

# OSD 读写源码分析

## 1 日志

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

```
1 $ git log --pretty=oneline
2 7fe91d5d5842e04be3b4f514d6dd990c54b29c76 (HEAD, tag: v18.2.1) 18.2.1
```

ceph中日志机制很高级，比如源码中是这样描述的：

```
1 dout(20) << __func__ << " alloc reply " << ctx->reply
2         << " result " << result << endl;
```

实际输出却是：

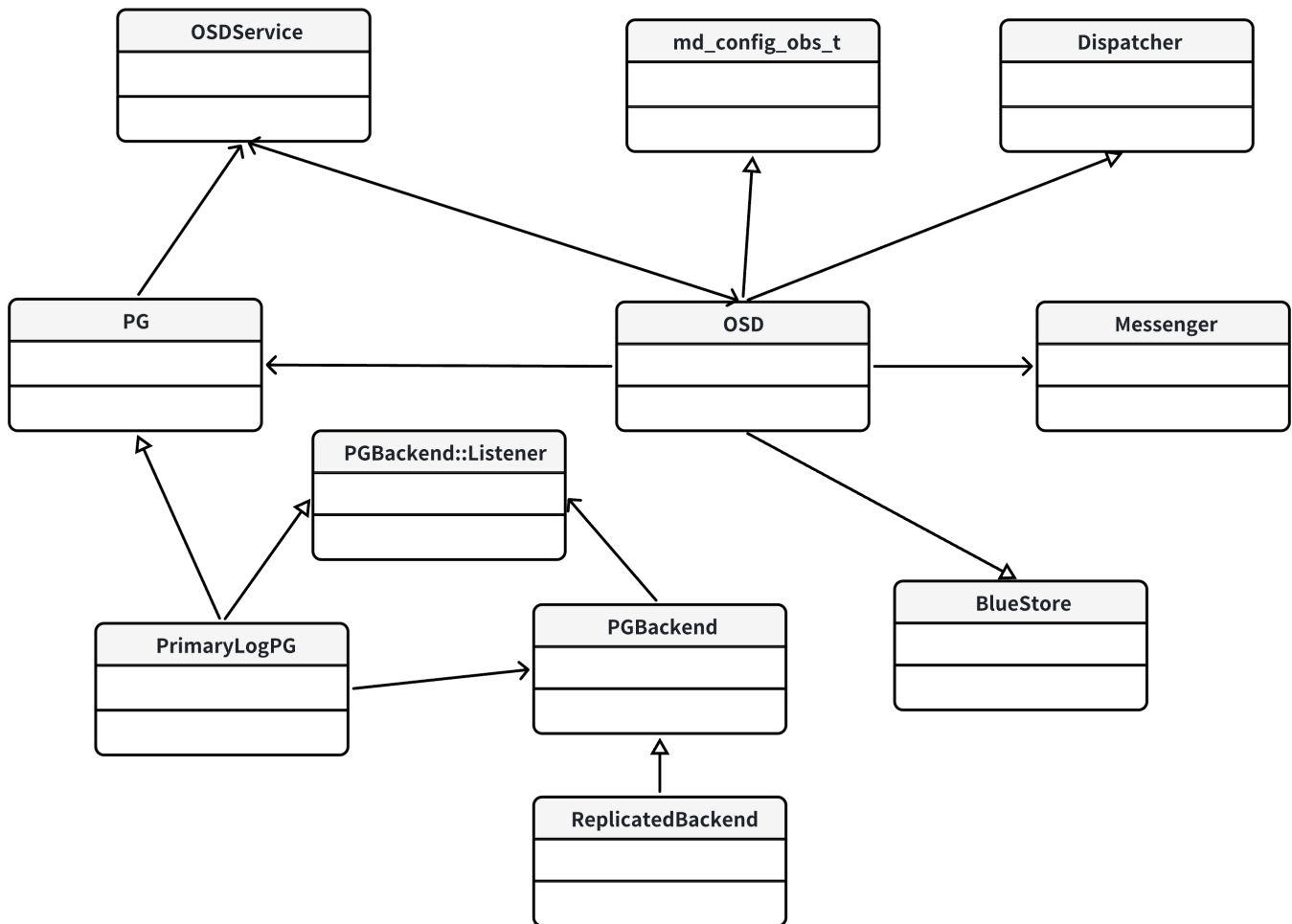
```
1      2024-01-17T05:39:56.991+0000 7f2018753700 20 osd.0 pg_epoch: 130
pg[4.16( v 130'48755 (129'46400,130'48755] local-lis/les=107/108 n=401
ec=45/24 lis/c=107/107 les/c/f=108/108/0 sis=107) [0,1,2] r=0 lpr=107
crt=130'48755 lcod 130'48754 mlcod 130'48754 active+clean] execute_ctx alloc
reply 0x55f568510000 result 0
```

在\_\_func\_\_前面，时间戳、线程号、日志级别、pg等，都是动态添加进去的。

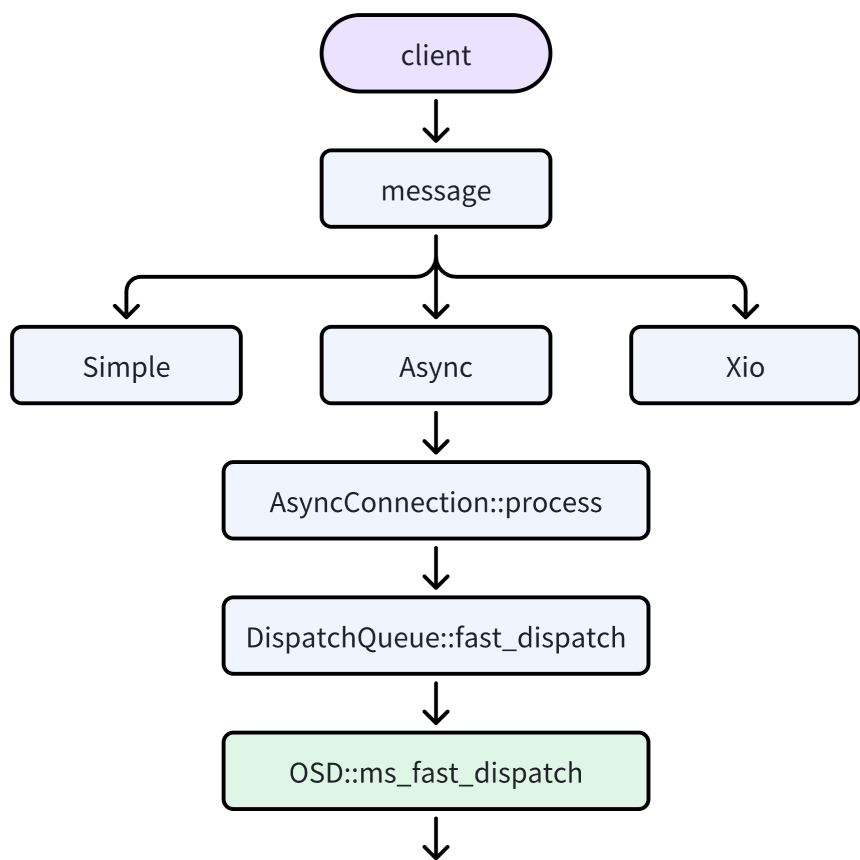
```
1 # 这个函数就添加了pg信息的日志前缀，具体是怎么做到的以后再看
2 template <typename T>
3 static ostream& _prefix(std::ostream *_dout, T *pg) {
4     return pg->gen_prefix(*_dout);
5 }
```

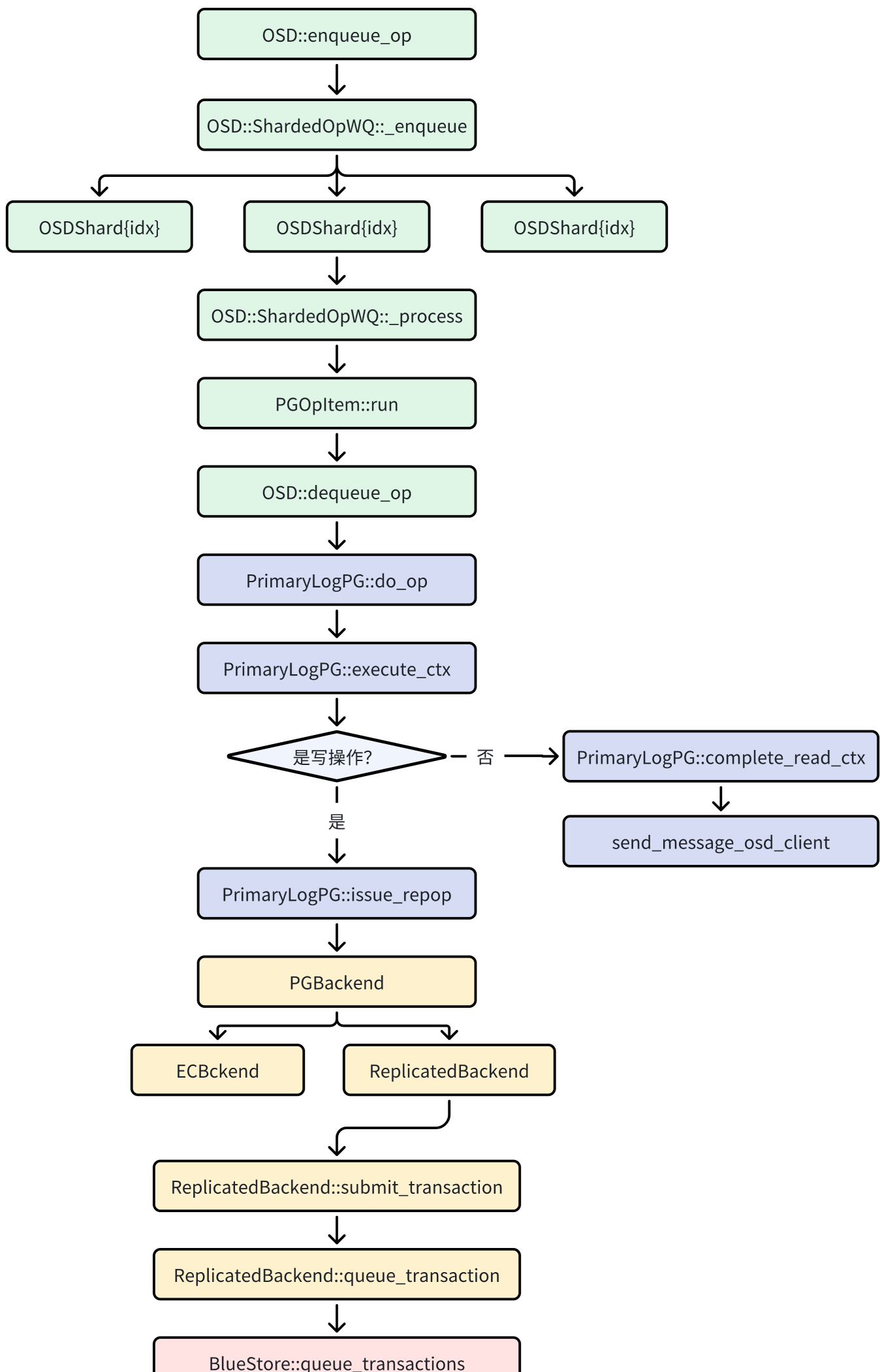
## 2 总览

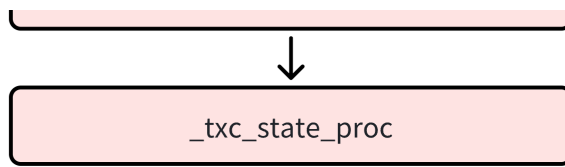
### 2.1 OSD读写类图



2.2 主要流程

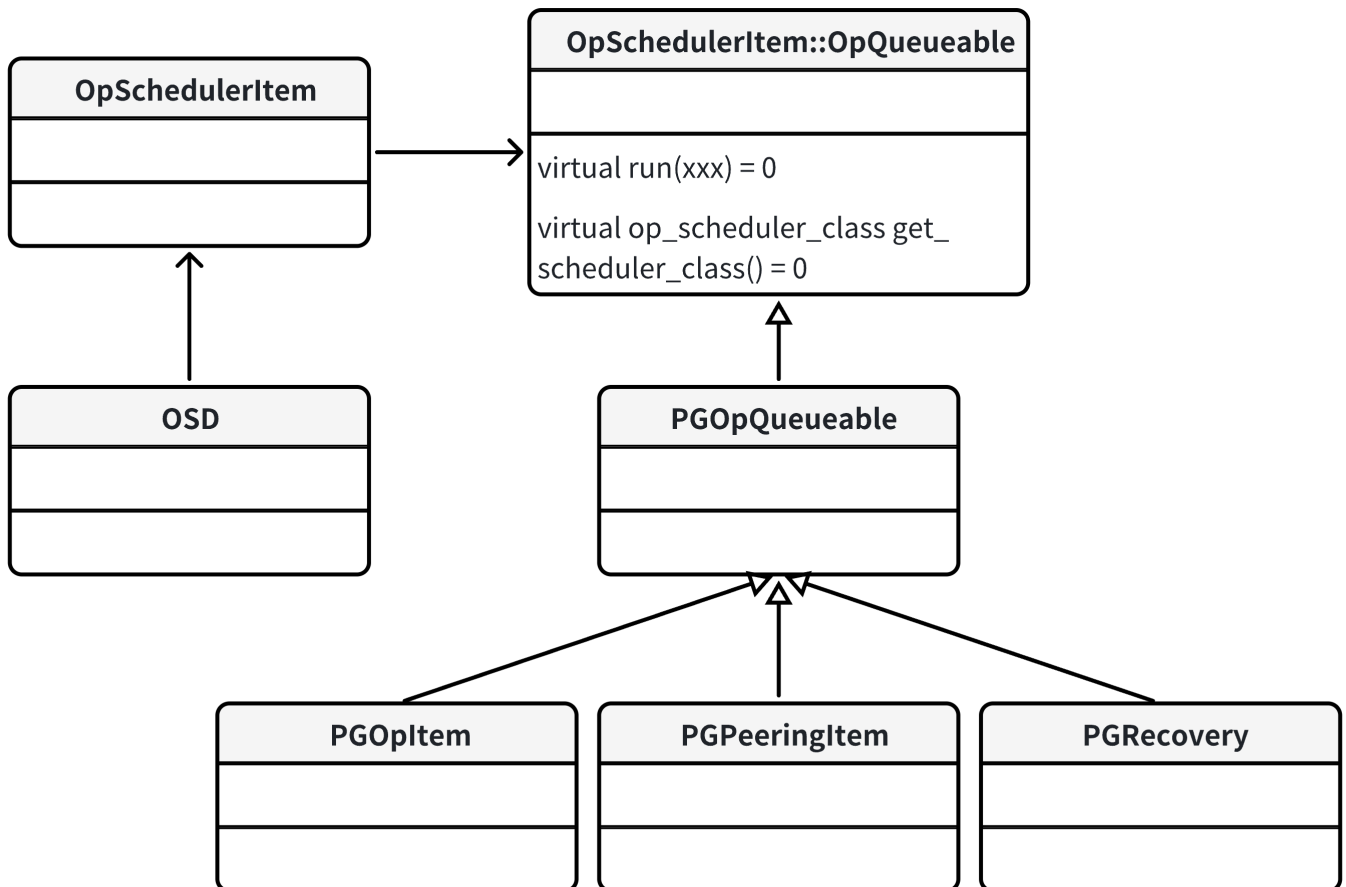






### 3 相关数据结构

#### OpSchedulerItem



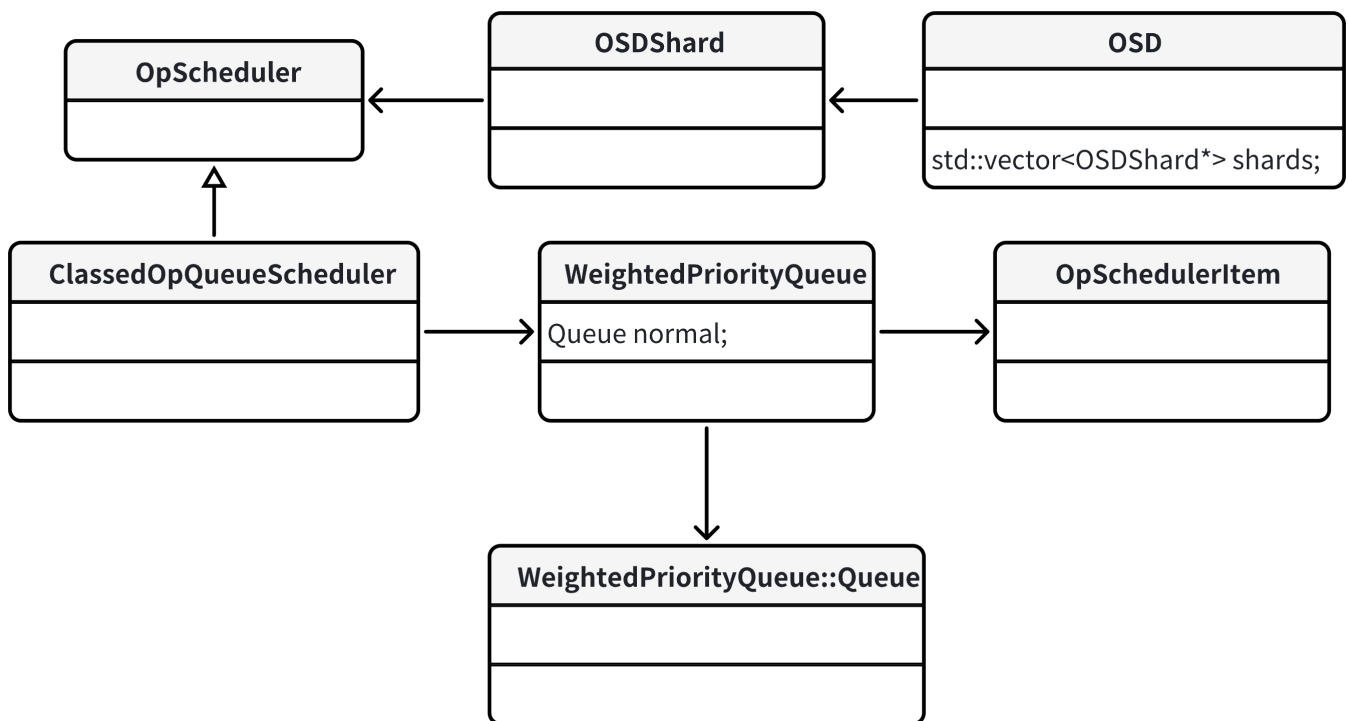
```
1 class OpSchedulerItem {
2     OpQueueable::Ref qitem;
3     int cost;
4     unsigned priority;
5     utime_t start_time;
6     uint64_t owner;
7     epoch_t map_epoch;
8
9     OpSchedulerItem(
10         OpQueueable::Ref &&item,
11         int cost,
12         unsigned priority,
```

```

13     utime_t start_time,
14     uint64_t owner,
15     epoch_t e)
16     : qitem(std::move(item)),
17       cost(cost),
18       priority(priority),
19       start_time(start_time),
20       owner(owner),
21       map_epoch(e) {}
22 };

```

## OSDShard



```

1 struct OSDShard {
2     /// priority queue
3     ceph::osd::scheduler::OpSchedulerRef scheduler;
4     OSDShard(
5         int id,
6         CephContext *cct,
7         OSD *osd);
8 };

```

```

1 class OpScheduler {

```

```

2 public:
3     // Enqueue op for scheduling
4     virtual void enqueue(OpSchedulerItem &&item) = 0;
5     ...
6 };
7
8 template <typename T>
9 class ClassedOpQueueScheduler final : public OpScheduler {
10     T queue;
11     ...
12 };

```

```

1 template <typename T, typename K>
2 class WeightedPriorityQueue : public OpQueue <T, K> {
3     ...
4     class Queue{};
5     Queue normal;
6 }

```

## 初始化

```

1 for (uint32_t i = 0; i < num_shards; i++) {
2     OSDShard *one_shard = new OSDShard(
3         i,
4         cct,
5         this);
6     shards.push_back(one_shard);
7 }
8 =====
9 OSDShard::OSDShard(
10     int id,
11     CephContext *cct,
12     OSD *osd)
13 : shard_id(id),
14   cct(cct),
15   osd(osd),
16   shard_name(string("OSDShard.") + stringify(id)),
17   sdata_wait_lock_name(shard_name + "::sdata_wait_lock"),
18   sdata_wait_lock{make_mutex(sdata_wait_lock_name)},
19   osdmap_lock{make_mutex(shard_name + "::osdmap_lock")},
20   shard_lock_name(shard_name + "::shard_lock"),
21   shard_lock{make_mutex(shard_lock_name)},
22   scheduler(ceph::osd::scheduler::make_scheduler(

```

```

23         cct, osd->whoami, osd->num_shards, id, osd->store->is_rotational(),
24         osd->store->get_type(), osd->monc)),
25         context_queue(sdata_wait_lock, sdata_cond){}
26 =====
27 OpSchedulerRef make_scheduler() {
28     return std::make_unique<
29         ClassedOpQueueScheduler<WeightedPriorityQueue<OpSchedulerItem, client>>>(
30         cct,
31         cct->_conf->osd_op_pq_max_tokens_per_priority,
32         cct->_conf->osd_op_pq_min_cost
33     );
34 }

```

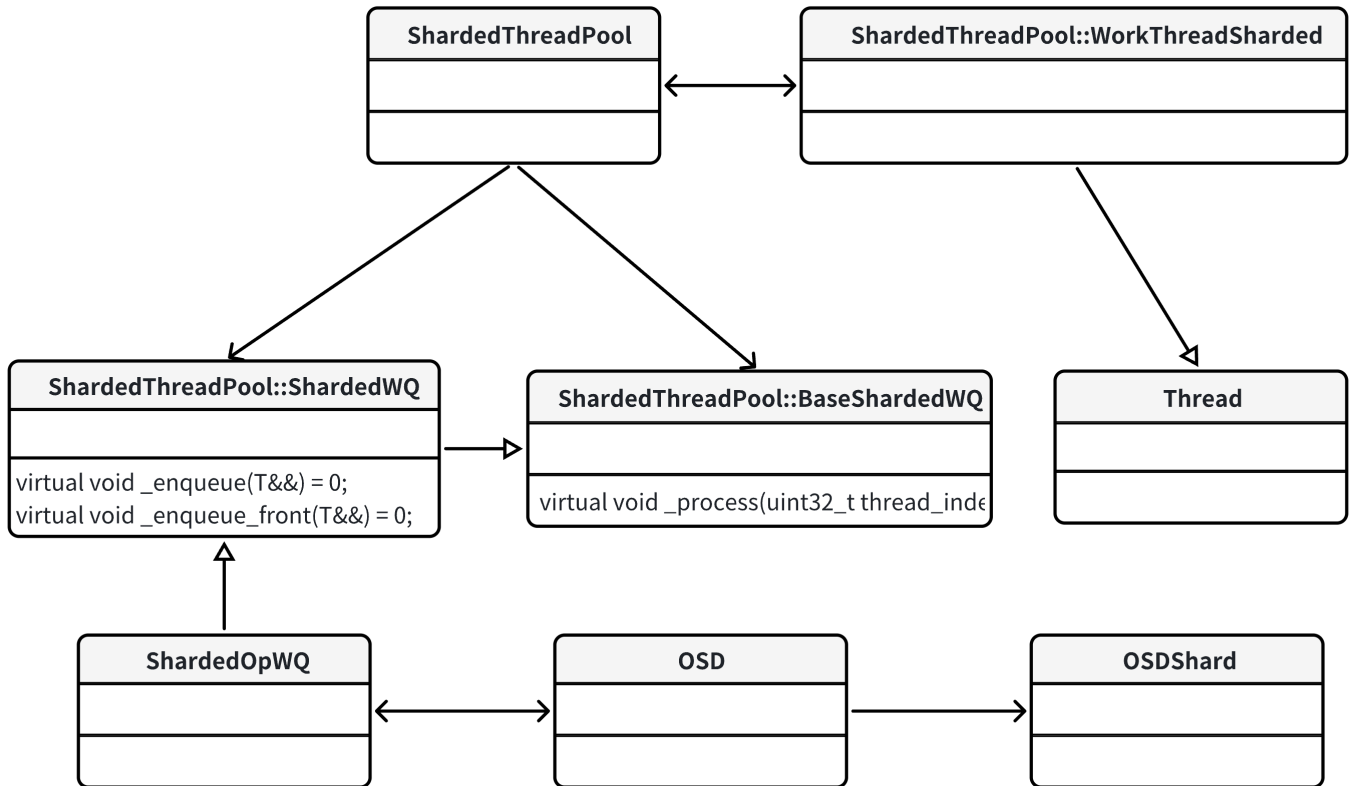
## 添加

```

1  =====调用=====
2  OSDShard* sdata = osd->shards[shard_index];
3  sdata->scheduler->enqueue(std::move(item));
4  =====ClassedOpQueueScheduler=====
5  void enqueue(OpSchedulerItem &&item) final {
6      unsigned priority = item.get_priority();
7      unsigned cost = item.get_cost();
8      queue.enqueue(item.get_owner(), priority, cost, std::move(item));
9  }
10 =====WeightedPriorityQueue=====
11 void enqueue(K cl, unsigned p, unsigned cost, T&& item) final {
12     normal.insert(p, cl, cost, std::move(item));
13 }
14 // normal是其内部类Queue的实例，就分析到这了，相应的，取出的操作：
15 T dequeue() override {
16     return normal.pop();
17 }

```

## ShardedOpWQ



## ShardedThreadPool

在Ceph中，ShardedThreadPool是一个特殊类型的线程池。普通的线程池中，如果任务之间有互斥性，那么正在处理该任务的两个线程有一个必须等待另一个处理完成后才能处理，从而导致线程的阻塞，性能下降。为了解决这个问题，ShardedThreadPool引入了"Shard"的概念。每个线程对应一个任务队列，所有需要顺序执行的任务都放在同一个线程的任务队列里，全部由该线程执行。这样，就可以避免因任务之间的互斥性导致的线程阻塞，从而提高性能。

## ShardedOpWQ

```

1  class ShardedOpWQ : public ShardedThreadPool::ShardedWQ<OpSchedulerItem>
2  {
3      OSD *osd;
4      bool m_fast_shutdown = false;
5  public:
6      ShardedOpWQ(OSD *o,
7                  ceph::timespan ti,
8                  ceph::timespan si,
9                  ShardedThreadPool* tp)
10         : ShardedThreadPool::ShardedWQ<OpSchedulerItem>(ti, si, tp),
11           osd(o) {}
12  void _process(uint32_t thread_index, ceph::heartbeat_handle_d *hb) override;
13  void _enqueue(OpSchedulerItem&& item) override;
14  void _enqueue_front(OpSchedulerItem&& item) override;
15  } op_shardedwq;
  
```



# 入队

```
1 op_shardedwq.queue(  
2   OpSchedulerItem(  
3     unique_ptr(new PGOpItem(pg, std::move(op))),  
4     cost, priority, stamp, owner, epoch));  
5 =====  
6 void OSD::ShardedOpWQ::_enqueue(OpSchedulerItem&& item) {  
7   uint32_t shard_index =  
8     item.get_ordering_token().hash_to_shard(osd->shards.size());  
9  
10  OSDShard* sdata = osd->shards[shard_index];  
11  
12  bool empty = true;  
13  {  
14    std::lock_guard l{sdata->shard_lock};  
15    empty = sdata->scheduler->empty();  
16    sdata->scheduler->enqueue(std::move(item));  
17  }  
18  
19  {  
20    std::lock_guard l{sdata->sdata_wait_lock};  
21    if (empty) {  
22      sdata->sdata_cond.notify_all();  
23    } else if (sdata->waiting_threads) {  
24      sdata->sdata_cond.notify_one();  
25    }  
26  }  
27 }
```

## 4 读写流程

### 4.1 从分发请求到完成io

首先要清楚的一些事情：

通信方式采用异步，存储模式使用主备复制，存储方式使用BlueStore

```
void OSD::ms_fast_dispatch(Message *m)  
{  
  // 如果该消息具有异步分发的特性或者消息类型不是CEPH_MSG_OSD_OP  
  if (m->get_connection()->has_features(CEPH_FEATUREMASK_RESEND_ON_SPLIT)) {  
    m->get_type() != CEPH_MSG_OSD_OP {  
      // 直接将该消息加入队列  
      enqueue_op(  
        static_cast<MOSDFastDispatchOp*>(m->get_spg()), std::move(op),  
        static_cast<MOSDFastDispatchOp*>(m->get_map_epoch()));  
    } else {  
      // 对于旧的客户端，这是一个MOSDOp（唯一没有明确spg_id的快分发消息）；  
      // 我们需要在保持交付顺序的同时将它们映射到spg。t  
      auto priv = m->get_connection()->get_priv();  
      // 尝试获取session  
      if (auto session = static_cast<Session*>(priv.get()); session) {  
        // 将操作请求加入等待映射的列表中  
        session->waiting_on_map.push_back(op);  
        // 获取最新OSD映射  
        OSDMapRef nextmap = service.get_nextmap_reserved();
```

旧版使用dispatch\_xxx

```
void OSD::dispatch_session_waiting(const ceph::ref_t<Session*>&  
session, OSDMapRef oadm)  
{  
  // 遍历等待映射的操作请求  
  auto i = session->waiting_on_map.begin();  
  while (i != session->waiting_on_map.end()) {  
    OpRequestRef op = &*i;  
    // 获取操作请求的消息  
    auto m = op->get_req("MOSDFastDispatchOp");  
    // 从等待映射的列表中移除操作请求  
    session->waiting_on_map.erase(i++);  
    // 减少操作请求的引用计数  
    op->put();  
    // 定义 PG ID  
    spg_t pgid;  
    // 如果操作请求的类型是 CEPH_MSG_OSD_OP
```

```
dispatch_session_waiting(session, nextmap);
--
新任务直接enqueue_op

// 创建一个操作请求并将其加入队列
void OSD::enqueue_op(pg_t pg, OpRequestRef& op, epoch_t epoch)
{
    if (PGRecoveryMsg::is_recovery_msg(op)) {
        // 是恢复消息
        op_shardedwq.queue(
            OpSchedulertem{
                unique_ptr<OpSchedulertem>:OpQueueable(new PGRecoveryMsg(pg,
                    std::move(op))),
                cost, priority, stamp, owner, epoch});
    } else {
        // 不是恢复消息
        op_shardedwq.queue(
            OpSchedulertem{
                unique_ptr<OpSchedulertem>:OpQueueable(new PGOpitem(pg, std::move(op))),
                cost, priority, stamp, owner, epoch});
    }
}
```

```
if (m->get_type() == CEPH_MSG_USO_OP) {
    pg_t actual_pgid = osdmap->raw_pg_to_pg(
        static_cast<const MOSDOp*>(m)->get_pgid());
    osdmap->get_primary_shard(actual_pgid, &pgid)
    } else {
        pgid = m->get_pgid();
    }
    // 将操作请求加入队列
    enqueue_op(pgid, std::move(op), m->get_map_epoch());
}
```

```
void ShardedThreadPool::queue(T& item) {
    _enqueue(std::move(item));
}
_enqueue 是纯虚函数
virtual void _enqueue(T& item) = 0;
```

```
// 将操作请求加入队列
void OSD::shardedOpWQ::enqueue(OpSchedulertem& item) {
    // 计算分片索引
    uint32_t shard_index =
        item.get_ordering_token().hash_to_shard(osd->shards.size());
    // 获取 OSD 分片数据
    OSDShard *sdata = osd->shards[shard_index];
    bool empty = true;
    {
        std::lock_guard l(sdata->shard_lock);
        // 判断消息是否已空
        empty = sdata->scheduled.empty();
        // 将操作请求加入队列
        sdata->scheduled.enqueue(std::move(item));
    }
    {
        std::lock_guard l(sdata->sdata_wait_lock);
        if (empty) { // 如果消息为空，则通知所有等待的线程
            sdata->sdata_cond.notify_all();
        } else if (sdata->waiting_threads) { // 否则，只通知一个等待的线程
            sdata->sdata_cond.notify_one();
        }
    }
}
```

```
void OSD::ShardedOpWQ::process(uint32_t thread_index, heartbeat_handle_d *hb) {
    // OSD "osd"
    // auto& sdata = osd->shards[thread_index];
    // PGRef pg = slot->pg;
    // ThreadPool::TPHandle tp_handle(osd->cct, hb, timeout_interval, suicide_interval);
    // OSDShard *sdata = slot->sdata;
    // auto qi = std::move(slot->process.front());
    qi.run(osd, sdata, pg, tp_handle);
    // qi是OpSchedulertem
}
```

```
void PGOpitem::run(
    OSD *osd,
    OSDShard *sdata,
    PGRef pg,
    ThreadPool::TPHandle &handle)
{
    osd->dequeue_op(pg, op, handle);
    pg->unlock();
}
```

```
void run(OSD *osd, OSDShard *sdata, PGRef pg, ThreadPool::TPHandle &handle) {
    qiitem = run(osd, sdata, pg, handle);
    // qiitem 是 OpQueueable，而OpQueueable是一个接口，当前场景中，其实现是
    PGOpitem
}
```

```
void OSD::dequeue_op(
    PGRef pg,
    OpRequestRef op,
    ThreadPool::TPHandle &handle)
{
    // 获取操作请求的消息
    const Message *m = op->get_req();
    // 获取当前时间
    utime_t now = ceph_clock_now();
    // 设置操作的截止时间
    op->set_dequeued_time(now);
    // 计算操作的延迟
    utime_t latency = now - m->get_recv_stamp();
    // 更新操作的延迟统计
    logger->tinc(_osd_before_dequeue_op_lat, latency);
    // 标记操作已经到达 pg
    op->mark_reached_pg();
    // 执行操作请求
    pg->do_request(op, handle);
}
```

```
int PrimaryLogPG::prepare_transaction(OpContext *ctx) {
    // 准备事务，将相关的操作封装为事务，并放入ctx->op_list中
    int result = do_osd_ops(ctx, ctx->ops);
    if (osd_snap == CEPH_NOSNAP) make_writable(ctx);
    finish_ctx(ctx);
    ctx->new_ops.exists ? pg_log_entry_t::MODIFY :
        pg_log_entry_t::DELETE;
    result;
}
```

```
int PrimaryLogPG::do_osd_ops(OpContext *ctx, vector<OSDOp*> ops) {
    PGTransaction *t = ctx->op_list.get(); // 取出事务
    case CEPH_OSD_OP_WRITEFULL; // 这是最复杂的，还有一个局部函数(CEPH_OSD_OP_WRITE)，唯一区别就是多了些偏移量的
    检查
    ++ctx->num_write;
    result = 0;
    {
        result = check_offset_and_length; // 检查偏移量和长度
        0, op.extent.length; // 偏移量为0，长度为操作的范围
        static_cast<Option>::size, in=osd_max_object_size, get_dpp(); // 对象的最大大小为osd的最大对象大小
        if (result < 0)
            break;
        if (pool.info.has_flag(pg_pool_t::FLAG_WRITE_FADVISE_DONTNEED))
            op.flags |= CEPH_OSD_OP_FLAG_FADVISE_DONTNEED;
        maybe_create_new_object(ctx);
        if (pool.info.is_erasure()) { // 如果是纠删码的
            t->truncate(soid, 0); // 数据对象到数据长度
        } else if (soid.exists && op.extent.length < soid.size) { // 否则，如果对象存在且操作长度小于对象信息的大小
            t->truncate(soid, op.extent.length); // 则截断对象到操作长度
        }
        if (op.extent.length) {
            t->write(soid, 0, op.extent.length, osd_op.indata, op.flags); // 写入数据到缓冲区
        }
        ctx->clean_regions.mark_data_region_dirty(0); // 标记数据区域为脏
        std::max(uint64_t, op.extent.length, op.size); // 数据区域的长度为操作长度和对象信息大小的较大者
        write_update_size_and_usage(ctx->delta_stats, op, ctx->modified_ranges; // 更新大小和使用情况
        0, op.extent.length, true);
        break;
    }
}
```

```
void PrimaryLogPG::do_request(OpRequestRef op,
    ThreadPool::TPHandle &handle)
{
    switch (msg_type) {
        case CEPH_MSG_OSD_OP:
            do_op(op);
            break;
        --
    }
}
```

```
void PrimaryLogPG::do_op(OpRequestRef op)
{
    execute_ctx(ctx);
}
```

```
void PrimaryLogPG::execute_ctx(OpContext *ctx)
{
    ctx->op_list.reset(new PGTransaction()); // 创建一个新的事务
    int result = prepare_transaction(ctx);
    bool pending_async_reads = ctx->pending_async_reads.empty();
    if (result == ENOPROGRESS) pending_async_reads {
        // come back later:
        if (pending_async_reads) {
            ceph_assert(pool.info.is_erasure());
            in_progress_async_reads.push_back(make_pair(op, ctx));
            ctx->start_async_reads(this);
        }
        return;
    }
    // 创建回复
    ctx->reply = new MOSDOpReply(m, result, get_osdmap_epoch(), 0, ignore_out_data);
    // 同步读
    if ((ctx->op_list.empty() || result == 0) && ctx->update_log_only) {
        if (result == 0)
            do_osd_op_effects(ctx, m->get_connection());
        complete_read_ctx(result, ctx);
        return;
    }
}
```

```
// 这个函数是测试的结果，打包发送给客户端
void PrimaryLogPG::complete_read_ctx(int result, OpContext *ctx)
{
    auto m = ctx->op->get_req->MOSDOp();
    ceph_assert(ctx->pending_async_reads.empty());
    for (auto p = ctx->ops->begin();
        p != ctx->ops->end(); ++p) {
        if (p->reply == 0 && !p->flags & CEPH_OSD_OP_FLAG_FAIL0K) {
            break;
        }
        ctx->bytes_read += p->outdata.length();
        ctx->reply->get_header().data_off = (ctx->data_off ? ctx->data_off : 0);
        MOSDOpReply *reply = ctx->reply;
        ctx->reply = nullptr;
        if (result == 0) {
            if (ctx->ignore_log_op_stats)
                log_op_stats(ctx->op, ctx->bytes_written, ctx->bytes_read);
            publish_stats_to_osd();
        }
        if (ctx->reply)
            reply->set_reply_versions(eversion_t(), ctx->obs->ol.user_version);
        else {
            reply->set_reply_versions(eversion_t(), ctx->user_at_version);
        }
        else if (result == -EIO) {
            reply->set_enevent_reply_versions(info.last_update, info.last_user_version);
        }
        reply->set_result(result);
        reply->add_flags(CEPH_OSD_FLAG_ACK | CEPH_OSD_FLAG_ONDISK);
        osd_send_message_osd_client(reply, m->get_connection());
        close_op_ctx(ctx);
    }
}
```

```
// 除此之外，如果是该操作，就是退步分支。
case CEPH_OSD_OP_SYNC_READ; // 同步读不支持异步读，从副本异步读后由副本读到
if (pool.info.is_erasure()) {
    result = -EOPNOTSUPP;
    break;
}
// fall through
case CEPH_OSD_OP_READ:
    ++ctx->num_read;
    tracepoint(osd, do_osd_op_pre_read, soid, old.name.c_str(),
        soid.snap_val, ol.size, ol.truncate_seq, op.extent.offset,
        op.extent.length, op.extent.truncate_size,
        op.extent.truncate_seq);
    if (op_finisher == nullptr) {
        if (ctx->data_off)
            ctx->data_off = op.extent.offset;
    }
    result = do_read(ctx, osd_op);
    else {
        result = op_finisher->execute();
    }
    break;
}
```

```
// 采用主备复制策略的不支持异步读
void ReplicatedBackend::objects_read_async(...)
{
    ceph_abort_msg("async read is not used by replica pool");
}
// 采用纠删码的可以
void ECBackend::objects_read_async(...)
{
    const subject_t &id;
    const list<pair<boost::tuple<uint64_t, uint64_t, uint32_t,
        pair<bufferlist*, Context*>>> &&to_read,
        Context*>, complete,
        bool fast_read)
    {
        objects_read_and_reconstruct(
            reads,
            fast_read,
            make_lambda_context(
                ...
            )
        );
    }
}
```

默认情况下，使用的  
的主备复制，不  
会走到退步分支

```
void PrimaryLogPG::OpContext::start_async_reads(PrimaryLogPG *pg)
{
    inflight_reads = 1;
    list<pair<boost::tuple<uint64_t, uint64_t, uint32_t,
        pair<bufferlist*, Context*>>> >> in;
    in.swap(pending_async_reads);
    pg->pgbackend->objects_read_async(
        obs->obs.ol.soid,
        in,
        new OnReadComplete(pg, this, pg->get_pool(), fast_read);
}
```

```
recovery_state.update_trim_to; // 更改需要修改的块的状态
ctx->register_on_commit;
```

```

{m, ctx, this}{
    if (ctx->op)
        log_op_stats["ctx->op, ctx->bytes_written, ctx->bytes_read"];

    if (m && ctx->sent_reply) {
        MOSDOpReply "reply = ctx->reply;
        ctx->reply = nullptr;
        reply->valid_flags |= CEPH_OSD_FLAG_ACK | CEPH_OSD_FLAG_ONDISK;
        dout(10) << "sending reply on " << m << " " << reply << endl;
        osd->send_message_osd_client(reply, m->get_connection());
        ctx->sent_reply = true;
        ctx->op->mark_commit_sent();
    }

    ctx->register_on_success(
        [ctx, this] {
            do_osd_op_effects(
                ctx,
                ctx->op ? ctx->op->get_req()>get_connection() :
                ConnectionRef{});
        });

    ctx->register_on_finish(
        [ctx] {
            delete ctx;
        });
}

```

register\_XXX

```

// 这里定义了一组回调函数，在事务执行到某一阶段时，调用它们，后续会讲
template <typename F>
void register_on_finish(F &&f) {
    on_finish.emplace_back(std::forward<F>(f));
}

template <typename F>
void register_on_success(F &&f) {
    on_success.emplace_back(std::forward<F>(f));
}

template <typename F>
void register_on_applied(F &&f) {
    on_applied.emplace_back(std::forward<F>(f));
}

template <typename F>
void register_on_commit(F &&f) {
    on_committed.emplace_back(std::forward<F>(f));
}

```

还是execute\_ctx中，继续

```

// 获取一个新的事务ID
ceph_tid_t rep_tid = osd->get_tid();

// 用于收集所有副本操作的信息
RepGather "repop = new_repop(ctx, rep_tid);

// 将操作发送到所有副本
issue_repop(repop, ctx);

// 评估副本操作的结果，检查操作是否成功，是否需要回滚等
eval_repop(repop);

```

issue\_repop

```

void PrimaryLogPG::issue_repop(RepGather "repop, OpContext "ctx)
{
    recovery_state_pre_submit_op(); // 预提交
    osd->log(
        ctx->at_version(),
        pgbackend->submit_transaction(); // 提交事务
    );
    osd->delta_stats(
        ctx->at_version(),
        std::move(ctx->op_t),
        recovery_state_get_pg_trim_to(),
        recovery_state_get_min_last_complete_ondisk(),
        std::move(ctx->log),
        ctx->updated_hset_history,
        on_all_commit,
        repop->rep_tid,
        ctx->reqid,
        ctx->op);
}

```

```

void PrimaryLogPG::eval_repop(RepGather "repop)
{
    // 检查所有的复制操作是否都已经提交
    if (repop->all_committed) {
        // 遍历所有已经提交的复制操作，并执行它们的回调函数
        for (auto p = repop->on_committed.begin();
            p != repop->on_committed.end();
            repop->on_committed.erase(p++)) {
            (*p)();
        }
        // 查找等待磁盘写入完成的复制操作，并发送重复的提交信息
        auto it = waiting_for_ondisk.find(repop->);
        if (it != waiting_for_ondisk.end()) {
            for (auto& i : it->second) {
                int return_code = repop->v;
                if (return_code == 0) {
                    return_code = std::get<2>(i);
                }
                osd->reply_op_error(std::get<0>(i), return_code, repop->v,
                    std::get<1>(i), std::get<3>(i));
            }
            waiting_for_ondisk.erase(it);
        }
        // 发布统计信息到OSD
        publish_stats_to_osd();
    }

    // 如果队列的第一个复制操作就是当前的复制操作，那么它会移除所有已经提交的复制操作，并执行它们的回调函数
    if (repop->queue.front() == repop) {
        RepGather "to_remove = nullptr;
        while (!repop->queue.empty() &&
            (to_remove = repop->queue.front())->all_committed) {
            repop->queue.pop_front();
            for (auto p = to_remove->on_success.begin();
                p != to_remove->on_success.end();
                to_remove->on_success.erase(p++)) {
                (*p)();
            }
            remove_repop(to_remove);
        }
    }
}

```

```

void ReplicatedBackend::submit_transaction(ox)
{
    // 根据提供的参数生成一个操作事务，并确定哪些对象被添加或删除
    vector<op_log_entry_t> log_entries(log_entries);
    ObjectStore::Transaction op_t(
        PGTransactionUpr t(std::move(t)),
        set<object_t> added, removed,
        generate_transaction(
            t,
            coll,
            log_entries,
            &op_t,
            &added,
            &removed,
            get_osdmap()>require_osd_release);
    // 添加要修改的osd
    op.waiting_for_commit.insert(
        parent->get_acting_recovery_backfill_shards().begin(),
        parent->get_acting_recovery_backfill_shards().end());
    // 发送-复制操作-信息
    issue_op(
        osd,
        at_version,
        tid,
        reqid,
        trim_to,
        min_last_complete_ondisk,
        added.size() ? "added begin() : object_t",
        log_entries,
        hset_history,
        &op,
        op_t);
    // 等待执行
    parent->queue_transactions(tls, op, op);
}

```

issue\_op

```

void ReplicatedBackend::issue_op(ox)
{
    if (parent->get_acting_recovery_backfill_shards().size() > 1) {
        for (const auto& shard : get_acting_recovery_backfill_shards()) {
            Message "w";
            w.r = generate_subop(osd, at_version, tid, reqid, pg_trim_to,
                min_last_complete_ondisk, new_temp_oid, discard_temp_oid, logs, hset_hist,
                op_t, shard, pinfo);
            // 发送操作信息
            get_parent()->send_message_osd_cluster(shard.oid, w, get_osdmap_epoch());
        }
    }
}

```

```

void PrimaryLogPG::queue_transactions(
    std::vector<ObjectStore::Transaction& ts, OpRequestRef op) override {
    osd->store->queue_transactions(ch, ts, op, NULL);
}

```

```

void BlueStore::queue_transactions(
    CollectionHandle& ch,
    vector<Transaction& ts,
    TrackedOpRef op,
    ThreadPool::TPHandle "handle)
{
    // 一大堆不知所云的检查
    // 事务处理，真正的写流程
    _tx_state_proc(tx);
}

```

```

map<object_t, pair<int, extent_map>> > os, ceph::noat, noat, to_read, on_completes;
}

```

```

void ECBackend::objects_read_and_reconstruct(
    const map<object_t, L>
    std::list<boost::tuple<uint64_t, uint64_t, uint32_t>>
    > &reads,
    bool fast_read,
    GenContextUpr& mapobject_t, pair<int, extent_map>> && &&f) {
    map<object_t, L> obj_want_to_read;
    set<int> want_to_read;
    get_want_to_read_shards(want_to_read);

    start_read_op(
        CEPH_MSG_PRIO_DEFAULT,
        obj_want_to_read,
        for_read_op,
        OpRequestRef(),
        fast_read, false);
}

```

```

void ECBackend::start_read_op(
    int priority,
    map<object_t, L> set-int>> &want_to_read,
    map<object_t, read_request_t> &to_read,
    OpRequestRef _op,
    bool do_redundant_reads,
    bool for_recovery) {
    ceph_tid_t tid = get_parent()->get_tid();
    ceph_assert(tid.to_read_map.count(tid));
    auto &op = tid.to_read_map.emplace(
        tid,
        ReadOp(
            priority,
            tid,
            do_redundant_reads,
            for_recovery,
            _op,
            std::move(want_to_read),
            std::move(to_read));
    if (!op) {
        op.trace = op->pg_trace;
        op.trace.event("start ec read");
    }
    do_read_op(op);
}

```

```

void ECBackend::do_read_op(ReadOp &op)
{
    //
    get_parent()->send_message_osd_cluster(m, get_osdmap_epoch());
}

```

## 4.2 写操作同步

4.1 写操作的最后，issue\_op使用send\_message\_osd\_cluster发送了同步消息，这节来观察这个操作的两端是怎么工作的

这函数有四个形式的重载，最终调用的都是Connection::send\_message，而send\_message实际又由AsyncConnection实现

```

1 void OSDService::send_message_osd_cluster(int peer, Message *m, epoch_t from_epoch)
2 {
3     OSDMapRef next_map = get_nextmap_reserved();
4     // service map is always newer/newest
5     ceph_assert(from_epoch <= next_map->get_epoch());

```

```

6
7  if (next_map->is_down(peer) ||
8      next_map->get_info(peer).up_from > from_epoch) {
9      m->put();
10     release_map(next_map);
11     return;
12 }
13 ConnectionRef peer_con;
14 if (peer == whoami) {
15     peer_con = osd->cluster_messenger->get_loopback_connection();
16 } else {
17     peer_con = osd->cluster_messenger->connect_to_osd(
18         next_map->get_cluster_addrs(peer), false, true);
19 }
20 maybe_share_map(peer_con.get(), next_map);
21 peer_con->send_message(m);
22 release_map(next_map);
23 }
24
25 void OSDService::send_message_osd_cluster(std::vector<std::pair<int,
    Message*>>& messages, epoch_t from_epoch)
26 {
27     OSDMapRef next_map = get_nextmap_reserved();
28     ceph_assert(from_epoch <= next_map->get_epoch());
29
30     for (auto& iter : messages) {
31         if (next_map->is_down(iter.first) ||
32             next_map->get_info(iter.first).up_from > from_epoch) {
33             iter.second->put();
34             continue;
35         }
36         ConnectionRef peer_con;
37         if (iter.first == whoami) {
38             peer_con = osd->cluster_messenger->get_loopback_connection();
39         } else {
40             peer_con = osd->cluster_messenger->connect_to_osd(
41                 next_map->get_cluster_addrs(iter.first), false, true);
42         }
43         maybe_share_map(peer_con.get(), next_map);
44         peer_con->send_message(iter.second);
45     }
46     release_map(next_map);
47 }
48
49 void send_message_osd_cluster(MessageRef m, Connection *con) {
50     con->send_message2(std::move(m));
51     /*

```

```

52     virtual int send_message2(MessageRef m)
53     {
54         return send_message(m.detach());
55     }
56     */
57 }
58 void send_message_osd_cluster(Message *m, const ConnectionRef& con) {
59     con->send_message(m);
60 }

```

send\_message之后的流程见<5.7 发送消息>

## 4.3 接收消息

4.1的流程开始是ms\_fast\_dispatch，这是谁调用的？

在<5.5 add\_fast\_dispatch\_tail>中，讲了add\_dispatcher\_tail的流程。类似的还有add\_dispatcher\_head

```

1  int OSD::init() {
2      cluster_messenger->add_dispatcher_head(this);
3  }
4  void add_dispatcher_head(Dispatcher *d) {
5      bool first = dispatchers.empty();
6      dispatchers.push_front(d);
7      if (d->ms_can_fast_dispatch_any())
8          fast_dispatchers.push_front(d);
9      if (first)
10         ready();
11 }

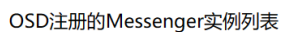
```

这里OSD将自身绑定到了一个fast\_dispatchers中，然后cluster\_messenger就能调用osd的ms\_fast\_dispatch

具体过程在<5.5 add\_fast\_dispatch\_tail>中说明了。

## 5 通信模块

### 5.1 OSD 通信类图

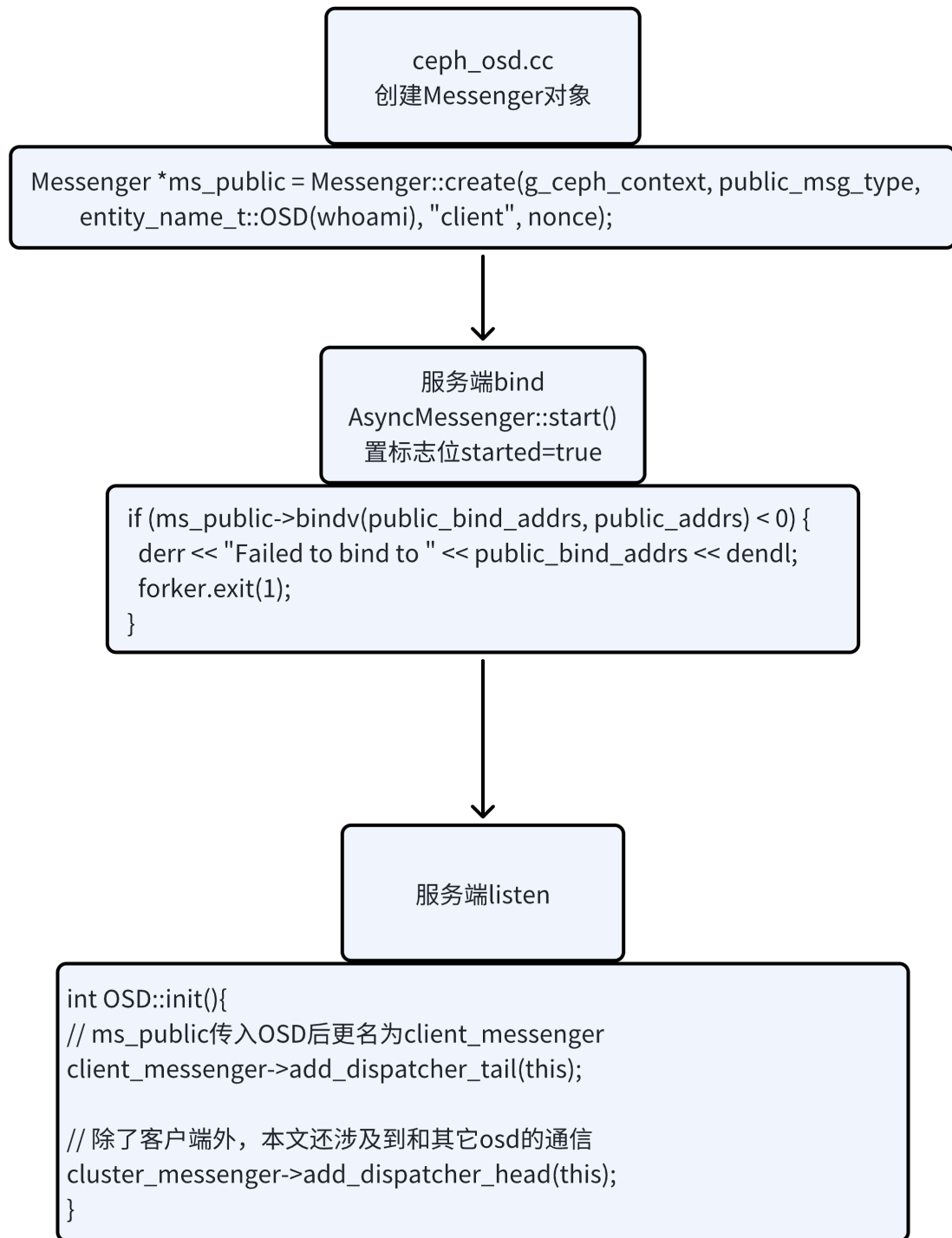


## 5.2 主流程

编程者联盟

1. **创建监听套接字**: 使用 `socket` 函数创建一个套接字。
2. **绑定本机IP和端口**: 使用 `bind` 函数将套接字绑定到特定的IP地址和端口上。
3. **监听**: 使用 `listen` 函数开始监听这个套接字。
4. **创建epoll实例**: 使用 `epoll_create` 函数创建一个epoll实例。
5. **注册文件描述符**: 使用 `epoll_ctl` 函数将监听套接字注册到epoll实例中, 监听 `EPOLLIN` 事件 (新的连接请求)。
6. **等待事件**: 使用 `epoll_wait` 函数等待注册在epoll实例上的文件描述符上的事件发生。
7. **处理事件**: 当 `epoll_wait` 返回时, 遍历所有就绪的事件, 对于每一个事件:
  - 如果是监听套接字上的事件, 那么使用 `accept` 函数接受新的连接, 然后使用 `epoll_ctl` 将新的连接套接字注册到epoll实例中, 监听 `EPOLLIN` 和 `EPOLLOUT` 事件。
  - 如果是连接套接字上的 `EPOLLIN` 事件, 那么读取数据, 处理请求。
  - 如果是连接套接字上的 `EPOLLOUT` 事件, 那么写入数据, 发送响应。
8. **关闭连接**: 当连接结束后, 使用 `close` 函数关闭连接套接字, 并使用 `epoll_ctl` 将其从epoll实例中删除。

ceph的通信模块采用事件驱动的模式, 先创建事件中心, 事件中心接收自定义的各类事件, 触发事件响应时调用相应的回调函数, 松散耦合。



## 5.3 create 流程

创建消息模块、初始化、启动事件监听，但是此时还没有事件加入。



```
Messenger *Messenger::create(CephContext *cct, const std::string &type,
    entity_name_t name, std::string name,
    uint64_t nonce)
{
    if (type == "random" || type.find("async") != std::string::npos)
        return new AsyncMessenger(cct, name, type, std::move(name), nonce);
    else
        return new SimplePolicyMessenger(cct, name, type, std::move(name), nonce);
    return nullptr;
}
```

```
class AsyncMessenger : public SimplePolicyMessenger {
private:
    NetworkStack *stack; // 用来起线程
    std::vector<Processor> processors; // 用来监听
    friend class Processor;
    DispatchQueue dispatch_queue; // 用来分发消息
    ceph::unordered_map<entity_name_t, AsyncConnect
    onRef> conns; // 保存已建立的连接
};

// 创建一个单例
auto single = &cct->lookup_or_create_singleton_object<StackSingleton>{
    "AsyncMessenger:NetworkStack:" + transport_type, true, cct};
single->ready(transport_type);
stack = single->stack.get();
// 启动！！
stack->start();
local_worker = stack->get_worker();
local_connection = ceph::make_ref<AsyncConnection>(cct, this, &dispatch_queue,
    local_worker, true, true);
init_local_connection();
reap_handler = new C_handle_reap(this);
unsigned processor_num = 1;
if (stack->support_local_listen_table())
    processor_num = stack->get_num_worker();
for (unsigned i = 0; i < processor_num; ++i)
    processors.push_back(new Processor(this, stack->get_worker(i), cct));
```

启动Worker中的线程  
加入监听事件  
等待并处理事件

```
// 注意stack是NetworkStack指针，指向PosixNetworkStack
void NetworkStack::start()
{
    ...
    for (Worker *worker : workers) {
        if (worker->is_init())
            continue;
        spawn_worker(add_thread(worker));
    }
    started = true;
    ...
    void spawn_worker(std::function<void (> &&func) override {
        threads.emplace_back(std::move(func));
    }
```

```
std::function<void (> NetworkStack::add_thread(Worker *w)
{
    return [this, w]() {
        rename_thread(w->id); // 重命名线程为当前id
        const unsigned EventMaxWaitUs = 2000000;
        w->center.set_owner(); // 加入监听事件，处理外部事件
        w->initialize();
        w->init_done();
        while (!w->done()) {
            ceph::timespan dur;
            int r = w->center.process_events(EventMaxWaitUs, &dur); // 等待并处理监听事件
            w->perf_logger->inc(&msgr_running_total_time, dur);
        }
        w->reset();
        w->destroy();
    };
}

void spawn_worker(std::function<void (> &&func) override {
    threads.emplace_back(std::move(func));
}
```

center.process\_events

```
int EventCenter::process_events(unsigned timeout_microseconds, ceph::timespan *working_dur)
{
    ...
    struct timeval tv;
    bool blocking = pollers.empty() && !external_num_events.load();
    if (blocking) // 轮询器有事件或有外部事件
        timeout_microseconds = 0;
    tv.tv_sec = timeout_microseconds / 1000000;
    tv.tv_usec = timeout_microseconds % 1000000;

    std::vector<FiredFileEvent> fired_events;
    numevents = driver->event_wait(fired_events, &tv); // 等待事件发生
    ...
}
```

event\_wait

```
int EventCenter::process_events(unsigned timeout_microseconds, ceph::timespan *working_dur)
{
    ...
    // numevents = driver->event_wait(fired_events, &tv);
    for (int event_id = 0; event_id < numevents; event_id++) {
        int rfired = 0;
        FileEvent *event;
        EventCallbackRef cb;
        event = _get_file_event(fired_events[event_id].fd);

        if (event->mask & fired_events[event_id].mask & EVENT_READABLE) {
            rfired = 1;
            cb = event->read_cb;
            cb->do_request(fired_events[event_id].fd); // 调用回调函数，这里是C_handle_notify->do_request
        }

        if (event->mask & fired_events[event_id].mask & EVENT_WRITEABLE) {
            if (rfired) event->read_cb = event->write_cb;
            cb = event->write_cb;
            cb->do_request(fired_events[event_id].fd); // 回调函数
        }
    }

    // 处理外部事件
    if (external_num_events.load()) {
        external_lock.lock();
        std::deque<EventCallbackRef> cur_process;
        cur_process.swap(external_events);
        external_num_events.store(0);
        external_lock.unlock();
        numevents += cur_process.size();
        while (!cur_process.empty()) {
            EventCallbackRef e = cur_process.front();
            Idout(cct, 30) << "func_" << e->do_request(0);
            e->do_request(0);
            cur_process.pop_front();
        }
    }
    return numevents;
}
```

```
class AsyncMessenger : public SimplePolicyMessenger {
private:
    NetworkStack *stack; // 用来起线程
    std::vector<Processor> processors; // 用来监听
    friend class Processor;
    DispatchQueue dispatch_queue; // 用来分发消息
    ceph::unordered_map<entity_name_t, AsyncConnect
    onRef> conns; // 保存已建立的连接
};
```

创建worker,每个Worker都会  
创建一个Epoll

```
class PosixNetworkStack : public NetworkStack {
    std::vector<std::thread> threads;
};

class NetworkStack {
    CephContext *cct;
    vector<Worker*> workers; // 存放worker
    ...
}
```

```
int EventCenter::init(int n, unsigned i, const std::string &t)
{
    ...
    // 这里创建了两个管道，分别赋值给send和receive
    int fds[2];
    notify_receive_fd = fds[0];
    notify_send_fd = fds[1];
    ...
}
```

```
void EventCenter::set_owner()
{
    owner = pthread_self();
    if (!global_centers) {
        global_centers = &cct->lookup_or_create_singleton_object<EventCenter::AssociatedCenters>{
            "AsyncMessenger:EventCenter:global_center:" + type, true};
        global_centers->centers[center_id] = this;
    }

    if (driver->need_wakeup()) {
        notify_handler = new C_handle_notify(this, cct);
        int r = create_file_event(notify_receive_fd, EVENT_READABLE, notify_handler);
    }

    notify_handler就是notify_receive_fd的事件回调函数，EventCenter::create_file_event将
    notify_receive_fd加入epoll的监听列表中，并且注册回调函数
}
```

```
int EventCenter::create_file_event(int fd, int mask, EventCallbackRef ctx)
{
    int r = 0;
    ...
    EventCenter::FileEvent *event = _get_file_event(fd);
    // driver是EpollDriver，封装了epoll操作，这里往其中添加监听事件
    r = driver->add_event(fd, event->mask, mask);
    event->mask = mask;
    if (mask & EVENT_READABLE) { event->read_cb = ctx; }
    if (mask & EVENT_WRITEABLE) { event->write_cb = ctx; }
    return 0;
}
```

create\_file\_event

```
int EpollDriver::event_wait(std::vector<FiredFileEvent> &fired_events, struct timeval *tp)
{
    int retval, numevents = 0;
    retval = epoll_wait(epfd, events, nevent,
        tvp ? (tvp->tv_sec*1000 + tvp->tv_usec/1000) : -1); // 如果tvp=0，那么就立即返回
    if (retval > 0) {
        numevents = retval;
        fired_events.resize(numevents);
        for (int event_id = 0; event_id < numevents; event_id++) {
            int mask = 0;
            struct epoll_event *e = &events[event_id];
            if (e->events & EPOLLIN) mask |= EVENT_READABLE;
            if (e->events & EPOLLOUT) mask |= EVENT_WRITEABLE;
            if (e->events & EPOLLERR) mask |= EVENT_READABLE|EVENT_WRITEABLE;
            if (e->events & EPOLLHUP) mask |= EVENT_READABLE|EVENT_WRITEABLE;
            fired_events[event_id].fd = e->data.fd;
            fired_events[event_id].mask = mask;
        }
        return numevents;
    }
}
```

event\_wait

```
void EventCenter::dispatch_event_external(EventCallbackRef e)
{
    external_events.push_back(e);
    num += external_num_events;
    if (num == 1 && !in_thread())
        wakeup();
}
```

基本流程涉及到的  
重要数据结构，全  
都用黄色框统一  
展示

```
struct StackSingleton {
    CephContext *cct;
    std::shared_ptr<NetworkStack> stack;
    StackSingleton(CephContext *c) : cct(c) {}
    void ready(std::string &type) {
        if (!stack)
            stack = NetworkStack::create(cct, type);
    }
};
```

```
std::shared_ptr<NetworkStack> NetworkStack::create(CephContext
    *c, const std::string &t)
{
    std::shared_ptr<NetworkStack> stack = nullptr;
    if (t == "posix")
        stack.reset(new PosixNetworkStack(c));
    unsigned num_workers = c->conf->ms_async_op_threads; // =3
    for (unsigned worker_id = 0; worker_id < num_workers; ++worker_id) {
        Worker *w = stack->create_worker(c, worker_id);
        int ret = w->center.init(EventNumber, worker_id, t);
        stack->workers.push_back(w);
    }
    return stack;
}
```

```
class PosixWorker : public Worker {
    ceph::NetHandler net;
    ...
};

class Worker {
    EventCenter center;
    ...
};
```

virtual Worker \*PosixNetworkStack::create\_worker(CephContext \*c,
 unsigned worker\_id) override {
 return new PosixWorker(c, worker\_id);
}

center.init

center.set\_owner

C\_handle\_notify

create\_file\_event

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

EventCenter (事件中心) 是一个处理事件的数据结  
构，相当于一个事件处理的容器。它本身并不真正去处  
理事件，通过回调函数的方式来处理事件。同样，如  
何获取需要处理的事件也不是事件中心来完成的，  
它只负责处理，具体对需要处理的事件的获取是通  
过EventDriver来完成的。EventDriver是一个接口类，  
其实现主要是由EpollDriver。  
EpollDriver封装了epoll的接口，事件驱动的接口主要依  
赖于epoll的方式，其中主要有三个函数：  
epoll\_create, epoll\_ctl, epoll\_wait

```
class EventCenter {
    EventDriver *driver;
    ...
};

class EpollDriver : public EventDriver {
    int epfd;
    struct epoll_event *events;
    CephContext *cct;
    int nevent;
    ...
};
```

virtual Worker \*PosixNetworkStack::create\_worker(CephContext \*c,
 unsigned worker\_id) override {
 return new PosixWorker(c, worker\_id);
}

center.init

center.set\_owner

C\_handle\_notify

create\_file\_event

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

```
class EventCenter {
    EventDriver *driver;
    ...
};

class EpollDriver : public EventDriver {
    int epfd;
    struct epoll_event *events;
    CephContext *cct;
    int nevent;
    ...
};
```

virtual Worker \*PosixNetworkStack::create\_worker(CephContext \*c,
 unsigned worker\_id) override {
 return new PosixWorker(c, worker\_id);
}

center.init

center.set\_owner

C\_handle\_notify

create\_file\_event

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

```
class EventCenter {
    EventDriver *driver;
    ...
};

class EpollDriver : public EventDriver {
    int epfd;
    struct epoll_event *events;
    CephContext *cct;
    int nevent;
    ...
};
```

virtual Worker \*PosixNetworkStack::create\_worker(CephContext \*c,
 unsigned worker\_id) override {
 return new PosixWorker(c, worker\_id);
}

center.init

center.set\_owner

C\_handle\_notify

create\_file\_event

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

event\_wait

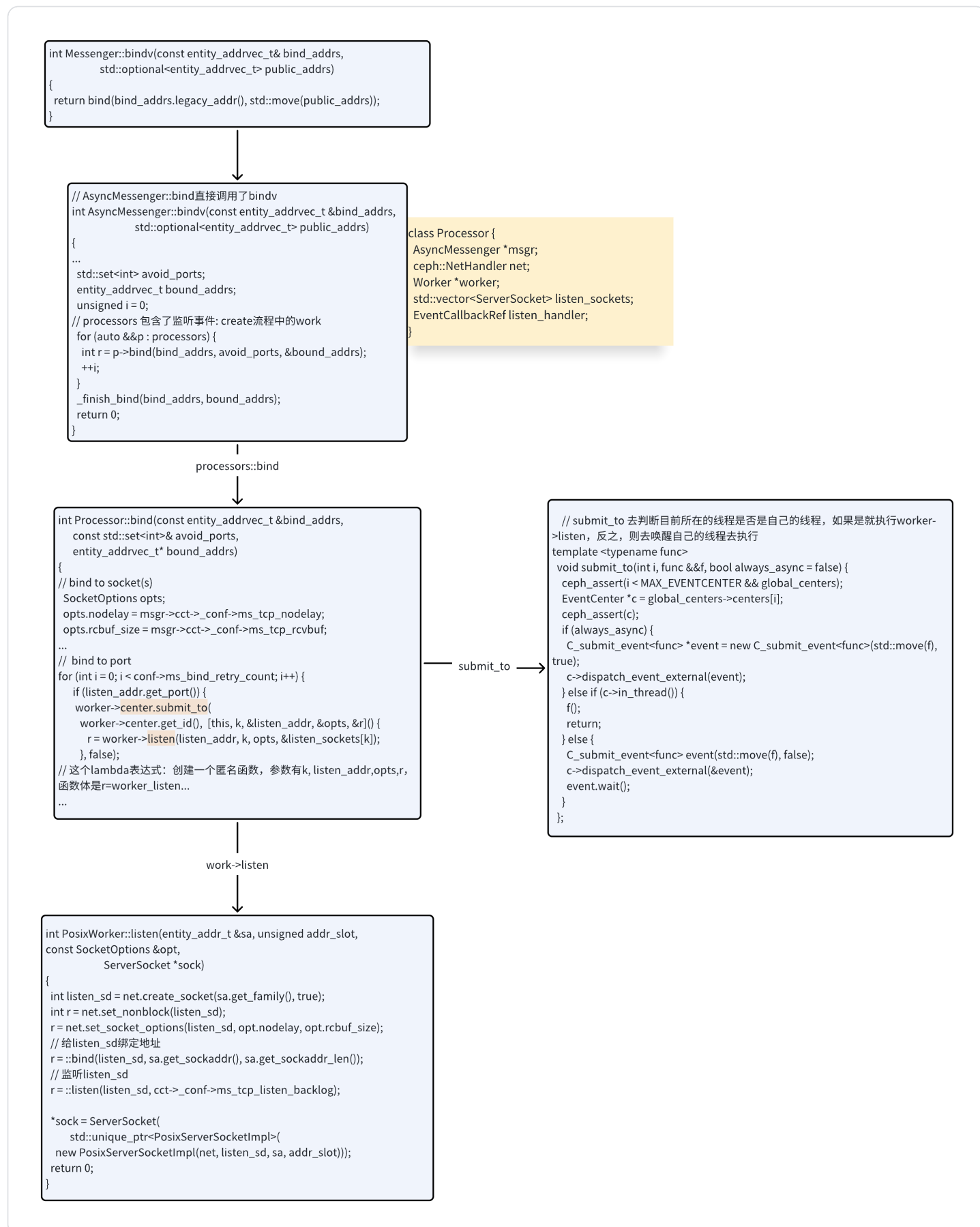
event\_wait

event\_wait

event\_wait

## 5.4 bindv 流程

### 绑定本机ip



## 5.5 add\_dispatcher\_tail 流程

```
void Messenger::add_dispatcher_tail(Dispatcher *d) {
    bool first = dispatchers.empty();
    dispatchers.push_back(d);
    if (d->ms_can_fast_dispatch_any())
        fast_dispatchers.push_back(d);
    if (first)
        ready();
}
```

```
void AsyncMessenger::ready()
{
    stack->ready();
    if (pending_bind) {
        int err = bindv(pending_bind_addr, saved_public_addr);
    }

    std::lock_guard l(lock);
    for (auto &&p : processors)
        p->start();
    dispatch_queue.start();
}
```

```
void Processor::start()
{
    // start thread
    worker->center.submit_to(worker->center.get_id(), [this]() {
        for (auto& listen_socket : listen_sockets) {
            if (listen_socket) {
                worker->center.create_file_event(listen_socket.fd(), EVENT_READABLE,
                    listen_handler);
            }
        }, false);
    }
    // center.create_file_event在5.3介绍过，功能是监听listen_socket.fd(), 回调函数是listen_handler，不清楚回调原理再去看看5.3的流程
}
```

```
class DispatchThread : public Thread {
    DispatchQueue *dq;
public:
    explicit DispatchThread(DispatchQueue *dq) : dq(dq) {}
    void *entry() override {
        dq->entry();
        return 0;
    }
} dispatch_thread;

class LocalDeliveryThread : public Thread {
    DispatchQueue *dq;
public:
    explicit LocalDeliveryThread(DispatchQueue *dq) : dq(dq) {}
    void *entry() override {
        dq->run_local_delivery();
        return 0;
    }
} local_delivery_thread;
```

```
void DispatchQueue::start()
{
    ceph_assert(!stop);
    ceph_assert(!dispatch_thread.is_started());
    // 启动两个线程
    dispatch_thread.create("ms_dispatch");
    local_delivery_thread.create("ms_local");
}
```

```
// listen_handler是一个EventCallbackRef对象，在构造中被初始化
Processor::Processor(AsyncMessenger *r, Worker *w, CephContext *c)
: mgr(r), net(c), worker(w),
  listen_handler(new C_processor_accept(this)) {}

class Processor::C_processor_accept : public EventCallback {
    Processor *pro;
public:
    explicit C_processor_accept(Processor *p): pro(p) {}
    void do_request(uint64_t id) override {
        pro->accept();
    }
}
```

还记得5.3中事件中心是怎么处理事件的吗？  
只要有事件到来就会执行回调函数对象中的do\_request

见5.6

```
void DispatchQueue::entry()
{
    while (true) {
        ...
        while (!lqueue.empty()) {
            ...
            const ref_t<Message> &m = qitem.get_message();
            if (stop) {
                ldout(cct, 10) << "stop flag set, discarding " << m << " " << *m << endl;
            } else {
                uint64_t msize = pre_dispatch(m);
                mgr->ms_deliver_dispatch(m);
                post_dispatch(m, msize);
            }
            if (stop)
                break;
        }
        cond.wait();
        ...
    }
}
```

```
void DispatchQueue::run_local_delivery()
{
    std::unique_lock l(local_delivery_lock);
    while (true) {
        if (stop_local_delivery)
            break;
        if (local_messages.empty()) {
            local_delivery_cond.wait();
            continue;
        }
        auto p = std::move(local_messages.front());
        local_messages.pop();
        l.unlock();
        const ref_t<Message> &m = p.first;
        int priority = p.second;
        fast_preprocess(m);
        if (can_fast_dispatch(m)) {
            fast_dispatch(m);
        } else {
            enqueue(m, priority, 0);
        }
        l.lock();
    }
}
```

```
void DispatchQueue::fast_dispatch(const ref_t<Message> &m)
{
    uint64_t msize = pre_dispatch(m);
    mgr->ms_fast_dispatch(m);
    post_dispatch(m, msize);
}
```

```
void ms_deliver_dispatch(const ceph::ref_t<Message> &m) {
    m->set_dispatch_stamp(ceph_clock_now());
    for (const auto &dispatcher : dispatchers) {
        if (dispatcher->ms_dispatch2(m))
            return;
    }
}
```

```
// Messenger.h
void ms_fast_dispatch(const ceph::ref_t<Message> &m) {
    m->set_dispatch_stamp(ceph_clock_now());
    for (const auto &dispatcher : fast_dispatchers) {
        if (dispatcher->ms_can_fast_dispatch2(m)) {
            dispatcher->ms_fast_dispatch2(m);
            return;
        }
    }
    ...
}
```

```
virtual bool ms_dispatch2(const MessageRef &m) {
    MessageRef mr(m);
    ms_dispatch(mr.get())
    return true;
}
virtual bool ms_dispatch(Message *m) {
    ceph_abort();
}
```

```
// Dispatcher.h
virtual void ms_fast_dispatch2(const MessageRef &m) {
    return ms_fast_dispatch(MessageRef(m).detach());
}
virtual void ms_fast_dispatch(Message *m) { ceph_abort(); }
```

OSD::ms\_dispatch

OSD::ms\_fast\_dispatch

## 5.7 发送消息



上层调用

```
int AsyncConnection::send_message(Message *m)
{
    if (!m->get_priority())
        m->set_priority(async_msgr->get_default_send_priority());

    m->get_header().src = async_msgr->get_myname();
    m->set_connection(this);

    protocol->send_message(m);
    return 0;
}
```

1. AsyncConnection::send\_message对消息做一层封装，然后使用protocol->send\_message(m);
2. protocol有两种实现：ProtocolV1 和 ProtocolV2，选择哪一个实现取决于建立连接时，对端支持的协议类型
3. 如果对端满足protocol->proto\_type == 2;那么使用V2，否则使用V1，如果是回环地址，也使用V1。
4. 最终都是通过connection->center->dispatch\_event\_external(connection->write\_handler);将消息递给事件中心，由事件中心去调度。

回调

回调函数是protocol->write\_event();  
此函数检查网络连接无误后调用 r = connection->\_try\_send();  
到此，就发送完成了

## 参考

[【ceph】 AsyncMessenger 网络模块总结--编辑中 - bdy - 博客园](#)

[【ceph】 CEPH源码解析:读写流程\\_ceph源码分析-CSDN博客](#)

[ceph的数据存储之路\(5\) -----osd数据处理 - 一只小江的个人空间 - OSCHINA - 中文开源技术交流社区](#)