SMB3.1.1 and beyond: Optimizing access from Linux Client to Samba, the Cloud and modern file servers

Steve French Principal Software Engineer Azure Storage - Microsoft



Legal Statement

- This work represents the views of the author(s) and does not necessarily reflect the views of Microsoft Corporation
- Linux is a registered trademark of Linus Torvalds.
- Other company, product, and service names may be trademarks or service marks of others.

Who am I?

- Steve French smfrench@gmail.com
- Author and maintainer of Linux cifs vfs (for accessing Samba, Windows and various SMB3/CIFS based NAS appliances)
- Also wrote initial SMB2 kernel client prototype
- Member of the Samba team, coauthor of SNIA CIFS Technical Reference, former SNIA CIFS Working Group chair
- Principal Software Engineer, Azure Storage: Microsoft

Outline

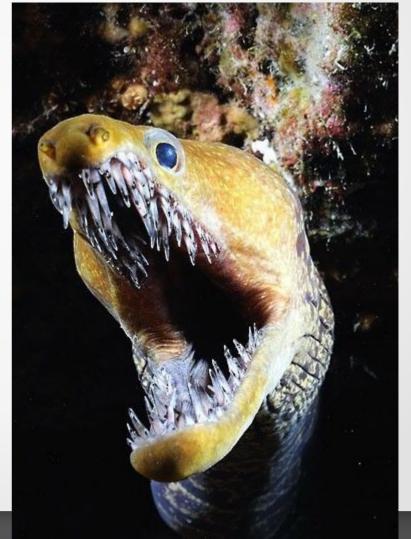
- General Linux File System Status Linux FS and VFS Activity
- What are the goals?
- Key Feature Status (add RDMA, compounding, handle caching, directory leasing)
 - SMB3.11
 - Handle caching and directory leases
 - Compounding
 - RDMA (see Long Li's talk)
 - CopyOffload
 - HA
 - Security Features/Encryption
 - Other optional SMB3 features
- Performance overview
- POSIX compatibility
 - Status of SMB3 POSIX Extensions
 - Alternatives
- Testing

A year ago ... and now ... kernel (including SMB3 client cifs.ko) improving

 13 months ago we had Linux version 4.11 ie "Fearless Coyote"



Three days ago we got 4.17 "Merciless Moray"



Discussions driving some of the FS development activity ?

- New mount API, new fsinfo API
- Many of the high priority, evolving storage features are critical:
 - Better support for faster storage
 - RDMA and low latency ways to access VERY high speed storage
 - NVMe
 - Faster (and cheaper) network adapters (10Gb \rightarrow 40Gb->100Gb ethernet ... and RDMA)
 - I/O priority
 - Now that statx (extended stat) is in, adding more metadata flags
 - Broadening use of copy offload (e.g. "copy_file_range" syscall)
 - In rsync, cp etc.
 - Shift to Cloud (longer latencies, object & file coexisting)

2018 Linux FS/MM summit (in April)

Great group of talented developers



Most Active Linux Filesystems this year

- 4357 kernel filesystem changesets in last year (since 4.12-rc4 kernel)! Continuing strong (up slightly)
 - FS activity: 5.75% of overall kernel changes (which are dominated by drivers). FS is watched carefully!
 - Kernel is now 17.17 million lines of source code (measured last week with sloccount tool)
- There are many Linux file systems (>50), but six (and the VFS layer itself) drive 70% of the activity
 - File systems represent about 5.1% of the overall kernel source code (876,000 lines of code)
- cifs.ko (cifs/smb3 client) among more active fs (#5 out of 60 and growing). More activity is good!
 - BTRFS 826 changesets (up)
 - VFS (overall fs mapping layer and common functions) 598 (down 13%)
 - XFS 524 (up slightly)
 - F2FS 357 (down 25%)
 - NFS client 276 (down over 40%!)
 - CIFS/SMB2/SMB3 client 250 (up 50%!). And speeding up! (70% in last 5 months)
 - cifs.ko is 47,690 lines of kernel code (not counting user space helpers and samba userspace tools)
 - Ext4 230 (flat)
 - NFS server 140 (down 7%). Linux NFS server is **MUCH** smaller than CIFS or NFS clients (or Samba).
 - And various other file systems ... Ceph 144 (down), GFS 130, AFS 120 ...
- NB: Samba is as active as all Linux file systems put together (>4000 changesets per year) broader in scope (by a lot) and also is user space not kernel. 100x larger than the NFS server in Linux!

What are the goals?

- Make SMB3 (SMB3.11 and followons) fastest, most secure general purpose way to access file data, whether in the cloud or on premises or from virtualized environments
- Implement all reasonable Linux/POSIX features so apps don't have to know running on SMB3 mounts (vs. local)
- Allow extensions so that as Linux evolves, and need for new features discovered, can quickly add them to Linux kernel client and Samba

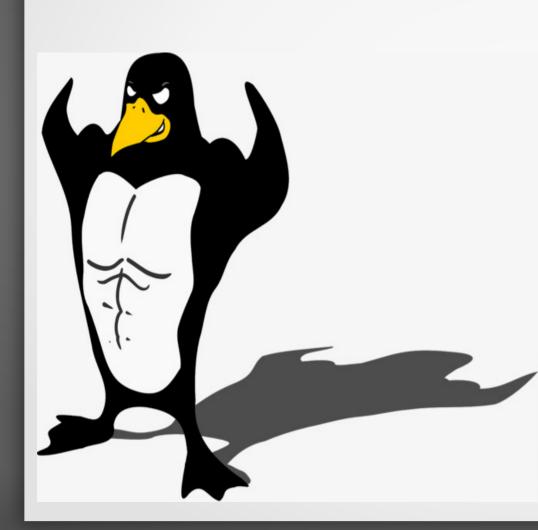


Exciting year!!

. . .

- Faster performance
- POSIX Extensions (finally)!
- SMB3.11, improved security
- LOTS of new features

Fixes and Features that were in progress last time ...



- Full SMB3.11 support!
- Statx (extended stat linux API returning additional metadata flags)
- Improved performance
- Improved POSIX compatibility (partial, in progress)
- ACLs and security improvements

35% more efficient mount & SMB3.11 works!

Filter:	smb2			 Expression. 	Clear Apply Save
No.	Time	Source	Destination	Protocol Len	gth Info
	4 0.000666558	172.16.194.1	172.16.194.128	SMB2	256 Negotiate Protocol Request
	5 0.002358268	172.16.194.128	172.16.194.1	SMB2	668 Negotiate Protocol Response
	7 0.002502467	172.16.194.1	172.16.194.128		192 Session Setup Request, NTLMSSP_NEGOTIATE
	8 0.003919218	172.16.194.128	172.16.194.1		382 Session Setup Response, Error: STATUS_MORE_PROCESSING_REQUIRED, NTL
1	9 0.004131694	172.16.194.1	172.16.194.128		454 Session Setup Request, NTLMSSP_AUTH, User: \testuser
1	0.007151312	172.16.194.128	172.16.194.1		144 Session Setup Response
_		172.16.194.1	172.16.194.128		188 Tree Connect Request Tree: \\172.16.194.128\IPC\$
		172.16.194.128	172.16.194.1		152 Tree Connect Response
		172.16.194.1	172.16.194.128		192 Tree Connect Request Tree: \\172.16.194.128\public
		172.16.194.128	172.16.194.1		152 Tree Connect Response
		172.16.194.1	172.16.194.128		200 Create Request File:
1	6 0.009128975	172.16.194.128	172.16.194.1		224 Create Response File: [unknown]
-		172.16.194.1	172.16.194.128		177 GetInfo Request FS_INFO/FileFsAttributeInformation File: [unknown]
1	8 0.009681622	172.16.194.128	172.16.194.1		164 GetInfo Response
		172.16.194.1	172.16.194.128		177 GetInfo Request FS_INFO/FileFsDeviceInformation File: [unknown]
_		172.16.194.128	172.16.194.1		152 GetInfo Response
2	1 0.010309488	172.16.194.1	172.16.194.128		177 GetInfo Request FS_INFO/FileFsSectorSizeInformation File: [unknown]
2	2 0.010566781	172.16.194.128	172.16.194.1		172 GetInfo Response
2	3 0.010721458	172.16.194.1	172.16.194.128		240 Ioctl Request FSCTL_DFS_GET_REFERRALS, File: \172.16.194.128\public
-		172.16.194.128	172.16.194.1		145 Ioctl Response, Error: STATUS_FS_DRIVER_REQUIRED
		172.16.194.1	172.16.194.128		176 GetInfo Request FILE_INFO/SMB2_FILE_ALL_INFO File: [unknown]
2	6 0.011595834	172.16.194.128	172.16.194.1	SMB2	248 GetInfo Response

▶ Frame 5: 668 bytes on wire (5344 bits), 668 bytes captured (5344 bits) on interface 0

Linux cooked capture

Internet Protocol Version 4, Src: 172.16.194.128, Dst: 172.16.194.1

> Transmission Control Protocol, Src Port: 445, Dst Port: 51128, Seq: 1, Ack: 189, Len: 600

NetBIOS Session Service

- SMB2 (Server Message Block Protocol version 2)
- SMB2 Header
- Negotiate Protocol Response (0x00)
 - StructureSize: 0x0041
 - Security mode: 0x01, Signing enabled Dialect: 0x0311

NegetieteContentCont

NegotiateContextCount: 2

Server Guid: e21779a0-c688-457d-86e9-dd2977809277

▶ Capabilities: 0x00000007, DFS, LEASING, LARGE MTU

Max Transaction Size: 8388608

And SMB3.11 encryption works ...

- "mount -t cifs //server/share /mnt -o vers=3.11,seal"
- Thanks Aurelien! | Filte

Time	Source		Protocol	Length	h Into				
2 602200520	107.0.0.1	Destination	Protocol	-					
3.692398538		127.0.0.1	SMB2		6 Negotiate Protocol Request				
3.699723875		127.0.0.1	SMB2		0 Negotiate Protocol Response				
					2 Session Setup Request, NTLMSSP_NEGOTIATE				
					2 Session Setup Response, Error: STATUS_MO				
					0 Session Setup Request, NTLMSSP_AUTH, Use 4 Session Setup Response				
					0 Encrypted SMB3				
					4 Encrypted SMB3				
					6 Encrypted SMB3				
					4 Encrypted SMB3				
		, , , , , , , , , , , , , , , , , , ,	captarea (
		127.0.0.1. Dst: 1	27.0.0.1						
				. Sea:	1. Ack: 189. Len: 272				
				,	-,,				
Server Messa	ge Block Protoco	l version 2)							
Header									
tiate Protoc	ol Response (0x00))							
ructureSize:	0x0041								
Security mode: 0x01, Signing enabled									
Dialect: 0x0311									
NegotiateContextCount: 2									
Server Guid: 0000000-0000-0000-0000000000000									
Capabilities: 0x00000007, DFS, LEASING, LARGE MTU									
			T						
		· · · · · · · · · · · · · · · · · · ·							
		0502a03e303ca00e	300c060a2b6	060104.					
5									
<i></i>			LITIES						
	3.699999132 3.700105072 3.704463585 3.704580849 3.704732834 3.704829715 3.712062928 33: 340 byte cooked captu et Protocol ission Contr S Session Se Server Messa Header tiate Protoc ructureSize: curity mode: alect: 0x031 gotiateConte: rver Guid: 00 pabilities: 0 x Transaction x Read Size: x Write Size rrent Time: No curity Blob: gotiateConte: gotiateConte:	cooked capture et Protocol Version 4, Src: 3 ission Control Protocol, Src S Session Service Server Message Block Protocol Header tiate Protocol Response (0x00 ructureSize: 0x0041 curity mode: 0x01, Signing er alect: 0x0311 gotiateContextCount: 2 rver Guid: 00000000-0000-0000 pabilities: 0x00000007, DFS, x Transaction Size: 8388608 x Read Size: 8388608 x Write Size: 8388608 rrent Time: Jun 4, 2018 21:00 ot Time: No time specified (6 curity Blob: 604806062b060105 gotiateContextOffset: 0x00d0 gotiate Context: SMB2_PREAUTH	3.699999132 127.0.0.1 127.0.0.1 3.700105072 127.0.0.1 127.0.0.1 3.704463585 127.0.0.1 127.0.0.1 3.704580849 127.0.0.1 127.0.0.1 3.704732834 127.0.0.1 127.0.0.1 3.704829715 127.0.0.1 127.0.0.1 3.712062928 127.0.0.1 127.0.0.1 33: 340 bytes on wire (2720 bits), 340 bytes of cooked capture et Protocol Version 4, Src: 127.0.0.1, Dst: 1: ission Control Protocol, Src Port: 445, Dst Portice Server Message Block Protocol version 2) Header tiate Protocol Response (0x00) ructureSize: 0x0041 curity mode: 0x01, Signing enabled alect: 0x0311 gotiateContextCount: 2 rver Guid: 00000000-0000-0000-00000000000 pabilities: 0x0000007, DFS, LEASING, LARGE MT x Transaction Size: 8388608 x Read Size: 8388608 rrent Time: Jun 4, 2018 21:04:23.161808000 CD ot Time: No time specified (0) curity Blob: 604806062b0601050502a03e303ca00e: gotiateContextOffset: 0x00d0	3.699999132 127.0.0.1 127.0.0.1 SMB2 3.700105072 127.0.0.1 127.0.0.1 SMB2 3.704463585 127.0.0.1 127.0.0.1 SMB2 3.704580849 127.0.0.1 127.0.0.1 SMB2 3.704732834 127.0.0.1 127.0.0.1 SMB2 3.704829715 127.0.0.1 127.0.0.1 SMB2 3.712062928 127.0.0.1 127.0.0.1 SMB2 33: 340 bytes on wire (2720 bits), 340 bytes captured (2 cooked capture et Protocol Version 4, Src: 127.0.0.1, Dst: 127.0.0.1 ission Control Protocol, Src Port: 445, Dst Port: 56698 S Session Service Server Message Block Protocol version 2) Header tiate Protocol Response (0x00) ructureSize: 0x0041 curity mode: 0x01, Signing enabled alect: 0x0311 gotiateContextCount: 2 rver Guid: 00000000-0000-0000-0000000000000 pabilities: 0x00000007, DFS, LEASING, LARGE MTU x Transaction Size: 8388608 x Read Size: 8388608 x Write Size: 8388608 rrent Time: Jun 4, 2018 21:04:23.161808000 CDT ot Time: No time specified (0) curity Blob: 604806062b0601050502a03e303ca00e300c060a2b6 gotiateContextOffset: 0x00d0 gotiateContext: SMB2_PREAUTH_INTEGRITY_CAPABILITIES	3.699999132 127.0.0.1 127.0.0.1 SMB2 36 3.700105072 127.0.0.1 127.0.0.1 SMB2 43 3.704463585 127.0.0.1 127.0.0.1 SMB2 14 3.704580849 127.0.0.1 127.0.0.1 SMB2 23 3.704732834 127.0.0.1 127.0.0.1 SMB2 20 3.704829715 127.0.0.1 127.0.0.1 SMB2 23 3.712062928 127.0.0.1 127.0.0.1 SMB2 20 33: 340 bytes on wire (2720 bits), 340 bytes captured (2720 b cooked capture et Protocol Version 4, Src: 127.0.0.1, Dst: 127.0.0.1 ission Control Protocol, Src Port: 445, Dst Port: 56698, Seq: S session Service Server Message Block Protocol version 2) Header tiate Protocol Response (0x00) ructureSize: 0x0041 curity mode: 0x01, Signing enabled alect: 0x0311 gotiateContextCount: 2 rver Guid: 00000000-0000-0000-0000000000000 pabilities: 0x0000007, DFS, LEASING, LARGE MTU x Transaction Size: 8388608 x Read Size: 8388608 rrent Time: Jun 4, 2018 21:04:23.161808000 CDT ot Time: No time specified (0) curity Blob: 604806062b0601050502a03e303ca00e300c060a2b060104 gotiateContextOffset: 0x00d0 gotiateContext1 SMB2 PREAUTH_INTEGRITY_CAPABILITIES				

Can load it as 'smb3' and even disable cifs

- Improving security: can disable cifs

root@smf-Thinkpad-P51

File Edit View Search Terminal Help
root@smf-Thinkpad-P51:~# modprobe smb3 disable_legacy_dialects=1
root@smf-Thinkpad-P51:~# mount -t cifs //localhost/scratch /mnt1 -o vers=1.0,username=testuser,
mount error(22): Invalid argument
Refer to the mount.cifs(8) manual page (e.g. man mount.cifs)
root@smf-Thinkpad-P51:~# dmesg
[294.844994] FS-Cache: Netfs 'cifs' registered for caching
[294.845081] Key type cifs.spnego registered
[294.845084] Key type cifs.idmap registered
[297.769583] CIES VES: mount with legacy_dialect_disabled

Tracing with the new ftrace is so easy ...

root@smf-Thinkp

File Edit View Search Terminal Help

root@smf-Thinkpad-P51:~# modprobe smb3
root@smf-Thinkpad-P51:~# trace-cmd start -e cifs
root@smf-Thinkpad-P51:~# mount -t cifs //localhost/test /mnt1 -o username=testuser,pass
root@smf-Thinkpad-P51:~# touch /mnt1/newfile
touch: cannot touch '/mnt1/newfile': Permission denied
root@smf-Thinkpad-P51:~# trace-cmd show

Current List of CIFS/SMB3 tracepoints and an example of detail for one

```
root@smf-Thinkpad-P51:/sys/kernel/debug/tracing/events/cifs# ls
              smb3_cmd_err smb3_flush_err smb3_open_err
enable
                                                               smb3_set_info_err
filter
              smb3_enter smb3_fsctl_err smb3_query_info_err smb3_write_done
smb3_close_err smb3_exit_done smb3_lock_err smb3_read_done
                                                              smb3 write err
smb3_cmd_done smb3_exit_err smb3_open_done smb3_read_err
root@smf-Thinkpad-P51:/sys/kernel/debug/tracing/events/cifs# cat smb3 fsctl err/
enable filter format hist id trigger
root@smf-Thinkpad-P51:/sys/kernel/debug/tracing/events/cifs# cat smb3 fsctl err/format
name: smb3 fsctl err
ID: 2554
format:
       field:unsigned short common_type;
                                           offset:0; size:2; signed:0;
                                           offset:2; size:1; signed:0;
       field:unsigned char common flags;
       field:unsigned char common preempt count;
                                                                 size:1: signed:0:
                                                   offset:3:
       field:int common pid; offset:4;
                                           size:4; signed:1;
       field:unsigned int xid; offset:8;
                                           size:4: signed:0:
       field:__u64 fid; offset:16;
                                           size:8; signed:0;
       field:__u32 tid; offset:24;
                                           size:4; signed:0;
       field:__u64 sesid; offset:32;
                                           size:8; signed:0;
       field:__u8 infclass; offset:40; size:1; signed:0;
       field:__u32 type; offset:44;
                                           size:4; signed:0;
       field:int rc; offset:48; size:4; signed:1;
```

print fmt: "xid=%u sid=0x%llx tid=0x%x fid=0x%llx class=%u type=0x%x rc=%d", REC->xid, REC->se EC->tid, REC->fid, REC->infclass, REC->type, REC->rc

Example output: tracing mount and touch (create file) failure

παι στι σγ σοι ττι σ --=> preempt-depth delav TASK-PID CPU# TIMESTAMP FUNCTION 1370.528512: smb3 enter: cifs mount: xid=0 mount.cifs-4557 [005] mount.cifs-4557 [005] 1370.528778: smb3 enter: cifs get smb ses: xid=1 mount.cifs-4557 [005] 1370.536041: smb3 cmd done: sid=0x0 tid=0x0 cmd=0 mid=0 sid=0xfb6289ac tid=0x0 cmd=1 mid=1 status=0xc0000016 rc=-5 mount.cifs-4557 [005] 1370.536324: smb3 cmd err: sid=0xfb6289ac tid=0x0 cmd=1 mid=2 mount.cifs-4557 [005] 1370.541155: smb3 cmd done: [005] mount.cifs-4557 1370.541181: smb3_exit_done: cifs_get_smb_ses: xid=1 mount.cifs-4557 [005] 1370.541183: smb3 enter: cifs setup ipc: xid=2 mount.cifs-4557 [005] 1370.541419: smb3 cmd done: sid=0xfb6289ac tid=0x92f0b9bb cmd=3 mid=3 mount.cifs-4557 [005] 1370.541588: smb3 cmd done: sid=0xfb6289ac tid=0x92f0b9bb cmd=11 mid=4 mount.cifs-4557 [005] 1370.541590: smb3 exit done: cifs setup ipc: xid=2 mount.cifs-4557 [005] 1370.541591: smb3 enter: cifs get tcon: xid=3 mount.cifs-4557 [005] 1370.541768: smb3 cmd done: sid=0xfb6289ac tid=0xb02df36d cmd=3 mid=5 mount.cifs-4557 1370.541873: smb3 cmd done: sid=0xfb6289ac tid=0xb02df36d cmd=11 mid=6 [005] mount.cifs-4557 1370.541874: smb3 exit done: cifs get tcon: xid=3 [005] mount.cifs-4557 [005] 1370.542069: smb3 cmd done: sid=0xfb6289ac tid=0xb02df36d cmd=5 mid=7 mount.cifs-4557 [005] 1370.542070: smb3 open done: xid=0 sid=0xfb6289ac tid=0xb02df36d fid=0xf976554e cr opts=0x0 des access=0x80 [005] 1370.542122: smb3 cmd done: sid=0xfb6289ac tid=0xb02df36d cmd=16 mid=8 mount.cifs-4557 sid=0xfb6289ac tid=0xb02df36d cmd=16 mid=9 mount.cifs-4557 [005] 1370.542140: smb3 cmd done: [005] sid=0xfb6289ac tid=0xb02df36d cmd=16 mid=10 mount.cifs-4557 1370.542159: smb3 cmd done: mount.cifs-4557 [005] 1370.542197: smb3 cmd err: sid=0xfb6289ac tid=0x92f0b9bb cmd=11 mid=11 status=0xc0000225 rc=-2 mount.cifs-4557 [005] 1370.542198: smb3 fsctl err: 1370.542200: smb3 exit done: mount.cifs-4557 [005] cifs mount: xid=0 cifs root iget: xid=4 mount.cifs-4557 [005] 1370.542259: smb3 enter: mount.cifs-4557 [005] 1370.542310: smb3 cmd done: sid=0xfb6289ac tid=0xb02df36d cmd=16 mid=12 mount.cifs-4557 [005] 1370.542317: smb3 exit done: cifs root iget: xid=4 1377.479938: smb3 enter: cifs atomic open: xid=5 touch-4562 [001] 1377.480702: smb3 cmd err: sid=0xfb6289ac tid=0xb02df36d cmd=5 mid=13 status=0xc0000022 rc=-13 touch-4562 [001] [001] 1377.480707: smb3 open err: xid=5 sid=0xfb6289ac tid=0xb02df36d cr opts=0x40 des access=0x40000080 rc=-13 touch-4562

Splice write fixed (also helps sendfile)

root@smf-Thinkpad-P51:~# gio copy /mnt1/trace.dat /mnt1/targe Transferred 7.2 MB out of 7.2 MB (7.2 MB/s) root@smf-Thinkpad-P51:~#

Statx (and cifs pseudoxattrs) and get/set real xattrs work

root@smf-Thinkpad-P51:/mnt1# getfattr file2 -d # file: file2 user.somexattr="somevalue" root@smf-Thinkpad-P51:/mnt1# ~/statx/test-statx file2 2M statx(file2) = 0 results=fdf Size: 0 Blocks: 0 IO Block: 16384 regular file Device: 00:38 Inode: 13107206 Links: 1 Access: (0755/-rwxr-xr-x) Uid: 0 Gid: 0 Modify: 2018-06-05 02:39:25.088837500-0500 Change: 2018-06-05 02:39:25.088837500-0500 Birth: 2018-05-31 18:06:01.644761500-0500 statx(2M) = 0results=fdf Size: 2097152 Blocks: 4096 IO Block: 16384 regular file Device: 00:38 Inode: 13107210 Links: 1 Access: (0755/-rwxr-xr-x) Uid: 0 Gid: 0 Modify: 2018-06-05 02:41:05.058102400-0500 Change: 2018-06-05 02:41:05.058102400-0500 Birth: 2018-06-05 02:41:05.054102300-0500 root@smf-Thinkpad-P51:/mnt1# getfattr 2M -n user.cifs.creationtime -e hex # file: 2M user.cifs.creationtime=0xdfff268fa0fcd301 root@smf-Thinkpad-P51:/mnt1# getfattr 2M -n user.cifs.dosattrib -e hex

file: 2M

<u>user cifs dosattrib=0x80000000</u>

SMB3/CIFS Fixes/Features by release

- 4.9 (37 changesets) December 11, 2016
 - - Various reconnect improvements (e.g. send echo ASAP to reconnect smb session/tcon quicker after socket reconnect
 - Uid/gid from special sid (new mount option "idsfromsid")
 - Can override number of credits (new mount option "max_credits")
 - Query file attributes or creation time via xattr (cifs.dosattrib, cifs.creationtime)
- 4.10 (17) February 9th, 2017 Bug Fixes
- 4.11 (51 changesets) April 30th, 2017
 - SMB3 reconnect improvements (including better persistent & durable handles). Much higher reliability now when server crashes or failsover while I/o in flight or cached. Lots of corner cases fixed (Thank you Germano!)
 - Server side copy works much better: Clone file range (and "cp -reflink" command) now support more common
 - "copychunk" copy offload style (had required less common "duplicate extents" support). Thank you Sachin!
 - SMB3 DFS support (Thank you Aurelien!)
 - SMB3 Encryption support (Thank you Pavel!)
 - Note that this allows mounts to the cloud: Azure shares often require encryption
- 4.12 (36 changesets) July 12th, 2017
 - Posix smb3 name mapping improvements
 - Improved aio support
 - Add support for enumerating snapshots (via ioctl to cifs.ko)
 - Bug fixes

SMB3/CIFS Features by release (cont)

- 4.13 (27 changesets) September 3rd, 2017
 - Change default dialect to SMB3 from CIFS
 - SMB3 support for "cifsacl" mount option (and mode emulation)
 - Bug fixes
- 4.14 (37 changesets) November 12th, 2017
 - Bug fixes (especially for SMB2.1/SMB3 validate negotiate)
 - Default dialect changed to multidialect (SMB2.1, SMB3, SMB3.02)
 - Added xattr support for SMB2/SMB3
- 4.15 (6 changesets) January 28, 2018
 - Minor bug fixes

SMB3/CIFS Features by release (cont)

- 4.16 (68 changesets) April 1
 - Add splice_write support
 - Add support for smbdirect (SMB3 rdma). Thanks Long Li!
- 4.17 (54 changesets) June 3
 - Bug fixes
 - Add signing support for smbdirect
 - Add support for SMB3.11 encryption, and preauth integrity
 - SMB3.11 dialect improvements (and no longer marked experimental)
- Linux next ie 4.18-rc (38 changesets)
 - RDMA and Direct I/O improvements (see Long Li's talk)
 - Bug fixes
 - SMB3 POSIX extensions (initial minimal set, open and negotiate context only. use 'posix' mnt parm)
 - Add "smb3" alias to cifs.ko ("insmod smb3")
 - Allow disabling less secure dialects through new module install parm (disable_legacy_dialects)
 - Add support for improved tracing (ftrace, trace-cmd)
 - Cache root file handle, reducing redundant opens, improving perf

Linux CIFS/SMB3 client bug status summary

- Bugzilla.kernel.org
 - 40 bugs mostly not serious/already fixed
- Bugzilla.samba.org
 - 53 bugs mostly not serious or already fixed
- Would love help to triage, and close out some of the bugs which are already fixed.

pdat

ndclas

SMB2/SMB3 Compounding

(Slides courtesy of Ronnie Sahlberg at RedHat who is doing great work improving this)

- Hard work is done by now. I.e. the separation of NBSS and SMB2 headers. Most of work is already merged into mainline now
- TODO: plumbing to operate on arrays of requests/responses that are all done in one one compound with an array of smb2 PDUs. Patches exist on the list for this.
- smb2 compounding is VERY flexible and there are a lot of places in cifs.ko where we will be able to use them to
 - improve performance
 - also make the client get slightly more posix like behavior from smb2.
- Once we have the compounding in, there are a HUGE number of places where we should switch to using compounding.

df

	ib2				Expres	sion
lo.	Time	Source	Destination		Length Info	
-	1 0.000000000	192.168.124.203	192.168.124.1	SMB2	198 Create Request File:	
	2 0.000864358	192.168.124.1	192.168.124.203	SMB2	222 Create Response File: [unknown]	
	4 0.001715177	192.168.124.203	192.168.124.1	SMB2	174 GetInfo Request FILE_INFO/SMB2_FILE_ALL_INFO File: [unknown]	
		192.168.124.1	192.168.124.203	SMB2	244 GetInfo Response	
	6 0.002746605	192.168.124.203	192.168.124.1	SMB2	158 Close Request File: [unknown]	
	7 0.002974102	192.168.124.1	192.168.124.203	SMB2	194 Close Response	
	8 0.003632539	192.168.124.203	192.168.124.1	SMB2	198 Create Request File:	
	9 0.004250306	192.168.124.1	192.168.124.203	SMB2	222 Create Response File: [unknown]	
		192.168.124.203	192.168.124.1	SMB2	174 GetInfo Request FILE_INFO/SMB2_FILE_FULL_EA_INFO File: [unknown]	
		192.168.124.1	192.168.124.203	SMB2	206 GetInfo Response	
		192.168.124.203	192.168.124.1	SMB2	158 Close Request File: [unknown]	
-		192.168.124.1	192.168.124.203	SMB2	194 Close Response	
		192.168.124.203	192.168.124.1	SMB2	<u>390 Create Request File: ;GetInfo Request FS_INFO/FileFsFullSizeInformation;Close R</u>	equest
	15 0.012183184	192.168.124.1	192.168.124.203	SMB2	454 Create Response File: [unknown];GetInfo Response;Close Response	
Fra	ame 14: 390 bytes	on wire (3120 bits). 390 bytes captured	(3120 bits) on interface 0	
), 390 bytes captured st: 52:54:00:55:3b:d4) on interface 0	
Et	nernet II, Srć: 5	52:54:00:cì:f8:ef, D	st: 52:54:00:55:3b:d4) on interface 0	
Et In	nernet II, Src: 5 Cernet Protocol V	52:54:00:cl̀:f8:ef, D /ersion 4, Src: 192.	st: 52:54:00:55:3b:d4 168.124.203, Dst: 192	.168.124.1) on interface 0 5, Ack: 887, Len: 324	
Et In Tra	nernet II, Src: 5 Cernet Protocol V	52:54:00:cÌ:f8:ef, D /ersion 4, Src: 192. bl Protocol, Src Por	st: 52:54:00:55:3b:d4 168.124.203, Dst: 192	.168.124.1	·	
Et In Tr Ne	nernet II, Srć: 5 cernet Protocol V ansmission Contro cBIOS Session Ser	52:54:00:cÌ:f8:ef, D /ersion 4, Src: 192. bl Protocol, Src Por	<pre>śt: 52:54:00:55:3b:d4 168.124.203, Dst: 192 t: 52458, Dst Port: 4</pre>	.168.124.1	·	
Et In Tr Ne SM	nernet II, Srć: 5 cernet Protocol V ansmission Contro cBIOS Session Ser	52:54:00:cÌ:f8:ef, D /ersion 4, Src: 192. Dl Protocol, Src Por vice	<pre>śt: 52:54:00:55:3b:d4 168.124.203, Dst: 192 t: 52458, Dst Port: 4</pre>	.168.124.1	·	
Et In Tr: Ne SM	nernet II, Src: 5 cernet Protocol V ansmission Contro BIOS Session Ser 32 (Server Messag SMB2 Header Create Request (52:54:00:cl:f8:ef, D /ersion 4, Src: 192. D Protocol, Src Por vice ge Block Protocol ve 0x05)	śt: 52:54:00:55:3b:d4 168.124.203, Dst: 192 t: 52458, Dst Port: 4 rsion 2)	.168.124.1	·	
Et In Tr SM	nernet II, Src: 5 cernet Protocol V ansmission Contro cBIOS Session Ser 32 (Server Messag SMB2 Header Create Request (32 (Server Messag	52:54:00:cl:f8:ef, D /ersion 4, Src: 192. D Protocol, Src Por vice Je Block Protocol ve	śt: 52:54:00:55:3b:d4 168.124.203, Dst: 192 t: 52458, Dst Port: 4 rsion 2)	.168.124.1	·	
Et In Tr: Ne SM	nernet II, Src: 5 cernet Protocol V ansmission Contro BIOS Session Ser 32 (Server Messag SMB2 Header Create Request (32 (Server Messag SMB2 Header	52:54:00:cl:f8:ef, D /ersion 4, Src: 192. ol Protocol, Src Por vice ge Block Protocol ve 0x05) ge Block Protocol ve	śt: 52:54:00:55:3b:d4 168.124.203, Dst: 192 t: 52458, Dst Port: 4 rsion 2)	.168.124.1	·	
Etl In Tr: Ne SM	nernet II, Srć: 5 ernet Protocol V ansmission Contro BIOS Session Ser 22 (Server Messag SMB2 Header Create Request (32 (Server Messag SMB2 Header GetInfo Request	52:54:00:cl:f8:ef, D Version 4, Src: 192. D) Protocol, Src Por Vice Block Protocol ve 0x05) ge Block Protocol ve (0x10)	<pre>śt: 52:54:00:55:3b:d4 168.124.203, Dst: 192 t: 52458, Dst Port: 4 rsion 2) rsion 2)</pre>	.168.124.1	·	
Et In Tr: Ne SM	nernet II, Srć: 5 cernet Protocol V ansmission Contro BIOS Session Ser 22 (Server Messag SMB2 Header Create Request (32 (Server Messag SMB2 Header GetInfo Request 32 (Server Messag 34 (Server Messag	52:54:00:cl:f8:ef, D /ersion 4, Src: 192. ol Protocol, Src Por vice ge Block Protocol ve 0x05) ge Block Protocol ve	<pre>śt: 52:54:00:55:3b:d4 168.124.203, Dst: 192 t: 52458, Dst Port: 4 rsion 2) rsion 2)</pre>	.168.124.1	·	
Et In Tr: Ne SM	nernet II, Srć: 5 cernet Protocol V ansmission Contro BIOS Session Ser 22 (Server Messag SMB2 Header Create Request (22 (Server Messag SMB2 Header GetInfo Request 22 (Server Messag SMB2 Header	52:54:00:cl:f8:ef, D /ersion 4, Src: 192. D Protocol, Src Por vice ge Block Protocol ve 0x05) ge Block Protocol ve (0x10) ge Block Protocol ve	<pre>śt: 52:54:00:55:3b:d4 168.124.203, Dst: 192 t: 52458, Dst Port: 4 rsion 2) rsion 2)</pre>	.168.124.1	·	
Et In Tr: Ne SM	nernet II, Srć: 5 cernet Protocol V ansmission Contro BIOS Session Ser 22 (Server Messag SMB2 Header Create Request (32 (Server Messag SMB2 Header GetInfo Request 32 (Server Messag 34 (Server Messag	52:54:00:cl:f8:ef, D /ersion 4, Src: 192. D Protocol, Src Por vice ge Block Protocol ve 0x05) ge Block Protocol ve (0x10) ge Block Protocol ve	<pre>śt: 52:54:00:55:3b:d4 168.124.203, Dst: 192 t: 52458, Dst Port: 4 rsion 2) rsion 2)</pre>	.168.124.1	·	

API

- You create an array of requests. One request at a time and set if they are related or not.
- The result is an array of iovectors, one vector per request.

First a CREATE at [0]

```
oparms.tcon = tcon;
oparms.desired_access = FILE_READ_ATTRIBUTES;
oparms.disposition = FILE_OPEN;
oparms.create_options = 0;
oparms.fid = &fid;
oparms.reconnect = false;
```

rc = SMB2_open_init(tcon, &rqst[0], &oplock, &oparms, &srch_path);
if (rc)
 goto qfs_exit;
smb2_set_next_command(&rqst[0]);

Then a QUERY INFO at [1]

rc = SMB2_query_info_init(tcon, &rqst[1], COMPOUND_FID, COMPOUND_FID, FS_FULL_SIZE_INFORMATION, SMB2_O_INFO_FILESYSTEM, 0, sizeof(struct smb2_fs_full_size_info)); if (rc) goto qfs_exit; smb2_set_next_command(&rqst[1]); smb2_set_related(&rqst[1]);

Finally a CLOSE at [2]

```
rc = SMB2_close_init(tcon, &rqst[2], COMPOUND_FID,
COMPOUND_FID);
if (rc)
goto qfs_exit;
smb2_set_related(&rqst[2]);
```

Send off the request

rsp_iov returns an array of 3 response vectors.

Better HA: Reconnect improvements

- Resilient and persistent handles are supported, and reconnect continues to improve
- Some remaining items:
 - Add lock sequence number
 - Fix EAGAIN rc which can occur for pending ops which overlap a reconnect
 - Reset credits on reconnect
 - Improve server to server failover
 - Allow alternate (failover) targets using DFS referrals
 - Witness protocol: server or share redirection

SMB3 and ACLs

 "cifsacl" mount option now supported for SMB3 for emulating mode bits via ACL

SMB3 Security Features

- SMB3.11 is no longer experimental, and works well
- SMB3.1.1 secure negotiate works (better than validate negotiate ioctl from SMB2.1 and SMB3)
- SMB3 and SMB3.11 Share Encryption works
 - AES128-CCM encryption algorithm is negotiated (AES128-GCM not supported yet for Linux client or Samba)

FSCTL passthrough ioctl ...

- Many interesting, useful features
 - Now we just need some python or C user space helpers to make them easier to use ...

Other Optional features

- statfs integration and new mount api integration
 - New API in Al Viro's tree
- IOCTLs e.g. to list alternate data streams
 - NB: Querying data in alternate data streams (e.g. for backup) requires disabling posix pathnames (due to conflict with ":")
- Clustering, Witness protocol integratio
- DFS reconnect to different DFS server
- Performance features (see next slides
- Other suggestions ...



Approach 3 – POSIX Extensions for SMB3!

See POSIX Extensions talk here!

root@Ubuntu-17-Virtual-Machine:~/cifs-2.6# cat /proc/mounts | grep cifs

//localhost/test-no-posix /mnt1 cifs rw,relatime,vers=3.1.1,cache=strict,username=testuser,domain=,uid=0,noforc euid,gid=0,noforcegid,addr=127.0.0.1,file mode=0755,dir mode=0755,soft,nounix,serverino,mapposix,rsize=1048576, wsize=1048576,echo interval=60,actimeo=1 0 0 //localhost/test /mnt **cifs** rw,relatime,vers=3.1.1,cache=strict,username=testuser,domain=,uid=0,noforceuid,gid=0 ,noforcegid,addr=127.0.0.1,file mode=0755,dir mode=0755,soft,posix,posix,paths,serverino,mapposix,rsize=1048576, wsize=1048576,echo interval=60,actimeo=1 0 0 root@Ubuntu-17-Virtual-Machine:~/cifs-2.6# cat /proc/fs/cifs/DebugData Display Internal CIFS Data Structures for Debugging CIFS Version 2.12 Features: dfs fscache lanman posix spnego xattr acl Active VFS Requests: 0 Servers: Number of credits: 16 Dialect 0x311 posix 1) Name: 127.0.0.1 Uses: 2 Capability: 0x300047 Session Status: 1 TCP status: 1 Local Users To Server: 1 SecMode: 0x1 Reg On Wire: 0 Shares: 0) IPC: \\127.0.0.1\IPC\$ Mounts: 1 DevInfo: 0x0 Attributes: 0x0 PathComponentMax: 0 Status: 1 type: 0 Share Capabilities: None Share Flags: 0x0 tid: 0x4f5511db Maximal Access: 0x1f00a9 1) \\localhost\test Mounts: 1 DevInfo: 0x20 Attributes: 0x1006f PathComponentMax: 255 Status: 1 type: DISK Share Capabilities: None Aligned, Partition Aligned, Share Flags: 0x0 tid: 0x8579c31d Optimal sector size: 0x200 Maximal Access: 0x1f01ff 2) \\localhost\test-no-posix Mounts: 1 DevInfo: 0x20 Attributes: 0x1006f PathComponentMax: 255 Status: 1 type: DISK Share Capabilities: None Aligned, Partition Aligned, Share Flags: 0x0 tid: 0x1813a493 Optimal sector size: 0x200 Maximal Access: 0x1f01ff

MIDs:

Mode bits on create and case sensitive!

```
root@Ubuntu-17-Virtual-Machine:/mnt# ~/create-4-files-with-mode-test
root@Ubuntu-17-Virtual-Machine:/mnt# cd /mnt1
root@Ubuntu-17-Virtual-Machine:/mnt1# ~/create-4-files-with-mode-test
root@Ubuntu-17-Virtual-Machine:/mnt1# ls /test /test-no-posix -la
/test:
total 12
drwxrwxrwx 3 root root 4096 May 31 16:55
drwxr-xr-x 32 root root 4096 May 31 16:46 ...
-rwx----- 1 testuser testuser 0 May 31 16:55 0700
-rwxrwx--- 1 testuser testuser 0 May 31 16:55 0770
-rwxrwxr-x 1 testuser testuser 0 May 31 16:55 0775
drwxr-xr-x 2 sfrench sfrench 4096 Mar 24 10:34 tmp
/test-no-posix:
total 8
drwxrwxrwx 2 root root 4096 May 31 16:55
drwxr-xr-x 32 root root 4096 May 31 16:46 ..
-rwxrw-r-- 1 testuser testuser 0 May 31 16:55 0700
-rwxrw-r-- 1 testuser testuser 0 May 31 16:55 0770
-rwxrw-r-- 1 testuser testuser 0 May 31 16:55 0775
root@Ubuntu-17-Virtual-Machine:/mnt1# mkdir UPPER
root@Ubuntu-17-Virtual-Machine:/mnt1# touch upper
root@Ubuntu-17-Virtual-Machine:/mnt1# cd /mnt
root@Ubuntu-17-Virtual-Machine:/mnt# mkdir UPPER
root@Ubuntu-17-Virtual-Machine:/mnt# touch upper
root@Ubuntu-17-Virtual-Machine:/mnt# ls /test /test-no-posix
/test:
0700 0770 0775 tmp upper UPPER
```

/test-no-posix: 0700 0770 0775 UPPER

Rename works with POSIX extensions!

root@Ubuntu-17-Virtual-Machine: ~ 🛛 😑 🗐 🧕	root@Ubuntu-17-Virtual-Machine: ~ 🔴 🗉 😒	
File Edit View Search Terminal Help	File Edit View Search Terminal Help	
<pre>root@Ubuntu-17-Virtual-Machine:~# ls /mnt-rename-test -la total 2052 drwxr-xr-x 2 root root 0 May 31 18:19 . Idrwxr-xr-x 34 root root 4096 May 31 18:13 t-rwxr-xr-x 1 root root 0 May 31 18:18 emptyfile -rwxr-xr-x 1 root root 0 May 31 18:19 emptyfile-posix /-rwxr-xr-x 1 root root 16 May 31 18:19 targetfile -rwxr-xr-x 1 root root 16 May 31 18:19 targetfile -rwxr-xr-x 1 root root 16 May 31 18:19 targetfile -rwar-xr-xr-x 1 root root 16 May 31 18:19 targetfile -rwar-xr-x 1 root root 16 May 31 18:19 targetfile -rwar-xr-x 1 root root 16 May 31 18:19 targetfile -rwar-xr-x 1 root root 16 May 31 18:19 targetfile -rwar-xr-x 1 root root 16 May</pre>	<pre>root@Ubuntu-17-Virtual-Machine:-# tail -f /mnt-rename-test/targetfile targetfile data tail: /mnt-rename-test/targetfile: No such file or directory tail: no files remaining root@Ubuntu-17-Virtual-Machine:-# [] ,soft,posix,posixpaths,serverino,mapposix,rsize=1048576,wsf he=strict,username=testuser,domain=,uid=0,noforceuid,gid=0, happosix,rsize=1048576,wsize=1048576,echo_interval=60,actime</pre>	,noforc
<pre>root@Ubuntu-17-Virtual-Machine:~# mv /mnt-rename-test/emptyfile /mnt-rename-test/targetfileo mv: cannot move '/mnt-rename-test/emptyfile' to '/mnt-rename-test/targetfile': Permission de nied</pre>		

root@Ubuntu-17-Virtual-Machine: ~

File Edit View Search Terminal Help

root@Ubuntu-17-Virtual-Machine:~# mount | grep rename

//localhost/rename-test on /mnt-rename-test type cifs (rw,relatime,vers=3.1.1,cache=strict,u
sername=testuser,domain=,uid=0,noforceuid,gid=0,noforcegid,addr=127.0.0.1,file_mode=0755,dir
_mode=0755,soft,posix,posixpaths,serverino,mapposix,rsize=1048576,wsize=1048576,echo_interva
l=60,actimeo=1)

<mark>root@Ubuntu-17-Virtual-Machine:</mark>~# mv /mnt-rename-test/emptyfile-posix /mnt-rename-test/targe tfile-posix

root@Ubuntu-17-Virtual-Machine:~#

SMB3 Performance – the Myth

 Googling NFS vs. SMB3 (or Samba) ... first result said:

"As you can see NFS offers a better performance and is unbeatable if the files are medium sized or small. If the files are large enough the timings of both methods get closer to each other. Linux and Mac OS owners should use NFS instead of SMB. Sadly Windows users are forced to use SMB ..."

Is NFS really always faster than Samba...





SMB3 to Samba is faster in many cases

- Localhost (network shouldn't be an issue. Default Ubuntu Samba server vs. NFS kernel server. Default parms. Comparing NFSv3, NFSv4.2 and cifs.ko (SMB3.02 dialect is default)
- fio with the read/write job file : SMB3 12.5% faster to Samba (than NFSv4.2 server) for random reads and SMB3 12.8% faster for writes
- For sequential: SMB3 31.8% faster for read, 31.2% faster for write (and not just because of stricter sync)
- Even simple DD command with large file i/o shows SMB3 much faster Linux to Linux for write than NFS

Just last night ... 1st test I tried SMB3 wins by 29% over NFS (defaults, localhost mounts)

-oot@smf-Thinkpad-P51:~/cifs-2.6# cat /proc/mounts | grep nfs

sd /proc/fs/nfsd nfsd rw,relatime 0 0 localhost:/nfsexport /mnt2 nfs4 rw.relatime.vers=4.2.rsize=1048576.wsize=1048576.namlen=255.hard.proto=tc 🖷 p,timeo=600,retrans=2,sec=sys,clientaddr=127.0.0.1,local lock=none,addr=127.0.0.1 0 0 root@smf-Thinkpad-P51:~/cifs-2.6# dd if=/dev/zero of=/mnt2/targetfile bs=10M count=100 100+0 records in 100+0 records out 1048576000 bytes (1.0 GB, 1000 MiB) copied, 1.83421 s, 572 MB/s root@smf-Thinkpad-P51:~/cifs-2.6# dd if=/dev/zero of=/mnt2/targetfile bs=10M count=100 100+0 records in 100+0 records out 1048576000 bytes (1.0 GB, 1000 MiB) copied, 1.67055 s, 628 MB/s root@smf-Thinkpad-P51:~/cifs-2.6# dd if=/dev/zero of=/mnt2/targetfile bs=10M count=100 100+0 records in 100+0 records out 1048576000 bytes (1.0 GB, 1000 MiB) copied, 1.80421 s, 581 MB/s root@smf-Thinkpad-P51:~/cifs-2.6# dd if=/dev/zero of=/mnt2/targetfile bs=10M count=100 100+0 records in 100+0 records out 1048576000 bytes (1.0 GB, 1000 MiB) copied, 1.80514 s, 581 MB/s root@smf-Thinkpad-P51:~/cifs-2.6# umount /mnt2 root@smf-Thinkpad-P51:~/cifs-2.6# mount | grep cifs root@smf-Thinkpad-P51:~/cifs-2.6# mount -t cifs //localhost/scratch /mnt2 -o username=sfrench,noperm Password for sfrench@//localhost/scratch: ********** root@smf-Thinkpad-P51:~/cifs-2.6# dd if=/dev/zero of=/mnt2/targetfile bs=10M count=100 100+0 records in 100+0 records out 1048576000 bytes (1.0 GB, 1000 MiB) copied, 0.834104 s, 1.3 GB/s root@smf-Thinkpad-P51:~/cifs-2.6# dd if=/dev/zero of=/mnt2/targetfile bs=10M count=100 100+0 records in 100+0 records out 1048576000 bytes (1.0 GB, 1000 MiB) copied, 1.76119 s, 595 MB/s root@smf-Thinkpad-P51:~/cifs-2.6# dd if=/dev/zero of=/mnt2/targetfile bs=10M count=100 100+0 records in 100+0 records out 1048576000 bytes (1.0 GB, 1000 MiB) copied, 1.76155 s, 595 MB/s root@smf-Thinkpad-P51:~/cifs-2.6# dd if=/dev/zero of=/mnt2/targetfile bs=10M count=100 100+0 records in 100+0 records out 1048576000 bytes (1.0 GB, 1000 MiB) copied, 1.78004 s, 589 MB/s root@smf-Thinkpad-P51:~/cifs-2.6# mount | grep cifs //localhost/scratch on /mnt2 type cifs (rw,relatime,vers=default,cache=strict,username=sfrench,domain=,ui d=0,noforceuid,gid=0,noforcegid,addr=127.0.0.1,file_mode=0755,dir_mode=0755,soft,nounix,serverino,mapposi x,noperm,rsize=1048576,wsize=1048576,echo interval=60,actimeo=1) root@smf-Thinkpad-P51:~/cifs-2.6# dd of=/dev/zero if=/mnt2/targetfile bs=10M count=100 100+0 records in 100+0 records out 1048576000 bytes (1.0 GB. 1000 MiB) copied. 0.244735 s. 4.3 GB/s

Maybe coincidence so lets try fio ... (at 1am!)

- Standard fio random read/write i/o job file, localhost Samba vs. NFS, using all defaults
- /mnt2: fio ~/fio/fio-rand-RW.job
- SMB3 20% faster than NFS for read, 21% for write

READ: bw=204MiB/s (214MB/s), 51.1MiB/s-51.1MiB/s (53.6MB/s-53.6MB/s), io=17.0GiB (19.3GB), run=90001-90001msec WRITE: bw=136MiB/s (143MB/s), 34.0MiB/s-34.1MiB/s (35.7MB/s-35.7MB/s), io=11.0GiB (12.9GB), run=90001-90001msec sfrench@smf-Thinkpad-P51:/mnt2\$ mount | grep mnt2 //localhost/scratch on /mnt2 type cifs (rw,relatime,vers=default,cache=none,username=sfrench,domain=,uid=0,noforceu =-0755 dir mode=0755 soft pount servering mapposix poperm rsize=2097152 wsize=2097152 echo interval=60 actimeo=1)

Run status group 0 (all jobs): READ: bw=170MiB/s (178MB/s), 42.5MiB/s-42.6MiB/s (44.6MB/s-44.7MB/s), io=14.0GiB (16.1GB), run=90001-90001msec WRITE: bw=113MiB/s (119MB/s), 28.3MiB/s-28.4MiB/s (29.7MB/s-29.7MB/s), io=9.97GiB (10.7GB), run=90001-90001msec sfrench@smf-Thinkpad-P51:/mnt2\$ mount | grep mnt2 localhost:/nfsexport on /mnt2 type nfs4 (rw,relatime,vers=4.2,rsize=1048576,wsize=1048576,namlen=255,hard,proto=tcp .0.0.1,local_lock=none,addr=127.0.0.1) sfrench@smf-Thinkpad-P51:/mnt2\$

SMB3 Performance WIP: great features ... but only if we implement them ...

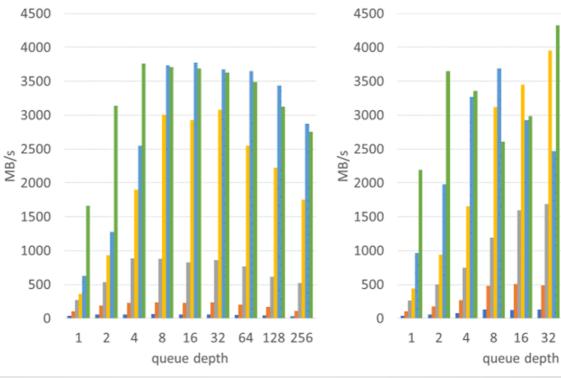
- Key Features
 - Compounding
 - Large file I/O
 - File Leases
 - Lease upgrades
 - Directory Leases
 - Handle caching
 - Crediting
 - I/O priority
 - Copy Offload
 - Multi-Channel
 - And optional RDMA
 - Linux specific protocol optimizations possible too ...

We have fun work to do ... (go to Long Li's talk to hear exciting improvements!)

- And not just for metadata heavy workloads
- But the SMB3 protocol is richer, more function that can help performance when implemented fully in client



• 85% IWarp



4K

■ 16K
■ 64K
■ 256K

■ 1M

64 128 256

Conclusion ... When is SMB3 good?

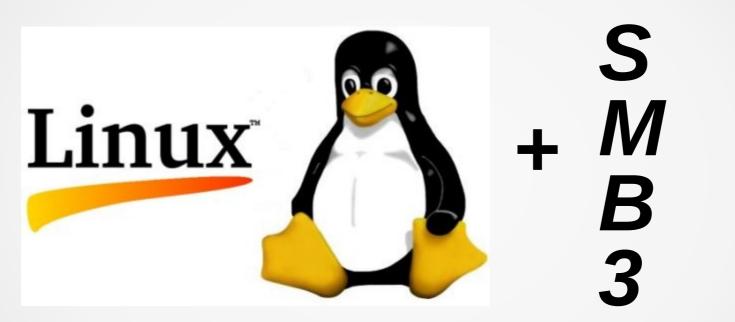
- When need nice security ...
- Workloads where performance with lots of large directories is not an obstacle (pending improvements to leasing and compounding in cifs.ko)
- Workloads which do not depend on case sensitivity (common unfortunately) and do not depend on advisory locking or delete of open files (more rare) ... (pending POSIX extensions being merged into Samba etc.)
- Where you can take advantage of smbdirect (RDMA)
- Where global namespace (DFS) helps
- Where rich features of SMB3 (snapshots, encrypted/compressed files, persistent handles) are helpful ...
- And of course ... to the cloud (Azure) and Macs and Windows and ... not just Samba

Testing ... testing ... testing

- See xfstesting page in cifs wiki https://wiki.samba.org/index.php/Xfstesting-cifs
- Easy to setup, exclude file for slow tests or failing ones
- XFSTEST status update
 - Bugzillas
 - Features in progress
 - Automating improvements

Thank you for your time

• Future is very bright!



Additional Resources to Explore for SMB3 and Linux

- https://msdn.microsoft.com/en-us/library/gg685446.aspx
 - In particular MS-SMB2.pdf at https://msdn.microsoft.com/en-us/library/cc246482.aspx
- https://wiki.samba.org/index.php/Xfstesting-cifs
- Linux CIFS client https://wiki.samba.org/index.php/LinuxCIFS
- Samba-technical mailing list and IRC channel
- And various presentations at http://www.sambaxp.org and Microsoft channel 9 and of course SNIA ... http://www.snia.org/events/storage-developer
- And the code:
 - https://git.kernel.org/cgit/linux/kernel/git/torvalds/linux.git/tree/fs/cifs
 - For pending changes, soon to go into upstream kernel see:
 - https://git.samba.org/?p=sfrench/cifs-2.6.git;a=shortlog;h=refs/heads/ for-next