

Ceph Monitor

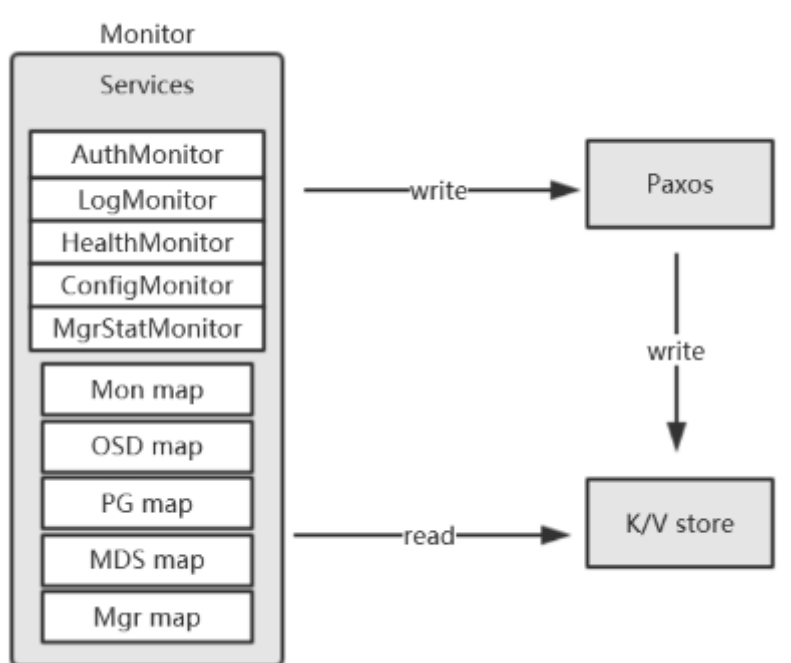
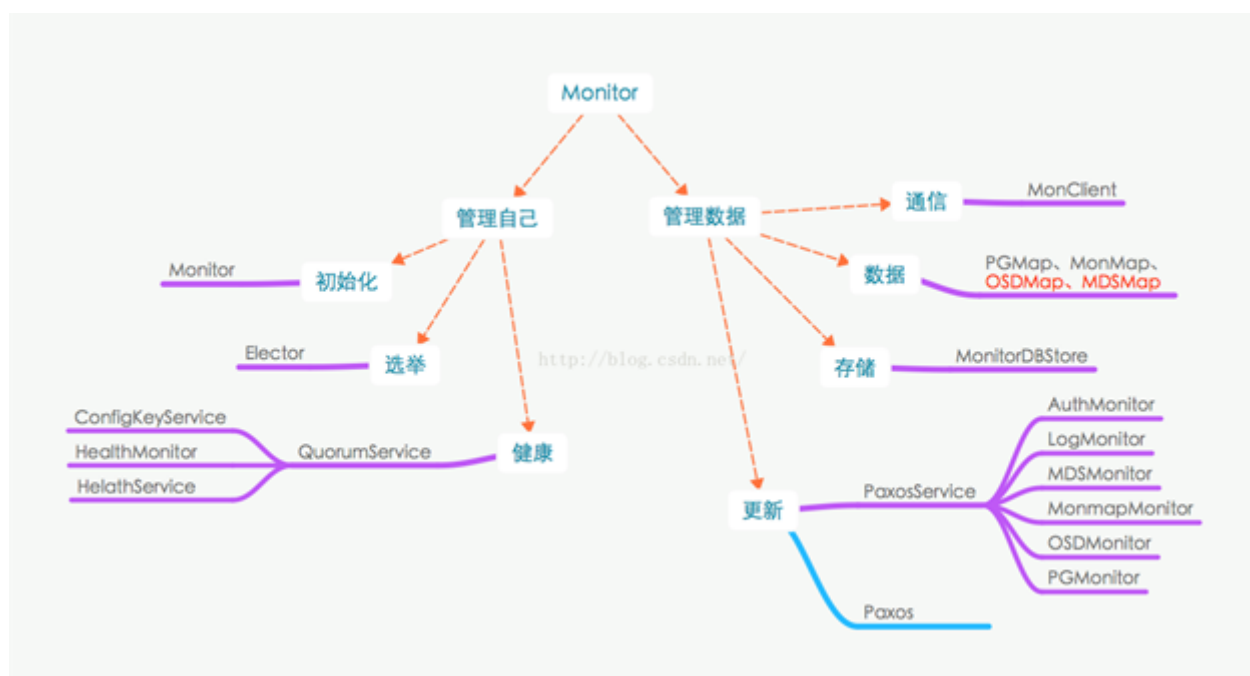
ceph版本: <https://github.com/ceph/ceph/tree/v18.2.1>

关于paxos的基本原理, 本文不再介绍, 网上有各路大神的分享

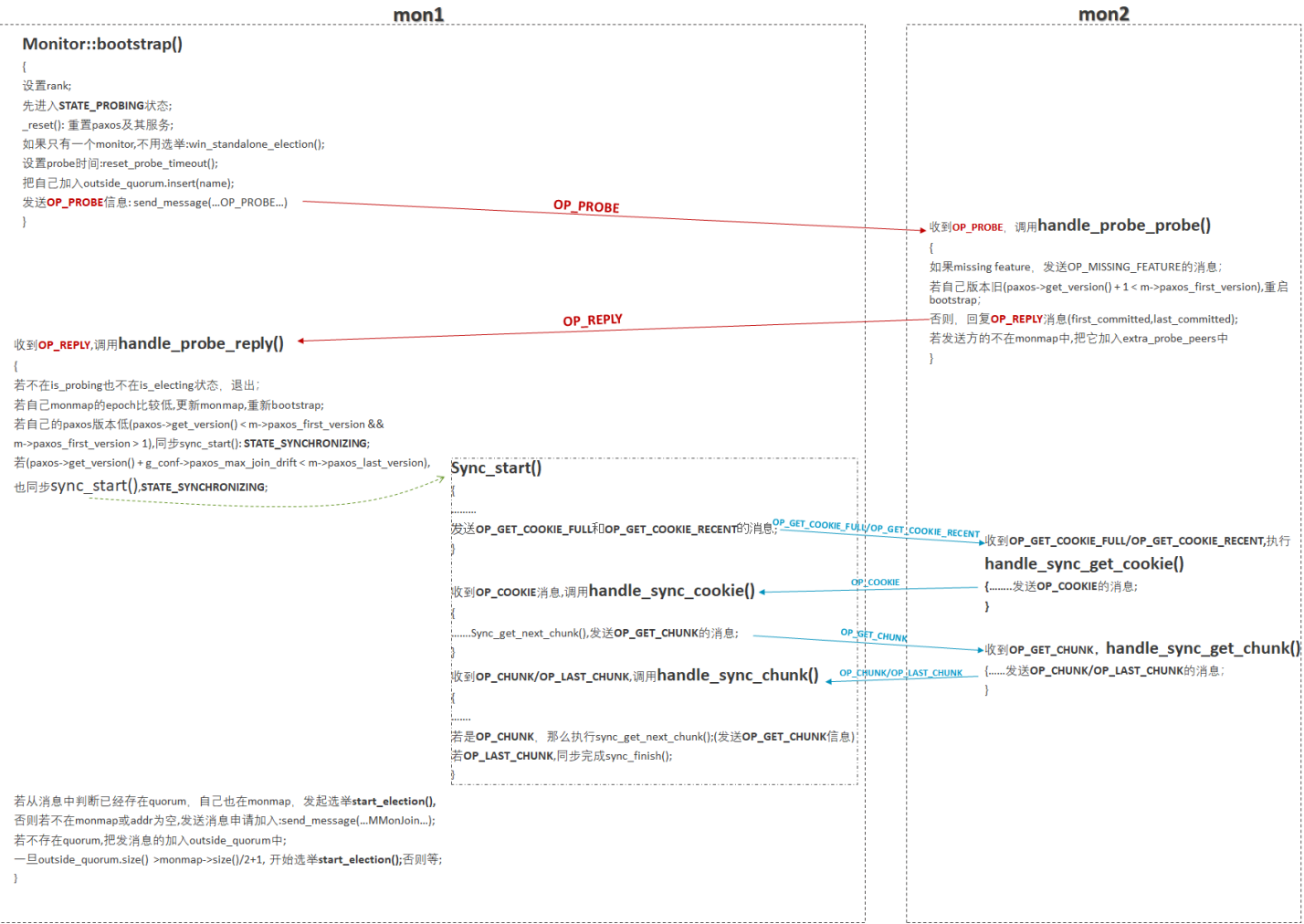
[官方介绍](#)

[可靠分布式系统-paxos的直观解释](#)

Mon



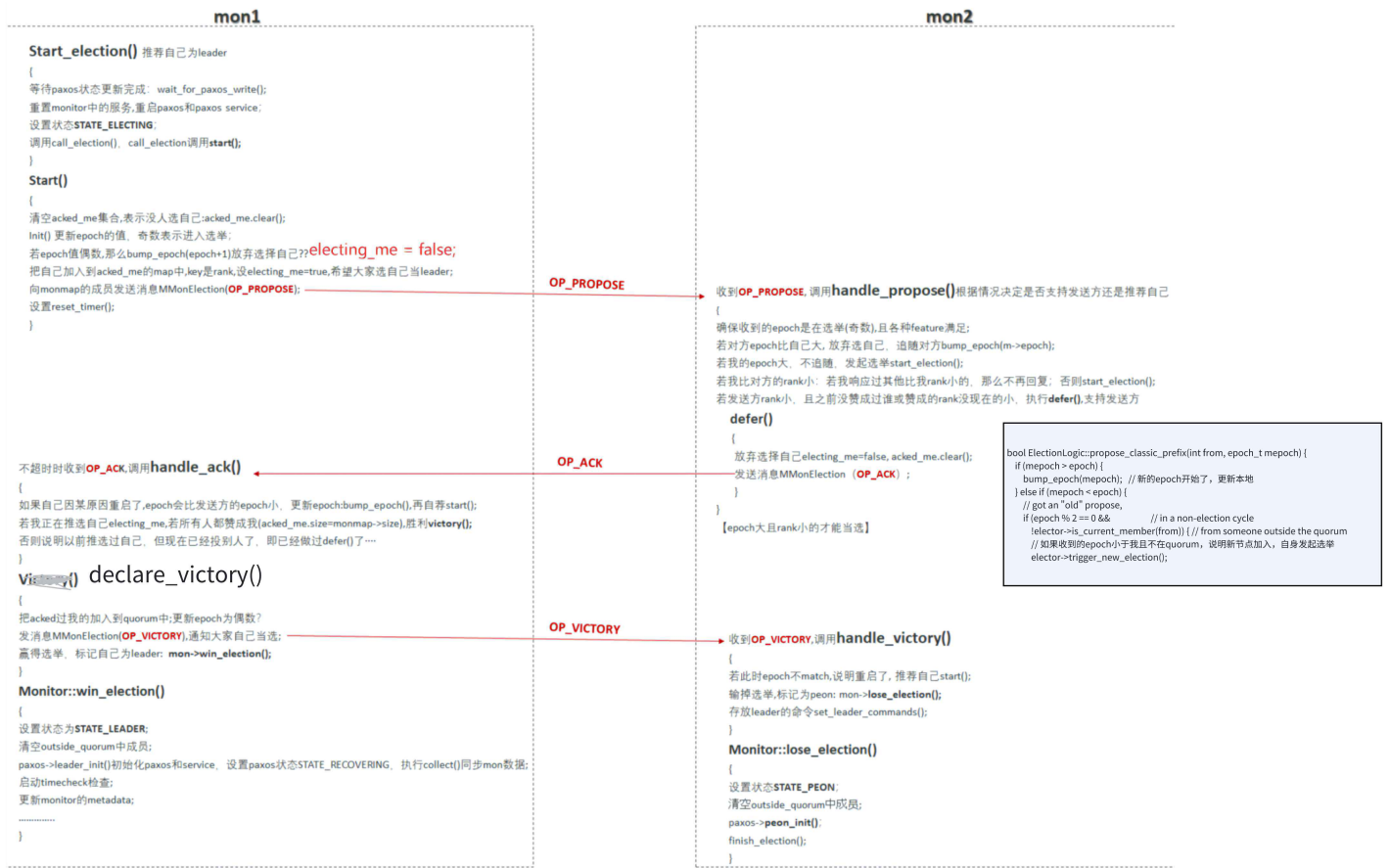
启动/同步： bootStrap



选举：election

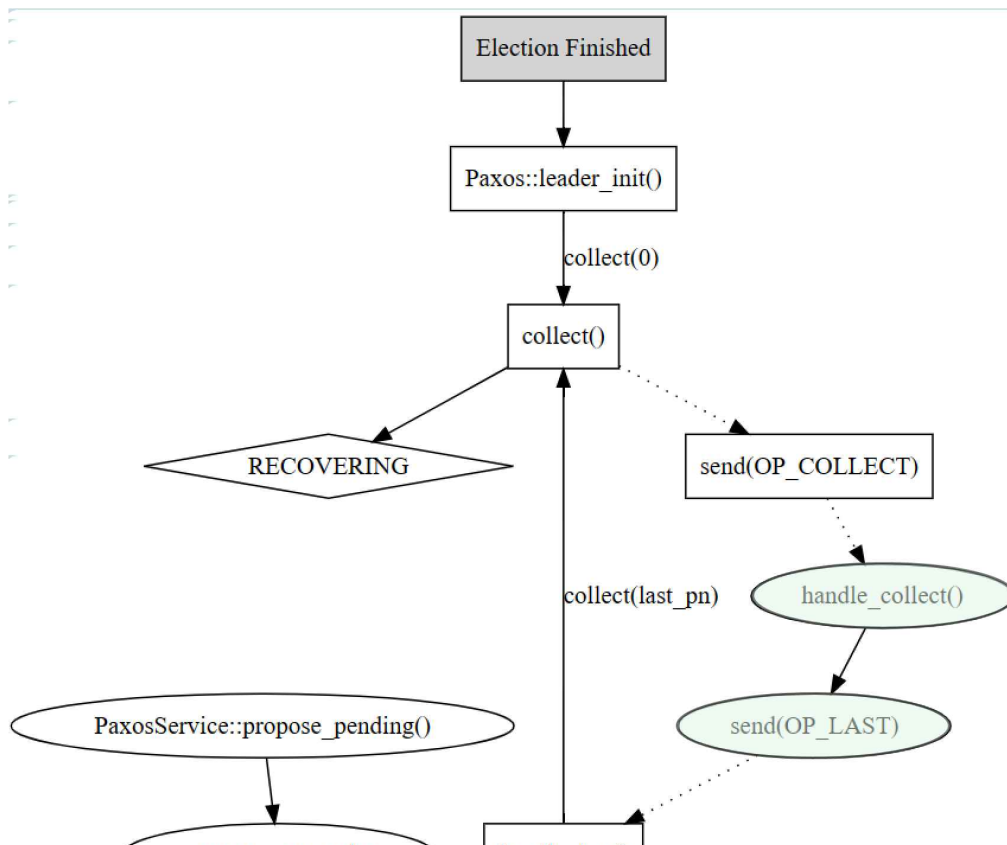
下图中的{若epoch是偶数,则...},实际上在init中, epoch会设置为奇数,也就是说 electing_me=false这段分支根本走不到(当时看的时候很矛盾,先设置为了false,后面又设置为 true,最终发现根本没有false的前提)

关于epoch的奇偶性对paxos的影响并不只是版本号,还表示一种状态,如果是奇数,表示当前正在进行选举;否则当前表示当前是除了选举之外的状态。而这也是针对ceph实现,更宽泛的说法, epoch为奇数表示当前正在商议提案,为偶数则表示已经商议完毕。对当前情况下的选举来说,就是指商议阶段。



一致性决策阶段

白色部分是leader的行为, 绿色部分是peon的行为





需要持久化的数据

名称	含义	其他
PN(Proposal Number)	提案编号，用于标识提案，维护提案顺序	Leader当选之后，会执行一次确定PN，在其为Leader的整个Phase 2共用一个PN
last_pn	上次当选leader后生成的PN	get_new_proposal_number选后，接着生成
accepted_pn	我接受过的PN，可能是别的leader提议的PN	peer根据这个值拒绝较小的
first_committed	本节点记录的第一个被commit的版本	更早的版本(日志)，本节点记录
last_committed	本节点记录的最后一次被commit的版本	往后的版本，未被commit，
uncommitted_v	本节点记录的未commit的版本，如果有，只能等于last_commit+1	ceph只允许有一个未commit
uncommitted_pn	未commit的版本对应的PN	与uncommitted_v，uncor一个事务中记录
uncommitted_value	未commit的版本的內容	与uncommitted_v，uncor一个事务中记录

leader_init

```
1 void Paxos::leader_init()
2 {
3     cancel_events();
4     // 清除之前提出的提议值
5     new_value.clear();
6
7     // 重置正在等待决议的提议
8     pending_proposal.reset();
9     // 重置commit的回调函数
10    reset_pending_committing_finishers();
11
12    // 根据法定人数设置状态为ACTIVE或RECOVERING
13    if (mon.get_quorum().size() == 1) {
14        state = STATE_ACTIVE;
15        return;
16    }
17
18    // 进入recovery状态
19    state = STATE_RECOVERING;
20    lease_expire = {};
21    // 收集其它节点信息，进行提议值选择和一致性决议
```

```
22     collect(0);
23 }
```

collect

```
1 void Paxos::collect(version_t oldpn) {
2
3     state = STATE_RECOVERING;
4     ceph_assert(mon.is_leader());
5
6     // 重置未提交提议的记录
7     uncommitted_v = 0;
8     uncommitted_pn = 0;
9     uncommitted_value.clear();
10    peer_first_committed.clear();
11    peer_last_committed.clear();
12
13    // 检查本地存储,获取未提交的提议值
14    if (get_store()->exists(get_name(), last_committed + 1)) {
15        version_t v = get_store()->get(get_name(), "pending_v");
16        version_t pn = get_store()->get(get_name(), "pending_pn");
17        if (v && pn && v == last_committed + 1) {
18            uncommitted_pn = pn;
19        } else {
20            uncommitted_pn = accepted_pn;
21        }
22        uncommitted_v = last_committed + 1;
23
24        get_store()->get(get_name(), last_committed + 1, uncommitted_value);
25        ceph_assert(uncommitted_value.length());
26    }
27
28    // pick new pn
29    accepted_pn = get_new_proposal_number(std::max(accepted_pn, oldpn));
30    accepted_pn_from = last_committed;
31    num_last = 1;
32
33    // send collect
34    for (auto p = mon.get_quorum().begin();
35         p != mon.get_quorum().end();
36         ++p) {
37        if (*p == mon.rank)
38            continue;
39    }
```

```

40         MMonPaxos *collect = new MMonPaxos(mon.get_epoch(),
        MMonPaxos::OP_COLLECT,
41                                     ceph_clock_now());
42         collect->last_committed = last_committed;
43         collect->first_committed = first_committed;
44         collect->pn = accepted_pn;
45         mon.send_mon_message(collect, *p);
46     }
47
48 // set timeout event
49     collect_timeout_event = mon.timer.add_event_after(
50         g_conf()->mon_accept_timeout_factor *
51         g_conf()->mon_lease,
52         new C_MonContext{&mon, [this](int r) {
53             if (r == -ECANCELED)
54                 return;
55             collect_timeout();
56             }});
57 }
58

```

每次collect, 都会生成新的accept_pn, get_new_proposal_number:

```

1 version_t Paxos::get_new_proposal_number(version_t gt) {
2     if (last_pn < gt)
3         last_pn = gt;
4
5     // update. make it unique among all monitors.
6     last_pn /= 100; // 将最后两位数字清零
7     last_pn++;      // 让前几位数字增1
8     last_pn *= 100; // 重新恢复最后两位为00
9     last_pn += (version_t)mon.rank; // 将rank作为最后两位
10    // 利用低两位来区分不同的mon, 这样可以保证mon生成的pn唯一
11    ...
12    return last_pn;
13 }

```

Peon: handle_collect

- 1 主要功能包括:
- 2 设置状态为RECOVERING表示正在进行恢复
- 3 更新副本节点lease过期时间
- 4 如果Leader提交的版本号太高则需要引导恢复

- 5 回复Leader本地最后提交和首提交信息
- 6 判断是否可以接受Leader的提议号,如果可以更新提议号状态
- 7 将本节点已提交但Leader未知的信息同步给Leader
- 8 如果本地有已提交但未确认的值,也发送给Leader
- 9 这段代码结合Paxos算法流程很清晰地实现了副本节点收到Leader **COLLECT**请求后的处理逻辑:
- 10 检查自己是否需要进行引导恢复
- 11 回传自己的状态信息
- 12 判断是否需要采纳Leader新的提议号
- 13 与Leader同步状态,将自己更高版本的信息传递给Leader

```
1 void Paxos::handle_collect(MonOpRequestRef op) {
2     op->mark_paxos_event("handle_collect");
3     auto collect = op->get_req<MMonPaxos>();
4     ceph_assert(mon.is_peon());
5     state = STATE_RECOVERING;
6
7     // 更新副本节点lease过期时间
8     reset_lease_timeout();
9     // 版本过于落后, 需要引导修复
10    if (collect->first_committed > last_committed + 1) {
11        op->mark_paxos_event("need to bootstrap");
12        mon.bootstrap();
13        return;
14    }
15
16    // reply
17    MMonPaxos *last = new MMonPaxos(mon.get_epoch(),
18                                    MMonPaxos::OP_LAST, ceph_clock_now());
19    last->last_committed = last_committed;
20    last->first_committed = first_committed;
21
22    version_t previous_pn = accepted_pn;
23
24    // can we accept this pn?
25    if (collect->pn > accepted_pn) {
26        // ok, accept it
27        accepted_pn = collect->pn;
28        accepted_pn_from = collect->pn_from;
29        auto t(std::make_shared<MonitorDBStore::Transaction>());
30        t->put(get_name(), "accepted_pn", accepted_pn);
31
32        JSONFormatter f(true);
33        t->dump(&f);
34        f.flush(*_dout);
35        *_dout << endl;
```



```

36 // 将接受的提案值通过事务写到存储
37     auto start = ceph::coarse_mono_clock::now();
38     get_store()->apply_transaction(t);
39     auto end = ceph::coarse_mono_clock::now();
40
41     } else {
42     }
43     last->pn = accepted_pn;
44     last->pn_from = accepted_pn_from;
45
46 // 本地已提交，但leader未知的信息，将leader缺少的记录给添加到last中
47     if (collect->last_committed < last_committed)
48         share_state(last, collect->first_committed, collect->last_committed);
49
50     bufferlist bl;
51     if (collect->last_committed <= last_committed &&
52         get_store()->exists(get_name(), last_committed + 1)) {
53         get_store()->get(get_name(), last_committed + 1, bl);
54         ceph_assert(bl.length() > 0);
55
56         last->values[last_committed + 1] = bl;
57
58         version_t v = get_store()->get(get_name(), "pending_v");
59         version_t pn = get_store()->get(get_name(), "pending_pn");
60         if (v && pn && v == last_committed + 1) {
61             last->uncommitted_pn = pn;
62         } else {
63             last->uncommitted_pn = previous_pn;
64         }
65
66         logger->inc(l_paxos_collect_uncommitted);
67     }
68
69 // send reply
70     collect->get_connection()->send_message(last);
71 }
72

```

handle_last

处理逻辑

- 1 更新收到回复节点的first_committed和last_committed状态

- 2 判断状态是否一致,不一致触发引导恢复
- 3 保存LAST消息中的已确认值
- 4 判断其他副本是否都接收当前提议号,是则进入ACTIVE状态
- 5 处理可能学习到的未确认值
- 6 如果还没有全体副本响应,继续等待
- 7 否则完成本轮提议,切换状态

```
1 void Paxos::handle_last(MonOpRequestRef op) {
2     op->mark_paxos_event("handle_last");
3     auto last = op->get_req<MMonPaxos>();
4     bool need_refresh = false;
5     int from = last->get_source().num();
6
7     if (!mon.is_leader()) {
8         return;
9     }
10
11
12     peer_first_committed[from] = last->first_committed;
13     peer_last_committed[from] = last->last_committed;
14
15 // 引导修复
16     if (last->first_committed > last_committed + 1) {
17         op->mark_paxos_event("need to bootstrap");
18         mon.bootstrap();
19         return;
20     }
21
22     ceph_assert(g_conf()->paxos_kill_at != 1);
23
24     // store any committed values if any are specified in the message
25     need_refresh = store_state(last);
26
27     ceph_assert(g_conf()->paxos_kill_at != 2);
28
29 // 更新每一个peer的信息
30     // is everyone contiguous and up to date?
31     for (auto p = peer_last_committed.begin();
32          p != peer_last_committed.end();
33          ++p) {
34         if (p->second + 1 < first_committed && first_committed > 1) {
35             op->mark_paxos_event("need to bootstrap");
36             mon.bootstrap();
37             return;
38         }
39     }
```

```

39         if (p->second < last_committed) {
40             // share committed values
41             MMonPaxos *commit = new MMonPaxos(mon.get_epoch(),
42                                                MMonPaxos::OP_COMMIT,
43                                                ceph_clock_now());
44             share_state(commit, peer_first_committed[p->first], p->second);
45             mon.send_mon_message(commit, p->first);
46         }
47     }
48
49 // 说明leader的pn落后, 开始新的提案
50     if (last->pn > accepted_pn) {
51         // cancel timeout event
52         mon.timer.cancel_event(collect_timeout_event);
53         collect_timeout_event = 0;
54
55         collect(last->pn);
56     } else if (last->pn == accepted_pn) {
57 // 增加已接受提案的计数
58         num_last++;
59         // did this person send back an accepted but uncommitted value?
60         if (last->uncommitted_pn) {
61             if (last->uncommitted_pn >= uncommitted_pn &&
62                 last->last_committed >= last_committed &&
63                 last->last_committed + 1 >= uncommitted_v) {
64                 uncommitted_v = last->last_committed + 1;
65                 uncommitted_pn = last->uncommitted_pn;
66                 uncommitted_value = last->values[uncommitted_v];
67             }
68         }
69
70 // 接收提案的节点达到了指定人数
71         if (num_last == mon.get_quorum().size()) {
72 // 取消超时事件, 并清理peer的状态信息
73             // cancel timeout event
74             mon.timer.cancel_event(collect_timeout_event);
75             collect_timeout_event = 0;
76             peer_first_committed.clear();
77             peer_last_committed.clear();
78
79             // almost...
80
81 // 如果当前节点收到了一个有效的未提交值, 设置状态为UPDATING_PREVIOUS
82             if (uncommitted_v == last_committed + 1 &&
83                 uncommitted_value.length()) {
84                 state = STATE_UPDATING_PREVIOUS;
85                 begin(uncommitted_value);

```

```

86             } else {
87 // 否则如果所有的节点都接受了提案编号且没有新的未提交值，延长租约
88             // active!
89             extend_lease();
90
91             need_refresh = false;
92             if (do_refresh()) {
93                 finish_round();
94             }
95         }
96     }
97 } else {
98     // no, this is an old message, discard
99 }
100
101 if (need_refresh)
102     (void)do_refresh();
103 }
104

```

propose_pending

```

1 void Paxos::propose_pending() {
2     ceph_assert(is_active());
3     ceph_assert(pending_proposal);
4
5     cancel_events();
6
7     bufferlist bl;
8     pending_proposal->encode(bl);
9
10    pending_proposal.reset();
11
12    committing_finishers.swap(pending_finishers);
13    state = STATE_UPDATING;
14    begin(bl);
15 }

```

begin

```

1 void Paxos::begin(bufferlist &v) {
2     ceph_assert(mon.is_leader());

```

```

3     ceph_assert(is_updating() || is_updating_previous());
4
5     // we must already have a majority for this to work.
6     ceph_assert(mon.get_quorum().size() == 1 ||
7                 num_last > (unsigned)mon.monmap->size() / 2);
8
9     // and no value, yet.
10    ceph_assert(new_value.length() == 0);
11
12    // accept it ourselves, 并设置新的提案
13    accepted.clear();
14    accepted.insert(mon.rank);
15    new_value = v;
16
17    // 第一个commit, 只有第一次提出提案的时候才会遇到
18    if (last_committed == 0) {
19        auto t(std::make_shared<MonitorDBStore::Transaction>());
20        // initial base case; set first_committed too
21        t->put(get_name(), "first_committed", 1);
22        decode_append_transaction(t, new_value);
23
24        bufferlist tx_bl;
25        t->encode(tx_bl);
26
27        new_value = tx_bl;
28    }
29
30    // 将编码后的值添加到存储事务中
31    auto t(std::make_shared<MonitorDBStore::Transaction>());
32    t->put(get_name(), last_committed + 1, new_value);
33
34    // note which pn this pending value is for.
35    t->put(get_name(), "pending_v", last_committed + 1);
36    t->put(get_name(), "pending_pn", accepted_pn);
37
38    // 将要发起的提案通过事务保存在本地
39    auto start = ceph::coarse_mono_clock::now();
40    get_store()->apply_transaction(t);
41    auto end = ceph::coarse_mono_clock::now();
42
43    logger->tinc(l_paxos_begin_latency, to_timespan(end - start));
44
45    ceph_assert(g_conf()->paxos_kill_at != 3);
46    // 集群只有一个mon, 就直接commit
47    if (mon.get_quorum().size() == 1) {
48        // we're alone, take it easy
49        commit_start();

```

```

50         return;
51     }
52
53 // 通知其它节点，处理这个提案
54 // ask others to accept it too!
55 for (auto p = mon.get_quorum().begin();
56      p != mon.get_quorum().end();
57      ++p) {
58     if (*p == mon.rank)
59         continue;
60
61     MMonPaxos *begin = new MMonPaxos(mon.get_epoch(), MMonPaxos::OP_BEGIN,
62                                     ceph_clock_now());
63     begin->values[last_committed + 1] = new_value;
64     begin->last_committed = last_committed;
65     begin->pn = accepted_pn;
66
67     mon.send_mon_message(begin, *p);
68 }
69
70 // set timeout event
71 accept_timeout_event = mon.timer.add_event_after(
72     g_conf()->mon_accept_timeout_factor * g_conf()->mon_lease,
73     new C_MonContext{&mon, [this](int r) {
74         if (r == -ECANCELED)
75             return;
76         accept_timeout();
77     }});
78 }
79

```

Peon: handle_begin

```

1 void Paxos::handle_begin(MonOpRequestRef op) {
2     op->mark_paxos_event("handle_begin");
3     auto begin = op->get_req<MMonPaxos>();
4
5 // can we accept this?
6     if (begin->pn < accepted_pn) {
7         op->mark_paxos_event("have higher pn, ignore");
8         return;
9     }
10    ceph_assert(begin->pn == accepted_pn);
11    ceph_assert(begin->last_committed == last_committed);
12

```

```

13     ceph_assert(g_conf()->paxos_kill_at != 4);
14
15     logger->inc(l_paxos_begin);
16
17 // set state.
18     state = STATE_UPDATING;
19     lease_expire = {}; // cancel lease
20
21 // yes.
22     version_t v = last_committed + 1;
23
24 // 这里和leader一样，先将数据放在事务中，然后提交事务
25     auto t(std::make_shared<MonitorDBStore::Transaction>());
26     t->put(get_name(), v, begin->values[v]);
27
28     // note which pn this pending value is for.
29     t->put(get_name(), "pending_v", v);
30     t->put(get_name(), "pending_pn", accepted_pn);
31
32     auto start = ceph::coarse_mono_clock::now();
33     get_store()->apply_transaction(t);
34     auto end = ceph::coarse_mono_clock::now();
35
36     logger->tinc(l_paxos_begin_latency, to_timespan(end - start));
37
38     ceph_assert(g_conf()->paxos_kill_at != 5);
39
40 // 回复接收提案的消息
41     MMonPaxos *accept = new MMonPaxos(mon.get_epoch(), MMonPaxos::OP_ACCEPT,
42                                     ceph_clock_now());
43     accept->pn = accepted_pn;
44     accept->last_committed = last_committed;
45     begin->get_connection()->send_message(accept);
46 }

```

hanle_accept

```

1 void Paxos::handle_accept(MonOpRequestRef op) {
2     op->mark_paxos_event("handle_accept");
3     auto accept = op->get_req<MMonPaxos>();
4     int from = accept->get_source().num();
5
6     if (accept->pn != accepted_pn) {

```

```

7         op->mark_paxos_event("have higher pn, ignore");
8         return;
9     }
10    if (last_committed > 0 &&
11        accept->last_committed < last_committed - 1) {
12        op->mark_paxos_event("old round, ignore");
13        return;
14    }
15    ceph_assert(accept->last_committed == last_committed ||    // not committed
16               accept->last_committed == last_committed - 1); // committed
17
18    ceph_assert(is Updating() || is Updating_previous());
19    ceph_assert(accepted.count(from) == 0);
20    accepted.insert(from);
21
22    ceph_assert(g_conf()->paxos_kill_at != 6);
23
24    // 只有通过数达到了quorum, 才能提交
25    if (accepted == mon.get_quorum()) {
26        op->mark_paxos_event("commit_start");
27        commit_start();
28    }
29 }

```

第二阶段: commit_start

```

1 void Paxos::commit_start() {
2
3     ceph_assert(g_conf()->paxos_kill_at != 7);
4
5     auto t(std::make_shared<MonitorDBStore::Transaction>());
6
7     // commit locally
8     t->put(get_name(), "last_committed", last_committed + 1);
9
10    // decode the value and apply its transaction to the store.
11    // this value can now be read from last_committed.
12    decode_append_transaction(t, new_value);
13
14    commit_start_stamp = ceph_clock_now();
15
16    // 提交事务, 与之前不同的是, 这里是异步调用
17    // 设置的回调函数是commit_finish
18    get_store()->queue_transaction(t, new C_Committed(this));
19

```



```

20     if (is_updating_previous())
21         state = STATE_WRITING_PREVIOUS;
22     else if (is_updating())
23         state = STATE_WRITING;
24     else
25         ceph_abort();
26     ++commits_started;
27
28     if (mon.get_quorum().size() > 1) {
29         // cancel timeout event
30         mon.timer.cancel_event(accept_timeout_event);
31         accept_timeout_event = 0;
32     }
33 }

```

commit_finish

```

1 void Paxos::commit_finish() {
2     utime_t end = ceph_clock_now();
3     logger->tinc(l_paxos_commit_latency, end - commit_start_stamp);
4
5     ceph_assert(g_conf()->paxos_kill_at != 8);
6
7     // 取消当前租约
8     lease_expire = {}; // cancel lease
9
10    last_committed++;
11    last_commit_time = ceph_clock_now();
12
13    // refresh first_committed; this txn may have trimmed.
14    first_committed = get_store()->get(get_name(), "first_committed");
15
16    _sanity_check_store();
17
18    // 通知所有人commit
19    for (auto p = mon.get_quorum().begin();
20         p != mon.get_quorum().end();
21         ++p) {
22        if (*p == mon.rank)
23            continue;
24
25
26        MMonPaxos *commit = new MMonPaxos(mon.get_epoch(),
27        MMonPaxos::OP_COMMIT,
28
29        ceph_clock_now());

```

```

28     commit->values[last_committed] = new_value;
29     commit->pn = accepted_pn;
30     commit->last_committed = last_committed;
31
32     mon.send_mon_message(commit, *p);
33 }
34
35     ceph_assert(g_conf()->paxos_kill_at != 9);
36
37 // 清空value, 准备下一次循环
38 // get ready for a new round.
39     new_value.clear();
40
41 // WRITING -> REFRESH
42 // among other things, this lets do_refresh() -> mon.bootstrap() ->
43 // wait_for_paxos_write() know that it doesn't need to flush the store
44 // queue. and it should not, as we are in the async completion thread now!
45     ceph_assert(is_writing() || is_writing_previous());
46     state = STATE_REFRESH;
47     ceph_assert(commits_started > 0);
48     --commits_started;
49
50 // 如果有新一轮的提案, 则延长租约
51 // 完成当前循环
52     if (do_refresh()) {
53         commit_proposal();
54         if (mon.get_quorum().size() > 1) {
55             extend_lease();
56         }
57
58         ceph_assert(g_conf()->paxos_kill_at != 10);
59
60         finish_round();
61     }
62 }

```

Peon: handle_commit

```

1 void Paxos::handle_commit(MonOpRequestRef op) {
2     op->mark_paxos_event("handle_commit");
3     auto commit = op->get_req<MMonPaxos>();
4
5     logger->inc(l_paxos_commit);
6
7     if (!mon.is_peon()) {

```

```
8         ceph_abort();
9         return;
10    }
11
12    op->mark_paxos_event("store_state");
13    store_state(commit);
14
15    (void)do_refresh();
16 }
17
```

参考

<http://bean-li.github.io/ceph-paxos/>

<http://bean-li.github.io/ceph-paxos-2/>

[ceph monitor----初始化和选举 - yimuxi - 博客园](#)

<https://www.cnblogs.com/yi-mu-xi/p/10364797.html>