descriptors without copying between kernel address space and user address space". [2] Using a technique that avoids copying of data blocks between kernel and user space is one way to achieve a speedup. Even if probably not measureable for a single *read()/write()* cycle, for Terabytes of data this saves a lot of *memcpy()* operations. The buffer that is used to hold the copied data blocks belongs to a pipe which needs to be created before calling *splice()*. This buffer is located inside the kernel and that's why no kernel-to-user-space copies are needed. We will see later how this pipe can be tuned to achieve more throughput.

As seen in appendix A.1, two *splice()* system calls are used to transfer one data block from one file to another. So, this technique equals the amount of system calls *dd* will be using, assuming the same block size.

The second technique to speed up file copying is to reduce the number of syscalls inside the copy loop and will be described in the next sections.

3 Using *splice()* with two CPU cores

The use of *splice()* alone often does not achieve much more performance than *dd* if using the same block size, since caching or scheduling effects throw away any performance gain of the zero-copy.

Even today's commodity PC hardware is equipped with multiple CPU cores. This allows to split the two *splice()* system calls to different cores, keeping the advantage of zero-copy between userland to kernel and back, as happening with a *read()/write()* cycle. That way, we effectively reduced the amount of syscalls inside the loop to 50% (one syscall vs. two). This technique is implemented in *odd* [7] using *fork()* and *sched_setaffinity()* to bind each process to its own CPU core. The speedup result can be found in Figure 1 and 2. Using more than two CPU cores does not make sense with this algorithm.

The experiment will also show that using *splice()* with cores without increasing the kernel pipe buffer size will give no good speedup. Since Linux kernel 2.6.35 [5] the `/proc/sys/fs/pipe-max-size` file controls the maximum possible pipe buffer size. Each process can then set its own pipe buffer size via *fcntl(fd, F_SETPIPE_SZ, value)*; which must not exceed that limit. Increasing the max pipe buffer size to 512MB or above shows really good results, whereas using statically block size of 4096 and/or a system that only has the same size for kernel pipe buffers will even be slower than normal *dd* operations. This is because the two processes are blocking each other too often due to a full or empty pipe.

Following good UNIX-tools practise, *odd* will not set the value inside `/proc/sys/fs/pipe-max-size`. This must be done by the administrator, but *odd* will try to increase its own pipe buffer size via the *fcntl()* call.

4 Using *mmap()*

An entirely different technique can be implemented by using the *mmap()* system call.

Instead of calling *read()/write()* on data of block size B , a target file with a *file hole* [6] of the same size F as the source file is created via *ftruncate()*. This allows us to *mmap()* a writable block of M bytes and call *read()* on the mapped address in chunks of size B . This reduces the amount of system calls used, if just a single CPU core is available.

$$S(M, B) = 2\frac{F}{M} + \frac{F}{M} \frac{M}{B} = 2\frac{F}{M} + \frac{F}{B} \quad (1)$$

$$S(B) = 2\frac{F}{B} \quad (2)$$

Formula 1 shows the amount of system calls used with the *mmap()* technique. For large *mmap* blocks M this comes closer to the optimum of F/B system calls (one call for each transferred block). Formula 2 shows the number of system calls for a *read()/write()* copy. A simplified version of the copy algorithm using *mmap()* can be found in appendix A.2.

5 Using *sendfile()*

”*sendfile()* copies data between one file descriptor and another. Because this copying is done within the kernel, *sendfile()* is more efficient than the combination of *read(2)* and *write(2)*, which would require transferring data to and from user space.” [3] Linux kernels since 2.6.33 [3] can use *sendfile()* with file as a target, not just with sockets. This is a straight forward optimization that also achieves an optimum of performance with

$$S(B) = \frac{F}{B} \quad (3)$$

system calls. In that sense, the name *sendfile()* is a bit misleading if used solely with files.

6 Other performance thoughts

The performance of the algorithms described above not only depends on block sizes¹ and *mmap* blocks. Copying a large file on the same disk kills any performance, and clearly the type of storage one is using (SSD vs. HDD vs. tape) is very important. I am also not using any real-time tricks to kill scheduling effects. Opening a file with the `O_NOATIME` flag should be considered, but probably only has an effect if copying a lot of files rather than a single large file.

¹The *stat()* system call reports the best file I/O blocksize in the `st_blksize` member.

POSIX.1-2001 offers *posix_fadvise()* [4] calls which might be of use with the `POSIX_FADV_SEQUENTIAL`² flag, but it was not measurable in the experiment.

7 The experiment

For the ease of use, `odd` mimics most of the `dd` options such as *bs*, *skip*, *seek*, *if* and *of*. Additionally it supports the *mmap=M* (use `mmap` algorithm with `M` MB chunks), *cores=N* (use `splice` on cores `N-2` and `N-1`) and *send* for the *sendfile()* algorithm.

Benchmarking file copying can be very frustrating. The expected result from the formulas not always hold in practise. Additionally much depends on the underlying file system, storage, scheduler and so on. Results from quite short copies (that take less than a minute) do not mean much. Too many scheduling, caching and FS effects make the results vary too much. In the long run however, Figure 1 and 2 show that `odd` has major speedups against `dd`. Both tools were used with the same block size of 4096. Real world copy operations however should consider use of larger block sizes. The USB flash drive from Figure 1 was the bottleneck, but still, `odd` has noticable impact. I run all copy operations multiple times and chose an average result. The file copy operations have had different source and target file systems and less RAM than the file to copy. Both target file systems where an EXT4. The experiment shows that not only the zero-copy and amount of system calls improves performance, but the overall performance also depends on the implementation of the system calls inside the kernel. *splice()* seems to have the most performing implementation.

By playing with different block sizes, it is also easy to discover that the reported best block size for IO operation should not lead to a decision about block sizes. Larger block sizes might force unaligned operations inside the file system layer, but the saved amount of syscalls due to the larger buffer has much more impact on performance. As a rule of thumb, larger block sizes produce better performance, but using an optimized algorithm still out-performs larger block sizes of the standard algorithm.

Please note that the `dd` statistics prints `mB` (millions of bytes) rather than `MB` (Megabytes) and it were therefore converted to allow comparison.

²This flag tells the FS-layer to increase the read ahead window for this file.

Figure 1: 8GB copy from raw USB flash to EXT4-file on HDD, bs=4096

algorithm	time/s	MB/s	speedup
dd	522.7	15.3	-
splice	422	18	19.3%
splice cores default pipe	436	17.5	16.6%
splice cores 512MB pipe	399	19.1	23.7%
mmap 1GB	437	17.5	16.4%
sendfile	426	17.9	18.5%

Figure 2: 4GB file copy from one HDD to another, EXT4, bs=4096

algorithm	time/s	MB/s	speedup
dd	80	51	-
splice	72	56.7	10%
splice cores default pipe	70	58	13%
splice cores 512MB pipe	n/a	n/a	n/a%
mmap 512MB	77	53	4%
sendfile	n/a	n/a	n/a

8 Outlook

Other Operating Systems than Linux provide similar system calls. Therefore it is likely that `odd` could be ported to BSD or Solaris. It could also be considered to use asynchronous IO (`aio_read()` etc.) from the POSIX.1-2001 standard, but the current glibc implementation does not use in-kernel AIO operations but simulates them using threads. An AIO implementation that is using the kernels async IO interface can be found at [8].

References

- [1] dd
On most Linux distributions, dd is part of the coreutils RPM package.
<http://www.gnu.org/software/coreutils>
- [2] splice manpage:
man 2 splice
- [3] sendfile manpage:
man 2 sendfile
- [4] posix_fadvise manpage:
man 2 posix_fadvise
- [5] fcntl manpage:
man 2 fcntl
- [6] *Advanced Programming in the UNIX Environment*
Addison-Wesley, 2005 , Second Edition, W.Richard Stevens, Stephen A. Rago p.65f
- [7] optimized dd:
<http://github.com/stealth/odd>
- [8] AIO implementation:
<http://stealth.openwall.net/misc/aio-0.51.tgz>

9 Appendix A.1

simplified splice copy algorithm:

```
1 pipe(p);
2 n = ddc->bs;
3 for (;ddc->b_out != ddc->count && !sigint;) {
4     if (n > ddc->count - ddc->b_out)
5         n = ddc->count - ddc->b_out;
6     r = splice(ifd, NULL, p[1], NULL, n, SPLICE_F_MORE);
7     if (r <= 0) {
8         ddc->saved_errno = errno;
9         break;
10    }
11    ++ddc->rec_in;
12    r = splice(p[0], NULL, ofd, NULL, r, SPLICE_F_MORE);
13    if (r <= 0) {
14        ddc->saved_errno = errno;
15        break;
16    }
17    ddc->b_out += r;
18    ++ddc->rec_out;
19 }
```

10 Appendix A.2

simplified mmap copy algorithm:

```
1 for (;ddc->b_out != ddc->count && !sigint;) {
2     n = ddc->mmap;
3     bs = ddc->bs;
4     if (n > ddc->count - ddc->b_out)
5         n = ddc->count - ddc->b_out;
6     if (bs > n)
7         bs = n;
8     addr = mmap(NULL, n, PROT_WRITE, MAP_SHARED, ofd,
9                ddc->b_out + ddc->skip);
10    if (addr == MAP_FAILED) {
11        ddc->saved_errno = errno;
12        break;
13    }
14    for (i = 0; i < n; i += r) {
15        if (i + bs > n)
16            bs = n - i;
17        r = read(ifd, addr + i, bs);
18        if (r <= 0) {
19            ddc->saved_errno = errno;
20            munmap(addr, n);
21            break;
22        }
23        ddc->b_out += r;
24        ++ddc->rec_in;
25    }
26    /* pass along the potential above break */
27    if (r <= 0)
28        break;
29    ++ddc->rec_out;
30    munmap(addr, n);
31 }
```