# 数据修复

ceph版本：https://github.com/ceph/ceph/tree/v18.2.1

PG日志记录了PG内对象的所有操作，包括创建、修改和删除等。每个日志条目（log entry）包含了执行操作时的元数据，例如对象的版本号（eversion）、操作类型、对象的元数据（如大小、时间戳等），这些日志条目按按照操作发生的顺序排列，形成了一个操作历史记录。

当Ceph执行数据恢复或同步操作时，它会参考PG日志来确定对象的当前状态。例如，在Recovery阶段，Ceph会使用PG日志来恢复缺失的对象副本，确保所有副本都是最新的。在这个过程中，PG日志充当了元数据变更的记录器，帮助Ceph维护数据的一致性和完整性。

Recovery是依据PG日志中的缺失记录来修复不一致的对象。Backfill是PG通过重新扫描所有的对象，对比发现缺失的对象，通过整体拷贝来修复。当一个OSD失效时间过长导致无法根据PG日志来修复，或者新加入的OSD导致数据迁移时，就会启动Backfill过程。

## 资源预约

如果有大量pg同时进行修复，那可能会影响到客户端读写的性能，需要对pg同时修复的数量进行限制，于是修复前需要进行资源的预约。

```
1  void request_reservation(
2      T item,                  ///< [in] reservation key
3      Context *on_reserved,    ///< [in] callback to be called on reservation
4      unsigned prio,           ///< [in] priority
5      Context *on_preempt = 0 ///< [in] callback to be called if we are
   preempted (optional)
6  ) {
7      std::lock_guard l(lock);
8      // 资源关键字、优先级、预留成功时的回调、可选的抢占回调
9      Reservation r(item, prio, on_reserved, on_preempt);
10
11     // 资源优先队列
12     queues[prio].push_back(r);
13     // item -> 资源位置
14     queue_pointers.insert(std::make_pair(item, std::make_pair(prio, --
   (queues[prio]).end())));
15     // 资源分配
16     do_queues();
17 }
```

**取消预约**

```cpp
void cancel_reservation(
    T item ///< [in] key for reservation to cancel
) {
    std::lock_guard l(lock);
    auto i = queue_pointers.find(item);
    if (i != queue_pointers.end()) {
        unsigned prio = i->second.first;
        const Reservation &r = *i->second.second;
        delete r.grant;
        delete r.preempt;
        queues[prio].erase(i->second.second);
        if (queues[prio].empty()) {
            queues.erase(prio);
        }
        queue_pointers.erase(i);
    } else {
        auto p = in_progress.find(item);
        if (p != in_progress.end()) {
            if (p->second.preempt) {
                preempt_by_prio.erase(std::make_pair(p->second.prio, p->second.item));
                delete p->second.preempt;
            }
            in_progress.erase(p);
        }
    }
    // 继续给队列中其他的请求分配
    do_queues();
}
```

# do_queues

做资源分配的函数

```cpp
void do_queues() {
    // 抢占或回收资源
    while (!preempt_by_prio.empty() &&
            (in_progress.size() > max_allowed ||
            preempt_by_prio.begin()->first < min_priority)) {
        preempt_one();
    }

}
```

```
 9      while (!queues.empty()) {
10          // choose highest priority queue
11          auto it = queues.end();
12          --it;
13          ceph_assert(!it->second.empty());
14          // 如果优先级过低
15          if (it->first < min_priority) {
16              break;
17          }
18
19          // 1. 正在进行的资源分配数量大于阈值
20          // 2. 存在可以抢占的资源
21          // 3. 优先级高于可抢占资源
22          if (in_progress.size() >= max_allowed &&
23              !preempt_by_prio.empty() &&
24              it->first > preempt_by_prio.begin()->first) {
25              preempt_one();
26          }
27          // 正在进行的资源分配数量大于阈值
28          if (in_progress.size() >= max_allowed) {
29              break; // no room
30          }
31          // 给资源
32          Reservation p = it->second.front();
33          // 从分配队列中取出待分配对象
34          queue_pointers.erase(p.item);
35          it->second.pop_front();
36          if (it->second.empty()) {
37              queues.erase(it);
38          }
39          // 分配资源
40          f->queue(p.grant);
41          p.grant = nullptr;
42          in_progress[p.item] = p;
43          if (p.preempt) {
44              s.insert(std::make_pair(p.prio, p.item));
45          }
46      }
47  }
```

# 数据修复

数据修复有两个过程，一个是Recovery（修复）一个是Backfill（回填），当数据无法通过日志修复时（Recovery），就通过向其它OSD寻找完整的数据，拷贝到自己这（Backfill），一般前者作为临时故障的修复，后者作为长时间故障的修复或用于集群变更时PG的迁移。
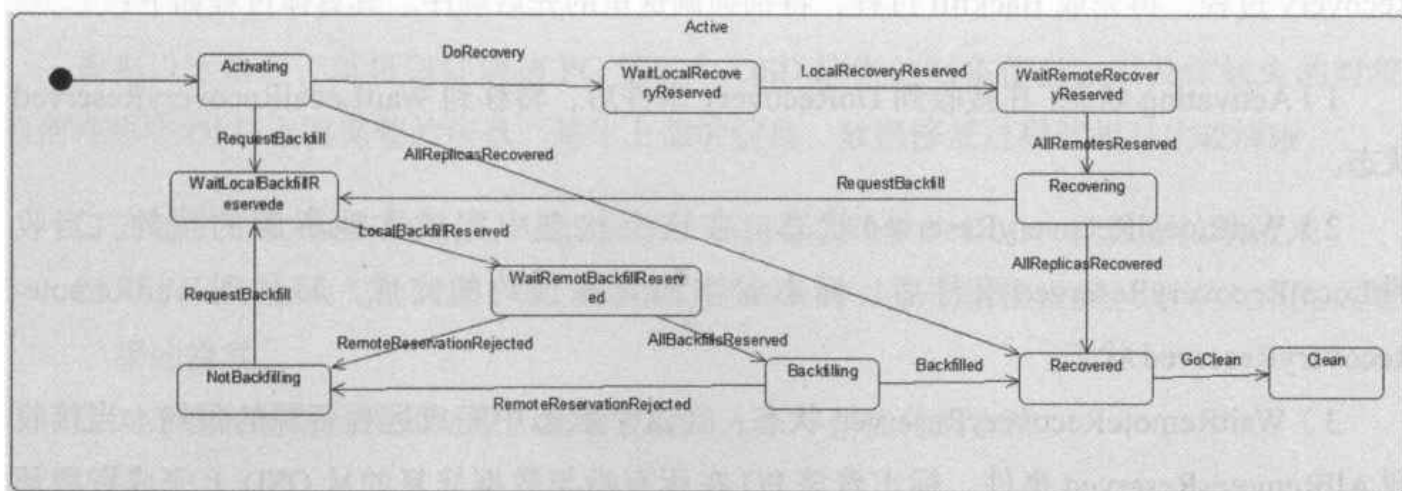
图 11-1 修复过程状态转换图

情况1：Activating状态，如果副本完整不需要修复，直接转换到Recovered状态，再到Clean状态

情况2：Activating状态，不需要Recovery时，就进行Backfill过程：

WaitLocalBackfillReserved状态申请资源，成功后，进入WaitRemoteBackfillReserved状态，**所有**副本资源预约成功后，主OSD进入Backfilling状态，完成修复进入Recovered状态。

异常：资源预约失败后进入NotBackfill状态，等待事件重新发起Backfill过程

情况3：Activating状态，进行Backfill过程：

WaitLocalRecoveryReserved状态申请资源，成功后转到WaitRemoteRecoveryReserved状态，当**所有**参与数据修复的OSD都预约到资源后，转到Recovering状态，该状态完成之后，会根据是否需要Backfill决定转入哪一个状态（见上图）

# Recovery过程

数据修复的依据是在 Peering 过程中产生的如下信息：

❑ 主副本上的缺失对象的信息保存在 pg_log 类的 pg_missing_t 结构中。

❑ 各从副本上的缺失对象信息保存在 OSD 对应的 peer_missing 中的 pg_missing_t 结构中。

❑ 缺失对象的位置信息保存在类 MissingLoc 中。

根据以上信息，就可以知道该 PG 里各个 OSD 缺失的对象信息，以及该缺失的对象目前在哪些 OSD 上有完整的信息。基于上面的信息，数据修复过程就相对比较清晰：

❑ 对于主 OSD 缺失的对象，随机选择一个拥有该对象的 OSD，把数据拉取过来。

❑ 对于 replica 缺失的对象，从主副本上把缺失的对象数据推送到从副本上来完成数据的修复。

❑ 对于比较特殊的快照对象，在修复时加入了一些优化的方法。

先修复主OSD上缺失或不一致的对象，然后修复从OSD上的对象。

当前Recovery一共有两种方式

1. **Pull** 主OSD自身存在待修复对象，由主OSD按照missing_loc选择合适的副本去拉取待修复的对象目标版本至本地，完成修复方式。

2. **Push** 主OSD知道从OSD存在待修复对象，主动推送每个待修复对象目标版本到相应从OSD，然后由其本地完成修复的方式。

主OSD自我修复过程中，可能有多个从OSD拥有待修复对象目标版本，处于负载均衡的目的可以随机选择副本，完成修复之后，开始修复各个从OSD损坏的对象，依靠peering阶段生成的missing列表，通过push的方式逐个完成从OSD的修复。

## 流程总览

```
 1  主OSD:
 2  OSD::do_recovery
 3  >> pg->start_recovery_ops
 4
 5  PrimaryLogPG::start_recovery_ops
 6  >> 1 recover_primary
 7  >> 2 recover_replicas
 8  >> 状态转变，见状态图，完成修复或进入backfill流程
 9
10  1
11  PrimaryLogPG::recover_primary
12  >> 获取缺失对象的最后一条日志
```

```
13  >> 1.1 处理"版本回退"操作，执行修复 recover_missing
14  >> 1.2 run_recovery_op 修复其它的对象
15
16  1.1
17  PrimaryLogPG::recover_missing
18  >> 修复已经删除的对象
19  >> 修复snap对象前，需要先递归修复head
20  >> 修复其它类型对象 pgbackend->recover_object
21
22  ReplicatedBackend::recover_object
23  >> prepare_pull 准备PullOp
24
25  1.2
26  ReplicatedBackend::run_recovery_op
27  >> send_pulls
28
29  2
30  PrimaryLogPG::recover_replicas
31  >> 遍历所有需要修复的OSD，准备PushOp
32  >> run_recovery_op
33
34  ReplicatedBackend::run_recovery_op
35  >> send_pushes
36
37  从OSD:
38  当主OSD send_pulls，从OSD收到
39  ReplicatedBackend::handle_pull
40  >> 封装对象数据，PushOp推送
41
42  当主OSD send_pushes，从OSD收到
43  ReplicatedBackend::handle_push
44  >> submit_push_data
45  >>>> submit_push_complete
46  >>>>>> ObjectStore::Transaction->clone_range // 将要修改的操作提交给事务，具体运行机
        制未分析
```

## 启动

```
1  void PGRecovery::run(
2    OSD *osd,
3    OSDShard *sdata,
4    PGRef& pg,
5    ThreadPool::TPHandle &handle)
6  {
7    osd->logger->tinc(
```

```
 8      l_osd_recovery_queue_lat,
 9      time_queued - ceph_clock_now());
10    osd->do_recovery(pg.get(), epoch_queued, reserved_pushes, priority, handle);
11    pg->unlock();
12  }
```

## do_recovery

```
 1  void OSD::do_recovery(
 2      PG *pg, epoch_t queued, uint64_t reserved_pushes, int priority,
 3      ThreadPool::TPHandle &handle) {
 4      uint64_t started = 0;
 5
 6      float recovery_sleep = get_osd_recovery_sleep();
 7      {
 8          // 休眠一定的时间再开始修复
 9          return;
10      }
11      // 开始修复
12      bool do_unfound = pg->start_recovery_ops(reserved_pushes, handle,
   &started);
13
14      if (do_unfound) { // 有些对象不知道能在哪个OSD上找到
15          PeeringCtx rctx;
16          rctx.handle = &handle;
17          pg->find_unfound(queued, rctx); // 去找那部分对象
18          dispatch_context(rctx, pg, pg->get_osdmap()); // 给其它OSD发送信息
19      }
20  }
```

## start_recovery_ops

```
 1  bool PrimaryLogPG::start_recovery_ops(
 2      uint64_t max,
 3      ThreadPool::TPHandle &handle,
 4      uint64_t *ops_started) {
 5      uint64_t &started = *ops_started;
 6      started = 0;
 7      bool work_in_progress = false;
 8      bool recovery_started = false;
 9  // 1 状态检查
10      ceph_assert(is_primary());
```

```
11        ceph_assert(is_peered());
12        ceph_assert(!recovery_state.is_deleting());
13
14        ceph_assert(recovery_queued);
15        recovery_queued = false;
16
17        if (!state_test(PG_STATE_RECOVERING) &&
18            !state_test(PG_STATE_BACKFILLING)) {
19            return have_unfound();
20        }
21
22 // 2 获取missing，missing是缺失的对象
23        const auto &missing = recovery_state.get_pg_log().get_missing();
24
25 // 2.1 unfound 是缺失但没有找到正确副本所在位置的对象
26        uint64_t num_unfound = get_num_unfound();
27
28 // 2.2 如果不缺失对象，设置 info.last_complete = info.last_update;
29        if (!recovery_state.have_missing()) {
30            recovery_state.local_recovery_complete();
31        }
32
33 // 3 主OSD没有missing，或者所有的missing都是unfound对象，就先修复副本
34        if (!missing.have_missing() ||
35            recovery_state.all_missing_unfound()) {
36            // 修复副本
37            started = recover_replicas(max, handle, &recovery_started);
38        }
39
40 // 3.1 表示启动修复的对象数量为0，修复主OSD上的对象
41        if (!started) {
42            started += recover_primary(max, handle);
43        }
44
45 // 3.2 如果仍然为0，且num_unfound发生变化，那么再次启动修复
46        if (!started && num_unfound != get_num_unfound()) {
47            // 修复副本
48            started = recover_replicas(max, handle, &recovery_started);
49        }
50
51        if (started || recovery_started)
52            work_in_progress = true;
53
54        bool deferred_backfill = false; // 推迟backfill
55
56 // 4 决定推迟backfill还是立即开始
57        if (recovering.empty() &&
```

```
58            state_test(PG_STATE_BACKFILLING) &&
59            !get_backfill_targets().empty() && started < max &&
60            missing.num_missing() == 0 &&
61            waiting_on_backfill.empty()) {
62            if (get_osdmap()->test_flag(CEPH_OSDMAP_NOBACKFILL)) {
63                deferred_backfill = true;
64            } else if (get_osdmap()->test_flag(CEPH_OSDMAP_NOREBALANCE) &&
   !is_degraded()) {
65                deferred_backfill = true;
66            } else if (!recovery_state.is_backfill_reserved()) {
67                // 如果backfill_reserved没有设置
68                /* DNMNOTE I think this branch is dead */
69                if (!backfill_reserving) {
70                    backfill_reserving = true;
71                    queue_peering_event(
72                        PGPeeringEventRef(
73                            std::make_shared<PGPeeringEvent>(
74                                get_osdmap_epoch(),
75                                get_osdmap_epoch(),
76                                PeeringState::RequestBackfill())));
77                }
78                deferred_backfill = true;
79            } else {
80                // 开始backfill过程
81                started += recover_backfill(max - started, handle,
   &work_in_progress);
82            }
83        }
84
85 // 5 if: 是否有正在进行的数据恢复操作，是否有正在进行的工作，是否有活跃的恢复操作，
86 // 是否有延迟的数据回填操作
87 // 返回值为true：表示需要继续进行数据修复，此时没有正在运行的恢复工作，且有未找到的对象
88    if (!recovering.empty() ||
89        work_in_progress || recovery_ops_active > 0 || deferred_backfill)
90        return !work_in_progress && have_unfound();
91
92    ceph_assert(recovering.empty());
93    ceph_assert(recovery_ops_active == 0);
94
95    int unfound = get_num_unfound();
96    if (unfound) {
97        return true;
98    }
99
100   if (missing.num_missing() > 0) {
101   // 这不应该发生
```

```cpp
102        osd->clog->error() << info.pgid << " Unexpected Error: recovery ending
   with "
103                          << missing.num_missing() << ": " << missing.get_items();
104        return false;
105    }
106
107    if (needs_recovery()) {
108    // 这不应该发生
109        osd->clog->error() << info.pgid
110                          << " Unexpected Error: recovery ending with missing
   replicas";
111        return false;
112    }
113
114 // 6 状态转换
115    if (state_test(PG_STATE_RECOVERING)) {
116        state_clear(PG_STATE_RECOVERING);
117        state_clear(PG_STATE_FORCED_RECOVERY);
118        if (needs_backfill()) {
119            queue_peering_event(
120                PGPeeringEventRef(
121                    std::make_shared<PGPeeringEvent>(
122                        get_osdmap_epoch(),
123                        get_osdmap_epoch(),
124                        PeeringState::RequestBackfill())));
125        } else {
126            state_clear(PG_STATE_FORCED_BACKFILL);
127            queue_peering_event(
128                PGPeeringEventRef(
129                    std::make_shared<PGPeeringEvent>(
130                        get_osdmap_epoch(),
131                        get_osdmap_epoch(),
132                        PeeringState::AllReplicasRecovered())));
133        }
134    } else { // backfilling
135        state_clear(PG_STATE_BACKFILLING);
136        state_clear(PG_STATE_FORCED_BACKFILL);
137        state_clear(PG_STATE_FORCED_RECOVERY);
138        queue_peering_event(
139            PGPeeringEventRef(
140                std::make_shared<PGPeeringEvent>(
141                    get_osdmap_epoch(),
142                    get_osdmap_epoch(),
143                    PeeringState::Backfilled())));
144    }
145
146    return false;
```

```
147 }
148
```

# recover_primary

```
1  uint64_t PrimaryLogPG::recover_primary(uint64_t max, ThreadPool::TPHandle
   &handle) {
2      ceph_assert(is_primary());
3
4      const auto &missing = recovery_state.get_pg_log().get_missing();
5
6      // look at log!
7      pg_log_entry_t *latest = 0;
8      unsigned started = 0;
9      int skipped = 0;
10
11 // 1 获取一个Handle，用于Push和Pull
12     PGBackend::RecoveryHandle *h = pgbackend->open_recovery_op();
13     map<version_t, hobject_t>::const_iterator p =
   missing.get_rmissing().lower_bound(recovery_state.get_pg_log().get_log().last_r
   equested);
14     while (p != missing.get_rmissing().end()) { // 遍历未被修复的对象
15         handle.reset_tp_timeout();
16         hobject_t soid;
17         version_t v = p->first;
18
19 // 2 lastest 是日志记录中保存的该缺失对象的最后的一条日志，soid为缺失的对象
20         auto it_objects = recovery_state.get_pg_log().get_log().objects.find(p-
   >second);
21         if (it_objects != recovery_state.get_pg_log().get_log().objects.end())
   {
22             latest = it_objects->second;
23             ceph_assert(latest->is_update() || latest->is_delete());
24             soid = latest->soid;
25         } else {
26             latest = 0;
27             soid = p->second;
28         }
29         const pg_missing_item &item = missing.get_items().find(p->second)-
   >second;
30         ++p;
31
32         hobject_t head = soid.get_head();
33         eversion_t need = item.need;
34
```

```
35  // 3 开始修复
36      if (latest) {
37          switch (latest->op) {
38          case pg_log_entry_t::CLONE:
39          /*
40          暂时取消了对这种特殊情况的处理，直到我们能够从旧的SnapSet中正确地构建一个准
    确的SnapSet。
41          */
42              break;
43  // 4 该记录类型为LOST_REVERT: 该revert操作为数据不一致时，管理员通过命令强行回退到指定版
    本
44  // reverting记录了回退的版本号
45          case pg_log_entry_t::LOST_REVERT: {
46
47  // 4.1 此条件表示: 日志记录显示当前已经拥有回退的版本,
48              if (item.have == latest->reverting_to) {
49              // 获取该对象的ObjectContext
50                  ObjectContextRef obc = get_object_context(soid, true);
51
52  // 4.1.1 如果检查对象当前的版本 obc->obs.oi.version 等于 latest->version，说明回退操
    作已完成
53                  if (obc->obs.oi.version == latest->version) {
54                      // I''m already reverting
55                  } else {
56  // 4.1.2 说明没有执行回退操作，直接修改对象的版本号为latest->version，关于为什么如此，
    见下文
57                      obc->obs.oi.version = latest->version;
58
59                      ObjectStore::Transaction t;
60                      bufferlist b2;
61                      obc->obs.oi.encode(
62                          b2,
63                          get_osdmap()->get_features(CEPH_ENTITY_TYPE_OSD,
    nullptr));
64                      ceph_assert(!pool.info.require_rollback());
65                      t.setattr(coll, ghobject_t(soid), OI_ATTR, b2);
66                      // 更新恢复状态
67                      recovery_state.recover_got(
68                          soid,
69                          latest->version,
70                          false,
71                          t);
72
73                      ++active_pushes;
74
75                      t.register_on_applied(new
    C_OSD_AppliedRecoveredObject(this, obc));
```

```cpp
 76                        t.register_on_commit(new C_OSD_CommittedPushedObject(
 77                            this,
 78                            get_osdmap_epoch(),
 79                            info.last_complete));
 80                        osd->store->queue_transaction(ch, std::move(t));
 81                        continue;
 82                    }
 83                } else {
```
// 4.2 需要拉取该reverting_to版本的对象，这里不做特殊处理，只是检查所有OSD是否拥有该版本
的对象，如果有就加入到missing_loc记录该版本的位置信息，由后续修复继续来完成
```cpp
 85                    eversion_t alternate_need = latest->reverting_to;
 86
 87                    set<pg_shard_t> good_peers;
 88                    for (auto p = recovery_state.get_peer_missing().begin();
 89                         p != recovery_state.get_peer_missing().end();
 90                         ++p) {
 91                        if (p->second.is_missing(soid, need) &&
 92                            p->second.get_items().at(soid).have ==
    alternate_need) {
 93                            good_peers.insert(p->first);
 94                        }
 95                    }
 96                    recovery_state.set_revert_with_targets(soid, good_peers);
 97                    /*
 98                    void PeeringState::set_revert_with_targets(
 99                        const hobject_t &soid,
100                        const set<pg_shard_t> &good_peers) {
101                        for (auto &&peer: good_peers) {
102                            missing_loc.add_location(soid, peer);
103                        }
104                    }
105                    */
106                }
107            } break;
108            }
109        }
110
```
// 5 如果当前soid没有在修复
```cpp
112        if (!recovering.count(soid)) {
113            if (recovering.count(head)) { // 或者head在修复
114                ++skipped;
115            } else {
116            // 修复
117                int r = recover_missing(soid, need,
    recovery_state.get_recovery_op_priority(), h);
118                switch (r) {
119                case PULL_YES:
```

```
120                  ++started;
121                  break;
122              case PULL_HEAD:
123                  ++started;
124              case PULL_NONE:
125                  ++skipped;
126                  break;
127              default:
128                  ceph_abort();
129              }
130              if (started >= max)
131                  break;
132          }
133      }
134
135      if (!skipped)
136          recovery_state.set_last_requested(v);
137  }
138
139 // 6 把PullOp或PushOp封装的消息发送出去
140     pgbackend->run_recovery_op(h, recovery_state.get_recovery_op_priority());
141     return started;
142 }
143
```

书上的例子：

例 11-1 日志修复过程。

PG 日志的记录如下：每个单元代表一条日志记录，分别为对象的名字和版本以及操作，版本的格式为（epoch,version）。灰色的部分代表本 OSD 上缺失的日志记录，该日志记录是从权威日志记录中拷贝过来的，所以当前该日志记录是连续完整的。

| obj2(1,3) modify | obj1(1,4) modify | obj2(1,5) modify | obj1(1,6) modify | obj1(1,7) modify | obj1(1,8) modify |

**情况 1：正常情况的修复。**

缺失的对象列表为 [obj1，obj2]。当前修复对象为 obj1。由日志记录可知：对象 obj1 被修改过三次，分别为版本 6,7,8。当前拥有的 obj1 对象的版本 have 值为 4，修复时只修复到最后修改的版本 8 即可。

**情况 2：最后一个操作为 LOST_REVERT 类型的操作。**

| obj2(1,3) modify | obj1(1,4) modify | obj2(1,5) modify | obj1(1,6) modify | obj1(1,7) modify | obj1(1,8)<br>lost_revert_<br>version = 8<br>prior_version=7<br>reverting_to=4 |
|---|---|---|---|---|---|
| | | | | | |

对于要修复的对象 obj1，最后一次操作为 LOST_REVERT 类型的操作，该操作当前版本 version 为 8，修改前的版本 prior_version 为 7，回退版本 reverting_to 为 4。

在这种情况下，日志显示当前已经有版本 4，检查对象 obj1 的实际版本，也就是 object_info 里保存的版本号：

1）如果该值是 8，说明最后一次 revert 操作成功，不需要做任何修复动作。

2）如果该值是 4，说明 LOST_REVERT 操作就没有执行。当然数据内容已经是版本 4 了，只需要修改 object_info 的版本为 8 即可。

如果回退的版本 reverting_to 不是版本 4，而是版本 6，那么最终还是需要把 obj1 的数据修复到版本 6 的数据。Ceph 在这里的处理，仅仅是检查其他 OSD 缺失的对象中是否有版本 6，如果有，就加入到 missing_loc 中，记录拥有该版本的 OSD 位置，待后续继续修复。

最后一种情况，如果是版本6，对象中并没有版本6的数据，将数据修复为版本6，当前版本设置为8。

# recover_missing

```
1  int PrimaryLogPG::recover_missing(
2    const hobject_t &soid, eversion_t v,
3    int priority,
4    PGBackend::RecoveryHandle *h) {
5
6  // 1 如果是unfound对象，无法修复
7    if (recovery_state.get_missing_loc().is_unfound(soid)) {
```

```
  8        return PULL_NONE;
  9    }
 10
 11  // 2 已删除的对象
 12    if (recovery_state.get_missing_loc().is_deleted(soid)) {
 13  // 2.1 开始修复操作
 14        start_recovery_op(soid);
 15        // 确保不在修复中
 16        ceph_assert(!recovering.count(soid));
 17        // 加入修复队列
 18        recovering.insert(make_pair(soid, ObjectContextRef()));
 19        epoch_t cur_epoch = get_osdmap_epoch();
 20        remove_missing_object(soid, v, new LambdaContext([=, this](int) {
 21            std::scoped_lock locker{*this};
 22            if (!pg_has_reset_since(cur_epoch)) {
 23              bool object_missing = false;
 24              for (const auto &shard : get_acting_recovery_backfill()) {
 25                    if (shard == pg_whoami)
 26                        continue;
 27                    if (recovery_state.get_peer_missing(shard).is_missing(soid))
{
 28                        object_missing = true;
 29                        break;
 30                    }
 31              }
 32              if (!object_missing) {
 33                    object_stat_sum_t stat_diff;
 34                    stat_diff.num_objects_recovered = 1;
 35                    if (scrub_after_recovery)
 36                        stat_diff.num_objects_repaired = 1;
 37                    on_global_recover(soid, stat_diff, true);
 38              } else {
 39                    auto recovery_handle = pgbackend->open_recovery_op();
 40                    pgbackend->recover_delete_object(soid, v, recovery_handle);
 41                    pgbackend->run_recovery_op(recovery_handle, priority);
 42              }
 43            }
 44        }));
 45        return PULL_YES;
 46    }
 47
 48  // is this a snapped object? if so, consult the snapset.. we may not need the
     entire object!
 49  // 3 快照对象
 50    ObjectContextRef obc;
 51    ObjectContextRef head_obc;
 52    if (soid.snap && soid.snap < CEPH_NOSNAP) {
```

```
53            // do we have the head?
54            hobject_t head = soid.get_head();
55            if (recovery_state.get_pg_log().get_missing().is_missing(head)) {
56                if (recovering.count(head)) {
57                    return PULL_NONE;
58                } else {
59                    int r = recover_missing(
60                        head,
   recovery_state.get_pg_log().get_missing().get_items().find(head)->second.need,
   priority,
61                        h);
62                    if (r != PULL_NONE)
63                        return PULL_HEAD;
64                    return PULL_NONE;
65                }
66            }
67            head_obc = get_object_context(
68                head,
69                false,
70                0);
71            ceph_assert(head_obc);
72        }
73    start_recovery_op(soid);
74    ceph_assert(!recovering.count(soid));
75    recovering.insert(make_pair(soid, obc));
76
77 // 4 修复对象
78    int r = pgbackend->recover_object(soid, v, head_obc, obc, h);
79    ceph_assert(r >= 0);
80    return PULL_YES;
81 }
```

## recover_object

pgbackend封装了不同类型的Pool的实现。ReplicatedBackend实现了replicate类型的PG相关的底层功能，ECbackend实现了Erasure code类型的PG相关的底层功能。

我们讨论基于副本的修复。

```
1 int ReplicatedBackend::recover_object(
2    const hobject_t &hoid,
3    eversion_t v,
4    ObjectContextRef head,
5    ObjectContextRef obc,
6    RecoveryHandle *_h) {
```

```
 7      RPGHandle *h = static_cast<RPGHandle *>(_h);
 8      if (get_parent()->get_local_missing().is_missing(hoid)) {
 9          ceph_assert(!obc);
10          // 把请求封装为PullOp
11          prepare_pull(v, hoid, head, h);
12      } else {
13          ceph_assert(obc);
14          // 把请求封装为PushOp
15          int started = start_pushes(hoid, obc, h);
16          if (started < 0) {
17              pushing[hoid].clear();
18              return started;
19          }
20      }
21      return 0;
22 }
```

# 1 pull

```
 1 struct PullOp {
 2     hobject_t soid;                              // 需要拉取的对象
 3     ObjectRecoveryInfo recovery_info;            // 对象修复的信息
 4     ObjectRecoveryProgress recovery_progress;    // 对象修复进度信息
 5 };
 6
 7 struct ObjectRecoveryInfo {
 8     hobject_t soid;
 9     eversion_t version;
10     uint64_t size;
11     object_info_t oi;
12     SnapSet ss;                          // 修复对象的快照信息
13
14     interval_set<uint64_t> copy_subset;    // 修复快照时，需要从其它OSD拷贝到本地的对
   象的区段集合
15     std::map<hobject_t, interval_set<uint64_t>> clone_subset;
16     // clone对象修复时，需要从本地拷贝来修复的区间
17     bool object_exist;
18 }
19
20 struct ObjectRecoveryProgress {
21     uint64_t data_recovered_to;      // data已修复的位置指针
22     std::string omap_recovered_to;   // omap已修复的位置指针
23     bool first;              // 是否是首次修复
24     bool data_complete;      // data是否修复完成
25     bool omap_complete;      // omap是否修复完成
```

```
26    bool error = false;
27 }
```

```
1  void ReplicatedBackend::prepare_pull(
2      eversion_t v,              // 要拉取对象的版本信息
3      const hobject_t& soid,     // 要拉取的对象
4      ObjectContextRef headctx,  // 拉取对象的ObjectContext信息
5      RPGHandle *h){             // 封装后保存的RecveryHandle
6  // 1 获取PG对象
7      const auto missing_iter = get_parent()-
   >get_local_missing().get_items().find(soid);
8      ceph_assert(missing_iter != get_parent()-
   >get_local_missing().get_items().end());
9      eversion_t _v = missing_iter->second.need;
10     ceph_assert(_v == v);
11
12 // 2 missing_loc 包含缺失对象的位置，peering_missing包含其它节点缺失对象的信息
13     const map<hobject_t, set<pg_shard_t>> &missing_loc(
14         get_parent()->get_missing_loc_shards());
15     const map<pg_shard_t, pg_missing_t> &peer_missing(
16         get_parent()->get_shard_missing());
17
18 // 3 查找soid所在的OSD集合
19     map<hobject_t, set<pg_shard_t>>::const_iterator q = missing_loc.find(soid);
20     ceph_assert(q != missing_loc.end());
21     ceph_assert(!q->second.empty());
22
23 // 4 选择一个特定的分片（OSD）作为拉取操作的目标
24     auto p = q->second.end();
25     if (cct->_conf->osd_debug_feed_pullee >= 0) {
26         for (auto it = q->second.begin(); it != q->second.end(); it++) {
27             if (it->osd == cct->_conf->osd_debug_feed_pullee) {
28                 p = it;
29                 break;
30             }
31         }
32     }
33 // 4.1 如果没有找到特定的pullee，可能是用户输入了错误的信息，随机选择一个目标
34     if (p == q->second.end()) {
35         // probably because user feed a wrong pullee
36         p = q->second.begin();
37         std::advance(p,ceph::util::generate_random_number<int>(0,
38             q->second.size() - 1));
39     }
40     ceph_assert(get_osdmap()->is_up(p->osd));
```

```cpp
41        pg_shard_t fromshard = *p;
42
43  // 5 确保选中的OSD上确实需要的对象
44        ceph_assert(peer_missing.count(fromshard));
45        const pg_missing_t &pmissing = peer_missing.find(fromshard)->second;
46        if (pmissing.is_missing(soid, v)) {
47            ceph_assert(pmissing.get_items().find(soid)->second.have != v);
48            v = pmissing.get_items().find(soid)->second.have;
49            ceph_assert(get_parent()->get_log().get_log().objects.count(soid) &&
50                        (get_parent()->get_log().get_log().objects.find(soid)-
    >second->op ==
51                         pg_log_entry_t::LOST_REVERT) &&
52                        (get_parent()->get_log().get_log().objects.find(soid)
53                            ->second->reverting_to ==v));
54        }
55
56        ObjectRecoveryInfo recovery_info;
57        ObcLockManager lock_manager;
58  // 6 如果是快照
59        if (soid.is_snap()) {
60            ceph_assert(!get_parent()-
    >get_local_missing().is_missing(soid.get_head()));
61            ceph_assert(headctx);
62            // check snapset
63            SnapSetContext *ssc = headctx->ssc;
64            ceph_assert(ssc);
65
66            recovery_info.ss = ssc->snapset;
67            calc_clone_subsets(
68                ssc->snapset, soid, get_parent()->get_local_missing(),
69                get_info().last_backfill,
70                recovery_info.copy_subset,
71                recovery_info.clone_subset,
72                lock_manager);
73            // FIXME: this may overestimate if we are pulling multiple clones in
    parallel...
74
75            ceph_assert(ssc->snapset.clone_size.count(soid.snap));
76            recovery_info.size = ssc->snapset.clone_size[soid.snap];
77            recovery_info.object_exist = missing_iter-
    >second.clean_regions.object_is_exist();
78        } else {
79  // 6.1 如果是head对象，拉取全部
80            // pulling head or unversioned object.
81            // always pull the whole thing.
82            recovery_info.copy_subset.insert(0, (uint64_t)-1);
83            assert(HAVE_FEATURE(parent->min_peer_features(), SERVER_OCTOPUS));
```

```
84          recovery_info.copy_subset.intersection_of(missing_iter-
    >second.clean_regions.get_dirty_regions());
85          recovery_info.size = ((uint64_t)-1);
86          recovery_info.object_exist = missing_iter-
    >second.clean_regions.object_is_exist();
87      }
88
89 // 7 创建PullOp对象，设置拉取操作的相关信息
90      h->pulls[fromshard].push_back(PullOp());
91      PullOp &op = h->pulls[fromshard].back();
92      op.soid = soid;
93
94      op.recovery_info = recovery_info;
95      op.recovery_info.soid = soid;
96      op.recovery_info.version = v;
97      op.recovery_progress.data_complete = false;
98      op.recovery_progress.omap_complete = !missing_iter-
    >second.clean_regions.omap_is_dirty();
99      op.recovery_progress.data_recovered_to = 0;
100     op.recovery_progress.first = true;
101
102     ceph_assert(!pulling.count(soid));
103     pull_from_peer[fromshard].insert(soid);
104     PullInfo &pi = pulling[soid];
105     pi.from = fromshard;
106     pi.soid = soid;
107     pi.head_ctx = headctx;
108     pi.recovery_info = op.recovery_info;
109     pi.recovery_progress = op.recovery_progress;
110     pi.cache_dont_need = h->cache_dont_need;
111     pi.lock_manager = std::move(lock_manager);
112 }
113
```

## 1.1 calc_clone_subsetes

用于修复快照对象，在此之前介绍两个概念

1.在SnapSet结构中，字段clone_overlap保存了clone对象和上一次clone对象的重叠部分（没有冲突部分）

```
1 struct SnapSet {
2     snapid_t seq;
3     std::vector<snapid_t> snaps;                          // descending
4     std::vector<snapid_t> clones;                         // ascending
```

```
5    std::map<snapid_t, interval_set<uint64_t>> clone_overlap; // overlap w/
  next newest
6    std::map<snapid_t, uint64_t> clone_size;
7    std::map<snapid_t, std::vector<snapid_t>> clone_snaps; // descending
8 }
```

## 2.clone_overlap



图 11-2    clone_overlap 示意图

snap3从snap2克隆过来，然后修改了区间3，4，其在对象中范围的offset和length为（4，8），（8，12）那么记录就为：

clone_overlap[3] = {(0, 4), (12, len(区间8-区间4))}

然后是calc_clone_subsets函数

```
1 void ReplicatedBackend::calc_clone_subsets(
2    SnapSet &snapset, const hobject_t &soid,
3    const pg_missing_t &missing,
4    const hobject_t &last_backfill,
5    interval_set<uint64_t> &data_subset,
6    map<hobject_t, interval_set<uint64_t>> &clone_subsets,
7    ObcLockManager &manager) {
8
9 // 1 获取快照大小，加入到data_subset中（虽然不知道这是在干嘛）
10    uint64_t size = snapset.clone_size[soid.snap];
11    if (size) data_subset.insert(0, size);
12
13    // any overlap with next older clone?
14    interval_set<uint64_t> cloning;
15    interval_set<uint64_t> prev;
16
17 // 2 往前查找完整的快照对象区间，添加到clone_subsets和cloning
18    if (size)
19        prev.insert(0, size);
20    for (int j = i - 1; j >= 0; j--) {
21        hobject_t c = soid;
22        c.snap = snapset.clones[j];
23        // 计算重叠区间
```

```cpp
24              prev.intersection_of(snapset.clone_overlap[snapset.clones[j]]);
25          if (!missing.is_missing(c) &&
26              c < last_backfill &&
27              get_parent()->try_lock_for_read(c, manager)) {
28              clone_subsets[c] = prev;
29              cloning.union_of(prev);
30              break;
31          }
32      }
33
34  // 同上，往后查找
35      // overlap with next newest?
36      interval_set<uint64_t> next;
37      if (size)
38          next.insert(0, size);
39      for (unsigned j = i + 1; j < snapset.clones.size(); j++) {
40          hobject_t c = soid;
41          c.snap = snapset.clones[j];
42          next.intersection_of(snapset.clone_overlap[snapset.clones[j - 1]]);
43          if (!missing.is_missing(c) &&
44              c < last_backfill &&
45              get_parent()->try_lock_for_read(c, manager)) {
46              clone_subsets[c] = next;
47              cloning.union_of(next);
48              break;
49          }
50      }
51      // 去重
52      data_subset.subtract(cloning);
53  }
```
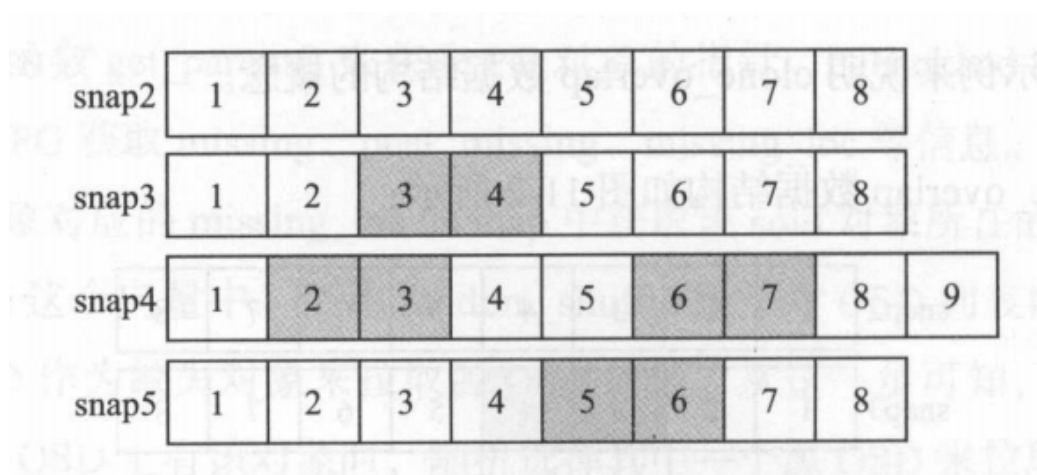
光看代码看不出是在做什么，书上给了一个例子：



图 11-3　快照对象修复计算示例图

snap3是克隆的snap2，snap4 克隆的snap3，snap5克隆的snap4，灰色区间表示clone后修改的区间，snap2、3、5都是完整的对象，要修复的对象是snap4，不同长度代表各个clone对象的size是不同的。

1.向前查找和snap4相同的区间，可以看到区间1，5，8和snap4相同，那么snap4就直接从本地拿到这三个区间。

2.向后查找，可以看到1，2，3，4，7，8未被修改，那么直接拿过来就行

3.去重，最后只有区间6需要从其他OSD上拷贝数据来修复。

## 2 push

获取actingbackfill的OSD列表，通过peering_missing查找缺失该对象的OSD，并发送信息给对方

```
1  int ReplicatedBackend::start_pushes(
2      const hobject_t &soid,
3      ObjectContextRef obc,
4      RPGHandle *h) {
5
6  // 1 用于存储需要数据的OSD
7      list<map<pg_shard_t, pg_missing_t>::const_iterator> shards;
8
9  // 遍历actingbackfill列表，统计需要信息的OSD
10     ceph_assert(get_parent()->get_acting_recovery_backfill_shards().size() >
   0);
11     for (set<pg_shard_t>::iterator i =
12            get_parent()->get_acting_recovery_backfill_shards().begin();
13       i != get_parent()->get_acting_recovery_backfill_shards().end();
14        ++i) {
15      if (*i == get_parent()->whoami_shard())
16          continue;
17      pg_shard_t peer = *i;
18      map<pg_shard_t, pg_missing_t>::const_iterator j =
19            get_parent()->get_shard_missing().find(peer);
20      ceph_assert(j != get_parent()->get_shard_missing().end());
21      if (j->second.is_missing(soid)) {
22          shards.push_back(j);
23      }
24     }
25
26     // If more than 1 read will occur ignore possible request to not cache
27     bool cache = shards.size() == 1 ? h->cache_dont_need : false;
28
29  // 发送信息
30     for (auto j : shards) {
31         pg_shard_t peer = j->first;
```

```
32          h->pushes[peer].push_back(PushOp());
33          int r = prep_push_to_replica(obc, soid, peer,
34                                        &(h->pushes[peer].back()), cache);
35          if (r < 0) {
36              // Back out all failed reads
37              for (auto k : shards) {
38                  pg_shard_t p = k->first;
39                  h->pushes[p].pop_back();
40                  if (p == peer)
41                      break;
42              }
43              return r;
44          }
45      }
46      return shards.size();
47  }
```

```
1  int ReplicatedBackend::prep_push_to_replica(
2      ObjectContextRef obc, const hobject_t &soid, pg_shard_t peer,
3      PushOp *pop, bool cache_dont_need) {
4      const object_info_t &oi = obc->obs.oi;
5      uint64_t size = obc->obs.oi.size;
6
7      map<hobject_t, interval_set<uint64_t>> clone_subsets;
8      interval_set<uint64_t> data_subset;
9
10     ObcLockManager lock_manager;
11 // 1 如果对象是一个快照
12     if (soid.snap && soid.snap < CEPH_NOSNAP) {
13         hobject_t head = soid;
14         head.snap = CEPH_NOSNAP;
15
16 // 1.1 尝试基于成功或先于当前快照的克隆来推送数据。这需要在本地获取对象的头部（head）和当
   前的快照集（SnapSet）
17         if (get_parent()->get_local_missing().is_missing(head)) {
18             return prep_push(obc, soid, peer, pop, cache_dont_need);
19         }
20 // 1.2 如果head存在，获取克隆数据子集和数据子集，使用这些信息来准备推送操作
21         SnapSetContext *ssc = obc->ssc;
22         ceph_assert(ssc);
23
24         pop->recovery_info.ss = ssc->snapset;
25         map<pg_shard_t, pg_missing_t>::const_iterator pm =
```

```
26                  get_parent()->get_shard_missing().find(peer);
27              ceph_assert(pm != get_parent()->get_shard_missing().end());
28              map<pg_shard_t, pg_info_t>::const_iterator pi =
29                  get_parent()->get_shard_info().find(peer);
30              ceph_assert(pi != get_parent()->get_shard_info().end());
31              calc_clone_subsets(
32                  ssc->snapset, soid,
33                  pm->second,
34                  pi->second.last_backfill,
35                  data_subset, clone_subsets,
36                  lock_manager);
37          } else if (soid.snap == CEPH_NOSNAP) {
38  // 2 如果不是快照，同样计算数据集，并使用这些信息来准备推送操作
39              SnapSetContext *ssc = obc->ssc;
40              ceph_assert(ssc);
41
42              calc_head_subsets(
43                  obc,
44                  ssc->snapset, soid, get_parent()->get_shard_missing().find(peer)-
    >second,
45                  get_parent()->get_shard_info().find(peer)->second.last_backfill,
46                  data_subset, clone_subsets,
47                  lock_manager);
48          }
49  // 3 构造PushOp
50          return prep_push(
51              obc,            // 对象上下文
52              soid,           // 对象ID
53              peer,           // 目标
54              oi.version,     // 对象版本
55              data_subset,    // 数据子集
56              clone_subsets,  // 克隆数据子集
57              pop,            // 推送操作指针
58              cache_dont_need,// 是否需要缓存
59              std::move(lock_manager)); // 锁管理器
60  }
61
62  prep_push:
63  >> int r = build_push_op(pi.recovery_info, pi.recovery_progress, &new_progress,
64          pop, &(pi.stat), cache_dont_need);
65
```

```
1  int ReplicatedBackend::build_push_op(const ObjectRecoveryInfo &recovery_info,
2                                       const ObjectRecoveryProgress &progress,
3                                       ObjectRecoveryProgress *out_progress,
```

```
 4                                                   PushOp *out_op,
 5                                                   object_stat_sum_t *stat,
 6                                                   bool cache_dont_need) {
 7      ObjectRecoveryProgress _new_progress;
 8      if (!out_progress)
 9          out_progress = &_new_progress;
10      ObjectRecoveryProgress &new_progress = *out_progress;
11      new_progress = progress;
12      eversion_t v = recovery_info.version;
13      object_info_t oi;
14
// 1 是第一次恢复，需要获取元数据信息
16      if (progress.first) {
17          int r = store->omap_get_header(ch, ghobject_t(recovery_info.soid),
   &out_op->omap_header);
18          r = store->getattrs(ch, ghobject_t(recovery_info.soid), out_op-
   >attrset);
19
// 2 解码对象信息，并检查本地版本号与请求的版本号是否一致。如果不一致，或者请求的版本号未
   知，则返回错误。
21          // Debug
22          try {
23              oi.decode(out_op->attrset[OI_ATTR]);
24          } catch (...) {
25              return -EINVAL;
26          }
27
28          // If requestor didn''t know the version, use ours
29          if (v == eversion_t()) {
30              v = oi.version;
31          } else if (oi.version != v) {
32              return -EINVAL;
33          }
34
// 3 标志为false，表示后续请求将包含版本信息
36          new_progress.first = false;
37      }
38      // 一旦我们提供了版本，随后的请求就会有它，所以在这一点上它必须是已知的。
39      ceph_assert(v != eversion_t());
40
// 4 上一步只获取了header，没有omap(key, value)信息，这一步获取
// 首先获取可用的数据块大小
43      uint64_t available = cct->_conf->osd_recovery_max_chunk;
44      if (!progress.omap_complete) {
45          ObjectMap::ObjectMapIterator iter =
46              store->get_omap_iterator(ch,
47                                       ghobject_t(recovery_info.soid));
```

```
48          ceph_assert(iter);
49
50 // 4.1 获取omap信息，添加到omap_entries,
51 // 一次获取的信息大小不能超过osd_recovery_max_chunk(available )
52          for (iter->lower_bound(progress.omap_recovered_to);
53               iter->valid();
54               iter->next()) {
55          if (!out_op->omap_entries.empty() &&
56              ((cct->_conf->osd_recovery_max_omap_entries_per_chunk > 0 &&
57               out_op->omap_entries.size() >= cct->_conf-
   >osd_recovery_max_omap_entries_per_chunk) ||
58               available <= iter->key().size() + iter->value().length()))
59              break;
60          out_op->omap_entries.insert(make_pair(iter->key(), iter->value()));
61
62          if ((iter->key().size() + iter->value().length()) <= available)
63              available -= (iter->key().size() + iter->value().length());
64          else
65              available = 0;
66      }
67      if (!iter->valid())
68          new_progress.omap_complete = true;
69      else
70          new_progress.omap_recovered_to = iter->key();
71  }
72
73 // 5 获取omap之后，现在要获取数据
74   if (available > 0) {
75       if (!recovery_info.copy_subset.empty()) {
76          interval_set<uint64_t> copy_subset = recovery_info.copy_subset;
77          map<uint64_t, uint64_t> m;
78 // 5.1 获取文件的映射信息
79          int r = store->fiemap(ch, ghobject_t(recovery_info.soid), 0,
80                               copy_subset.range_end(), m);
81          if (r >= 0) {
82 // 5.2 保存下需要推送的数据区间
83              interval_set<uint64_t> fiemap_included(std::move(m));
84              copy_subset.intersection_of(fiemap_included);
85          } else {
86          // copy_subset和empty interval_set的交集无论如何都是空的
87              copy_subset.clear();
88          }
89 // 5.3 将数据更新到data_included
90          out_op->data_included.span_of(copy_subset,
   progress.data_recovered_to,
91                                       available);
92
```

```
 93            if (out_op->data_included.empty() ||
 94                out_op->data_included.range_end() == copy_subset.range_end())
 95                new_progress.data_recovered_to =
    recovery_info.copy_subset.range_end();
 96            else
 97                new_progress.data_recovered_to = out_op-
    >data_included.range_end();
 98        }
 99    } else {
100        out_op->data_included.clear();
101    }
102
103    auto origin_size = out_op->data_included.size();
104    bufferlist bit;
// 6 将数据读到缓冲区bit，说实话不知道为什么又要把数据转一遍
106    int r = store->readv(ch, ghobject_t(recovery_info.soid),
107                        out_op->data_included, bit,
108                        cache_dont_need ? CEPH_OSD_OP_FLAG_FADVISE_DONTNEED :
    0);
109
110    if (cct->_conf->osd_debug_random_push_read_error &&
111        (rand() % (int)(cct->_conf->osd_debug_random_push_read_error * 100.0))
    == 0) {
112        r = -EIO;
113    }
// 7 数据块大小发生变化
115    if (out_op->data_included.size() != origin_size) {
116        new_progress.data_complete = true;
117    }
// 8 将数据追加到out_op->data，是data_included不能直接追加吗
119    out_op->data.claim_append(bit);
120    // 检验数据完整性，数据校验和等
121    if (progress.first && !out_op->data_included.empty() &&
122        out_op->data_included.begin().get_start() == 0 &&
123        out_op->data.length() == oi.size && oi.is_data_digest()) {
124        uint32_t crc = out_op->data.crc32c(-1);
125        if (oi.data_digest != crc) {
126            return -EIO;
127        }
128    }
129
// 9 恢复进度完成
131    if (new_progress.is_complete(recovery_info)) {
132        new_progress.data_complete = true;
133        if (stat) {
134            stat->num_objects_recovered++;
135            if (get_parent()->pg_is_repair())
```

```
136            stat->num_objects_repaired++;
137        }
138    } else if (progress.first && progress.omap_complete) {
139        new_progress.omap_complete = false;
140    }
141
142    if (stat) {
143        stat->num_keys_recovered += out_op->omap_entries.size();
144        stat->num_bytes_recovered += out_op->data.length();
145        get_parent()->get_logger()->inc(l_osd_rbytes, out_op->omap_entries.size() + out_op->data.length());
146    }
147
148    get_parent()->get_logger()->inc(l_osd_push);
149    get_parent()->get_logger()->inc(l_osd_push_outb, out_op->data.length());
150
151 // 10 更新out_op信息，注意out_op其实就是PushOp的引用
152    // send
153    out_op->version = v;
154    out_op->soid = recovery_info.soid;
155    out_op->recovery_info = recovery_info;
156    out_op->after_progress = new_progress;
157    out_op->before_progress = progress;
158    return 0;
159 }
```

## run_recovery_op

```
 1 void ReplicatedBackend::run_recovery_op(
 2     PGBackend::RecoveryHandle *_h,
 3     int priority) {
 4     RPGHandle *h = static_cast<RPGHandle *>(_h);
 5 // 这里的pushes和pulls就是上文中的PushOp和PullOp集合
 6     send_pushes(priority, h->pushes);
 7     send_pulls(priority, h->pulls);
 8     send_recovery_deletes(priority, h->deletes);
 9     delete h;
10 }
```

## handle_pull

当主OSD将对象推送给从OSD后，从OSD需要调用handle_push来实现数据的写入；同样，当主OSD给从OSD发起拉取请求后，需要handle_pull处理对应的请求。

```cpp
1  void ReplicatedBackend::handle_pull(pg_shard_t peer, PullOp &op, PushOp
   *reply) {
2      const hobject_t &soid = op.soid;
3      struct stat st;
4  // 1 验证对象是否存在
5      int r = store->stat(ch, ghobject_t(soid), &st);
6
7      if (r != 0) {
8  // 不存在返回空值
9          prep_push_op_blank(soid, reply);
10     } else {
11         ObjectRecoveryInfo &recovery_info = op.recovery_info;
12         ObjectRecoveryProgress &progress = op.recovery_progress;
13
14 // 2 第一次修复，且全部拷贝
15         if (progress.first && recovery_info.size == ((uint64_t)-1)) {
16             recovery_info.size = st.st_size;
17             if (st.st_size) {
18                 interval_set<uint64_t> object_range;
19                 object_range.insert(0, st.st_size);
20                 // 添加数据
21                 recovery_info.copy_subset.intersection_of(object_range);
22             } else {
23                 recovery_info.copy_subset.clear();
24             }
25             assert(recovery_info.clone_subset.empty());
26         }
27 // 3 构建PushOp操作
28         r = build_push_op(recovery_info, progress, 0, reply);
29         if (r < 0) prep_push_op_blank(soid, reply);
30     }
31 }
```

## recover_replicas

```cpp
1  uint64_t PrimaryLogPG::recover_replicas(uint64_t max, ThreadPool::TPHandle
   &handle,
2                                          bool *work_started) {
3      uint64_t started = 0;
4
5  // 1 获取一个Handle，用于Push和Pull
6      PGBackend::RecoveryHandle *h = pgbackend->open_recovery_op();
7
8      ceph_assert(!get_acting_recovery_backfill().empty());
9      std::vector<std::pair<unsigned int, pg_shard_t>> replicas_by_num_missing,
```

```cpp
10          async_by_num_missing;
11      replicas_by_num_missing.reserve(get_acting_recovery_backfill().size() - 1);
// 2 遍历所有正在修复的OSD
13      for (auto &p : get_acting_recovery_backfill()) {
14          if (p == get_primary()) {
15              continue;
16          }
17          auto pm = recovery_state.get_peer_missing().find(p);
18          ceph_assert(pm != recovery_state.get_peer_missing().end());
19          auto nm = pm->second.num_missing();
20          if (nm != 0) {
21              if (is_async_recovery_target(p)) {
22                  async_by_num_missing.push_back(make_pair(nm, p));
23              } else {
24                  replicas_by_num_missing.push_back(make_pair(nm, p));
25              }
26          }
27      }
28      // 排序函数
29      auto func = [](const std::pair<unsigned int, pg_shard_t> &lhs,
30                     const std::pair<unsigned int, pg_shard_t> &rhs) {
31          return lhs.first < rhs.first;
32      };
33
// 3 按缺失数量排序副本，从小到大
35      std::sort(replicas_by_num_missing.begin(), replicas_by_num_missing.end(),
   func);
36      std::sort(async_by_num_missing.begin(), async_by_num_missing.end(), func);
37      replicas_by_num_missing.insert(replicas_by_num_missing.end(),
38                                     async_by_num_missing.begin(),
   async_by_num_missing.end());
39
// 4 遍历副本
41      for (auto &replica : replicas_by_num_missing) {
42          pg_shard_t &peer = replica.second;
43          ceph_assert(peer != get_primary());
44          auto pm = recovery_state.get_peer_missing().find(peer);
45          ceph_assert(pm != recovery_state.get_peer_missing().end());
46          size_t m_sz = pm->second.num_missing();
47
48          // oldest first!
49          const pg_missing_t &m(pm->second);
50
// 4.1 遍历副本的缺失对象列表
52          for (map<version_t, hobject_t>::const_iterator p =
   m.get_rmissing().begin();
53               p != m.get_rmissing().end() && started < max;
```

```
54            ++p) {
55                handle.reset_tp_timeout();
56                const hobject_t soid(p->second);
57
58                if (recovery_state.get_missing_loc().is_unfound(soid)) {
59                    continue;
60                }
61
62 // 4.2 对象版本比副本的last_backfill新，说明数据已经存在副本OSD上了
63                const pg_info_t &pi = recovery_state.get_peer_info(peer);
64                if (soid > pi.last_backfill) {
65                    if (!recovering.count(soid)) {
66                        ceph_abort();
67                    }
68                    continue;
69                }
70
71                if (recovering.count(soid)) {
72                    continue;
73                }
74
75 // 4.3 如果对象已被删除，准备删除操作
76                if (recovery_state.get_missing_loc().is_deleted(soid)) {
77                    map<hobject_t, pg_missing_item>::const_iterator r =
   m.get_items().find(soid);
78                    started += prep_object_replica_deletes(soid, r->second.need,
   h, work_started);
79                    continue;
80                }
81
82 // 4.4 如果对象是快照，并且快照的头部对象也在缺失列表中
83                if (soid.is_snap() &&
84                    recovery_state.get_pg_log().get_missing().is_missing(
85                        soid.get_head())) {
86                    continue;
87                }
88
89                if (recovery_state.get_pg_log().get_missing().is_missing(soid)) {
90                    continue;
91                }
92 // 4.5 准备推送
93                map<hobject_t, pg_missing_item>::const_iterator r =
   m.get_items().find(soid);
94                started += prep_object_replica_pushes(soid, r->second.need, h,
   work_started);
95        }
96    }
```

```
97
98  // 5 推送
99      pgbackend->run_recovery_op(h, recovery_state.get_recovery_op_priority());
100     return started;
101 }
```

## handle_push

```
1  void ReplicatedBackend::handle_push(
2      pg_shard_t from, const PushOp &pop, PushReplyOp *response,
3      ObjectStore::Transaction *t, bool is_repair) {
4  /*
5  从 PushOp 结构体中提取数据和进度信息。
6  first 表示是否是恢复过程的开始。
7  complete 表示数据和对象映射（omap）是否已经完全恢复。
8  clear_omap 表示是否需要清除对象映射。
9  data_zeros 是一个区间集合，用于记录数据中的零区间。
10 z_offset 和 z_length 分别表示数据恢复的起始偏移量和长度。
11 */
12     bufferlist data;
13     data = pop.data;
14     bool first = pop.before_progress.first;
15     bool complete = pop.after_progress.data_complete &&
16                     pop.after_progress.omap_complete;
17
18     bool clear_omap = !pop.before_progress.omap_complete;
19     interval_set<uint64_t> data_zeros;
20     uint64_t z_offset = pop.before_progress.data_recovered_to;
21     uint64_t z_length = pop.after_progress.data_recovered_to -
   pop.before_progress.data_recovered_to;
22     if (z_length)
23         data_zeros.insert(z_offset, z_length);
24
25 // 1 构造响应对象
26     response->soid = pop.recovery_info.soid;
27
28 // 2 将数据提交给对象存储，应用更新
29     submit_push_data(pop.recovery_info, first, complete,
30         clear_omap, true, // must be replicate
31         data_zeros, pop.data_included, data, pop.omap_header,
32         pop.attrset, pop.omap_entries, t);
33
34     if (complete) {
35         if (is_repair) {
36             get_parent()->inc_osd_stat_repaired();
```

```
37             }
38  // 3 通知父对象本地修复完成
39         get_parent()->on_local_recover(
40             pop.recovery_info.soid,
41             pop.recovery_info,
42             ObjectContextRef(), // ok, is replica
43             false,
44             t);
45     }
46  }
```

# Backfill过程

Recovery过程通过日志修复对象，日志无法修复就通过Backfill过程直接拷贝数据。

## 数据结构

last_backfill 是backfill过程中，修复进程的指针。last_backfill初始化为 MIN对象，用来记录Backfiil过程中已修复的对象，last_backfill随着修复的进程不断推进，如果对象小于等于last_backfill，就是已经完成修复的对象，如果对象大于last_backfill，需要进一步判断是否需要修复。

如下，last_backfill从MIN推进，此时obj2对象以及obj1已经修复完成，之后的对象待修复。

| Bacakfill对象列表 | MIN | obj1(1,0) | obj2(1,1) | obj3(1,4) | obj4(1,5) | obj5(1, |
|---|---|---|---|---|---|---|
| | | | last_backfill | | | |

```
1  struct BackfillInterval {
2      eversion_t version;
3      std::map<hobject_t,eversion_t> objects;
4      hobject_t begin;
5      hobject_t end;
6  }
7  BackfillInterval backfill_info;
8  std::map<pg_shard_t, BackfillInterval> peer_backfill_info;
```

backfill_info在primary上记录backfill的进度，peer_backfill_info（replica 的 backfill_info的集合）记录其余副本的进度。

**earliset_peer_backfill() 和 earlist_backfill()**

这两个函数非常像，但是返回结果的含义完全不同：

```cpp
hobject_t PrimaryLogPG::earliest_peer_backfill() const {
    hobject_t e = hobject_t::get_max();
    for (const pg_shard_t &peer : get_backfill_targets()) {
        const auto iter = peer_backfill_info.find(peer);
        ceph_assert(iter != peer_backfill_info.end());
        e = std::min(e, iter->second.begin);
    }
    return e;
}

hobject_t PeeringState::earliest_backfill() const {
    hobject_t e = hobject_t::get_max();
    for (const pg_shard_t &bt : get_backfill_targets()) {
        const pg_info_t &pi = get_peer_info(bt);
        e = std::min(pi.last_backfill, e);
    }
    return e;
}
```

带peer的返回的是参与backfill的OSD中，最小的begin值。

而earliest_backfill返回的是参与backfill的OSD中，最小的last_backfill值。

## 流程

1. 首先，check指向所有副本的earliest_backfill()，然后判断check与primary上backfill_info.begin 的大小关系。

1.1 check < backfill_info.begin 判定为多余对象，加入to_move中。

1.2 check >= backfill_info.begin 继续判断是否需要backfill

2. 如果check == backfill_info.begin 且 backfill_info.begin == peer_backfill_info.begin （replica上 的对象和primary一致），进一步判断：如果版本号一致则不用backfill；不一致则加入 need_ver_targs，表示需要修复。

3. 如果2不成立，表示当前check对象和primary不一致，进一步判断：

3.1 如果replica的last_backfill < backfill_info.begin 表示对象缺失，加入missing_targs；

3.2 如果replica的last_backfill >= backfill_info.begin 表示replica上的该对象已经做过backfill，加入 skip_targs跳过对象

4. 更新指针 last_backfill_started = backfill_info.begin; backfill_info.pop_front(); 情况2中涉及到 的replica全部更新peer_backfill_info.pop_front();

# 示例

如下图：每一行是一个BackfillInterval，所有replica的BackfillInterval集合被称为
peer_backfill_info。

| OSD 0<br>info.lastbackfill =<br>hobject() | obj4(1,1)<br>pbi[0].begin | obj5(1,4) | obj6(1,10) |
|---|---|---|---|
| OSD 1<br>info.lastbackfill =<br>hobject() | | obj5(1,3)<br>pbi[1].begin | |
| OSD 2<br>info.lastbackfill =<br>hobject() | obj4(1,1)<br>pbi[2].begin | | obj6(1,4) |
| OSD 3<br>info.lastbackfill = obj5 | | ~~obj5(1,4)~~<br>从pbi[3]中剔除 | obj6(1,1)<br>pbi[3].begin |
| OSD 4<br>info.lastbackfill =<br>hobject() | | obj5(1,4)<br>pbi[4].begin | obj6(1,10) |
| OSD 5<br>（主） | | obj5(1,4)<br>backfill_info.begin | obj6(1,10) |
| | | | |

初始时，lastbackfill初始化为hobject()，是个空对象，而last_backfill_started = earliest_backfill()
也是一个空对象

上图中有一个特殊的OSD 3，它可能是之前做过backfill（听起来很怪，但不管出于什么原因，姑且当
它backfill的进度目前超过了Primary），此时它的last_backfill并不是MIN对象，而是obj5。

**第一次**recover_backfill：

check指向earliest_peer_backfill()

首先，发现check（obj4）＜backfill_info.begin，说明有冗余对象，查找和check相同的对象：OSD
0和OSD 2，加入to_move队列中准备删除。一直循环查找，直到冗余对象全部加入to_move队列，上
图的情况，循环一次即可。

同时，被删除的对象会出队

```
1  peer_backfill_info.pop_front();  // 也就是peer_backfill_info.begin前移
2  last_backfill_started = check
```

| | | | |
|---|---|---|---|
| OSD 0<br>info.lastbackfill =<br>hobject() | ~~obj4(1,1)~~ | obj5(1,4)<br>pbi[0].begin | obj6(1,10) |
| OSD 1<br>info.lastbackfill =<br>hobject() | | obj5(1,3)<br>pbi[1].begin | |
| OSD 2<br>info.lastbackfill =<br>hobject() | ~~obj4(1,1)~~ | | obj6(1,4)<br>pbi[2].begin |
| OSD 3<br>info.lastbackfill = obj5 | | | obj6(1,1)<br>pbi[3].begin |
| OSD 4<br>info.lastbackfill =<br>hobject() | | obj5(1,4)<br>pbi[4].begin | obj6(1,10) |
| OSD 5<br>(主) | | obj5(1,4)<br>backfill_info.begin | obj6(1,10) |
| | last_backfill_started | | |

**第二次循环** check = earliest_peer_backfill() ,此时一定大于等于backfill_info.begin，此时按OSD一个个对比：

1. 对于begin对象等于check对象的OSD: 0, 1, 4。判断版本是否正确

1.1 如果不正确就加入修复队列need_ver_targs

1.2 正确则加入keep_ver_targs，不需要修复

2. 对于begin对象不等于check对象的OSD: 2, 3。根据last_backfill判断其修复进度

2.1 对于OSD 2而言，其last_backfill小于backfill_info.begin，判定为缺失了obj5对象，加入missing列表

2.1 对于OSD 3而言，其last_backfill等于backfill_info.begin，判断其已经做过backfill，加入skip队列

下一步，need_ver_targs和keep_ver_targs中的begin全部前移

```
1  peer_backfill_info.pop_front();
2  last_backfill_started = backfill_info.begin
```

| | | |
|---|---|---|
| OSD 0<br>info.lastbackfill =<br>hobject() | ~~obj5(1,4)~~ | obj6(1,10)<br>pbi[0].begin |
| OSD 1<br>info.lastbackfill =<br>hobject() | ~~obj5(1,3)~~ | pbi[1].begin |
| OSD 2<br>info.lastbackfill =<br>hobject() | | obj6(1,4)<br>pbi[2].begin |
| OSD 3<br>info.lastbackfill = obj5 | | obj6(1,1)<br>pbi[3].begin |
| OSD 4<br>info.lastbackfill =<br>hobject() | ~~obj5(1,4)~~ | obj6(1,10)<br>pbi[4].begin |
| OSD 5<br>(主) | ~~obj5(1,4)~~ | obj6(1,10)<br>backfill_info.begin |
| | last_backfill_started | |

假设最大操作数是1，而刚刚我们循环了两次，第一次找了一列to_move对象，第二次处理了一列obj5对象。

1. 对于to_move对象的处理，不算一次"操作"，因为删除是异步的，且消耗资源较少。

2. 对于obj5，整理出了四种操作类型，分别对应队列：keep_ver_targs, need_ver_targs, skip_targs, missing_targs

每一次循环中，如果本次循环处理的不是to_move，那么会处理需要修复的对象：

```
// 构造PushOp做backfill
if (!need_ver_targs.empty() || !missing_targs.empty()) {
    ObjectContextRef obc = get_object_context(backfill_info.begin, false);
    prep_backfill_object_push(backfill_info.begin, obj_v, obc, all_push, h);
    ops ++; // 操作数
}
```

如果达到最大操作数后退出循环，并首先处理to_move队列，通知replica删除对象

```
pg_shard_t peer = to_remove[i].get<2>();
```

```
2    reqs[peer] = new MOSDPGBackfillRemove(spg_t(info.pgid.pgid, peer.shard),
     get_osdmap_epoch());
3    for (auto p : reqs) {
4        osd->send_message_osd_cluster(p.first.osd, p.second, get_osdmap_epoch());
5    }
6    // 通知过后，紧接着将准备好的PushOp发出去
7    pgbackend->run_recovery_op(h, recovery_state.get_recovery_op_priority());
```

完成backfill后，更新replica的last_backfill：

```
1    hobject_t new_last_backfill = recovery_state.earliest_backfill();
2    recovery_state.update_peer_last_backfill(bt, new_last_backfill);
```

### 第二次recover_backfill

| | | |
|---|---|---|
| OSD 0<br>info.lastbackfill = obj5 | ~~obj5(1,4)~~ | obj6(1,10)<br>pbi[0].begin |
| OSD 1<br>info.lastbackfill = obj5 | ~~obj5(1,3)~~ | pbi[1].begin |
| OSD 2<br>info.lastbackfill = obj5 | | obj6(1,4)<br>pbi[2].begin |
| OSD 3<br>info.lastbackfill = obj5 | | obj6(1,1)<br>pbi[3].begin |
| OSD 4<br>info.lastbackfill = obj5 | ~~obj5(1,4)~~ | obj6(1,10)<br>pbi[4].begin |
| OSD 5<br>(主) | ~~obj5(1,4)~~ | obj6(1,10)<br>backfill_info.begin |
| | last_backfill_started | |

此时check = earliest_peer_backfill()等于obj6。和obj5一样的处理思路：

1. 对于和check相等的OSD：0 2 3 4：0，4版本正确，不用修复，2，3版本错误，需要修复

2. 和check不相等的OSD 1（此时它已经是空对象了），判断它的last_backfill <
   backfill_info.begin，判定为缺失对象，加入missing中。

同样的，它们的begin会再次前移，直到backfill_info为空或者达到最大并发处理数

## 源码

```
1   uint64_t PrimaryLogPG::recover_backfill(
2       uint64_t max,
3       ThreadPool::TPHandle &handle, bool *work_started) {
4       ceph_assert(!get_backfill_targets().empty());
5
6   // 1 初始化backfill区间
7       // Initialize from prior backfill state
8       if (new_backfill) {
9           // on_activate() was called prior to getting here
10          ceph_assert(last_backfill_started ==
    recovery_state.earliest_backfill());
11          new_backfill = false;
12
13          // initialize BackfillIntervals
14          for (set<pg_shard_t>::const_iterator i =
    get_backfill_targets().begin();
15              i != get_backfill_targets().end();
16              ++i) {
17          peer_backfill_info[*i].reset(
18              recovery_state.get_peer_info(*i).last_backfill);
19          }
20          backfill_info.reset(last_backfill_started);
21
22          backfills_in_flight.clear();
23          // 保存着需要删除的对象
24          pending_backfill_updates.clear();
25      }
26
27  // 2 更新backfill_info.begin, update_range更新需要进行Backfill操作的对象列表
28      backfill_info.begin = last_backfill_started;
29      update_range(&backfill_info, handle);
30
31      unsigned ops = 0;
32      vector<boost::tuple<hobject_t, eversion_t, pg_shard_t>> to_remove;
33      set<hobject_t> add_to_stat;
34
35  // 3 去掉不需要backfill的区间
36      for (set<pg_shard_t>::const_iterator i = get_backfill_targets().begin();
37          i != get_backfill_targets().end();
```

```
38                  ++i) {
39          peer_backfill_info[*i].trim_to(
40              std::max(
41                  recovery_state.get_peer_info(*i).last_backfill,
42                  last_backfill_started));
43      }
44      backfill_info.trim_to(last_backfill_started);
45
46      PGBackend::RecoveryHandle *h = pgbackend->open_recovery_op();
47      while (ops < max) {
48  // 4 当前区间完成backfill后，添加新的backfill区间
49          if (backfill_info.begin <= earliest_peer_backfill() &&
50              !backfill_info.extends_to_end() && backfill_info.empty()) {
51              hobject_t next = backfill_info.end;
52              backfill_info.reset(next);
53              backfill_info.end = hobject_t::get_max();
54              update_range(&backfill_info, handle);
55              backfill_info.trim();
56          }
57
58          bool sent_scan = false;
59          for (set<pg_shard_t>::const_iterator i =
    get_backfill_targets().begin();
60              i != get_backfill_targets().end();
61              ++i) {
62              pg_shard_t bt = *i;
63              BackfillInterval &pbi = peer_backfill_info[bt];
64
65  // 5 当前backfillInterval没有需要回填的区间，对相应副本发起扫描操作，更新数据
66              if (pbi.begin <= backfill_info.begin &&
67                  !pbi.extends_to_end() && pbi.empty()) {
68                  epoch_t e = get_osdmap_epoch();
69                  MOSDPGScan *m = new MOSDPGScan(
70                      MOSDPGScan::OP_SCAN_GET_DIGEST, pg_whoami, e,
    get_last_peering_reset(),
71                      spg_t(info.pgid.pgid, bt.shard),
72                      pbi.end, hobject_t());
73
74                  if (cct->_conf->osd_op_queue == "mclock_scheduler") {
75                      /* This guard preserves legacy WeightedPriorityQueue
    behavior for
76                       * now, but should be removed after Reef */
77                      m->set_priority(recovery_state.get_recovery_op_priority());
78                  }
79                  osd->send_message_osd_cluster(bt.osd, m, get_osdmap_epoch());
80                  ceph_assert(waiting_on_backfill.find(bt) ==
    waiting_on_backfill.end());
```

```
 81                 waiting_on_backfill.insert(bt);
 82                 sent_scan = true;
 83             }
 84         }
 85
 86         // Count simultaneous scans as a single op and let those complete
 87  // 6 获取OSd的对象列表后，对比当前主OSD的对象列表来进行修复
 88         if (sent_scan) {
 89             ops++;
 90             start_recovery_op(hobject_t::get_max()); // XXX: was pbi.end
 91             break;
 92         }
 93
 94         if (backfill_info.empty() && all_peer_done()) {
 95             break;
 96         }
 97
 98         // Get object within set of peers to operate on and
 99         // the set of targets for which that object applies.
100         hobject_t check = earliest_peer_backfill();
101  // 7 check指向当前OSD中最小的需要进行Backfill操作的对象
102         if (check < backfill_info.begin) {
103  // 7.1 冗余对象，加入到to_remove队列中
104             set<pg_shard_t> check_targets;
105             for (set<pg_shard_t>::const_iterator i =
    get_backfill_targets().begin();
106                  i != get_backfill_targets().end();
107                  ++i) {
108                 pg_shard_t bt = *i;
109                 BackfillInterval &pbi = peer_backfill_info[bt];
110                 if (pbi.begin == check)
111                     check_targets.insert(bt);
112             }
113             ceph_assert(!check_targets.empty());
114
115             for (set<pg_shard_t>::iterator i = check_targets.begin();
116                  i != check_targets.end();
117                  ++i) {
118                 pg_shard_t bt = *i;
119                 BackfillInterval &pbi = peer_backfill_info[bt];
120                 ceph_assert(pbi.begin == check);
121
122                 to_remove.push_back(boost::make_tuple(check,
    pbi.objects.begin()->second, bt));
123                 pbi.pop_front();
124             }
125
```

```
126                 last_backfill_started = check;
127             } else {
// 7.2 check等于backfill_info_begin，判断对象是否需要backfill
129             eversion_t &obj_v = backfill_info.objects.begin()->second;
130             vector<pg_shard_t> need_ver_targs, missing_targs, keep_ver_targs,
    skip_targs;
131             for (set<pg_shard_t>::const_iterator i =
    get_backfill_targets().begin();
132                  i != get_backfill_targets().end();
133                  ++i) {
134                 pg_shard_t bt = *i;
135                 BackfillInterval &pbi = peer_backfill_info[bt];
136                 // Find all check peers that have the wrong version
137                 if (check == backfill_info.begin && check == pbi.begin) {
138                     if (pbi.objects.begin()->second != obj_v) {
139                         need_ver_targs.push_back(bt);
140                     } else {
141                         keep_ver_targs.push_back(bt);
142                     }
143                 } else {
144                     const pg_info_t &pinfo = recovery_state.get_peer_info(bt);
145                     if (backfill_info.begin > pinfo.last_backfill)
146                         missing_targs.push_back(bt);
147                     else
148                         skip_targs.push_back(bt);
149                 }
150             }
151
152             if (!keep_ver_targs.empty()) {
153                 // These peers have version obj_v
154                 // assert(!waiting_for_degraded_object.count(check));
155             }
// 8 对于keep_ver_targs中的OSD，不需要进行Backfill操作
// 对于need_ver_targs和missing_targs中的OSD，需要进行Backfill操作，准备PushOp
158             if (!need_ver_targs.empty() || !missing_targs.empty()) {
159                 ObjectContextRef obc = get_object_context(backfill_info.begin,
    false);
160                 ceph_assert(obc);
161                 if (obc->get_recovery_read()) {
162                     vector<pg_shard_t> all_push = need_ver_targs;
163                     all_push.insert(all_push.end(), missing_targs.begin(),
    missing_targs.end());
164
165                     handle.reset_tp_timeout();
166                     int r = prep_backfill_object_push(backfill_info.begin,
    obj_v, obc, all_push, h);
167                     if (r < 0) {
```

```
168                        *work_started = true;
169                            break;
170                    }
171                    ops++;
172                } else {
173                        *work_started = true;
174                        break;
175                }
176            }
177
178            last_backfill_started = backfill_info.begin;
179            add_to_stat.insert(backfill_info.begin);
180            backfill_info.pop_front();
181            vector<pg_shard_t> check_targets = need_ver_targs;
182            check_targets.insert(check_targets.end(), keep_ver_targs.begin(),
    keep_ver_targs.end());
183            for (vector<pg_shard_t>::iterator i = check_targets.begin();
184                 i != check_targets.end();
185                 ++i) {
186                pg_s    hard_t bt = *i;
187                BackfillInterval &pbi = peer_backfill_info[bt];
188                pbi.pop_front();
189            }
190        }
191    }
192
193    for (set<hobject_t>::iterator i = add_to_stat.begin();
194         i != add_to_stat.end();
195         ++i) {
196        ObjectContextRef obc = get_object_context(*i, false);
197        ceph_assert(obc);
198        pg_stat_t stat;
199        add_object_context_to_pg_stat(obc, &stat);
200        pending_backfill_updates[*i] = stat;
201    }
202
203 // 9 冗余对象
204    map<pg_shard_t, MOSDPGBackfillRemove *> reqs;
205    for (unsigned i = 0; i < to_remove.size(); ++i) {
206        handle.reset_tp_timeout();
207        const hobject_t &oid = to_remove[i].get<0>();
208        eversion_t v = to_remove[i].get<1>();
209        pg_shard_t peer = to_remove[i].get<2>();
210        MOSDPGBackfillRemove *m;
211        auto it = reqs.find(peer);
212        if (it != reqs.end()) {
213            m = it->second;
```

```
214          } else {
215              m = reqs[peer] = new MOSDPGBackfillRemove(
216                  spg_t(info.pgid.pgid, peer.shard),
217                  get_osdmap_epoch());
218              if (cct->_conf->osd_op_queue == "mclock_scheduler") {
219                  m->set_priority(recovery_state.get_recovery_op_priority());
220              }
221          }
222          m->ls.push_back(make_pair(oid, v));
223
224          if (oid <= last_backfill_started)
225              pending_backfill_updates[oid]; // add empty stat!
226      }
227      for (auto p : reqs) {
228          osd->send_message_osd_cluster(p.first.osd, p.second,
    get_osdmap_epoch());
229      }
230
231 // 10 发送PushOp
232      pgbackend->run_recovery_op(h, recovery_state.get_recovery_op_priority());
233
234 // 11 处理已完成backfill的对象，更新backfill位置和状态
235      hobject_t next_backfill_to_complete = backfills_in_flight.empty() ?
    backfill_pos : *(backfills_in_flight.begin());
236      hobject_t new_last_backfill = recovery_state.earliest_backfill();
237      for (map<hobject_t, pg_stat_t>::iterator i =
238              pending_backfill_updates.begin();
239          i != pending_backfill_updates.end() &&
240          i->first < next_backfill_to_complete;
241          pending_backfill_updates.erase(i++)) {
242      ceph_assert(i->first > new_last_backfill);
243      recovery_state.update_complete_backfill_object_stats(
244          i->first,
245          i->second);
246      new_last_backfill = i->first;
247      }
248
249      ceph_assert(!pending_backfill_updates.empty() ||
250              new_last_backfill == last_backfill_started);
251      if (pending_backfill_updates.empty() &&
252          backfill_pos.is_max()) {
253          ceph_assert(backfills_in_flight.empty());
254          new_last_backfill = backfill_pos;
255          last_backfill_started = backfill_pos;
256      }
257
258 // 12 更新各个OSD的backfill位置
```

```
259        for (set<pg_shard_t>::const_iterator i = get_backfill_targets().begin();
260             i != get_backfill_targets().end();
261             ++i) {
262          pg_shard_t bt = *i;
263          const pg_info_t &pinfo = recovery_state.get_peer_info(bt);
264
265          if (new_last_backfill > pinfo.last_backfill) {
266            recovery_state.update_peer_last_backfill(bt, new_last_backfill);
267            epoch_t e = get_osdmap_epoch();
268            MOSDPGBackfill *m = NULL;
// 12.1 如果值为MAX，说明backfill完成，发送OP_BACKFILL_FINISH
270            if (pinfo.last_backfill.is_max()) {
271              m = new MOSDPGBackfill(
272                MOSDPGBackfill::OP_BACKFILL_FINISH,
273                e,
274                get_last_peering_reset(),
275                spg_t(info.pgid.pgid, bt.shard));
276              // Use default priority here, must match sub_op priority
277              start_recovery_op(hobject_t::get_max());
278            } else {
// 12.2 否则发送OP_BACKFILL_PROGRESS
280              m = new MOSDPGBackfill(
281                MOSDPGBackfill::OP_BACKFILL_PROGRESS,
282                e,
283                get_last_peering_reset(),
284                spg_t(info.pgid.pgid, bt.shard));
285              // Use default priority here, must match sub_op priority
286            }
287            m->last_backfill = pinfo.last_backfill;
288            m->stats = pinfo.stats;
289
290            if (cct->_conf->osd_op_queue == "mclock_scheduler") {
291              /* This guard preserves legacy WeightedPriorityQueue behavior
    for
292               * now, but should be removed after Reef */
293              m->set_priority(recovery_state.get_recovery_op_priority());
294            }
// 12.3 发送消息
296            osd->send_message_osd_cluster(bt.osd, m, get_osdmap_epoch());
297          }
298        }
299
300      if (ops)
301        *work_started = true;
302      return ops;
303    }
```

# 对象排序方法

做backfill时，即使各个OSD中相同版本号不同，也能保证各个OSD内对象的相对顺序一致。

```
1   // 首先对对象计算出对象的唯一hash值
2   // 传入参数中一共就两个成员变量：对象名，快照序号
3     explicit hobject_t(const sobject_t &o) :
4       oid(o.oid), snap(o.snap), max(false), pool(POOL_META) {
5       set_hash(std::hash<sobject_t>()(o));
6     }
7
8     namespace std {
9   template<> struct hash<hobject_t> {
10    size_t operator()(const hobject_t &r) const {
11      static rjhash<uint64_t> RJ;
12      return RJ(r.get_hash() ^ r.snap);
13    }
14  };
15  }
16
17  inline uint64_t rjhash64(uint64_t key) {
18    key = (~key) + (key << 21); // key = (key << 21) - key - 1;
19    key = key ^ (key >> 24);
20    key = (key + (key << 3)) + (key << 8); // key * 265
21    key = key ^ (key >> 14);
22    key = (key + (key << 2)) + (key << 4); // key * 21
23    key = key ^ (key >> 28);
24    key = key + (key << 31);
25    return key;
26  }
27
28  // 然后根据hash值，计算出用于排序的序号
29    void build_hash_cache() {
30      nibblewise_key_cache = _reverse_nibbles(hash);
31      hash_reverse_bits = _reverse_bits(hash);
32    }
33
34  uint32_t reverse_bits(uint32_t v) {
35    if (v == 0)
36      return v;
37
38    /* reverse bits
39     * swap odd and even bits
40     */
41    v = ((v >> 1) & 0x55555555) | ((v & 0x55555555) << 1);
42    /* swap consecutive pairs */
```

```
43    v = ((v >> 2) & 0x33333333) | ((v & 0x33333333) << 2);
44    /* swap nibbles ... */
45    v = ((v >> 4) & 0x0F0F0F0F) | ((v & 0x0F0F0F0F) << 4);
46    /* swap bytes */
47    v = ((v >> 8) & 0x00FF00FF) | ((v & 0x00FF00FF) << 8);
48    /* swap 2-byte long pairs */
49    v = ( v >> 16               ) | ( v                << 16);

50    return v;
51 }
52
53 // 对象放入backfill的map中后，根据以下规则排列，一般到reverse_bits就能计算出结果
54   auto operator<=>(const hobject_t &rhs) const noexcept {
55     auto cmp = max <=> rhs.max;              // 比较max与rhs.max，相等返回0，不相等返回比较结果
56     if (cmp != 0) return cmp;                              // 不相等返回比较结果
57     cmp = pool <=> rhs.pool;                // 比较pool
58     if (cmp != 0) return cmp;
59     cmp = get_bitwise_key() <=> rhs.get_bitwise_key();    // reverse_bits
60     if (cmp != 0) return cmp;
61     cmp = nspace <=> rhs.nspace;            // 命名空间
62     if (cmp != 0) return cmp;
63     if (!(get_key().empty() && rhs.get_key().empty())) {
64       cmp = get_effective_key() <=> rhs.get_effective_key(); // 自定义key和对象名
65       if (cmp != 0) return cmp;
66     }
67     cmp = oid <=> rhs.oid;                   // 对象id
68     if (cmp != 0) return cmp;
69     return snap <=> rhs.snap;               // 快照信息
70   }
71
```