## btrfs: fix deadlock with concurrent chunk allocations involving syste…

**Browse files**

...m chunks

When a task attempting to allocate a new chunk verifies that there is not
currently enough free space in the system space_info and there is another
task that allocated a new system chunk but it did not finish yet the
creation of the respective block group, it waits for that other task to
finish creating the block group. This is to avoid exhaustion of the system
chunk array in the superblock, which is limited, when we have a thundering
herd of tasks allocating new chunks. This problem was described and fixed
by commit eafa4fd ("btrfs: fix exhaustion of the system chunk array
due to concurrent allocations").

However there are two very similar scenarios where this can lead to a
deadlock:

1) Task B allocated a new system chunk and task A is waiting on task B
   to finish creation of the respective system block group. However before
   task B ends its transaction handle and finishes the creation of the
   system block group, it attempts to allocate another chunk (like a data
   chunk for an fallocate operation for a very large range). Task B will
   be unable to progress and allocate the new chunk, because task A set
   space_info->chunk_alloc to 1 and therefore it loops at
   btrfs_chunk_alloc() waiting for task A to finish its chunk allocation
   and set space_info->chunk_alloc to 0, but task A is waiting on task B
   to finish creation of the new system block group, therefore resulting
   in a deadlock;

2) Task B allocated a new system chunk and task A is waiting on task B to
   finish creation of the respective system block group. By the time that
   task B enter the final phase of block group allocation, which happens
   at btrfs_create_pending_block_groups(), when it modifies the extent
   tree, the device tree or the chunk tree to insert the items for some
   new block group, it needs to allocate a new chunk, so it ends up at
   btrfs_chunk_alloc() and keeps looping there because task A has set
   space_info->chunk_alloc to 1, but task A is waiting for task B to
   finish creation of the new system block group and release the reserved
   system space, therefore resulting in a deadlock.

In short, the problem is if a task B needs to allocate a new chunk after
it previously allocated a new system chunk and if another task A is
currently waiting for task B to complete the allocation of the new system
chunk.

Unfortunately this deadlock scenario introduced by the previous fix for
the system chunk array exhaustion problem does not have a simple and short
fix, and requires a big change to rework the chunk allocation code so that
chunk btree updates are all made in the first phase of chunk allocation.
And since this deadlock regression is being frequently hit on zoned
filesystems and the system chunk array exhaustion problem is triggered
in more extreme cases (originally observed on PowerPC with a node size
of 64K when running the fallocate tests from stress-ng), revert the
changes from that commit. The next patch in the series, with a subject
of "btrfs: rework chunk allocation to avoid exhaustion of the system
chunk array" does the necessary changes to fix the system chunk array
exhaustion problem.

⑂ master
🏷 v6.1  …  v5.14-rc2

🧑 **fdmanana** authored and **kdave** committed on Jul 7, 2021  parent 5f93e77    commit 1cb3db1cf383a3c7dbda1aa0ce748b0958759947

Showing **3 changed files** with **1 addition** and **69 deletions**.

Split    Unified

⌄ ⬥ 58 ▮▮▮▮  fs/btrfs/block-group.c  ⎘

```
3377  3377       */
3378  3378       void check_system_chunk(struct btrfs_trans_handle *trans, u64 type)
3379  3379       {
3380      -           struct btrfs_transaction *cur_trans = trans->transaction;
3381  3380           struct btrfs_fs_info *fs_info = trans->fs_info;
3382  3381           struct btrfs_space_info *info;
3383  3382           u64 left;
3392  3391           lockdep_assert_held(&fs_info->chunk_mutex);
3393  3392
3394  3393           info = btrfs_find_space_info(fs_info, BTRFS_BLOCK_GROUP_SYSTEM);
3395      -   again:
3396  3394           spin_lock(&info->lock);
3397  3395           left = info->total_bytes - btrfs_space_info_used(info, true);
3398  3396           spin_unlock(&info->lock);
3411  3409
3412  3410           if (left < thresh) {
3413  3411               u64 flags = btrfs_system_alloc_profile(fs_info);
```

```
3414    -                        u64 reserved = atomic64_read(&cur_trans->chunk_bytes_reserved);
3415    -
3416    -                        /*
3417    -                         * If there's not available space for the chunk tree (system
3418    -                         * space) and there are other tasks that reserved space for
3419    -                         * creating a new system block group, wait for them to complete
3420    -                         * the creation of their system block group and release excess
3421    -                         * reserved space. We do this because:
3422    -                         *
3423    -                         * *) We can end up allocating more system chunks than necessary
3424    -                         *    when there are multiple tasks that are concurrently
3425    -                         *    allocating block groups, which can lead to exhaustion of
3426    -                         *    the system array in the superblock;
3427    -                         *
3428    -                         * *) If we allocate extra and unnecessary system block groups,
3429    -                         *    despite being empty for a long time, and possibly forever,
3430    -                         *    they end not being added to the list of unused block groups
3431    -                         *    because that typically happens only when deallocating the
3432    -                         *    last extent from a block group - which never happens since
3433    -                         *    we never allocate from them in the first place. The few
3434    -                         *    exceptions are when mounting a filesystem or running scrub,
3435    -                         *    which add unused block groups to the list of unused block
3436    -                         *    groups, to be deleted by the cleaner kthread.
3437    -                         *    And even when they are added to the list of unused block
3438    -                         *    groups, it can take a long time until they get deleted,
3439    -                         *    since the cleaner kthread might be sleeping or busy with
3440    -                         *    other work (deleting subvolumes, running delayed iputs,
3441    -                         *    defrag scheduling, etc);
3442    -                         *
3443    -                         * This is rare in practice, but can happen when too many tasks
3444    -                         * are allocating blocks groups in parallel (via fallocate())
3445    -                         * and before the one that reserved space for a new system block
3446    -                         * group finishes the block group creation and releases the space
3447    -                         * reserved in excess (at btrfs_create_pending_block_groups()),
3448    -                         * other tasks end up here and see free system space temporarily
3449    -                         * not enough for updating the chunk tree.
3450    -                         *
3451    -                         * We unlock the chunk mutex before waiting for such tasks and
3452    -                         * lock it again after the wait, otherwise we would deadlock.
3453    -                         * It is safe to do so because allocating a system chunk is the
3454    -                         * first thing done while allocating a new block group.
3455    -                         */
3456    -                        if (reserved > trans->chunk_bytes_reserved) {
3457    -                                const u64 min_needed = reserved - thresh;
3458    -
3459    -                                mutex_unlock(&fs_info->chunk_mutex);
3460    -                                wait_event(cur_trans->chunk_reserve_wait,
3461    -                                   atomic64_read(&cur_trans->chunk_bytes_reserved) <=
3462    -                                   min_needed);
3463    -                                mutex_lock(&fs_info->chunk_mutex);
3464    -                                goto again;
3465    -                        }
3466  3412
3467  3413                        /*
3468  3414                         * Ignore failure to create system chunk. We might end up not
3477  3423                        ret = btrfs_block_rsv_add(fs_info->chunk_root,
3478  3424                                                  &fs_info->chunk_block_rsv,
3479  3425                                                  thresh, BTRFS_RESERVE_NO_FLUSH);
3480    -                        if (!ret) {
3481    -                                atomic64_add(thresh, &cur_trans->chunk_bytes_reserved);
      3426  +                        if (!ret)
3482  3427                                trans->chunk_bytes_reserved += thresh;
3483    -                        }
3484  3428                }
3485  3429        }
3486  3430
```

∨ ⊕ 5 ▣▣▣▣▣ fs/btrfs/transaction.c ⟱

```
260   260        void btrfs_trans_release_chunk_metadata(struct btrfs_trans_handle *trans)
261   261        {
262   262                struct btrfs_fs_info *fs_info = trans->fs_info;
263     -                struct btrfs_transaction *cur_trans = trans->transaction;
264   263
265   264                if (!trans->chunk_bytes_reserved)
266   265                        return;
269   268
270   269                btrfs_block_rsv_release(fs_info, &fs_info->chunk_block_rsv,
271   270                                        trans->chunk_bytes_reserved, NULL);
272     -                atomic64_sub(trans->chunk_bytes_reserved, &cur_trans->chunk_bytes_reserved);
273     -                cond_wake_up(&cur_trans->chunk_reserve_wait);
274   271                trans->chunk_bytes_reserved = 0;
275   272        }
276   273
386   383                spin_lock_init(&cur_trans->dropped_roots_lock);
387   384                INIT_LIST_HEAD(&cur_trans->releasing_ebs);
388   385                spin_lock_init(&cur_trans->releasing_ebs_lock);
389     -                atomic64_set(&cur_trans->chunk_bytes_reserved, 0);
390     -                init_waitqueue_head(&cur_trans->chunk_reserve_wait);
391   386                list_add_tail(&cur_trans->list, &fs_info->trans_list);
392   387                extent_io_tree_init(fs_info, &cur_trans->dirty_pages,
393   388                                    IO_TREE_TRANS_DIRTY_PAGES, fs_info->btree_inode);
```

∨ ⊕ 7 ▣▣▣▣▣ fs/btrfs/transaction.h ⟱

```
96    96
97    97                spinlock_t releasing_ebs_lock;
```

```
 98   98            struct list_head releasing_ebs;
 99        -
100        -        /*
101        -         * The number of bytes currently reserved, by all transaction handles
102        -         * attached to this transaction, for metadata extents of the chunk tree.
103        -         */
104        -        atomic64_t chunk_bytes_reserved;
105        -        wait_queue_head_t chunk_reserve_wait;
106   99    };
107  100
108  101    #define __TRANS_FREEZABLE        (1U << 0)
```

**0 comments on commit** `1cb3db1`