

一致性检查

概念

现实中，我们的计算机组件都是不完美的，都有可能产生静默数据错误，即便通过数据修复，每个副本仍然可能存在本地数据错误，副本之间也仍然可能会有数据一致性问题。

Scrub 通过对比各个对象副本的数据和元数据，实现副本的一致性检查。它不是一种实时的纠错机制，而是一种事后的检查机制。

按照扫描的内容分为两种方式：

浅扫描：Scrub，通过检查元数据来检查数据的一致性，如果元数据一致就通过检查

深扫描：deep-Scrub，进一步检查对象的数据内容是否一致，要扫描磁盘上的所有数据，比较耗时，占用的资源也更多。

Scrub有在线扫描（不影响系统正常运行）和离线扫描（需要系统暂停业务）两种方式，一般情况下我们使用在线扫描，这个过程不会中断系统，但是被扫描的对象会被锁定暂停访问，直到对象完成Scrub操作后才能解锁。

扫描粒度

如何确定Scrub的扫描粒度，

1. 以对象为粒度，集群中往往存在大量的对象，逐个对象执行scrub的效率太低了，且对象的生命周期不固定，对象可能对是被删除和创建。
2. 以PG为粒度，集群中PG的数量比较固定，生命周期也更稳定，OSD的读写等都围绕PG进行组织，所以使用PG作为粒度是比较合理的
3. 2仍然有问题，如果一个PG中含有大量的对象，这种锁粒度会导致PG的写请求被长时间阻塞，于是有了新的方案：针对当前PG，一次只扫描部分对象，分批次扫描完整个PG。这个方案的可取之处在于：PG内的对象顺序基本保持一致，新对象的加入不会打乱旧对象之间的相对顺序，也就是说即使有对象创建或删除，也能保证同一个对象不被重复扫描，而未被扫描的对象，会被安排在下一轮的扫描中被检查。

ScrubMap

保存校验对象以及相应校验数据

```
1 struct ScrubMap {  
2     struct object {
```

```

3      std::map<std::string, ceph::buffer::ptr, std::less<>> attrs; // 对象的属性
4      uint64_t size; // 对象的大小
5      __u32 omap_digest; // 对象映射的CRC32C校验和
6      __u32 digest; // 数据的CRC32C校验和
7      bool digest_present : 1; // 表示数据校验和是否存在
8      bool omap_digest_present : 1; // 表示omap校验和是否存在
9      ...
10 };
11 ...
12 std::map<hobject_t, object> objects; // 需要校验的对象 -> 校验信息的映射
13 };

```

调度

调度总流程

```

1 # 注册定时任务, C_Tick_WithoutOSDLock回调函数是tick_without_osd_lock
2 OSD::init
3 >> tick_timer_without_osd_lock.add_event_after(get_tick_interval(),
4         new C_Tick_WithoutOSDLock(this));
5
6 OSD::tick_without_osd_lock
7 >> OSD::sched_scrub
8
9 OSD::sched_scrub
10 >> Scrub::schedule_result_t ScrubQueue::select_pg_and_scrub
11 >>>> Scrub::schedule_result_t ScrubQueue::select_from_group
12 >>>>> Scrub::schedule_result_t OSDService::initiate_a_scrub
13 >>>>>> pg::sched_scrub
14
15 pg::sched_scrub
16 >> osd->queue_for_scrub(this, Scrub::scrub_prio_t::low_priority);

```

OSD::sched_scrub

控制启动时机

```

1 void OSD::sched_scrub() {
2 // 1 获取调度器
3     auto &scrub_scheduler = service.get_scrub_services();
4 // 2 有被锁定的对象

```

```

5     if (auto blocked_pgs = scrub_scheduler.get_blocked_pgs_count();
6         blocked_pgs > 0) {
7         // 此OSD管理的某些PG在清除过程中被锁定的对象阻止。
8         // 这意味着我们现在可能没有清理所需的资源。
9     }
10
11 // 3 没有可用资源进行Scrub
12     if (!scrub_scheduler.can_inc_scrubs()) {
13         return;
14     }
15     /*
16     scrub_scheduler.can_inc_scrubs():
17         if (scrubs_local + scrubs_remote < conf()->osd_max_scrubs) {
18             return true;
19         }
20     */
21
22     // 如果有一个PG正在尝试保留清理副本资源
23     // 我们应该等待，不要开始新的Scrub
24     // bool is_reserving_now() const { return a_pg_is_reserving; }
25     if (scrub_scheduler.is_reserving_now()) {
26         return;
27     }
28
29     Scrub::ScrubPreconds env_conditions;
30 // 4 如果处于恢复期间，判断是否允许在回复期间进行Scrub
31     if (service.is_recovery_active() && !cct->_conf-
32         >osd_scrub_during_recovery) {
33         if (!cct->_conf->osd_repair_during_recovery) {
34             return;
35         }
36         env_conditions.allow_requested_repair_only = true;
37     }
38     if (g_conf()->subsys.should_gather<ceph_subsys_osd, 20>()) {
39         auto all_jobs = scrub_scheduler.list_registered_jobs();
40     }
41 // 5 选择一个PG开始检查
42     auto was_started = scrub_scheduler.select_pg_and_scrub(env_conditions);
43 }
44

```

```

1 Scrub::schedule_result_t ScrubQueue::select_pg_and_scrub(
2     Scrub::ScrubPreconds& preconds) {
3

```

```

4   utime_t now_is = time_now();
5   // 1 检查是否在允许的时间段内,
6   // 检查now是否在配置中的区间: [begin_time, end_time)或[end_time, begin_time)
7   preconds.time_permit = scrub_time_permit(now_is);
8   // 2 检查负载是否允许
9   preconds.load_is_low = scrub_load_below_threshold();
10  preconds.only_deadlined = !preconds.time_permit || !preconds.load_is_low;
11
12  std::unique_lock lck{jobs_lock};
13
14  // 3 检查处罚列表, 如果某个作业已到截止时间, 或已被更新, 就移除
15  scan_penalized(restore_penalized, now_is);
16  restore_penalized = false;
17
18  // 4 remove the 'updated' flag from all entries
19  std::for_each(to_scrub.begin(),
20               to_scrub.end(),
21               [](const auto& jobref) -> void { jobref->updated = false; });
22
23  // 5 将那些在之前的Scrub尝试中失败的PG移动到处罚列表中
24  move_failed_pgs(now_is);
25
26  // 从两个列表中收集所有有效且成熟的作业副本
27  auto to_scrub_copy = collect_ripe_jobs(to_scrub, now_is);
28  auto penalized_copy = collect_ripe_jobs(penalized, now_is);
29  lck.unlock(); // 解锁互斥锁
30
31  // 6 首先尝试从常规队列中选择一个 PG 进行Scrub
32  auto res = select_from_group(to_scrub_copy, preconds, now_is);
33
34  // 7 如果常规队列中没有准备好的作业, 且处罚队列不为空, 则尝试从处罚队列中选择
35  if (res == Scrub::schedule_result_t::none_ready && !penalized_copy.empty()) {
36      res = select_from_group(penalized_copy, preconds, now_is);
37      restore_penalized = true; // 设置恢复处罚作业的标志为 true
38  }
39
40  return res;
41 }

```

penalized 直译为“处罚”

ceph中指的是某些 PG 因为特定的问题或错误而被标记为需要特别处理的状态。这些问题可能包括但不限于数据损坏、副本不一致、或者在之前的Scrub过程中发现的错误。处罚是一种机制, 用于确保这些有问题的 PGs能够被优先处理, 以便修复或恢复它们的状态。

处罚的 PG 可能会被暂时从正常的Scrub队列中移除，并放入一个特殊的处罚列表中。这样，系统就可以在处理正常Scrub任务的同时，特别关注这些需要额外注意的 PG。处罚列表中的 PGs 会在一定条件下被重新评估，例如在负载较低或者在特定的时间窗口内，以便进行必要的Scrub或修复操作。在函数中，处罚的 PGs 会在特定条件下被考虑进行Scrub，例如当常规的Scrub队列中没有其他 PG 准备好进行Scrub时。这种机制有助于确保 Ceph 集群中的数据完整性和可靠性。

```
1 Scrub::schedule_result_t ScrubQueue::select_from_group(
2     ScrubQContainer& group,
3     const Scrub::ScrubPreconds& preconds,
4     utime_t now_is)
5 {
6     for (auto& candidate : group) {
7         if (preconds.only_deadlined && (candidate->schedule.deadline.is_zero() ||
8             candidate->schedule.deadline >= now_is)) {
9             continue;
10        }
11        switch (
12            osd_service.initiate_a_scrub(candidate->pgid,
13                preconds.allow_requested_repair_only)) {
14            case Scrub::schedule_result_t::scrub_initiated:
15                return Scrub::schedule_result_t::scrub_initiated;
16            // 其它的case都是return none_ready, 省略不看了
17        }
18    }
19    return Scrub::schedule_result_t::none_ready;
20 }
21
```

```
1 Scrub::schedule_result_t OSDService::initiate_a_scrub(spg_t pgid,
2     bool allow_requested_repair_only) {
3
4     // 我们有一个候选人要淘汰。我们需要一些PG信息来了解是否允许Scrub
5
6     PGRef pg = osd->lookup_lock_pg(pgid);
7     if (!pg) {
8         return Scrub::schedule_result_t::no_such_pg;
9     }
10
11    // This has already started, so go on to the next scrub job
12    if (pg->is_scrub_queued_or_active()) {
13        pg->unlock();
14    }
15
16    // ... (rest of the function code) ...
17
18    return Scrub::schedule_result_t::scrub_initiated;
19 }
20
```

```

14         return Scrub::schedule_result_t::already_started;
15     }
16     // Skip other kinds of scrubbing if only explicitly requested repairing is
    allowed
17     if (allow_requested_repair_only && !pg->get_planned_scrub().must_repair) {
18         pg->unlock();
19         return Scrub::schedule_result_t::preconditions;
20     }
21
22     auto scrub_attempt = pg->sched_scrub();
23     pg->unlock();
24     return scrub_attempt;
25 }

```

pg::sched_scrub

设置Scrub任务参数，并完成本地资源预约

```

1 Scrub::schedule_result_t PG::sched_scrub() {
2     using Scrub::schedule_result_t;
3     ceph_assert(ceph_mutex_is_locked(_lock));
4     ceph_assert(m_scrubber);
5     // 1 状态检测
6     if (is_scrub_queued_or_active()) {
7         return schedule_result_t::already_started;
8     }
9     if (!is_primary() || !is_active() || !is_clean()) {
10        return schedule_result_t::bad_pg_state;
11    }
12    if (state_test(PG_STATE_SNAPTRIM) || state_test(PG_STATE_SNAPTRIM_WAIT)) {
13        return schedule_result_t::bad_pg_state;
14    }
15
16    // 2 决定Scrub的深浅
17    auto updated_flags = validate_scrub_mode();
18    if (!updated_flags) { // 无法开始Scrub
19        return schedule_result_t::preconditions;
20    }
21    // 3 本地资源预约
22    if (!m_scrubber->reserve_local()) {
23        return schedule_result_t::no_local_resources;
24    }
25
26    m_planned_scrub = *updated_flags;
27

```

```

28     state_clear(PG_STATE_REPAIR);
29
30     m_scrubber->set_op_parameters(m_planned_scrub);
31
32 // 4 加入Scrub队列，等待调度
33     osd->queue_for_scrub(this, Scrub::scrub_prio_t::low_priority);
34     return schedule_result_t::scrub_initiated;
35 }

```

[illegible]

```

33                                     time_for_deep,
34                                     has_deep_errors,
35                                     m_planned_scrub);
36     } else {
37 // 5 验证周期性Scrub
38         ceph_assert(!m_planned_scrub.must_deep_scrub);
39         upd_flags = validate_periodic_mode(allow_deep_scrub,
40                                           try_to_auto_repair,
41                                           allow_shallow_scrub,
42                                           time_for_deep,
43                                           has_deep_errors,
44                                           m_planned_scrub);
45         if (!upd_flags) { // 取消周期性scrub
46             return std::nullopt;
47         }
48     }
49     upd_flags->need_auto = false;
50     return upd_flags;
51 }

```

validate_initiated_scrub和validate_periodic_mode内部会确认Scrub类型，依据下面这张表：

1	Periodic	type		none		no-scrub		no-scrub+no-deep		no-deep
2	-----									
3	-----									
4	periodic			shallow		x		x		shallow
5	-----									
6	periodic + t.f.deep			deep		deep		x		shallow
7	-----									
8	initiated			shallow		shallow		shallow		shallow
9	-----									
10	init. + t.f.deep			deep		deep		shallow		shallow
11	-----									
12	initiated deep			deep		deep		deep		deep
13	-----									
14	"periodic"	- 周期性Scrub			- if !must_scrub && !must_deep_scrub;					
15	"initiated deep"	- 手动发起的深层Scrub			- if must_scrub && must_deep_scrub;					
16	"initiated"	- 手动发起的浅层Scrub			- if must_scrub && !must_deep_scrub;					
17	"t.f.deep"	- triggered deep scrub			触发深层Scrub，意味着系统检测到了需要进行深层Scrub					

这张表展示在什么操作类型下，遇到什么PG状态/配置会发起怎样程度的Scrub，比如 `init. + t.f.deep` 操作遇到none，执行深层Scrub，而遇到no-deep状态执行浅层Scrub。

执行

通过比较对象各个副本上的元数据和数据，来完成元数据和数据的校验。

状态机

Scrub状态机的转换是这样的一个套路：

1.调用OSD层的queue_for_srcub_xxx

```
1 m_osds->queue_for_scrub_resched(m_pg, Scrub::scrub_prio_t::high_priority);
```

2.封装成PGScrub任务，要转什么状态，就用对应的实现，比如这里是PGScrubResched

```
1 void OSDService::queue_for_scrub_resched(PG *pg, Scrub::scrub_prio_t
  with_priority) {
2     queue_scrub_event_msg<PGScrubResched>(pg, with_priority);
3 }
```

3.进一步封装成OPItem，并加入OSD的任务队列

```
1 enqueue_back(OpSchedulerItem(
2     unique_ptr<OpSchedulerItem::OpQueueable>(msg), get_scrub_cost(),
3     pg->scrub_requeue_priority(with_priority), ceph_clock_now(), 0,
    epoch));
```

4.从任务队列取出后，调用run，比如这里是PGScrubResched实现

```
1 void PGScrubResched::run(OSD* osd,
2     OSDShard* sdata,
3     PGRef& pg,
4     ThreadPool::TPHandle& handle)
5 {
6     pg->scrub_send_scrub_resched(epoch_queued, handle);
7     pg->unlock();
8 }
```

5.这个scrub_send_scrub_xxx 也不知道用的什么抽象的技术，反正是看不懂的

```

1 void scrub_send_scrub_resched(epoch_t queued, ThreadPool::TPHandle& handle)
2 {
3     forward_scrub_event(&ScrubPgIF::send_scrub_resched, queued,
4         "InternalSchedScrub");
5 }
6 void PG::forward_scrub_event(ScrubAPI fn, epoch_t epoch_queued,
7     std::string_view desc) {
8     ceph_assert(m_scrubber);
9     if (is_active()) {
10         ((*m_scrubber).*fn)(epoch_queued);
11     } else {
12         // pg might be in the process of being deleted
13     }
14 }

```

6.实际调用的是send_scrub_xxx函数（看上一段代码的第三行）

```

1 void PgScrubber::send_scrub_resched(epoch_t epoch_queued) {
2     if (is_message_relevant(epoch_queued)) {
3         // 抛出InternalSchedScrub事件
4         m_fsm->process_event(InternalSchedScrub{});
5     }
6 }

```

7.最后投递的事件被当前状态接收，完成状态转换

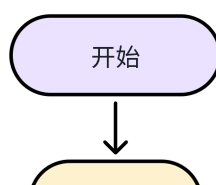
```

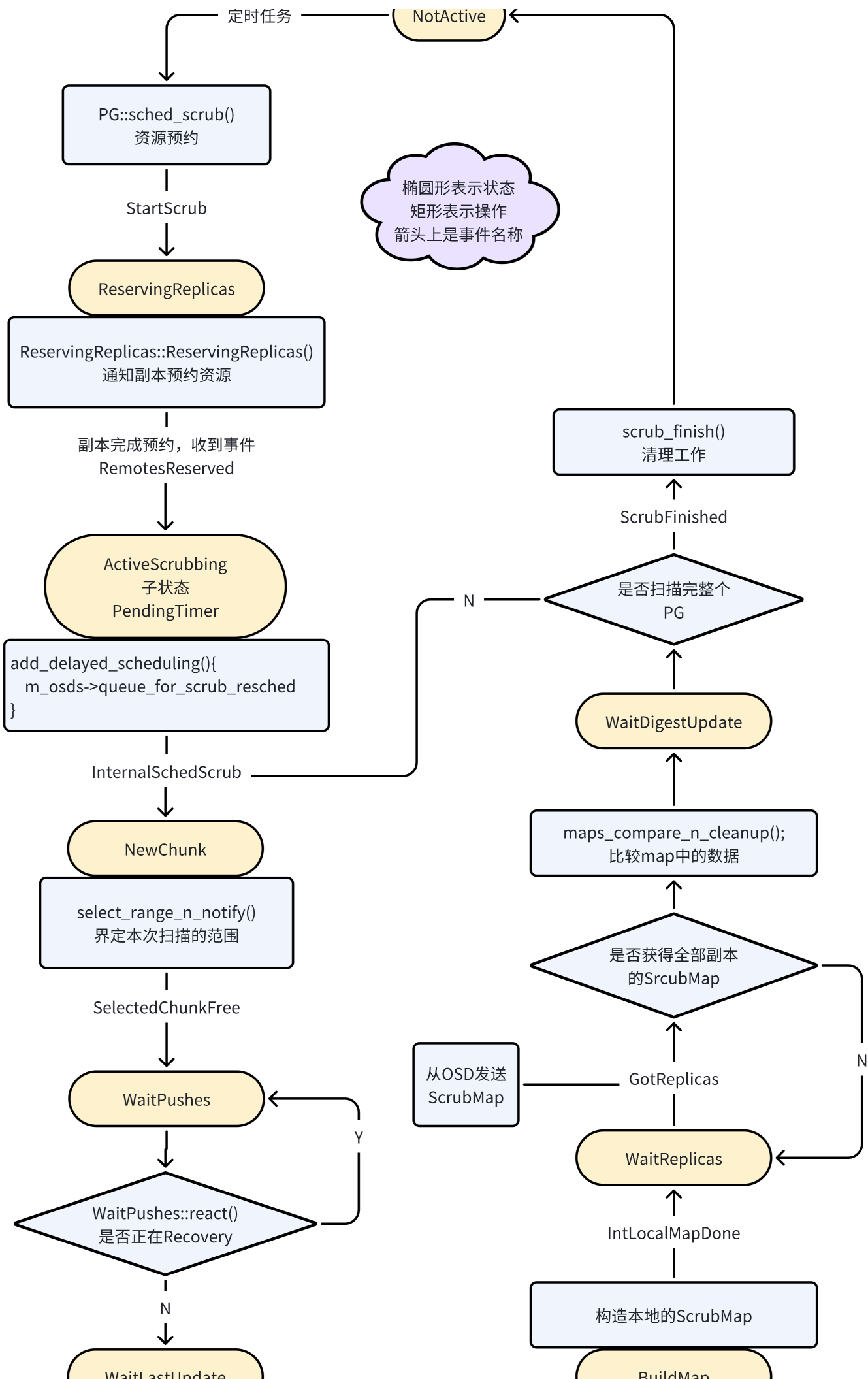
1 struct PendingTimer : sc::state<PendingTimer, ActiveScrubbing>, NamedSimply {
2     // 当收到InternalSchedScrub事件时，自动转到状态NewChunk
3     using reactions = mpl::list<sc::transition<InternalSchedScrub, NewChunk>>;
4 };

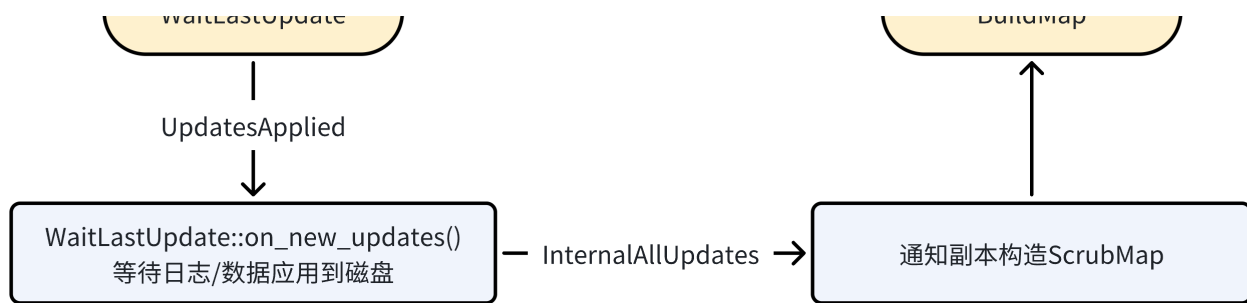
```

说实话是理解不了的，为什么切换状态要涉及这么多层，还要重新获得线程资源，为了方便被随时抢占吗。。。

总体流程







资源预留

```

1 void PGScrub::run() ==> 抛出StartScrub事件
2
3 sc::result NotActive::react(const StartScrub&)
4 {
5     DECLARE_LOCALS;
6     scrbr->set_scrub_begin_time();
7     return transit<ReservingReplicas>();
8 }
9
10 ReservingReplicas::ReservingReplicas(my_context ctx) : my_base(ctx)
11     , NamedSimply(context<ScrubMachine>().m_scrbr, "ReservingReplicas")
12 {
13     DECLARE_LOCALS;
14     // 预留资源副本, 防止 OSD 在资源预留过程中启动另一个扫描操作
15     scrbr->set_reserving_now(); // a_pg_is_reserving = true;
16
17     // 通知副本预留资源
18     scrbr->reserve_replicas();
19 }
20
21 副本预约完成后, 收到RemotesReserved, 切换状态
22 sc::transition<RemotesReserved, ActiveScrubbing>
  
```

ActiveScrubbing

```

1 ActiveScrubbing::ActiveScrubbing(my_context ctx)
2     : my_base(ctx)
3     , NamedSimply(context<ScrubMachine>().m_scrbr, "ActiveScrubbing")
4 {
5     DECLARE_LOCALS; // 'scrbr' & 'pg_id' aliases
6     scrbr->on_init();
7 }
8
  
```

```

9 void PgScrubber::on_init() {
10     // 断言确保没有其他活跃的扫描操作正在进行
11     ceph_assert(!is_scrub_active());
12
13     // 重置 PG 中已扫描对象的计数
14     m_pg->reset_objects_scrubbed();
15
16     // 重置预处理数据
17     preemption_data.reset();
18
19     // 设置扫描操作的起始时间间隔
20     m_interval_start = m_pg->get_history().same_interval_since;
21
22     // 创建负责执行实际扫描操作的后端对象
23     m_be = std::make_unique<ScrubBackend>(
24         *this, // 当前 PgScrubber 实例
25         *m_pg, // PG 实例
26         m_pg_whoami, // PG 的标识符
27         m_is_repair, // 是否执行修复操作
28         m_is_deep ? scrub_level_t::deep : scrub_level_t::shallow, // 扫描级别
29         m_pg->get_actingset() // PG 的活动副本集合
30     );
31
32     // 创建一个新的存储对象
33     {
34         ObjectStore::Transaction t;
35         cleanup_store(&t); // 清理现有存储
36         m_store.reset(
37             Scrub::Store::create(m_pg->osd->store, &t, m_pg->info.pgid, m_pg-
>coll) // 创建新的存储实例
38         );
39         m_pg->osd->store->queue_transaction(m_pg->ch, std::move(t), nullptr);
40     } // 将事务加入队列
41
42     // 设置扫描操作的起始对象
43     m_start = m_pg->info.pgid.pgid.get_hobj_start();
44
45     // 将扫描操作的状态设置为活跃
46     m_active = true;
47
48     // 增加会话计数器
49     ++m_sessions_counter;
50
51     // 将扫描操作的状态和会话计数器发布到 OSD
52     m_pg->publish_stats_to_osd();
53 }

```

然后，它有一个默认子状态，决定是否立即开始Scrub

```
1 struct ActiveScrubbing
2     : sc::state<ActiveScrubbing, ScrubMachine, PendingTimer>, NamedSimply {}
3
4 PendingTimer::PendingTimer(my_context ctx)
5     : my_base(ctx)
6     , NamedSimply(context<ScrubMachine>().m_scrbr, "Act/PendingTimer")
7 {
8     scrbr->add_delayed_scheduling();
9 }
10
11 立即或推迟后调用:
12 m_osds->queue_for_scrub_resched(m_pg, Scrub::scrub_prio_t::high_priority);
13 ==> 切换到NewChunk状态
```

NewChunk

```
1 NewChunk::NewChunk(my_context ctx)
2     : my_base(ctx), NamedSimply(context<ScrubMachine>().m_scrbr,
3     "Act/NewChunk") {
4
5     DECLARE_LOCALS; // 'scrbr' & 'pg_id' aliases
6     scrbr->get_preemptor().adjust_parameters();
7     scrbr->select_range_n_notify();
8 }
9
10 void PgScrubber::select_range_n_notify() {
11 // 1 范围界定
12     if (select_range()) {
13         // 抛出SelectedChunkFree事件，转换到WaitPushes状态
14         m_osds->queue_scrub_chunk_free(m_pg,
15         Scrub::scrub_prio_t::low_priority);
16     } else {
17         m_osds->queue_scrub_chunk_busy(m_pg,
18         Scrub::scrub_prio_t::low_priority);
19     }
20 }
21
22 sc::result NewChunk::react(const SelectedChunkFree &) {
23     DECLARE_LOCALS;
```

```

22     scrbr->set_subset_last_update(scrbr->search_log_for_updates());
23     return transit<WaitPushes>();
24 }

```

WaitPushes

```

1 WaitPushes::WaitPushes(my_context ctx)
2     : my_base(ctx), NamedSimply(context<ScrubMachine>().m_scrbr,
  "Act/WaitPushes") {
3     post_event(ActivePushesUpd{});
4 }
5
6 sc::result WaitPushes::react(const ActivePushesUpd &) {
7     DECLARE_LOCALS;
8     if (!scrbr->pending_active_pushes()) {
9 // 表示没有正在进行修复，跳转到下一状态
10         return transit<WaitLastUpdate>();
11     }
12 // 如果PG正在修复，忽略这次事件，保持WaitPushes
13     return discard_event();
14 }
15

```

WaitLastUpdate

```

1 struct WaitLastUpdate : sc::state<WaitLastUpdate, ActiveScrubbing>,
2     NamedSimply {
3     ...
4     using reactions =
5         mpl::list<sc::in_state_reaction<UpdatesApplied, WaitLastUpdate,
6             &WaitLastUpdate::on_new_updates>>;
7     ...
8 };
9
10 WaitLastUpdate::WaitLastUpdate(my_context ctx)
11     : my_base(ctx), NamedSimply(context<ScrubMachine>().m_scrbr,
  "Act/WaitLastUpdate") {
12     post_event(UpdatesApplied{});
13 }
14
15 void WaitLastUpdate::on_new_updates(const UpdatesApplied &) {
16     DECLARE_LOCALS;

```

```

17
18 // 等待日志应用完毕 return last_applied >= m_subset_last_update;
19     if (scrbr->has_pg_marked_new_updates()) {
20         post_event(InternalAllUpdates{});
21     } else {
22         // will be requeued by op_applied
23     }
24 }
25
26 sc::result WaitLastUpdate::react(const InternalAllUpdates &) {
27     DECLARE_LOCALS;
28 // 通知副本准备好map
29     scrbr->get_replicas_maps(scrbr->get_preemptor().is_preemptable());
30 // 状态转换
31     return transit<BuildMap>();
32 }

```

续

```

1 void PgScrubber::get_replicas_maps(bool replica_can_preempt) {
2
3     m_primary_scrubmap_pos.reset();
4
5     for (const auto &i : m_pg->get_actingset()) {
6
7         if (i == m_pg_whoami)
8             continue;
9         m_maps_status.mark_replica_map_request(i);
10        _request_scrub_map(i,
11                            m_subset_last_update,
12                            m_start,
13                            m_end,
14                            m_is_deep,
15                            replica_can_preempt);
16    }
17 }
18
19 void PgScrubber::_request_scrub_map(pg_shard_t replica,
20                                     eversion_t version,
21                                     hobject_t start,
22                                     hobject_t end,
23                                     bool deep,
24                                     bool allow_preemption) {
25     ceph_assert(replica != m_pg_whoami);
26

```



```

27     auto repscrubop = new MOSDRepScrub(spg_t(m_pg->info.pgid.pgid,
        replica.shard),
28         version, get_osdmap_epoch(), m_pg->get_last_peering_reset(),
29         start, end, deep, allow_preemption, m_flags.priority,
30         m_pg->ops_blocked_by_scrub());
31 // 后续见从副本的流程
32     m_osds->send_message_osd_cluster(replica.osd, repscrubop,
        get_osdmap_epoch());
33 }
34

```

BuildMap

```

1 BuildMap::BuildMap(my_context ctx)
2     : my_base(ctx), NamedSimply(context<ScrubMachine>().m_scrbr,
        "Act/BuildMap") {
3     DECLARE_LOCALS;
4
5     if (scrbr->get_preemptor().was_preempted()) {
6         // we were preempted, either directly or by a replica
7         scrbr->mark_local_map_ready();
8         post_event(IntBmPreempted{});
9     } else {
10 // 构造本地的ScrubMap
11     auto ret = scrbr->build_primary_map_chunk();
12     if (ret == -EINPROGRESS) {
13         // must wait for the backend to finish. No specific event provided.
14         // build_primary_map_chunk() has already requeued us.
15     } else if (ret < 0) {
16         post_event(InternalError{});
17     } else {
18         // the local map was created
19         post_event(IntLocalMapDone{});
20     }
21 }
22 }
23
24 sc::result BuildMap::react(const IntLocalMapDone &) {
25     DECLARE_LOCALS;
26     scrbr->mark_local_map_ready();
27     return transit<WaitReplicas>();
28 }
29

```

续

```
1 int PgScrubber::build_primary_map_chunk() {
2     epoch_t map_building_since = m_pg->get_osdmap_epoch();
3     // 见文末
4     auto ret = build_scrub_map_chunk(m_be->get_primary_scrubmap(),
5         m_primary_scrubmap_pos, m_start, m_end, m_is_deep);
6
7     if (ret == -EINPROGRESS) {
8         // 重新安排另一轮请求后端收集清理数据
9         // 抛出InternalSchedScrub事件
10        m_osds->queue_for_scrub_resched(m_pg,
11            Scrub::scrub_prio_t::low_priority);
12    }
13    return ret;
14 }
15
```

从副本准备和发送ScrubMap

主OSD在WaitLastUp阶段发送了一个MOSDRepScrub信息给从OSD，用于获取对象的校验信息。

```
1 void PrimaryLogPG::do_request(
2     OpRequestRef& op,
3     ThreadPool::TPHandle &handle)
4 {
5     case MSG_OSD_REP_SCRUB:
6         replica_scrub(op, handle);
7         break;
8 }
9
10 void PG::replica_scrub(OpRequestRef op, ThreadPool::TPHandle &handle) {
11     ceph_assert(m_scrubber);
12     m_scrubber->replica_scrub_op(op);
13 }
```

```
1 void PgScrubber::replica_scrub_op(OpRequestRef op) {
2     op->mark_started();
3     auto msg = op->get_req<MOSDRepScrub>();
```

```

4
5 // 这是过时的请求
6     if (msg->map_epoch < m_pg->info.history.same_interval_since) {
7         return;
8     }
9
10 // 意思是当前的清理操作尚未完成时，就收到了新的请求，重置scrub状态
11     if (is_queued_or_active()) {
12         scrub_clear_state(); // 转为NotActive
13     }
14
15     // 确保状态机处于非活动状态
16     m_fsm->assert_not_active();
17
18     replica_scrubmap = ScrubMap{};
19     replica_scrubmap_pos = ScrubMapBuilder{};
20
21     m_replica_min_epoch = msg->min_epoch;
22     m_start = msg->start;
23     m_end = msg->end;
24     m_max_end = msg->end;
25     m_is_deep = msg->deep;
26     m_interval_start = m_pg->info.history.same_interval_since;
27     m_replica_request_priority = msg->high_priority
28                               ? Scrub::scrub_prio_t::high_priority
29                               : Scrub::scrub_prio_t::low_priority;
30     m_flags.priority = msg->priority ? msg->priority : m_pg-
>get_scrub_priority();
31
32     preemption_data.reset();
33     preemption_data.force_preemptability(msg->allow_preemption);
34
35     replica_scrubmap_pos.reset();
36
37     set_queued_or_active();
38
39 // 最终的行为是调用 send_start_replica
40     m_osds->queue_for_rep_scrub(m_pg,
41                               m_replica_request_priority,
42                               m_flags.priority,
43                               m_current_token);
44 }
45
46
47 void PgScrubber::send_start_replica(epoch_t epoch_queued,
48                                     Scrub::act_token_t token) {
49     if (is_primary()) {

```

```

50         return;
51     }
52
53 // 检查消息是否过时 return epoch_to_verify >= m_pg->get_same_interval_since();
54     if (check_interval(epoch_queued) && is_token_current(token)) {
55         if (pending_active_pushes())
56             m_fsm->process_event(StartReplica{});
57         else
58             m_fsm->process_event(StartReplicaNoWait{});
59     }
60 }
61
62
63 struct NotActive : sc::state<NotActive, ScrubMachine>, NamedSimply {
64     using reactions = mpl::list<
65         sc::transition<StartReplica, ReplicaWaitUpdates>,
66         sc::transition<StartReplicaNoWait, ActiveReplica>>;
67 };
68
69 ReplicaWaitUpdates 状态会等待 pushes操作结束, 然后转到 ActiveReplica

```

```

1 ActiveReplica::ActiveReplica(my_context ctx)
2     : my_base(ctx)
3     , NamedSimply(context<ScrubMachine>().m_scrbr, "ActiveReplica")
4 {
5     DECLARE_LOCALS;
6     scrbr->on_replica_init();
7     post_event(SchedReplica{});
8 }
9
10 sc::result ActiveReplica::react(const SchedReplica&)
11 {
12     DECLARE_LOCALS;
13 // 被抢占
14     if (scrbr->get_preemptor().was_preempted()) {
15         scrbr->send_preempted_replica();
16         scrbr->replica_handling_done();
17         return transit<NotActive>();
18     }
19 // 构造ScrubMap, 并发送给主OSD
20     auto ret_init = scrbr->build_replica_map_chunk();
21     if (ret_init != -EINPROGRESS) {
22         return transit<NotActive>();
23     }
24

```

```

25     return discard_event();
26 }
27
28
29
30 int PgScrubber::build_replica_map_chunk() {
31     ceph_assert(m_be);
32
33 // 构造ScrubMap, 同主OSD流程
34     auto ret = build_scrub_map_chunk(replica_scrubmap,
35                                     replica_scrubmap_pos,
36                                     m_start,
37                                     m_end,
38                                     m_is_deep);
39
40
41     auto required_fixes = m_be->replica_clean_meta(replica_scrubmap,
42                                                  m_end.is_max(),
43                                                  m_start,
44                                                  get_snap_mapper_accessor());
45     auto reply = prep_replica_map_msg(PreemptionNoted::no_preemption);
46     replica_handling_done();
47     // 发给主OSD
48     // 消息类型是MOSDRepScrubMap
49     send_replica_map(reply);
50
51     return ret;
52 }

```

WaitReplicas

回到主OSD, 收到从OSD的消息后会抛出GotReplicas事件

```

1 WaitReplicas::WaitReplicas(my_context ctx)
2     : my_base(ctx), NamedSimply(context<ScrubMachine>().m_scrbr,
3     "Act/WaitReplicas") {
4     post_event(GotReplicas{});
5 }
6
7 sc::result WaitReplicas::react(const GotReplicas &) {
8     DECLARE_LOCALS;
9
10    if (!all_maps_already_called && scrbr->are_all_maps_available()) {
11        all_maps_already_called = true;

```

```

12
13         // were we preempted?
14         if (scrbr->get_preemptor().disable_and_test()) {
15             return transit<PendingTimer>();
16         } else {
17 // 比较map数据
18             scrbr->maps_compare_n_cleanup();
19             return transit<WaitDigestUpdate>(); // 切状态
20         }
21     } else {
22         return discard_event();
23     }
24 }

```

WaitDigestUpdate

```

1 WaitDigestUpdate::WaitDigestUpdate(my_context ctx)
2     : my_base(ctx)
3     , NamedSimply(context<ScrubMachine>().m_scrbr, "Act/WaitDigestUpdate")
4 {
5     DECLARE_LOCALS;
6     post_event(DigestUpdate{});
7 }
8
9 sc::result WaitDigestUpdate::react(const DigestUpdate&)
10 {
11     DECLARE_LOCALS;
12     scrbr->on_digest_updates();
13     return discard_event();
14 }
15
16 void PgScrubber::on_digest_updates() {
17     if (num_digest_updates_pending > 0) {
18         return;
19     }
20
21     if (m_end.is_max()) { // 扫描完毕
22 // 抛出ScrubFinished
23         m_osds->queue_scrub_is_finished(m_pg);
24     } else {
25         // go get a new chunk (via "requeue")
26         preemption_data.reset();
27 // 抛出NewChunk, 继续扫描当前PG中剩下的对象
28         m_osds->queue_scrub_next_chunk(m_pg, m_pg->is_scrub_blocking_ops());
29     }

```

```

30 }
31
32 sc::result WaitDigestUpdate::react(const ScrubFinished&)
33 {
34     DECLARE_LOCALS;
35
36     scrbr->set_scrub_duration();
37     scrbr->scrub_finish();
38     return transit<NotActive>();
39 }
40

```

scrub_finish

这个函数主要做清理工作，重置清理计划，更新修复状态等，最重要的是如果发现错误，会启动数据修复。

```

1 void PgScrubber::scrub_finish() {
2     if (has_error) {
3         m_pg->queue_peering_event(PGPeeringEventRef(
4             std::make_shared<PGPeeringEvent>(get_osdmap_epoch(),
5                                               get_osdmap_epoch(),
6                                               PeeringState::DoRecovery())));
7     }
8 }

```

补充

更具体的细节，比如怎么构建ScrubMap，怎么根据ScrubMap做检查。

1 build_scrub_map_chunk

构建指定范围内，所有对象的校验信息，并保存在SrcubMap中

```

1
2 int PgScrubber::build_scrub_map_chunk(ScrubMap &map,
3                                       ScrubMapBuilder &pos,
4                                       hobject_t start,
5                                       hobject_t end,
6                                       bool deep) {
7     // start
8     while (pos.empty()) {
9

```

```

10         pos.deep = deep;
11         map.valid_through = m_pg->info.last_update;
12
13 // 1.1 获取指定范围内的对象
14         vector<ghobject_t> rollback_obs;
15         pos.ret = m_pg->get_pgbackend()->objects_list_range(start,
16                                                             end,
17                                                             &pos.ls,
18                                                             &rollback_obs);
19         if (pos.ret < 0) {
20             return pos.ret;
21         }
22         if (pos.ls.empty()) {
23             break;
24         }
25         m_pg->_scan_rollback_obs(rollback_obs);
26         pos.pos = 0;
27         return -EINPROGRESS;
28     }
29
30 // 1.2 扫描对象，并构建map
31     while (!pos.done()) {
32         int r = m_pg->get_pgbackend()->be_scan_list(map, pos);
33         if (r == -EINPROGRESS) {
34             return r;
35         }
36     }
37
38     // finish
39     ceph_assert(pos.done());
40     repair_oinfo_oid(map);
41
42     return 0;
43 }

```

1.1 objects_list_range

```

1 int PGBackend::objects_list_range(
2     const hobject_t &start,
3     const hobject_t &end,
4     vector<hobject_t> *ls,
5     vector<ghobject_t> *gen_obs) {
6     ceph_assert(ls);
7     vector<ghobject_t> objects;
8     int r;

```



```

9      if (HAVE_FEATURE(parent->min_upacting_features(),
10                      OSD_FIXED_COLLECTION_LIST)) {
11  // 这里是用BlueStore的函数，寻找指定范围中的对象
12      r = store->collection_list(
13          ch,
14          ghobject_t(start, ghobject_t::NO_GEN, get_parent()-
15                    >whoami_shard().shard),
16          ghobject_t(end, ghobject_t::NO_GEN, get_parent()-
17                    >whoami_shard().shard),
18          INT_MAX,
19          &objects,
20          NULL);
21  } else {
22      // 旧的接口
23      r = store->collection_list_legacy(
24          ch,
25          ghobject_t(start, ghobject_t::NO_GEN, get_parent()-
26                    >whoami_shard().shard),
27          ghobject_t(end, ghobject_t::NO_GEN, get_parent()-
28                    >whoami_shard().shard),
29          INT_MAX,
30          &objects,
31          NULL);
32  }
33  ls->reserve(objects.size());
34  for (vector<ghobject_t>::iterator i = objects.begin();
35       i != objects.end();
36       ++i) {
37      // 如果是临时对象，跳过
38      if (i->is_pgmeta() || i->hobj.is_temp()) {
39          continue;
40      }
41      // 没有生成号的加入ls，有生成号的加入gen，这个gen_obs在调用方是回滚对象
42      if (i->is_no_gen()) {
43          ls->push_back(i->hobj);
44      } else if (gen_obs) {
45          gen_obs->push_back(*i);
46      }
47  }
48  return r;
49 }

```

```

1 void PG::_scan_rollback_obs(const vector<ghobject_t> &rollback_obs) {

```

```

2     ObjectStore::Transaction t;
3     eversion_t trimmed_to =
    recovery_state.get_last_rollback_info_trimmed_to_applied();
4     for (vector<ghobject_t>::const_iterator i = rollback_obs.begin();
5         i != rollback_obs.end();
6         ++i) {
7         if (i->generation < trimmed_to.version) {
8             t.remove(coll, *i);
9         }
10    }
11    if (!t.empty()) {
12        osd->store->queue_transaction(ch, std::move(t), NULL);
13    }
14 }

```

1.2 be_scan_list

```

1  int PGBackend::be_scan_list(
2      ScrubMap &map,
3      ScrubMapBuilder &pos) {
4      ceph_assert(!pos.done());
5      ceph_assert(pos.pos < pos.ls.size());
6      hobject_t &poid = pos.ls[pos.pos];
7
8      // 获取对象的状态信息
9      struct stat st;
10     int r = store->stat(ch,
11         ghobject_t(poid, ghobject_t::NO_GEN, get_parent()-
12 >whoami_shard().shard),
13         &st, true);
14
15     if (r == 0) {
16         ScrubMap::object &o = map.objects[poid];
17         o.size = st.st_size;
18         ceph_assert(!o.negative);
19         // 获取对象的属性，用于浅扫描
20         store->getattrs(ch,
21             ghobject_t(poid, ghobject_t::NO_GEN, get_parent()-
22 >whoami_shard().shard),
23             o.attrs);
24         // 如果需要深扫描
25         if (pos.deep) {
26             r = be_deep_scrub(poid, map, pos, o);
27         }
28     }
29 }

```

```
27     pos.next_object();
28     return 0;
29 }
```

这个函数一共四步：读数据，读header，读omap，计算digest（校验信息）

```
1  int ReplicatedBackend::be_deep_scrub(
2      const hobject_t &poid,
3      ScrubMap &map,
4      ScrubMapBuilder &pos,
5      ScrubMap::object &o) {
6      int r;
7      uint32_t fadvise_flags = CEPH_OSD_OP_FLAG_FADVISE_SEQUENTIAL |
8                               CEPH_OSD_OP_FLAG_FADVISE_DONTNEED |
9                               CEPH_OSD_OP_FLAG_BYPASS_CLEAN_CACHE;
10
11      utime_t sleeptime;
12      sleeptime.set_from_double(cct->_conf->osd_debug_deep_scrub_sleep);
13      if (sleeptime != utime_t()) {
14          sleeptime.sleep();
15      }
16
17      ceph_assert(poid == pos.ls[pos.pos]);
18      if (!pos.data_done()) {
19          if (pos.data_pos == 0) {
20              pos.data_hash = bufferhash(-1);
21          }
22
23          const uint64_t stride = cct->_conf->osd_deep_scrub_stride;
24
25          bufferlist bl;
26          // 1 一次读取stride长度的数据，如果读取除了完整的stride，说明可能还需要读取的块
27          // 返回-EINPROGRESS，同时上层会再次调用到这里，循环直到读取完
28          r = store->read(
29              ch,
30              ghobject_t(
31                  poid, ghobject_t::NO_GEN, get_parent()->whoami_shard().shard),
32                  pos.data_pos,
33                  stride, bl,
34                  fadvise_flags);
35          if (r < 0) {
36              o.read_error = true;
37              return 0;
38          }
39          if (r > 0) {
```

```

40         pos.data_hash << bl;
41     }
42     pos.data_pos += r;
43     if (static_cast<uint64_t>(r) == stride) {
44         return -EINPROGRESS;
45     }
46     // done with bytes
47     pos.data_pos = -1;
48     o.digest = pos.data_hash.digest();
49     o.digest_present = true;
50 }
51
52 // 2 omap header
53 if (pos.omap_pos.empty()) {
54     pos.omap_hash = bufferhash(-1);
55
56     bufferlist hdrbl;
57     r = store->omap_get_header(
58         ch,
59         ghobject_t(
60             poid, ghobject_t::NO_GEN, get_parent()->whoami_shard().shard),
61         &hdrbl, true);
62     if (r == -EIO) {
63         o.read_error = true;
64         return 0;
65     }
66     if (r == 0 && hdrbl.length()) {
67         bool encoded = false;
68         pos.omap_hash << hdrbl;
69     }
70 }
71
72 // 3 omap
73 ObjectMap::ObjectMapIterator iter = store->get_omap_iterator(
74     ch,
75     ghobject_t(
76         poid, ghobject_t::NO_GEN, get_parent()->whoami_shard().shard));
77 ceph_assert(iter);
78 if (pos.omap_pos.length()) {
79     iter->lower_bound(pos.omap_pos);
80 } else {
81     iter->seek_to_first();
82 }
83 int max = g_conf()->osd_deep_scrub_keys;
84 while (iter->status() == 0 && iter->valid()) {
85     pos.omap_bytes += iter->value().length();
86     ++pos.omap_keys;

```

```

87         --max;
88         // fixme: we can do this more efficiently.
89         bufferlist bl;
90         encode(iter->key(), bl);
91         encode(iter->value(), bl);
92         pos.omap_hash << bl;
93
94         iter->next();
95
96         if (iter->valid() && max == 0) {
97             pos.omap_pos = iter->key();
98             return -EINPROGRESS;
99         }
100        if (iter->status() < 0) {
101
102            o.read_error = true;
103            return 0;
104        }
105    }
106
107    if (pos.omap_keys > cct->_conf-
>osd_deep_scrub_large_omap_object_key_threshold ||
108        pos.omap_bytes > cct->_conf-
>osd_deep_scrub_large_omap_object_value_sum_threshold) {
109
110        o.large_omap_object_found = true;
111        o.large_omap_object_key_count = pos.omap_keys;
112        o.large_omap_object_value_size = pos.omap_bytes;
113        map.has_large_omap_object_errors = true;
114    }
115    // 4 计算digest
116    o.omap_digest = pos.omap_hash.digest();
117    o.omap_digest_present = true;
118
119    // Sum up omap usage
120    if (pos.omap_keys > 0 || pos.omap_bytes > 0) {
121        map.has_omap_keys = true;
122        o.object_omap_bytes = pos.omap_bytes;
123        o.object_omap_keys = pos.omap_keys;
124    }
125
126    // done!
127    return 0;
128 }
129

```

2 maps_compare_n_cleanup

在上文中，等到所有的副本都计算好了digest，就调用这个函数进行比较

```
1 void PgScrubber::maps_compare_n_cleanup() {
2     // 更新已扫描对象的数量
3     m_pg->add_objects_scrubbed_count(m_be-
4     >get_primary_scrubmap().objects.size());
5
6     // 比较map并获取需要修复的对象
7     auto required_fixes = m_be->scrub_compare_maps(m_end.is_max(),
8     get_snap_mapper_accessor());
9
10    if (!required_fixes.inconsistent_objs.empty()) { // 有不一致的对象
11        if (state_test(PG_STATE_REPAIR)) {
12            // 如果PG正在修复中，丢弃
13        } else {
14            // 如果不是修复状态，执行修复操作，处理不一致的对象
15            persist_scrub_results(std::move(required_fixes.inconsistent_objs));
16        }
17    }
18
19    // 应用快照映射器的修复
20    apply_snap_mapper_fixes(required_fixes.snap_fix_list);
21
22    auto chunk_err_counts = m_be->get_error_counts();
23    m_shallow_errors += chunk_err_counts.shallow_errors;
24    m_deep_errors += chunk_err_counts.deep_errors;
25
26    // 更新起始位置，为下一轮清理做准备
27    m_start = m_end;
28
29    // 运行回调，可能用于通知其他组件清理完成
30    run_callbacks();
31
32    requeue_waiting();
33 }
```

```
1 objs_fix_list_t ScrubBackend::scrub_compare_maps(
2     bool max_reached,
3     SnapMapReaderI &snaps_getter) {
4
5     ceph_assert(m_scrubber.is_primary());
6 }
```

```

7 // 1 构建权威map
8 // insert(this_chunk->received_maps[m_pg_whoami])
9     m_cleaned_meta_map.insert(my_map());
10 // 1.1 合并所有的map
11 // 将this_chunk->received_maps中的对象ID全添加到authoritative_set
12     merge_to_authoritative_set();
13 // 1.2 收集统计信息
14     omap_checks();
15 // 1.3 更新权威map，通过各个map间的共同信息，计算出真正的权威map
16     update_authoritative();
17
18 // 5 返回一个包含不一致对象的列表
19     auto for_meta_scrub = clean_meta_map(m_cleaned_meta_map, max_reached);
20
21 // 6 检查快照元数据
22     scrub_snapshot_metadata(for_meta_scrub);
23 // 7 返回一个包含不一致对象列表，和快照扫描结果的对象
24     return objs_fix_list_t{std::move(this_chunk->m_inconsistent_objs),
25                             scan_snaps(for_meta_scrub, snaps_getter)};
26 }

```

计算权威对象的过程很乱，但是遵循以下原则：

1. 没有状态错误，静默数据错误，OI属性正确，对象大小正确等
2. 如果各项指标正确，优先选择校验信息更多的副本
3. 如果存在满足上述条件的多个副本，优先选择主OSD对应的副本，否则随机选择