

## Ceph - Bug #1356

### OSD crashes during recovery with OSDMap::decode(ceph::buffer::list&)

08/04/2011 07:45 AM - Wido den Hollander

<b>Status:</b>	Can't reproduce	<b>Start date:</b>	08/04/2011
<b>Priority:</b>	Normal	<b>Due date:</b>	
<b>Assignee:</b>	Sage Weil	<b>% Done:</b>	0%
<b>Category:</b>	OSD	<b>Estimated time:</b>	0.00 hour
<b>Target version:</b>	v0.34	<b>Spent time:</b>	5.50 hours
<b>Source:</b>		<b>Reviewed:</b>	
<b>Tags:</b>		<b>Affected Versions:</b>	
<b>Backport:</b>		<b>ceph-qa-suite:</b>	
<b>Regression:</b>	No	<b>Pull request ID:</b>	
<b>Severity:</b>	3 - minor	<b>Crash signature:</b>	

#### Description

Hi,

Like I said some time ago, I've been seeing these kind of crashes lately.

I just tried to start my cluster (40 OSD's) up again and as always it started bouncing around, during this bouncing I saw a couple of OSD's going down with:

```
(gdb) bt
#0 0x00007f8afb7967bb in raise () from /lib/libpthread.so.0
#1 0x000000000057cdc3 in reraise_fatal (signum=2382) at global/signal_handler.cc:59
#2 0x000000000057d38c in handle_fatal_signal (signum=<value optimized out>) at global/signal_handler.cc:106
#3 <signal handler called>
#4 0x00007f8afa3c9a75 in raise () from /lib/libc.so.6
#5 0x00007f8afa3cd5c0 in abort () from /lib/libc.so.6
#6 0x00007f8afac7f8e5 in __gnu_cxx::__verbose_terminate_handler() () from /usr/lib/libstdc++.so.6
#7 0x00007f8afac7dd16 in ?? () from /usr/lib/libstdc++.so.6
#8 0x00007f8afac7dd43 in std::terminate() () from /usr/lib/libstdc++.so.6
#9 0x00007f8afac7de3e in __cxa_throw () from /usr/lib/libstdc++.so.6
#10 0x000000000049c3a6 in ceph::buffer::list::iterator::advance (this=0x7fff9327f4c0, len=2, dest=0x7fff9327f56c "\377\177") at ./include/buffer.h:315
#11 ceph::buffer::list::iterator::copy (this=0x7fff9327f4c0, len=2, dest=0x7fff9327f56c "\377\177") at ./include/buffer.h:369
#12 0x00000000005630c8 in OSDMap::decode(ceph::buffer::list&) ()
#13 0x000000000052f571 in OSD::get_map (this=0x1f1aca0, epoch=8974) at osd/OSD.cc:3492
#14 0x000000000053a5ce in OSD::init (this=0x1f1aca0) at osd/OSD.cc:555
#15 0x000000000049a6ca in main (argc=<value optimized out>, argv=<value optimized out>) at cosd.cc:298
(gdb)
```

The logging was very low since it kills the nodes even more if I increase the level, but what I do have is:

```
2011-08-04 16:28:37.243041 7f8afb7720 journal read_entry 417271808 : seq 1523780 483 bytes
2011-08-04 16:28:37.243138 7f8afb7720 journal read_entry 417280000 : seq 1523781 2547 bytes
2011-08-04 16:28:37.243256 7f8afb7720 journal read_entry 417288192 : seq 1523782 483 bytes
2011-08-04 16:28:37.243351 7f8afb7720 journal read_entry 417296384 : seq 1523783 3677 bytes
2011-08-04 16:28:37.243445 7f8afb7720 journal read_entry 417304576 : seq 1523784 483 bytes
2011-08-04 16:28:37.248665 7f8afb7720 journal read_entry 417312768 : seq 1523785 502916 bytes
2011-08-04 16:28:37.253398 7f8afb7720 journal kernel version is 2.6.39
2011-08-04 16:28:37.253994 7f8afb7720 journal _open /dev/data/journal0 fd 12: 1996488704 bytes,
block size 4096 bytes, directio = 1
*** Caught signal (Aborted) **
```

```
in thread 0x7f8afbbb7720
ceph version 0.31 (commit:9019c6ce64053ad515a493e912e2e63ba9b8e278)
1: /usr/bin/cosd() [0x57d154]
2: (()+0xf8f0) [0x7f8afb7968f0]
3: (gsignal()+0x35) [0x7f8afa3c9a75]
4: (abort()+0x180) [0x7f8afa3cd5c0]
5: (__gnu_cxx::__verbose_terminate_handler()+0x115) [0x7f8afac7f8e5]
6: (()+0xcad16) [0x7f8afac7dd16]
7: (()+0xcad43) [0x7f8afac7dd43]
8: (()+0xcae3e) [0x7f8afac7de3e]
9: (ceph::buffer::list::iterator::copy(unsigned int, char*)+0x156) [0x49c3a6]
10: (OSDMap::decode(ceph::buffer::list&)+0x78) [0x5630c8]
11: (OSD::get_map(unsigned int)+0x221) [0x52f571]
12: (OSD::init()+0x47e) [0x53a5ce]
13: (main()+0x25ea) [0x49a6ca]
14: (__libc_start_main()+0xfd) [0x7f8afa3b4c4d]
15: /usr/bin/cosd() [0x497cd9]
```

There is no real way to reproduce it, if I start this particular OSD again it will probably go one, but it could also be that it crashes, you never know.

#### Related issues:

Duplicated by Ceph - Bug #1486: osd: 0-length meta/pginfo\_\* files

Resolved

09/01/2011

#### History

##### #1 - 08/04/2011 09:31 AM - Sage Weil

- Assignee set to Sage Weil

- Target version set to v0.33

##### #2 - 08/04/2011 11:16 AM - Sage Weil

do you have a core file? in frame 10, can you look at bt->\_len ?

also, can you attach the osdmap\_full\_8974 from this node to the bug? should be something like current/meta/osdmap\_full\_8974\_head.

thanks!

##### #3 - 08/04/2011 11:47 AM - Wido den Hollander

The osdmap is attached, picked it from the mon out of osdmap\_full.

I took a look at the core dump (/core.atom2.3816 on atom2.ceph.widodh.nl), bt->\_len=4 if I'm correct.

##### #4 - 08/04/2011 11:50 AM - Wido den Hollander

- File osdmap.8974\_0 added

I also attached the osdmap from osd.9.

##### #5 - 08/04/2011 11:52 AM - Wido den Hollander

- File 8974.osdmap added

I'm doing something very wrong here... Attached is osdmap 8974 from the monitor (again).

##### #6 - 08/04/2011 12:58 PM - Sage Weil

This looks like it's crashing during cosd start-up... can you confirm what's going on here? Basically,

- osds are crashing
- someone is restarting them
- some of them crash while restarting

Is that right?

#### #7 - 08/04/2011 01:02 PM - Wido den Hollander

Yes, that's like it's going. In the process where they are all starting again they start bouncing up and down, in this process some start to crash.

#### #8 - 08/04/2011 01:15 PM - Wido den Hollander

I just started the OSD again with debug osd and filestore on 20, got a different backtrace:

```
(gdb) bt
#0 0x00007f8c28c0b7bb in raise () from /lib/libpthread.so.0
#1 0x00000000057cdc3 in reraise_fatal (signum=7730) at global/signal_handler.cc:59
#2 0x00000000057d38c in handle_fatal_signal (signum=<value optimized out>) at global/signal_handler.cc:106
#3 <signal handler called>
#4 0x00007f8c2783ea75 in raise () from /lib/libc.so.6
#5 0x00007f8c278425c0 in abort () from /lib/libc.so.6
#6 0x00007f8c280f48e5 in __gnu_cxx::__verbose_terminate_handler() () from /usr/lib/libstdc++.so.6
#7 0x00007f8c280f2d16 in ?? () from /usr/lib/libstdc++.so.6
#8 0x00007f8c280f2d43 in std::terminate() () from /usr/lib/libstdc++.so.6
#9 0x00007f8c280f2e3e in __cxa_throw () from /usr/lib/libstdc++.so.6
#10 0x00000000049c3a6 in ceph::buffer::list::iterator::advance (this=0x7fff80ddefc0, len=4, dest=0x7fff80dded68 "81W") at ./include/buffer.h:315
#11 ceph::buffer::list::iterator::copy (this=0x7fff80ddefc0, len=4, dest=0x7fff80dded68 "81W") at ./include/buffer.h:369
#12 0x000000000643f91 in void decode<unsigned int, PG::Interval>(std::map<unsigned int, PG::Interval, std::less<unsigned int>, std::allocator<std::pair<unsigned int const, PG::Interval> > >&, ceph::buffer::list::iterator&) ()
#13 0x0000000000627c7c in PG::read_state (this=0x229fa90, store=<value optimized out>) at osd/PG.cc:2415
#14 0x000000000051eb82 in OSD::load_pgs (this=0x21d7c20) at osd/OSD.cc:1095
#15 0x000000000053a7db in OSD::init (this=0x21d7c20) at osd/OSD.cc:568
#16 0x000000000049a6ca in main (argc=<value optimized out>, argv=<value optimized out>) at cosd.cc:298
(gdb)
```

The logs show:

```
2011-08-04 22:08:06.412484 7f8c2902c720 osd9 9263 _open_lock_pg 2.935
2011-08-04 22:08:06.412584 7f8c2902c720 osd9 9263 _get_pool 2 7 -> 8
2011-08-04 22:08:06.412771 7f8c2902c720 osd9 9263 pg[2.935( DNE empty n=0 ec=0 les/c 0/0 0/0/0) [] r=0 mlcod 0 '0 inactive] enter Initial
2011-08-04 22:08:06.412939 7f8c2902c720 osd9 9263 pg[2.935( DNE empty n=0 ec=0 les/c 0/0 0/0/0) [] r=0 mlcod 0 '0 inactive] enter NotTrimming
2011-08-04 22:08:06.413057 7f8c2902c720 filestore(/var/lib/ceph/osd.9) collection_getattr /var/lib/ceph/osd.9/current/2.935_head 'info'
2011-08-04 22:08:06.413352 7f8c2902c720 filestore(/var/lib/ceph/osd.9) collection_getattr /var/lib/ceph/osd.9/current/2.935_head 'info' = 309
2011-08-04 22:08:06.413463 7f8c2902c720 filestore(/var/lib/ceph/osd.9) read meta/pginfo_2.935/0 0~0
2011-08-04 22:08:06.413556 7f8c2902c720 filestore(/var/lib/ceph/osd.9) lfn_get cid=meta oid=pginfo_2.935/0 pat hname=/var/lib/ceph/osd.9/current/meta/pginfo_2.935_0 lfn=pginfo_2.935_0 is_lfn=0
2011-08-04 22:08:06.413838 7f8c2902c720 filestore(/var/lib/ceph/osd.9) FileStore::read meta/pginfo_2.935/0 0~2 348/2348
2011-08-04 22:08:06.414138 7f8c2902c720 filestore(/var/lib/ceph/osd.9) collection_getattr /var/lib/ceph/osd.9/current/2.935_head 'ondisklog'
2011-08-04 22:08:06.414314 7f8c2902c720 filestore(/var/lib/ceph/osd.9) collection_getattr /var/lib/ceph/osd.9/current/2.935_head 'ondisklog' = 17
2011-08-04 22:08:06.414501 7f8c2902c720 osd9 9263 pg[2.935( empty n=0 ec=2 les/c 9259/9263 8998/8998/8912) []
```

```

r=0 mlcod 0'0 inactive] read_log 0~0
2011-08-04 22:08:06.414634 7f8c2902c720 osd9 9263 pg[2.935( empty n=0 ec=2 les/c 9259/9263 8998/8998/8912) []
r=0 mlcod 0'0 inactive] read_log done
2011-08-04 22:08:06.414965 7f8c2902c720 osd9 9263 pg[2.935( empty n=0 ec=2 les/c 9259/9263 8998/8998/8912) []/
[9,33] r=0 mlcod 0'0 inactive] handle_backlog_loaded
2011-08-04 22:08:06.415087 7f8c2902c720 osd9 9263 pg[2.935( empty n=0 ec=2 les/c 9259/9263 8998/8998/8912) []/
[9,33] r=0 mlcod 0'0 inactive] exit Initial 0.002316 0 0.000000
2011-08-04 22:08:06.415231 7f8c2902c720 osd9 9263 pg[2.935( empty n=0 ec=2 les/c 9259/9263 8998/8998/8912) []/
[9,33] r=0 mlcod 0'0 inactive] enter Reset
2011-08-04 22:08:06.415345 7f8c2902c720 osd9 9263 load_pgs loaded pg[2.935( empty n=0 ec=2 les/c 9259/9263 899
8/8998/8912) []/[9,33] r=0 mlcod 0'0 inactive] log(0'0,0'0)
2011-08-04 22:08:06.415462 7f8c2902c720 osd9 9263 _open_lock_pg 0.937
2011-08-04 22:08:06.415562 7f8c2902c720 osd9 9263 _get_pool 0 7 -> 8
2011-08-04 22:08:06.415757 7f8c2902c720 osd9 9263 pg[0.937( DNE empty n=0 ec=0 les/c 0/0 0/0/0) [] r=0 mlcod 0
'0 inactive] enter Initial
2011-08-04 22:08:06.415917 7f8c2902c720 osd9 9263 pg[0.937( DNE empty n=0 ec=0 les/c 0/0 0/0/0) [] r=0 mlcod 0
'0 inactive] enter NotTrimming
2011-08-04 22:08:06.416035 7f8c2902c720 filestore(/var/lib/ceph/osd.9) collection_getattr /var/lib/ceph/osd.9/
current/0.937_head 'info'
2011-08-04 22:08:06.416320 7f8c2902c720 filestore(/var/lib/ceph/osd.9) collection_getattr /var/lib/ceph/osd.9/
current/0.937_head 'info' = 309
2011-08-04 22:08:06.416425 7f8c2902c720 filestore(/var/lib/ceph/osd.9) read meta/pginfo_0.937/0 0~0
2011-08-04 22:08:06.416513 7f8c2902c720 filestore(/var/lib/ceph/osd.9) lfn_get cid=meta oid=pginfo_0.937/0 pat
hname=/var/lib/ceph/osd.9/current/meta/pginfo_0.937_0 lfn=pginfo_0.937_0 is_lfn=0
2011-08-04 22:08:06.416748 7f8c2902c720 filestore(/var/lib/ceph/osd.9) FileStore::read meta/pginfo_0.937/0 0~0
/0
*** Caught signal (Aborted) **
in thread 0x7f8c2902c720
ceph version 0.31 (commit:9019c6ce64053ad515a493e912e2e63ba9b8e278)
1: /usr/bin/cosd() [0x57d154]
2: (()+0xf8f0) [0x7f8c28c0b8f0]
3: (gsignal()+0x35) [0x7f8c2783ea75]
4: (abort()+0x180) [0x7f8c278425c0]
5: (__gnu_cxx::__verbose_terminate_handler()+0x115) [0x7f8c280f48e5]
6: (()+0xcad16) [0x7f8c280f2d16]
7: (()+0xcad43) [0x7f8c280f2d43]
8: (()+0xcae3e) [0x7f8c280f2e3e]
9: (ceph::buffer::list::iterator::copy(unsigned int, char*)+0x156) [0x49c3a6]
10: (void decode<unsigned int, PG::Interval>(std::map<unsigned int, PG::Interval, std::less<unsigned int>, st
d::allocator<std::pair<unsigned int const, PG::Interval> > >&, ceph::buffer::list::iterator&)+0x31) [0x643f91]
11: (PG::read_state(ObjectStore*)+0x2cc) [0x627c7c]
12: (OSD::load_pgs()+0x272) [0x51eb82]
13: (OSD::init()+0x68b) [0x53a7db]
14: (main()+0x25ea) [0x49a6ca]
15: (__libc_start_main()+0xfd) [0x7f8c27829c4d]
16: /usr/bin/cosd() [0x497cd9]

```

I think `/var/lib/ceph/osd.9/current/meta/pginfo_0.937_0` is the problem here, it has a size of 0 bytes:

```

root@atom2:~# stat /var/lib/ceph/osd.9/current/meta/pginfo_0.937_0
  File: `/var/lib/ceph/osd.9/current/meta/pginfo_0.937_0'
  Size: 0          Blocks: 0          IO Block: 4096   regular empty file
Device: 20h/32d   Inode: 357          Links: 1
Access: (0644/-rw-r--r--)  Uid: (   0/   root)   Gid: (   0/   root)
Access: 2011-06-30 13:52:06.438052123 +0200
Modify: 2011-08-04 15:26:11.882911939 +0200
Change: 2011-08-04 15:26:11.882911939 +0200
root@atom2:~#

```

**#9 - 08/04/2011 01:22 PM - Sage Weil**

can you look in the snap\_\* directories and see if that pglog file is 0 in those too?

**#10 - 08/04/2011 01:48 PM - Wido den Hollander**

Just to sum up what we said on IRC:

osd.9 has 47 pginfo\* files which are empty in the current dir and the two snapshots.

Even worse, osd.4 has a 148 empty pginfo\* files!

I then tried to start osd.4 with full debugging and got:

```
2011-08-04 22:45:50.495007 7f3a36395700 filestore(/var/lib/ceph/osd.4) sync_entry waiting for max_interval 5.0
00000
2011-08-04 22:45:50.495209 7f3a34b92700 filestore(/var/lib/ceph/osd.4) flusher_entry start
2011-08-04 22:45:50.495297 7f3a34b92700 filestore(/var/lib/ceph/osd.4) flusher_entry sleeping
2011-08-04 22:45:50.495662 7f3a3d3c2720 osd4 0 boot
2011-08-04 22:45:50.495824 7f3a3d3c2720 filestore(/var/lib/ceph/osd.4) read meta/osd_superblock/0 0~0
2011-08-04 22:45:50.495969 7f3a3d3c2720 filestore(/var/lib/ceph/osd.4) lfn_get cid=meta oid=osd_superblock/0 p
athname=/var/lib/ceph/osd.4/current/meta/osd_superblock_0 lfn=osd_superblock_0 is_lfn=0
2011-08-04 22:45:50.496272 7f3a3d3c2720 filestore(/var/lib/ceph/osd.4) FileStore::read meta/osd_superblock/0 0
~123/123
2011-08-04 22:45:50.496433 7f3a3d3c2720 osd4 0 read_superblock sb(a235c24e-5dc8-5cad-b84c-816eadd854bb osd4 e8
974 [1,8974] lci=[3,8974])
2011-08-04 22:45:50.496563 7f3a3d3c2720 osd4 0 get_map 8974 - loading and decoding 0xee6040
2011-08-04 22:45:50.496660 7f3a3d3c2720 filestore(/var/lib/ceph/osd.4) read meta/osdmap.8974/0 0~0
2011-08-04 22:45:50.496743 7f3a3d3c2720 filestore(/var/lib/ceph/osd.4) lfn_get cid=meta oid=osdmap.8974/0 path
name=/var/lib/ceph/osd.4/current/meta/osdmap.8974_0 lfn=osdmap.8974_0 is_lfn=0
2011-08-04 22:45:50.496914 7f3a3d3c2720 filestore(/var/lib/ceph/osd.4) FileStore::read(meta/osdmap.8974/0): op
en error error 2: No such file or directory
*** Caught signal (Aborted) **
in thread 0x7f3a3d3c2720
ceph version 0.31 (commit:9019c6ce64053ad515a493e912e2e63ba9b8e278)
1: /usr/bin/cosd() [0x57d154]
2: (()+0xf8f0) [0x7f3a3cfa18f0]
3: (gsignal()+0x35) [0x7f3a3bbd4a75]
4: (abort()+0x180) [0x7f3a3bbd85c0]
5: (__gnu_cxx::__verbose_terminate_handler()+0x115) [0x7f3a3c48a8e5]
6: (()+0xcad16) [0x7f3a3c488d16]
7: (()+0xcad43) [0x7f3a3c488d43]
8: (()+0xcae3e) [0x7f3a3c488e3e]
9: (ceph::buffer::list::iterator::copy(unsigned int, char*)+0x156) [0x49c3a6]
10: (OSDMap::decode(ceph::buffer::list&)+0x78) [0x5630c8]
11: (OSD::get_map(unsigned int)+0x221) [0x52f571]
12: (OSD::init()+0x47e) [0x53a5ce]
13: (main()+0x25ea) [0x49a6ca]
14: (__libc_start_main()+0xfd) [0x7f3a3bbbfc4d]
15: /usr/bin/cosd() [0x497cd9]
```

Notice the fact that it's about the same osdmap, 8974!

**#11 - 08/08/2011 09:18 AM - Sage Weil**

- Target version changed from v0.33 to v0.34

**#12 - 08/08/2011 09:25 AM - Sage Weil**

- translation missing: en.field\_position set to 21

**#13 - 08/08/2011 03:30 PM - Sage Weil**

Wido den Hollander wrote:

Notice the fact that it's about the same osdmap, 8974!

Does that osdmap file size match the one on the monitor? md5sum too?

This looks like a low level corruption somewhere in the filestore or btrfs, since the osd *never* creates a pginfo without writing data to it. Either the interaction with the btrfs clone is broken, or btrfs itself is screwing up.

Which kernel version are you running?

**#14 - 08/09/2011 03:57 AM - Wido den Hollander**

Sage Weil wrote:

Wido den Hollander wrote:

Notice the fact that it's about the same osdmap, 8974!

Does that osdmap file size match the one on the monitor? md5sum too?

You might have missed it:

```
2011-08-04 22:45:50.496743 7f3a3d3c2720 filestore(/var/lib/ceph/osd.4) lfn_get cid=meta oid=osdmap.8974/0 path
name=/var/lib/ceph/osd.4/current/meta/osdmap.8974_0 lfn=osdmap.8974_0 is_lfn=0
2011-08-04 22:45:50.496914 7f3a3d3c2720 filestore(/var/lib/ceph/osd.4) FileStore::read(meta/osdmap.8974/0): op
en error error 2: No such file or directory
```

The whole osdmap is missing on that particular OSD, nothing to compare it to.

This looks like a low level corruption somewhere in the filestore or btrfs, since the osd *never* creates a pginfo without writing data to it. Either the interaction with the btrfs clone is broken, or btrfs itself is screwing up.

FYI, the filesystem has never been filled to a 100%, the average was about 20%.

Which kernel version are you running?

2.6.39.2 on all the OSD's.

**#15 - 08/10/2011 03:29 AM - Wido den Hollander**

I couldn't resist the urge to try, so I removed all the empty pginfo files and their corresponding *head directories in the snap* and current directories. osd.4 and osd.9 are up and running again, not sure what they will do on the long run.

**#16 - 08/10/2011 07:34 AM - Sage Weil**

Great! You were lucky that it didn't affect all replicas for any of the PGs.

**#17 - 08/10/2011 07:48 AM - Wido den Hollander**

The thing is, none of the PGs were assigned to these OSDs....

**#18 - 08/10/2011 01:08 PM - Sage Weil**

- *Status changed from New to Can't reproduce*

Let's keep an eye out for future appearances of 0-byte pginfo files, as per our conversation on irc yesterday.

**Files**

---

osdmap.8974_0	32.5 KB	08/04/2011	Wido den Hollander
8974.osdmap	32.5 KB	08/04/2011	Wido den Hollander