

io_uring

Status Update within Samba

Stefan Metzmacher <metze@samba.org>

Samba Team / SerNet

2023-09-20

<https://samba.org/~metze/presentations/2023/SDC/>

Topics

- ▶ What is io-uring?
- ▶ io-uring for Samba
- ▶ Performance research, prototyping and ideas
- ▶ The road to upstream
- ▶ Future Improvements
- ▶ Questions? Feedback!

- ▶ I gave a similar talk at the storage developer conference 2020:
 - ▶ See <https://samba.org/~metze/presentations/2020/SDC/>
 - ▶ It explains the milestones and design up to Samba 4.13 (in detail)
- ▶ I gave a similar talk at the storage developer conference 2021:
 - ▶ See <https://samba.org/~metze/presentations/2021/SDC/>
 - ▶ It explains the milestones and updates up to Samba 4.15 (in detail)
- ▶ I gave a similar talk at the SambaXP conference 2023:
 - ▶ See <https://samba.org/~metze/presentations/2023/SambaXP/>
 - ▶ It explains the milestones and updates up to Samba 4.19 (in detail)

What is io-uring? (Part 1)

- ▶ Linux 5.1 introduced a new scalable AIO infrastructure
 - ▶ It's designed to avoid syscalls as much as possible
 - ▶ kernel and userspace share mmap'ed rings:
 - ▶ submission queue (SQ) ring buffer
 - ▶ completion queue (CQ) ring buffer
 - ▶ See "[Ringing in a new asynchronous I/O API](#)" on LWN.NET
- ▶ This can be nicely integrated with our async tevent model
 - ▶ It may delegate work to kernel threads
 - ▶ It seems to perform better compared to our userspace threadpool
 - ▶ It can also inline non-blocking operations

- ▶ Between userspace and filesystem (available from 5.1):
 - ▶ IORING_OP_READV, IORING_OP_WRITEV and IORING_OP_FSYNC
 - ▶ Supports buffered and direct io
 - ▶ IORING_OP_FSETXATTR, IORING_OP_FGETXATTR (from 5.19)
 - ▶ IORING_OP_GETDENTS, under discussion, but seems to be tricky
 - ▶ IORING_OP_FADVISE (from 5.6)
- ▶ Path based syscalls with async impersonation (from 5.6)
 - ▶ IORING_OP_OPENAT2, IORING_OP_STATX
 - ▶ Using IORING_REGISTER_PERSONALITY for impersonation
 - ▶ IORING_OP_UNLINKAT, IORING_OP_RENAMEAT (from 5.10)
 - ▶ IORING_OP_MKDIRAT, IORING_OP_SYMLINKAT, IORING_OP_LINKAT (from 5.15)
 - ▶ IORING_OP_SETXATTR, IORING_OP_GETXATTR (from 5.19)

- ▶ Between userspace and socket (and also filesystem) (from 5.8)
 - ▶ IORING_OP_SENDMSG, IORING_OP_RECVMSG
 - ▶ Improved MSG_WAITALL support (5.12, backported to 5.11, 5.10)
 - ▶ Maybe using IOSQE_ASYNC in order to avoid inline memcpy
 - ▶ IORING_OP_SPLICE, IORING_OP_TEE
 - ▶ IORING_OP_SENDMSG_ZC, zero copy with an extra completion (from 6.1)
 - ▶ IORING_OP_GET_BUF, under discussion to replace IORING_OP_SPLICE

- ▶ With Samba 4.12 we added "io_uring" vfs module
 - ▶ For now it only implements SMB_VFS_PREAD,PWRITE,FSYNC_SEND/RECV
 - ▶ It has less overhead than our pthreadpool default implementations
 - ▶ I was able to speed up a smbclient 'get largefile /dev/null'
 - ▶ Using against smbd on loopback
 - ▶ The speed changes from 2.2GBytes/s to 2.7GBytes/s
- ▶ The improvement only happens by avoiding context switches
 - ▶ But the data copying still happens:
 - ▶ From/to a userspace buffer to/from the filesystem/page cache
 - ▶ The data path between userspace and socket is completely unchanged
 - ▶ For both cases the cpu is mostly busy with memcpy

Performance research (SMB2 Read)

- ▶ In October 2020 I was able to do some performance research
 - ▶ With 100Gbit/s interfaces and two NUMA nodes per server.
- ▶ At that time I focussed on the SMB2 Read performance only
 - ▶ We had limited time on the given hardware
 - ▶ We mainly tested with fio.exe on a Windows client
 - ▶ Linux kernel 5.8.12 on the server
- ▶ More verbose details can be found here:
 - ▶ <https://lists.samba.org/archive/samba-technical/2020-October/135856.html>

smbclient ORING_OP_SENDMSG/SPICE (network)

4 connections, ~11 GBytes/s, smb2 8.6% cpu, with 4 io.wqework threads (pipe to socket) at ~20% cpu each.

smbclient is the bottleneck here too

```

printing file 1086.dat of size 2897233996 on <dev>vll112731512.0 KiloBytes/sec [average 4760083.8 KiloBytes/sec]
printing file 1087.dat of size 2897233996 on <dev>vll112731512.0 KiloBytes/sec [average 4760083.8 KiloBytes/sec]
printing file 1088.dat of size 2897233996 on <dev>vll11292521.7 KiloBytes/sec [average 7190686.8 KiloBytes/sec]
printing file 1089.dat of size 2897233996 on <dev>vll1252822.7 KiloBytes/sec [average 2920089.8 KiloBytes/sec]
printing file 1090.dat of size 2897233996 on <dev>vll1252984.5 KiloBytes/sec [average 3204892.8 KiloBytes/sec]
printing file 1091.dat of size 2897233996 on <dev>vll1226268.3 KiloBytes/sec [average 2746039.8 KiloBytes/sec]
printing file 1092.dat of size 2897233996 on <dev>vll1202626.2 KiloBytes/sec [average 2320983.8 KiloBytes/sec]
printing file 1093.dat of size 2897233996 on <dev>vll1202626.2 KiloBytes/sec [average 2320983.8 KiloBytes/sec]
printing file 1094.dat of size 2897233996 on <dev>vll1202626.2 KiloBytes/sec [average 2320983.8 KiloBytes/sec]
printing file 1095.dat of size 2897233996 on <dev>vll1202626.2 KiloBytes/sec [average 2320983.8 KiloBytes/sec]
printing file 1096.dat of size 2897233996 on <dev>vll1202626.2 KiloBytes/sec [average 2320983.8 KiloBytes/sec]
printing file 1097.dat of size 2897233996 on <dev>vll1202626.2 KiloBytes/sec [average 2320983.8 KiloBytes/sec]
printing file 1098.dat of size 2897233996 on <dev>vll1202626.2 KiloBytes/sec [average 2320983.8 KiloBytes/sec]
printing file 1099.dat of size 2897233996 on <dev>vll1202626.2 KiloBytes/sec [average 2320983.8 KiloBytes/sec]

```

```

top - 02:41:58 up 37 days, 17:46, 1 user, load average: 3.31, 4.32, 3.55
task: 897 total, 8 running, 892 sleeping, 8 stopped, 8 zombie
cpus: 0.0 us, 0.0 sy, 0.0 id, 0.0 wa, 0.0 st, 99.8 sr, 0.0 hi, 0.0 si, 0.0 st, 0.0 dt
MiB Mem 18424.4 total, 12713.7 free, 8813.5 buff, 8895.4 cache
top user 100.0% total, 727.8 user, 282.8 sys, 11204.0 wall Mem

```

pid	ppid	ps	NI	U	St	VSZ	USS	SMEM	CSM	CMEM	COMMAND	
102818	root	28	0	0	0	12000	2812	1304	0	0	0.126160	smtp
498030	root	28	0	0	0	37684	3810	99	0	0	0.96140	smtp
498123	root	28	0	0	0	37684	3810	99	0	0	0.96140	smtp
498386	root	28	0	0	0	37684	3810	99	0	0	0.96140	smtp
498706	root	28	0	0	0	37684	3810	99	0	0	0.96140	smtp
499219	root	28	0	0	0	37684	3810	99	0	0	0.96140	smtp
500332	root	28	0	0	0	37684	3810	99	0	0	0.96140	smtp
501038	root	28	0	0	0	37684	3810	99	0	0	0.96140	smtp
501743	root	28	0	0	0	37684	3810	99	0	0	0.96140	smtp

```

top - 02:41:59 up 37 days, 17:46, 1 user, load average: 1.11, 1.09, 0.92
task: 897 total, 8 running, 889 sleeping, 8 stopped, 8 zombie
cpus: 0.0 us, 0.0 sy, 0.0 id, 0.0 wa, 0.0 st, 99.8 sr, 0.0 hi, 0.0 si, 0.0 st, 0.0 dt
MiB Mem 18424.4 total, 17749.5 free, 8662.5 buff, 8682.0 cache
top user 100.0% total, 727.8 user, 282.8 sys, 11204.0 wall Mem

```

pid	ppid	ps	NI	U	St	VSZ	USS	SMEM	CSM	CMEM	COMMAND	
102818	root	28	0	0	0	0	0	0	0	0	0.126160	io.wqework
102819	root	28	0	0	0	0	0	0	0	0	0.126160	io.wqework
102820	root	28	0	0	0	0	0	0	0	0	0.126160	io.wqework
102821	root	28	0	0	0	0	0	0	0	0	0.126160	io.wqework
102822	root	28	0	0	0	0	0	0	0	0	0.126160	io.wqework
102823	root	28	0	0	0	0	0	0	0	0	0.126160	io.wqework
102824	root	28	0	0	0	0	0	0	0	0	0.126160	io.wqework

name	type	path	perms	inode	size	uid	gid	mode	ctime	mtime	atime	crtime
..	dir	/	drwxr-xr-x	2	4096	root	root	0755	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00
.	dir	/	drwxr-xr-x	2	4096	root	root	0755	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00
..	dir	/	drwxr-xr-x	2	4096	root	root	0755	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00
.	dir	/	drwxr-xr-x	2	4096	root	root	0755	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00
..	dir	/	drwxr-xr-x	2	4096	root	root	0755	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00
.	dir	/	drwxr-xr-x	2	4096	root	root	0755	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00

smbclient ORING_OP_SENDMSG/SPICE (loopback)

8 connections, ~22 GBytes/s, smb2 22% cpu, with 4 io.wqework threads (pipe to socket) at ~22% cpu each.

smbclient is the bottleneck here too, it triggers the memory copy done by copy_user_enhanced_fast_string()

```

printing file 1086.dat of size 2897233996 on <dev>vll1252822.7 KiloBytes/sec [average 2920089.8 KiloBytes/sec]
printing file 1087.dat of size 2897233996 on <dev>vll1252822.7 KiloBytes/sec [average 2920089.8 KiloBytes/sec]
printing file 1088.dat of size 2897233996 on <dev>vll1292521.7 KiloBytes/sec [average 7190686.8 KiloBytes/sec]
printing file 1089.dat of size 2897233996 on <dev>vll1252822.7 KiloBytes/sec [average 2920089.8 KiloBytes/sec]
printing file 1090.dat of size 2897233996 on <dev>vll1252984.5 KiloBytes/sec [average 3204892.8 KiloBytes/sec]
printing file 1091.dat of size 2897233996 on <dev>vll1226268.3 KiloBytes/sec [average 2746039.8 KiloBytes/sec]
printing file 1092.dat of size 2897233996 on <dev>vll1202626.2 KiloBytes/sec [average 2320983.8 KiloBytes/sec]
printing file 1093.dat of size 2897233996 on <dev>vll1202626.2 KiloBytes/sec [average 2320983.8 KiloBytes/sec]
printing file 1094.dat of size 2897233996 on <dev>vll1202626.2 KiloBytes/sec [average 2320983.8 KiloBytes/sec]
printing file 1095.dat of size 2897233996 on <dev>vll1202626.2 KiloBytes/sec [average 2320983.8 KiloBytes/sec]
printing file 1096.dat of size 2897233996 on <dev>vll1202626.2 KiloBytes/sec [average 2320983.8 KiloBytes/sec]
printing file 1097.dat of size 2897233996 on <dev>vll1202626.2 KiloBytes/sec [average 2320983.8 KiloBytes/sec]
printing file 1098.dat of size 2897233996 on <dev>vll1202626.2 KiloBytes/sec [average 2320983.8 KiloBytes/sec]
printing file 1099.dat of size 2897233996 on <dev>vll1202626.2 KiloBytes/sec [average 2320983.8 KiloBytes/sec]

```

name	type	path	perms	inode	size	uid	gid	mode	ctime	mtime	atime	crtime
..	dir	/	drwxr-xr-x	2	4096	root	root	0755	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00
.	dir	/	drwxr-xr-x	2	4096	root	root	0755	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00
..	dir	/	drwxr-xr-x	2	4096	root	root	0755	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00
.	dir	/	drwxr-xr-x	2	4096	root	root	0755	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00
..	dir	/	drwxr-xr-x	2	4096	root	root	0755	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00
.	dir	/	drwxr-xr-x	2	4096	root	root	0755	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00	2018-12-31 12:10:00

More loopback testing on brand new hardware

- ▶ Recently I re-did the loopback read tests IORING_OP_SENDFMSG/SPLICE (from /dev/shm/)
 - ▶ 1 connection, ~10-13 GBytes/s, smbd 7% cpu, with 4 iou-wrk threads at 7%-50% cpu.
 - ▶ 4 connections, 24-30 GBytes/s, smbd 18% cpu, with 16 iou-wrk threads at 3%-35% cpu.
- ▶ I also implemented SMB2 writes with IORING_OP_RECVMSG/SPLICE (tested to /dev/null)
 - ▶ 1 connection, ~7-8 GBytes/s, smbd 5% cpu, with 3 io-wrk threads at 1%-20% cpu.
 - ▶ 4 connections, ~10 GBytes/s, smbd 15% cpu, with 12 io-wrk threads at 1%-20% cpu.
- ▶ I tested with a Linux Kernel 5.13
 - ▶ In both cases the bottleneck is clearly on the smbclient side
 - ▶ We could apply similar changes to smbclient and add true multichannel support
 - ▶ It seems that the filesystem->pipe->socket path is much better optimized



SAMBA⁺

Stefan Metzmacher

io_uring (15/21)

SerNet

The road to upstream (TEVENT_FD_ERROR)

- ▶ We need support for TEVENT_FD_ERROR in order to monitor errors
 - ▶ When using IORING_OP_SEND,RECVMSG we still want to notice errors
 - ▶ This is the main merge request:
 - ▶ https://gitlab.com/samba-team/samba/-/merge_requests/2793
 - ▶ This merge request converts Samba to use TEVENT_FD_ERROR:
 - ▶ https://gitlab.com/samba-team/samba/-/merge_requests/2885
 - ▶ (It also simplifies other places in the code without io_uring)



SAMBA⁺

Stefan Metzmacher

io_uring (16/21)

SerNet

API glue to tevent:

```
void samba_io_uring_ev_register(void);

const struct samba_io_uring_features *samba_io_uring_system_features(void);

struct samba_io_uring *samba_io_uring_ev_context_get_ring(struct tevent_context *ev);

const struct samba_io_uring_features *samba_io_uring_get_features(
    const struct samba_io_uring *ring);

ev = tevent_context_init_byname(mem_ctx, "samba_io_uring_ev");
```

- ▶ samba_io_uring abstraction factored out of vfs_io_uring:
 - ▶ samba_io_uring_ev_hybrid tevent backend (glued on epoll backend)
 - ▶ It means every layer getting the tevent_context can use io_uring
 - ▶ No #ifdef's just checking if the required features are available



SAMBA+

Stefan Metzmacher

io_uring (17/21)

SerNet

generic submission/completion api:

```
void samba_io_uring_completion_prepare(struct samba_io_uring_completion *completion,
    void (*completion_fn)(struct samba_io_uring_completion *completion,
        void *completion_private,
        const struct io_uring_cqe *cqe),
    void *completion_private);

void samba_io_uring_submission_prepare(struct samba_io_uring_submission *submission,
    void (*submission_fn)(struct samba_io_uring *ring,
        struct samba_io_uring_submission *submission,
        void *submission_private),
    void *submission_private,
    struct samba_io_uring_completion *completion);

struct io_uring_sqe *samba_io_uring_submission_sqe(struct samba_io_uring_submission *
    submission);

size_t samba_io_uring_queue_submissions(struct samba_io_uring *ring,
    struct samba_io_uring_submission *submission);
```

- ▶ Using it ...
 - ▶ convert vfs_io_uring
 - ▶ use it in smb2_server.c
 - ▶ In future use it in other performance critical places too.



SAMBA+

Stefan Metzmacher

io_uring (18/21)

SerNet

- ▶ Refactoring of smb2_server.c
 - ▶ add optional IORING_OP_SENDMSG, IORING_OP_RECVMSG support
- ▶ There are structural problems with splice from a file
 - ▶ I had a discussion with the Linux developers about it:
 - ▶ The page content from the page cache may change unexpectedly
 - ▶ <https://lists.samba.org/archive/samba-technical/2023-February/thread.html#137945>
 - ▶ We may not be able to use IORING_OP_SENDMSG/SPLICE by default
 - ▶ Maybe IORING_OP_RECVMSG/SPLICE is possible
- ▶ With IORING_OP_SENDMSG_ZC only 1 copy is used:
 - ▶ It is able to avoid copying to the socket
 - ▶ We get an extra completion once the buffers are not needed anymore
 - ▶ Only with real hardware, not on loopback in an upstream kernel
 - ▶ A custom kernel loopback gives ~7.5 GBytes/s instead of ~3.5 GBytes/s
 - ▶ With a noop vfs module, we get ~18 GBytes/s instead of ~6 GBytes/s

Future Improvements

- ▶ Patches are slowly getting prepared for master
 - ▶ Some preparations are already in or pending merge requests
 - ▶ We even have basic automated ci testing in place now
 - ▶ But changes need to be checked for performance regressions
- ▶ We can use io_uring deep inside of the smbclient code
 - ▶ The low layers can just use `samba_io_uring_ev_context_get_ring()`
 - ▶ And use it if available without changing the whole stack

- ▶ Stefan Metzmacher, metze@samba.org
- ▶ <https://www.sernet.com>
- ▶ <https://samba.plus>

→ SerNet/SAMBA+ sponsor booth

Slides: <https://samba.org/~metze/presentations/2023/SDC/>



Stefan Metzmacher

io_uring (21/21)

SerNet