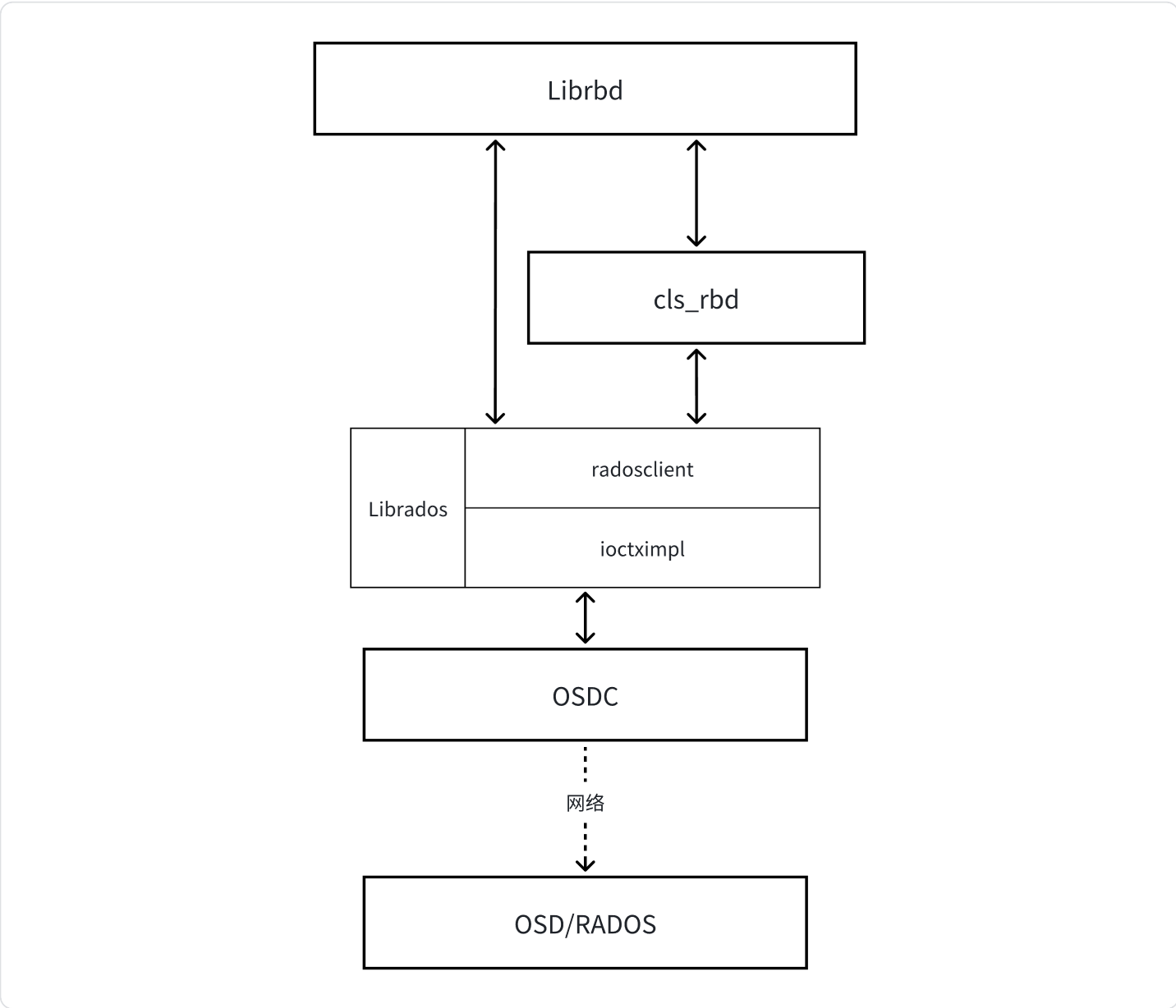


Librbd

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

RBD（RADOS Block Device）是由Ceph提供的存储结果，对于虚拟机来说，rbd设备就是本机的磁盘，对于ceph来说，rbd设备是分散在各个OSD节点的数据。

Librbd模块实现了RBD接口，基于librados实现对rbd的基本操作。



cls_rbd是一个扩展模块，实现了RBD的元数据相关操作，RBD的数据直接通过librados访问，最底层由OSDC完成数据的发送。

Cls

Cls是允许用户自定义对象的操作接口和实现方法，为用户提供了一种比较方便的接口扩展方式，比如cls_rbd就被rbd用来做元数据的操作。

类classHandler用来管理所有的扩展模块，函数register_class用来注册模块：

```
1
2 class ClassHandler {
3     // 注册的模块名 —> 模块元数据信息
4     std::map<std::string, ClassData> classes;
5
6     struct ClassData {
7         enum Status {
8             CLASS_UNKNOWN,
9             CLASS_MISSING,      // missing
10            CLASS_MISSING_DEPS, // missing dependencies
11            CLASS_INITIALIZING, // calling init() right now
12            CLASS_OPEN,         // initialized, usable
13        } status = CLASS_UNKNOWN;
14
15        std::string name;
16        ClassHandler *handler = nullptr;
17        void *handle = nullptr;
18
19        std::map<std::string, ClassMethod> methods_map; // 注册的方法
20        std::map<std::string, ClassFilter> filters_map; // 注册过滤方法
21
22        std::set<ClassData *> dependencies;           /* our dependencies */
23        std::set<ClassData *> missing_dependencies; /* only missing
24        dependencies */
25        };
26        struct ClassMethod {
27            const std::string name;
28            using func_t = std::variant<cls_method_cxx_call_t, cls_method_call_t>;
29            // 函数指针类型
30            func_t func; // 函数指针
31            int flags = 0;
32            ClassData *cls = nullptr; // 所属模块的ClassDate指针
33        };
34    };
```

注册一个模块

```
1 int cls_register(const char *name, cls_handle_t *handle) {
2     ClassHandler::ClassData *cls =
```

```

3         ClassHandler::get_instance().register_class(name);
4         *handle = (cls_handle_t)cls;
5         return (cls != NULL);
6     }

```

注册一个模块所属的函数

```

1 int cls_register_method(cls_handle_t hclass, const char *method,
2                         int flags,
3                         cls_method_call_t class_call, cls_method_handle_t
4                         *handle) {
5     if (!(flags & (CLS_METHOD_RD | CLS_METHOD_WR)))
6         return -EINVAL;
7     ClassHandler::ClassData *cls = (ClassHandler::ClassData *)hclass;
8     cls_method_handle_t hmethod = (cls_method_handle_t)cls->
9     >register_method(method, flags, class_call);
10    if (handle)
11        *handle = hmethod;
12    return (hmethod != NULL);
13 }

```

执行方法

```

1 int PrimaryLogPG::do_osd_ops(OpContext *ctx, vector<OSDOp> &ops):
2 case CEPH_OSD_OP_CALL: {
3     string cname, mname;
4     bufferlist indata;
5
6     bp.copy(op.cls.class_len, cname);
7     bp.copy(op.cls.method_len, mname);
8     bp.copy(op.cls.indata_len, indata);
9
10    ClassHandler::ClassData *cls;
11    result = ClassHandler::get_instance().open_class(cname, &cls);
12
13    ClassHandler::ClassMethod *method = cls->get_method(mname);
14
15    int flags = method->get_flags();
16    if (flags & CLS_METHOD_WR)
17        ctx->user_modify = true;
18
19    bufferlist outdata;
20

```

```

21     int prev_rd = ctx->num_read;
22     int prev_wr = ctx->num_write;
23     result = method->exec((cls_method_context_t)&ctx, indata, outdata);
24
25     op.extent.length = outdata.length();
26     osd_op.outdata.claim_append(outdata);
27 } break;

```

cls_rbd

cls_rbd.cc

注册了rbd模块和自定义的方法

```

1 cls_register("rbd", &h_class);
2 cls_register_cxx_method(h_class, "create",
3                          CLS_METHOD_RD | CLS_METHOD_WR,
4                          create, &h_create);

```

上面注册了create方法，其实现也在cls_rbd.cc中，这一类方法有三个参数：操作上下文，输入缓冲区，输出缓冲区。

```

1 int create(cls_method_context_t hctx, bufferlist *in, bufferlist *out)

```

注册完后，在Librados层的ioctx中会封装为CEPH_OSD_OP_CALL的操作

```

1 void create_image(librados::ObjectWriteOperation *op, uint64_t size,
2                  uint8_t order, uint64_t features,
3                  const std::string &object_prefix, int64_t data_pool_id)
4 {
5     bufferlist bl;
6     encode(size, bl);
7     encode(order, bl);
8     encode(features, bl);
9     encode(object_prefix, bl);
10    encode(data_pool_id, bl);
11
12    op->exec("rbd", "create", bl);
13 }
14
15 void librados::ObjectOperation::exec(const char *cls, const char *method,

```

```

16             bufferlist& inbl)
17 {
18     ceph_assert(impl);
19     ::ObjectOperation *o = &impl->o;
20     o->call(cls, method, inbl);
21 }

```

```

1 void call(const char *cname, const char *method, ceph::buffer::list &indata) {
2     add_call(CEPH_OSD_OP_CALL, cname, method, indata, NULL, NULL, NULL);
3 }
4
5 void add_call(int op, std::string_view cname, std::string_view method,
6               const ceph::buffer::list &indata,
7               fu2::unique_function<void(boost::system::error_code,
8                                         const ceph::buffer::list&) &&> f) {
9     OSDOp& osd_op = add_op(op);
10
11     set_handler([f = std::move(f)](boost::system::error_code ec,
12                                     int,
13                                     const ceph::buffer::list& bl) mutable {
14         std::move(f)(ec, bl);
15     });
16
17     osd_op.op.cls.class_len = cname.size();
18     osd_op.op.cls.method_len = method.size();
19     osd_op.op.cls.indata_len = indata.length();
20     osd_op.indata.append(cname.data(), osd_op.op.cls.class_len);
21     osd_op.indata.append(method.data(), osd_op.op.cls.method_len);
22     osd_op.indata.append(indata);
23 }
24
25

```

OSD收到操作后，通过exec接口调用具体的函数实现

```

1 int PrimaryLogPG::do_osd_ops(OpContext *ctx, vector<OSDOp> &ops):
2 case CEPH_OSD_OP_CALL:
3     result = method->exec((cls_method_context_t)&ctx, indata, outdata);

```

Librbd

使用示例, 完整版见examples/librbd/hello_word.cc

```
1 #include <iostream>
2 #include <rados/librados.hpp>
3 #include <rbd/librbd.hpp>
4 #include <sstream>
5 #include <string>
6
7 int main(int argc, const char **argv) {
8     ret = rados.connect();
9     // 创建pool
10     ret = rados.pool_create(pool_name);
11     ret = rados.ioctx_create(pool_name, io_ctx);
12
13     std::string name = "librbd_test";
14     uint64_t size = 2 << 20;
15     int order = 0;
16     librbd::RBD rbd;
17     librbd::Image image;
18     // 创建image
19     ret = rbd.create(io_ctx, name.c_str(), size, &order);
20     // 打开image
21     ret = rbd.open(io_ctx, image, name.c_str(), NULL);
22
23     int TEST_IO_SIZE = 512;
24     char test_data[TEST_IO_SIZE + 1];
25     int i;
26
27     for (i = 0; i < TEST_IO_SIZE; ++i) {
28         test_data[i] = (char)(rand() % (126 - 33) + 33);
29     }
30     test_data[TEST_IO_SIZE] = '\0';
31
32     size_t len = strlen(test_data);
33     ceph::bufferlist bl;
34     bl.append(test_data, len);
35
36     // 写数据
37     ret = image.write(0, len, bl);
38
39     ceph::bufferlist bl_r;
40     int read;
41     // 读数据
```

```

42     read = image.read(0, TEST_IO_SIZE, bl_r);
43
44     std::string bl_res(bl_r.c_str(), read);
45
46     int res = memcmp(bl_res.c_str(), test_data, TEST_IO_SIZE);
47 // 关闭image
48     image.close();
49
50     ret = rbd.remove(io_ctx, name.c_str());
51
52     ret = EXIT_SUCCESS;
53
54     int delete_ret = rados.pool_delete(pool_name);
55
56     rados.shutdown();
57
58     return ret;
59 }
60

```

创建RBD

internal.cc

librbd.cc中librbd::create有很多个版本，最终都是调用internal.cc下的create

```

1  int create(IoCtx &io_ctx, const std::string &image_name,
2             const std::string &image_id, uint64_t size,
3             ImageOptions &opts,
4             const std::string &non_primary_global_image_id,
5             const std::string &primary_mirror_uuid,
6             bool skip_mirror_enable) {
7 // 1. 如果image_id为空会生成一个新的id
8     std::string id(image_id);
9     if (id.empty()) {
10         id = util::generate_image_id(io_ctx);
11     }
12
13     CephContext *cct = (CephContext *)io_ctx.cct();
14     uint64_t option;
15 // 2. 不支持这两种选项
16     if (opts.get(RBD_IMAGE_OPTION_FLATTEN, &option) == 0) {
17         return -EINVAL;
18     }
19

```

```

19     if (opts.get(RBD_IMAGE_OPTION_CLONE_FORMAT, &option) == 0) {
20         return -EINVAL;
21     }
22
23 // 3. 获取image格式
24     uint64_t format;
25     if (opts.get(RBD_IMAGE_OPTION_FORMAT, &format) != 0)
26         format = cct->_conf.get_val<uint64_t>("rbd_default_format");
27     bool old_format = format == 1;
28
29 // 4. 检查image是否已存在, 如果已存在就返回错误
30     int r = detect_format(io_ctx, image_name, NULL, NULL);
31
32 // 5. 获取order, 决定了RBD对象的大小
33 // image的order表示其数据会被分成2^order个部分
34     uint64_t order = 0;
35     if (opts.get(RBD_IMAGE_OPTION_ORDER, &order) != 0 || order == 0) {
36         order = cct->_conf.get_val<uint64_t>("rbd_default_order");
37     }
38     // 验证order
39     r = image::CreateRequest<>::validate_order(cct, order);
40
41 // 6. 创建, 如果是旧格式, 且设置了RBD_FORCE_ALLOW_V1(略)使用create_v1
42     if (old_format) {
43         r = create_v1(io_ctx, image_name.c_str(), size, order);
44     } else {
45 // 6.1 否则使用AsioEngine创建image
46         AsioEngine asio_engine(io_ctx);
47
48         ConfigProxy config{cct->_conf};
49         api::Config<>::apply_pool_overrides(io_ctx, &config);
50
51         uint32_t create_flags = 0U;
52         uint64_t mirror_image_mode = RBD_MIRROR_IMAGE_MODE_JOURNAL;
53         if (skip_mirror_enable) {
54             create_flags = image::CREATE_FLAG_SKIP_MIRROR_ENABLE;
55         } else if (opts.get(RBD_IMAGE_OPTION_MIRROR_IMAGE_MODE,
56                             &mirror_image_mode) == 0) {
57             create_flags = image::CREATE_FLAG_FORCE_MIRROR_ENABLE;
58         }
59
60         C_SaferCond cond;
61         image::CreateRequest<> *req = image::CreateRequest<>::create(
62             config, io_ctx, image_name, id, size, opts, create_flags,
63             static_cast<cls::rbd::MirrorImageMode>(mirror_image_mode),
64             non_primary_global_image_id, primary_mirror_uuid,
65             asio_engine.get_work_queue(), &cond);

```



```

66         req->send();
67
68         r = cond.wait();
69     }
70
71     int r1 = opts.set(RBD_IMAGE_OPTION_ORDER, order);
72     ceph_assert(r1 == 0);
73     return r;
74 }
75

```

CreateRequest

```

1  template <typename ImageCtxT = ImageCtx>
2  class CreateRequest {
3  public:
4      static CreateRequest *create(const ConfigProxy& config, IoCtx &iocx,
5                                   const std::string &image_name,
6                                   const std::string &image_id, uint64_t size,
7                                   const ImageOptions &image_options,
8                                   uint32_t create_flags,
9                                   cls::rbd::MirrorImageMode mirror_image_mode,
10                                  const std::string &non_primary_global_image_id,
11                                  const std::string &primary_mirror_uuid,
12                                  asio::ContextWQ *op_work_queue,
13                                  Context *on_finish) {
14          return new CreateRequest(config, iocx, image_name, image_id, size,
15                                   image_options, create_flags,
16                                   mirror_image_mode, non_primary_global_image_id,
17                                   primary_mirror_uuid, op_work_queue, on_finish);
18      }
19
20      static int validate_order(CephContext *cct, uint8_t order);
21
22      void send();
23  }
24
25  template<typename I>
26  void CreateRequest<I>::send() {
27
28      int r = validate_features(m_cct, m_features);
29
30      r = validate_order(m_cct, m_order);
31
32      r = validate_striping(m_cct, m_order, m_stripe_unit, m_stripe_count);

```

```

33
34     if (((m_features & RBD_FEATURE_OBJECT_MAP) != 0) &&
35         (!validate_layout(m_cct, m_size, m_layout))) {
36         complete(-EINVAL);
37         return;
38     }
39
40     validate_data_pool();
41 }
42
43 template <typename I>
44 void CreateRequest<I>::validate_data_pool() {
45     m_data_io_ctx = m_io_ctx;
46     if ((m_features & RBD_FEATURE_DATA_POOL) != 0) {
47         librados::Rados rados(m_io_ctx);
48         int r = rados.ioctx_create(m_data_pool.c_str(), m_data_io_ctx);
49         m_data_pool_id = m_data_io_ctx.get_id();
50         m_data_io_ctx.set_namespace(m_io_ctx.get_namespace());
51     }
52
53     if (!m_config.get_val<bool>("rbd_validate_pool")) {
54         add_image_to_directory();
55         return;
56     }
57
58     auto ctx = create_context_callback<
59         CreateRequest<I>, &CreateRequest<I>::handle_validate_data_pool>(this);
60     auto req = ValidatePoolRequest<I>::create(m_data_io_ctx, ctx);
61     req->send();
62 }

```

ValidatePoolRequest

```

1 class ValidatePoolRequest {
2 public:
3     static ValidatePoolRequest* create(librados::IoCtx& io_ctx,
4                                         Context *on_finish) {
5         return new ValidatePoolRequest(io_ctx, on_finish);
6     }
7
8     ValidatePoolRequest(librados::IoCtx& io_ctx, Context *on_finish);
9
10    void send();
11 }
12

```

```

13 template <typename I>
14 void ValidatePoolRequest<I>::send() {
15     read_rbd_info();
16 }
17
18 template <typename I>
19 void ValidatePoolRequest<I>::read_rbd_info() {
20     ldout(m_cct, 5) << endl;
21
22     auto comp = create_rados_callback<
23         ValidatePoolRequest<I>,
24         &ValidatePoolRequest<I>::handle_read_rbd_info>(this);
25
26     librados::ObjectReadOperation op;
27     op.read(0, 0, nullptr, nullptr);
28
29     m_out_bl.clear();
30     int r = m_io_ctx.aio_operate(RBD_INFO, comp, &op, &m_out_bl);
31     ceph_assert(r == 0);
32     comp->release();
33 }

```

异步回调处理 AioCompletion

1. create_rados_callback

提供多种重载，应对不同的回调方式，可以看到均返回AioCompletion* 类型，在上一个函数中，调用的是第三个重载版本。

```

1 librados::AioCompletion *create_rados_callback(Context *on_finish) {
2     return create_rados_callback<Context, &Context::complete>(on_finish);
3 }
4
5 template <typename T>
6 librados::AioCompletion *create_rados_callback(T *obj) {
7     return librados::Rados::aio_create_completion(
8         obj, &detail::rados_callback<T>);
9 }
10
11 template <typename T, void(T::*MF)(int)>
12 librados::AioCompletion *create_rados_callback(T *obj) {
13     return librados::Rados::aio_create_completion(
14         obj, &detail::rados_callback<T, MF>);
15 }
16

```

```

17 // 这个版本带一个完成后删除对象的选项
18 template <typename T, Context*(T::*MF)(int*), bool destroy=true>
19 librados::AioCompletion *create_rados_callback(T *obj) {
20     return librados::Rados::aio_create_completion(
21         obj, &detail::rados_state_callback<T, MF, destroy>);
22 }

```

1.1 detail::rados_callback

这个函数有两个重载，而本质上都是在调用传入的函数

```

1 template <typename T>
2 void rados_callback(rados_completion_t c, void *arg) {
3     reinterpret_cast<T*>(arg)->complete(rados_aio_get_return_value(c));
4 }
5
6 template <typename T, void(T::*MF)(int)>
7 void rados_callback(rados_completion_t c, void *arg) {
8     T *obj = reinterpret_cast<T*>(arg);
9     int r = rados_aio_get_return_value(c);
10    (obj->*MF)(r);
11 }

```

`rados_state_callback` 和 `rados_callback` 不同的是它会根据执行结果和 `destroy` 标志选择删除对象

```

1 template <typename T, Context*(T::*MF)(int*), bool destroy>
2 void rados_state_callback(rados_completion_t c, void *arg) {
3     T *obj = reinterpret_cast<T*>(arg);
4     int r = rados_aio_get_return_value(c);
5     Context *on_finish = (obj->*MF)(r);
6     if (on_finish != nullptr) {
7         on_finish->complete(r);
8         if (destroy) {
9             delete obj;
10        }
11    }
12 }

```

1.2 aio_create_completion

传入参数和回调方法，返回 `AioCompletion` 对象

```

1 librados::AioCompletion *librados::Rados::aio_create_completion(void *cb_arg,
2                                                                    callback_t
   cb_complete)
3 {
4     AioCompletionImpl *c;
5     int r = rados_aio_create_completion2(cb_arg, cb_complete, (void**)&c);
6     ceph_assert(r == 0);
7     return new AioCompletion(c);
8 }
9
10 extern "C" int rados_aio_create_completion2(void *cb_arg,
11                                              rados_callback_t cb_complete,
12                                              rados_completion_t *pc)
13 {
14     librados::AioCompletionImpl *c = new librados::AioCompletionImpl;
15     if (cb_complete) {
16         c->set_complete_callback(cb_arg, cb_complete);
17     }
18     *pc = c;
19     return 0;
20 }
21
22 int set_complete_callback(void *cb_arg, rados_callback_t cb) {
23     std::scoped_lock l{lock};
24     callback_complete = cb;
25     callback_complete_arg = cb_arg;
26     return 0;
27 }

```

2. 回调类的发送

得到AioCompletion对象后，传入执行操作的函数

```

1 // AioCompletion* comp = funcxxx
2 // librados::ObjectReadOperation op
3 m_io_ctx.aio_operate(RBD_INFO, comp, &op, &m_out_bl);

```

```

1 int librados::IoCtx::aio_operate(const std::string& oid, AioCompletion *c,
2                                   librados::ObjectReadOperation *o,
3                                   bufferlist *pbl)
4 {
5     object_t obj(oid);
6     return io_ctx_impl->aio_operate_read(obj, &o->impl->o, c->pc, 0, pbl);

```

```

7 }
8
9 int librados::IoCtxImpl::aio_operate_read(const object_t &oid,
10                                           ::ObjectOperation *o,
11                                           AioCompletionImpl *c,
12                                           int flags,
13                                           bufferlist *pbl,
14                                           const blkin_trace_info *trace_info) {
15     Context *oncomplete = new C_aio_Complete(c);
16
17     c->is_read = true;
18     c->io = this;
19
20     Objecter::Op *objecter_op = objecter->prepare_read_op(
21         oid, oloc,
22         *o, snap_seq, pbl, flags | extra_op_flags,
23         oncomplete, &c->objver, nullptr, 0, &trace);
24
25     objecter->op_submit(objecter_op, &c->tid);
26     trace.event("rados operate read submitted");
27
28     return 0;
29 }

```

上面代码中能看到，AioCompletionImpl又被封装了一层，变成了C_aio_Complete。

下一步，oncomplete被塞到Op中：

```

1 Op *prepare_read_op(...) {
2     ...
3     // 这里，oncomplete摇身一变成了onack,参与Op的构造
4     Op *o = new Op(oid, oloc, std::move(op.ops), flags | global_op_flags |
5                   CEPH_OSD_FLAG_READ, onack, objver,
6                   data_offset, parent_trace);
7     ...
8     return o;
9 }
10
11 // onack在这里又换了一个名字，叫fin，然后fin赋值给了onfinish
12 Op(const object_t& o, const object_locator_t& ol, osdc_opvec&& _ops,
13     int f, Context* fin, version_t *ov, int *offset = nullptr,
14     ZTracer::Trace *parent_trace = nullptr) :
15     target(o, ol, f),
16     ops(std::move(_ops)),
17     out_bl(ops.size(), nullptr),
18     out_handler(ops.size()),

```

```

19     out_rval(ops.size(), nullptr),
20     out_ec(ops.size(), nullptr),
21     onfinish(fin),
22     objver(ov),
23     data_offset(offset) { ... }

```

AioCompletion全部都是按照上面类似的流程，封装到Op中并发送到服务端。

3. 执行回调*

当收到服务端响应后，网络模块将消息转到ms_dispatch

```

1 bool Objecter::ms_dispatch(Message *m) {
2 ...
3     switch (m->get_type()) {
4         case CEPH_MSG_OSD_OPREPLY:
5             handle_osd_op_reply(static_cast<MOSDOpReply *>(m));
6             return true;
7 ...
8 }

```

处理响应时，调用之前设置好的回调函数

```

1 void Objecter::handle_osd_op_reply(MOSDOpReply *m) {
2     if (op->has_completion()) {
3         num_in_flight--;
4         onfinish = std::move(op->onfinish);
5         op->onfinish = nullptr;
6     }
7     if (Op::has_completion(onfinish)) {
8         Op::complete(std::move(onfinish), osdcode(rc), rc);
9     }
10 }
11
12 // 这里用到了C++17的多个特性，不展开了
13 // 总之在当前情景，它实际调用了onfinish->complete(r)
14
15 static void complete(decltype(onfinish)&& f, boost::system::error_code ec,
16                     int r) {
17     std::visit([ec, r](auto&& arg) {
18         if constexpr (std::is_same_v<std::decay_t<decltype(arg)>,
19                             Context*>) {
20             arg->complete(r);

```

```

21         } else if constexpr (std::is_same_v<std::decay_t<decltype(arg)>,
22                               fu2::unique_function<OpSig>>) {
23             std::move(arg)(ec);
24         } else {
25             arg->defer(std::move(arg), ec);
26         }
27     }, std::move(f));
28 }

```

ofinish是C_aio_Complete类型，再来看

```

1 // C_aio_Completion并没有complete函数，在其父类Context中
2 virtual void complete(int r) {
3     finish(r);
4     delete this;
5 }
6 virtual void finish(int r) = 0;
7
8
9 // C_aio_Completion
10 void librados::IoCtxImpl::C_aio_Complete::finish(int r) {
11     c->lock.lock();
12     ...
13     c->complete = true;
14     c->cond.notify_all();
15
16 // 将CB_AioComplete加入到异步队列中，等待执行
17 // 再深入真太难看了，鬼知道队列中什么逻辑，这里当个黑盒先吧，反正是会调用CB_AioComplete
18     if (c->callback_complete ||
19         c->callback_safe) {
20         boost::asio::defer(c->io->client->finish_strand, CB_AioComplete(c));
21     }
22     ...
23     c->put_unlock();
24 }
25
26 // CB_AioComplete类中operator()执行的就是之前设置的回调函数
27 struct CB_AioComplete {
28     AioCompletionImpl *c;
29
30     explicit CB_AioComplete(AioCompletionImpl *cc) : c(cc) {
31         c->_get();
32     }
33
34     void operator()() {

```



```

35     rados_callback_t cb_complete = c->callback_complete;
36     void *cb_complete_arg = c->callback_complete_arg;
37     if (cb_complete)
38         cb_complete(c, cb_complete_arg);
39
40     rados_callback_t cb_safe = c->callback_safe;
41     void *cb_safe_arg = c->callback_safe_arg;
42     if (cb_safe)
43         cb_safe(c, cb_safe_arg);
44
45     c->lock.lock();
46     c->callback_complete = NULL;
47     c->callback_safe = NULL;
48     c->cond.notify_all();
49     c->put_unlock();
50 }
51 };

```

4. 整体流程

这个类就是这样，不断的执行一个操作并设置一个回调，直到finish

```

1  1. 构造函数 (ValidatePoolRequest<I>::ValidatePoolRequest):
2      - 创建一个ValidatePoolRequest类
3
4  2. send:
5      - 启动验证过程。调用read_rbd_info读取RBD信息。
6
7  3. read_rbd_info:
8      - 创建一个RADOS回调函数，设置回调 handle_read_rbd_info。
9      - 执行异步读取操作，读取RBD_INFO对象。
10
11 4. handle_read_rbd_info:
12     - 检查异步读取操作的返回值r。
13     - r>=0, 比较读取到的bufferlist (m_out_bl) 与预定义的字符串OVERWRITE_VALIDATED和
      VALIDATE。
14     - - m_out_bl == OVERWRITE_VALIDATED, 表示池已经被验证，直接调用 finish。
15     - - m_out_bl == VALIDATE, 表示快照已经成功创建，调用 overwrite_rbd_info。
16     - 如果r < 0 && r!= -ENOENT, 表示读取失败，调用finish并传递错误代码。
17     - 若以上条件不满足，调用create_snapshot来创建一个新的快照。
18
19 5. create_snapshot 和 handle_create_snapshot:
20     - 创建一个快照，设置回调 handle_create_snapshot
21     = 处理快照创建的结果。成功，调用 write_rbd_info
22

```

```

23 6. write_rbd_info 和 handle_write_rbd_info:
24   - 写新RBD的信息到RBD_INFO对象, 设置回调 handle_write_rbd_info
25   = 处理写操作的结果。成功, 调用 remove_snapshot
26
27 7. remove_snapshot 和 handle_remove_snapshot :
28   - 移除之前创建的验证快照, 设置回调 handle_remove_snapshot
29   = 处理快照移除的结果。成功, 调用 overwrite_rbd_info
30
31 8. overwrite_rbd_info 和 handle_overwrite_rbd_info :
32   - 覆盖RBD信息, 设置回调 handle_overwrite_rbd_info
33   = 处理覆盖操作的结果。无论是否成功, 都finish。
34
35 9. finish :
36   - delete this;
37   - m_on_finish->complete(r)

```

CreateRequest 续

ValidatePoolRequest::finish中的回调函数, 是之前CreateRequest<I>::validate_data_pool设置的。但是这里的机制又不一样了(真该死啊)。

create_context_callback

```

1 auto ctx = create_context_callback<
2     CreateRequest<I>, &CreateRequest<I>::handle_validate_data_pool>(this);
3 auto req = ValidatePoolRequest<I>::create(m_data_io_ctx, ctx);
4 req->send();

```

提供多种重载版本, 满足任何需求

```

1 // 基本回调适配器
2 template <typename T, void (T::*MF)(int) = &T::complete>
3 Context *create_context_callback(T *obj) {
4     return new detail::C_CallbackAdapter<T, MF>(obj);
5 }
6
7 // 状态回调适配器 提供删除对象参数
8 template <typename T, Context *(T::*MF)(int *), bool destroy = true>
9 Context *create_context_callback(T *obj) {
10     return new detail::C_StateCallbackAdapter<T, MF, destroy>(obj);
11 }
12
13 // 带引用计数的回调适配器

```

```

14 template <typename T, void (T::*MF)(int) = &T::complete>
15 Context *create_context_callback(T *obj, RefCountedPtr refptr) {
16     return new detail::C_RefCallbackAdapter<T, MF>(obj, refptr);
17 }
18
19 // 带引用计数的状态回调适配器
20 template <typename T, Context *(T::*MF)(int *)>
21 Context *create_context_callback(T *obj, RefCountedPtr refptr) {
22     return new detail::C_RefStateCallbackAdapter<T, MF>(obj, refptr);
23 }
24
25 // 下面两种太抽象了算了以后碰到再说
26 template <typename T, void (T::*MF)(int) = &T::complete, typename R>
27 typename std::enable_if<not std::is_base_of<RefCountedPtr, R>::value, Context
    *>::type
28 create_context_callback(T *obj, R *refptr) {
29     return new detail::C_CallbackAdapter<T, MF>(obj);
30 }
31
32 template <typename T, Context *(T::*MF)(int *), typename R, bool destroy =
    true>
33 typename std::enable_if<not std::is_base_of<RefCountedPtr, R>::value, Context
    *>::type
34 create_context_callback(T *obj, R *refptr) {
35     return new detail::C_StateCallbackAdapter<T, MF, destroy>(obj);
36 }

```

每一种模板返回的东西都完全不一样，在rbd_create场景中，使用的是第一个适配器

```

1 template <typename T, void (T::*MF)(int)>
2 class C_CallbackAdapter : public Context {
3     T *obj;
4
5 public:
6     C_CallbackAdapter(T *obj) : obj(obj) {
7     }
8
9 protected:
10     void finish(int r) override {
11         (obj->*MF)(r);
12     }
13 };
14 // 其它适配器虽然名字唬人，但其实做的事情类似，引用计数的就是多了个引用，状态的就是多考虑
    了下状态

```

这些适配器的父类是Context，之前提到它有一个虚函数和一个纯虚函数：

```
1 virtual void complete(int r) {
2     finish(r);
3     delete this;
4 }
5 virtual void finish(int r) = 0;
```

到这里context_callback的回调机制就分析完了，相比于rados_callback简单很多

complete

CreateRequest中有一个m_on_finish回调指针，由上层指定，在internal.cc的代码中，有这么一段

```
1 C_SaferCond cond;
2 image::CreateRequest<> *req = image::CreateRequest<>::create(
3     config, io_ctx, image_name, id, size, opts, create_flags,
4     static_cast<cls::rbd::MirrorImageMode>(mirror_image_mode),
5     non_primary_global_image_id, primary_mirror_uuid,
6     asio_engine.get_work_queue(), &cond);
7 req->send();
8
9 r = cond.wait();
```

m_on_finish被设置为cond，如果CreateRequest完成了操作，不管失败与否，都会调用complete函数

```
1 template<typename I>
2 void CreateRequest<I>::complete(int r) {
3     ldout(m_cct, 10) << "r=" << r << endl;
4
5     m_data_io_ctx.close();
6     auto on_finish = m_on_finish;
7     delete this;
8     on_finish->complete(r);
9 }
10
11 class C_SaferCond : public Context {
12     void complete(int r) override {
13         std::lock_guard l(lock);
14         done = true;
15         rval = r;
```

```
16     cond.notify_all();
17 }
18 }
```

整体流程

现在思路清晰了

1. CreateRequest执行操作，如果需要和外部交互的话，将操作交给其他类，同时设置回调 context_callback，自己就不管事了，等着内部类finish回调就好，如果失败，调用complete结束操作。
2. 要和外部交互的话，就设置rados_callback回调函数，等待异步回调，直到finish。
3. 这个类中，几乎全是按照这个模式走创建RBD的流程。

```
1  1. 构造函数 (CreateRequest<I>::CreateRequest):
2    - 设置配置、I/O 上下文、图像名称和ID、大小、图像选项、创建标志、镜像模式、
3    - 非主镜像全局ID、主镜像UUID、操作工作队列和回调。
4
5  2. send:
6    - 启动创建image的流程
7    - validate_features 检查请求的特性是否受支持
8    - validate_order 确保order在区间[12, 25]，一个条带对象大小为2^order
9    - validate_stripping 验证条带大小和数量是否有效
10   - validate_layout 验证图像大小是否与对象映射兼容
11   - validate_data_pool 验证pool，
12   - - 根据配置信息跳过验证pool，直接调用 add_image_to_directory
13   - - 走ValidatePoolRequest流程(见上文)
14   - - 设置context_callback类型回调，回调函数是 handle_validate_data_pool
15
16  3. handle_validate_data_pool:
17    - 处理pool验证结果。成功，调用 add_image_to_directory
18
19  4. add_image_to_directory 和 handle_add_image_to_directory:
20    - IoCtx写操作，将新image添加到RBD目录中。设置rados_callback类型的回调
21    = 验证结果。成功则调用 create_id_object
22
23  9. create_id_object 和 handle_create_id_object:
24    - IoCtx写操作，创建image的ID对象。设置rados_callback类型的回调
25    = 验证结果。成功则调用 negotiate_features
26
27  10. negotiate_features 和 handle_negotiate_features:
28    - IoCtx读操作，获取服务器的特性集，设置rados_callback类型回调
29    - 也可以根据配置，跳过这一步 直接调用 create_image
30    = 和服务器的特性集比较，不一致则更新，目的在于使用服务器支持的特性创建image
```

```

31     = 代码中没有例外情况，使用服务器支持的特性，调用 create_image
32
33 11. create_image 和 handle_create_image:
34     - IoCtx写操作，创建镜像，设置rados_callback类型回调
35     = 成功则调用 set_stripe_unit_count
36
37 12. set_stripe_unit_count 和 handle_set_stripe_unit_count:
38     - 如果unit count都没被设置或者都为默认值，调用 object_map_resize
39     - 否则需要设置，IoCtx写操作，在服务端改掉默认配置，设置rados_callback类型的回调
40     = 成功则调用 object_map_resize
41
42 13. object_map_resize 和 handle_object_map_resize:
43     - IoCtx写操作，调整object_map的大小，rados_callback
44     = 成功则 fetch_mirror_mode
45
46 14. fetch_mirror_mode 和 handle_fetch_mirror_mode:
47     - IoCtx读操作，获取镜像模式
48     = 成功 journal_create
49
50 15. journal_create 和 handle_journal_create:
51     - 创建一个新的日志请求，用于在RBD镜像中启用基于日志的快照和复制功能
52     - 创建context_callback类型回调，回调函数是 handle_journal_create
53     - 进入流程: librbd::journal::CreateRequest(见下文)
54     = 创建成功 mirror_image_enable
55
56 16. mirror_image_enable 和 handle_mirror_image_enable 回调:
57     - 启用镜像模式
58     - 创建context_callback类型回调，回调函数是 handle_mirror_image_enable
59     - 进入流程mirror::EnableRequest(见下文)
60     = 成功则 complete
61
62 17. complete:
63     - delete this
64     - on_finish->complete(r);
65
66 18. 清理操作
67 如果出现错误，会调用清理操作，从“高层”到“低”的顺序为：
68 journal_remove
69 remove_object_map
70 remove_header_object
71 remove_id_object
72 remove_from_dir

```

IoCtx写操作，在 [OSDC](#) 中介绍过

librbd::journal::CreateRequest

之前的部分已经非常详细的描述这类代码的运行框架，之后仅描述流程，不再涉及实现细节

```
1 1. 构造函数 (CreateRequest<I>::CreateRequest):
2   - 初始化I/O上下文、镜像ID、日志顺序、展开宽度、对象池名称、
3   - 标签类别、标签数据、客户端ID、操作工作队列和完成回调。
4
5 2. 发送请求 (CreateRequest<I>::send):
6   - 启动创建日志的流程。
7   - 确保 order在区间[12, 64], m_splay_width==0 // 这个参数没听过，据说是为了控制日志
   的分散程度，减少负载之类的，没具体了解
8   - 调用函数get_pool_id
9
10 3. get_pool_id
11   - 获取当前pool的id
12   - 调用 create_journal
13
14 4. create_journal 和 handle_create_journal
15   - 创建一个Journaler对象，用于管理镜像的日志，并启动日志创建过程
16   - 设置context_callback类型回调，回调函数是 handle_create_journal
17   - Journaler->create 使用IoCtx写操作，创建一个Journaler对象，并执行上层给的回调
18   = 若创建成功，调用 allocate_journal_tag
19
20 5. allocate_journal_tag 和 handle_register_client
21   - 为日志分配标签
22   - context_callback类型回调，函数为handle_register_client
23   - m_journaler->allocate_tag 同样将操作给到IoCtx层，完成后回调
24   = 成功则 register_client
25
26 6. register_client 和 handle_register_client
27   - 在日志中注册客户端，套路不变，m_journaler->register_client
28   = 成功，shut_down_journaler
29
30 7. shut_down_journaler 和 handle_journaler_shutdown
31   - 关闭Journaler对象，套路同上，m_journaler->shut_down
32   = 成功，complete(0)
33
34 8. complete
35   - m_on_finish->complete(r)
36   - delete this
37
38 9. remove_journal 和 handle_remove_journal
39   - 在创建好Journaler之后的所有失败处理，由这部分移除
40   - 创建RemoveRequest对象，进行移除操作
41   - complete(result)
```

这里主要是流程控制，具体创建操作是靠Journaler类实现，Journaler类又通过IoCtx执行写操作

mirror::EnableRequest

源码中给出了流程，那我就仔细分析了。

代码框架前面已经写几遍了，无非就是Op——>handle_Op——>NextOp的循环，直到finish调用上层回调。

这个流程的目的是启用镜像复制

```
1  <start>
2  |
3  v
4  GET_MIRROR_IMAGE * * * * * * *          获取镜像的复制状态
5  | * (on error)
6  v (skip if not needed) *
7  GET_TAG_OWNER * * * * * * *          获取标签所有者
8  | *
9  v (skip if not needed) *
10 OPEN_IMAGE *          打开镜像
11 | *
12 v (skip if not needed) *
13 CREATE_PRIMARY_SNAPSHOT * * *          创建初始快照
14 | *
15 v (skip if not opened) *
16 CLOSE_IMAGE *          关闭镜像
17 | *
18 v (skip if not needed) *
19 ENABLE_NON_PRIMARY_FEATURE *          启用非主特性，这个特性保证其只读
20 | *
21 v (skip if not needed) *
22 IMAGE_STATE_UPDATE * * * * * * *      更新镜像状态
23 | *
24 v *
25 <finish> < * * * * * * * * *
```

至此，rbd镜像的创建流程分析完了

Open

```
1  int RBD::open(IoCtx& io_ctx, Image& image, const char *name,
2               const char *snap_name)
3  {
```



```

4     ImageCtx *ictx = new ImageCtx(name, "", snap_name, io_ctx, false);
5
6     int r = ictx->state->open(0);
7
8     image.ctx = (image_ctx_t) ictx;
9
10    return 0;
11 }
12
13 struct ImageCtx {
14     ImageState<ImageCtx> *state;
15 }
16
17 template <typename I>
18 int ImageState<I>::open(uint64_t flags) {
19     C_SaferCond ctx;
20     open(flags, &ctx);
21
22     int r = ctx.wait();
23     return r;
24 }
25
26 template <typename I>
27 void ImageState<I>::open(uint64_t flags, Context *on_finish) {
28     CephContext *cct = m_image_ctx->cct;
29
30     m_lock.lock();
31     ceph_assert(m_state == STATE_UNINITIALIZED);
32     m_open_flags = flags;
33
34     Action action(ACTION_TYPE_OPEN);
35     action.refresh_seq = m_refresh_seq;
36     execute_action_unlock(action, on_finish);
37 }
38
39 template <typename I>
40 void ImageState<I>::send_open_unlock() {
41     ceph_assert(ceph_mutex_is_locked(m_lock));
42     CephContext *cct = m_image_ctx->cct;
43
44     m_state = STATE_OPENING;
45
46     Context *ctx = create_context_callback<ImageState<I>,
47         &ImageState<I>::handle_open>(this);
48     image::OpenRequest<I> *req =
49         image::OpenRequest<I>::create(m_image_ctx, m_open_flags, ctx);

```

```

50     m_lock.unlock();
51     req->send();
52 }

```

```

1  <start>
2  |
3  | (v1)
4  | -----> V1_DETECT_HEADER 尝试检测RBD镜像的v1格式头部
5  |         |
6  |         \-----\
7  | (v2)                                           |
8  \-----> V2_DETECT_HEADER                       /
9      检测v2格式头部                               |
10         |                                       |
11         v                                       |
12         获取镜像的ID或名称呼                               |
13         V2_GET_ID|NAME                               |
14         |                                       |
15         v (skip if have name)                               |
16         在正常位置找不到镜像名称，尝试在垃圾箱中检索                               |
17         V2_GET_NAME_FROM_TRASH                               |
18         |                                       |
19         v                                       |
20         获取镜像的初始元数据，如大小、对象前缀和特性。                               |
21         V2_GET_INITIAL_METADATA                               |
22         |                                       |
23         v                                       |
24         获取条带单元大小、数量                               |
25         V2_GET_STRIPE_UNIT_COUNT (skip if                               |
26         |                                       disabled) |
27         v                                       |
28         获取镜像创建时的时间戳                               |
29         V2_GET_CREATE_TIMESTAMP                               |
30         |                                       |
31         v                                       |
32         获取镜像的访问和修改时间戳                               |
33         V2_GET_ACCESS_MODIFY_TIMESTAMP                               |
34         |                                       |
35         v                                       |
36         获取镜像所在pool信息                               |
37         V2_GET_DATA_POOL -----> REFRESH 使用获取的元数据刷新镜像状态
38         |
39         v
40         初始化插件注册表
41         INIT_PLUGIN_REGISTRY

```

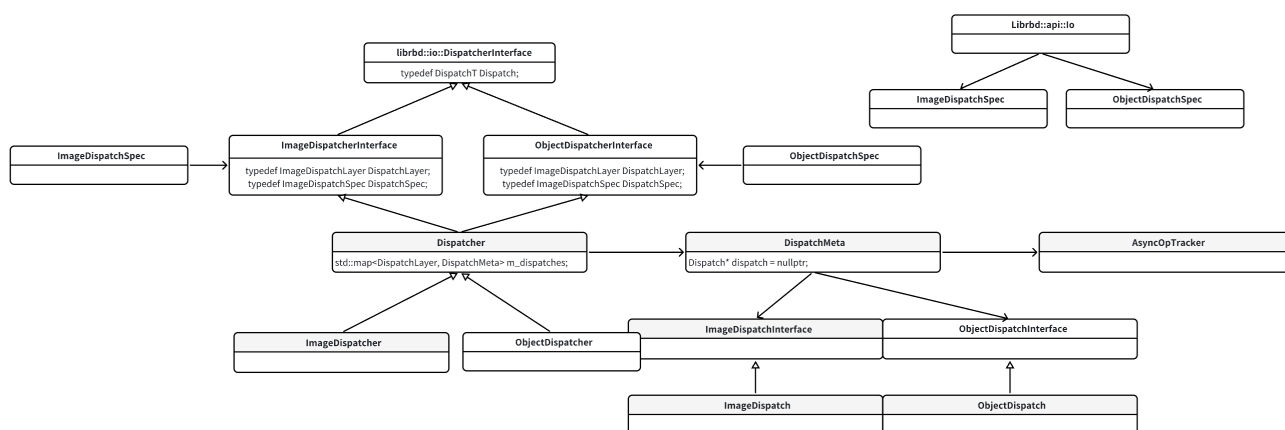
```

42      |
43      v
44      初始化缓存
45      INIT_CACHE
46      |
47      v
48      注册镜像更改监视
49      REGISTER_WATCH (skip if
50                      |      read-only)
51                      v
52                      设置快照, 没有指定快照就跳过
53      SET_SNAP (skip if no snap)
54      |
55      v
56      <finish>
57      ^
58      |
59      (on error)
60      * * * * * > CLOSE -----/

```

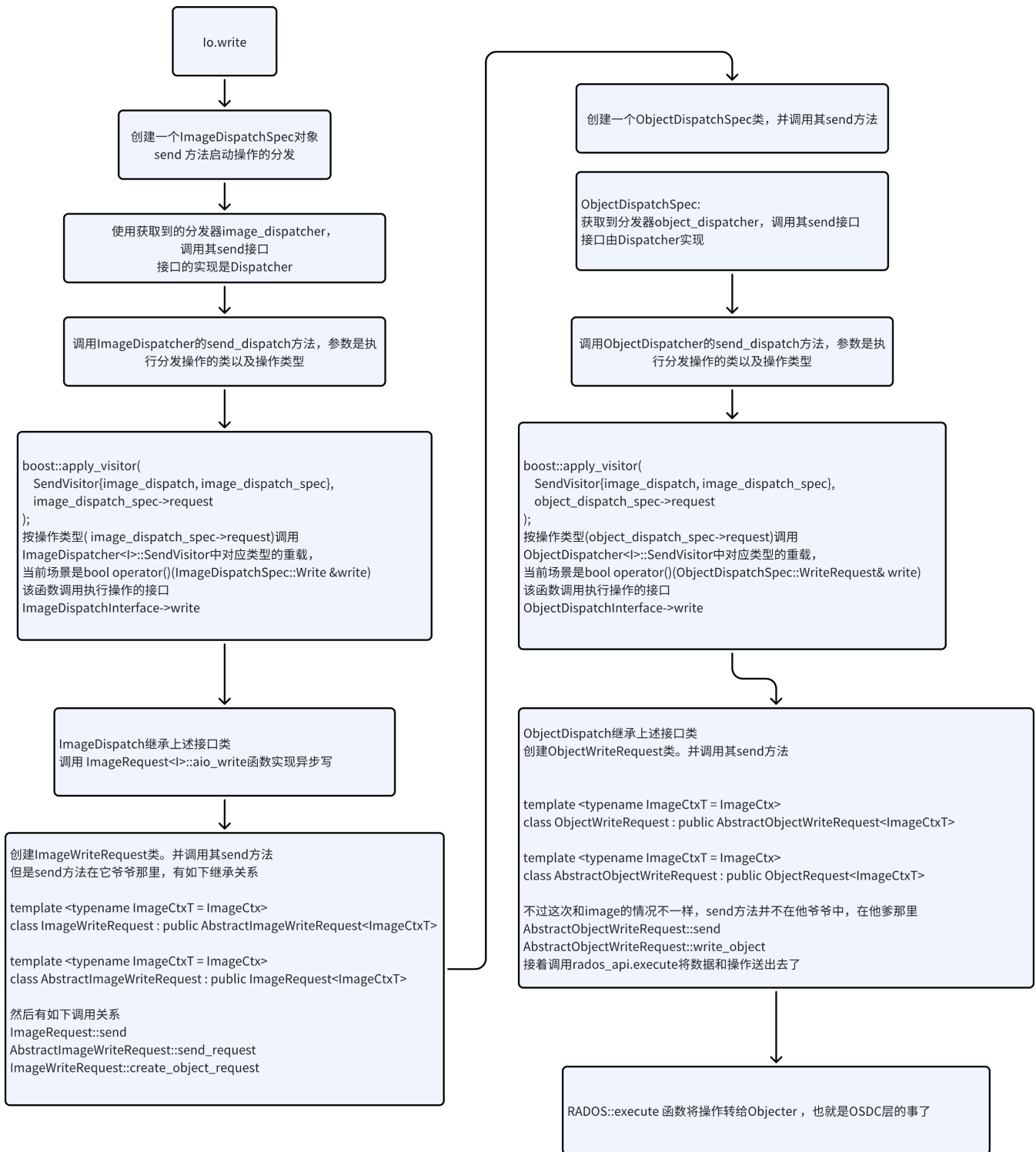
写数据

类图



总流程

站在顶层来看，分发的过程比较简单：



API调用

```

1 ssize_t Image::write(uint64_t ofs, size_t len, bufferlist &bl) {
2     ImageCtx *ictx = (ImageCtx *)ctx;
3
4     int r = api::Io<>::write(*ictx, ofs, len, bufferlist{bl}, 0);
5
6     return r;
7 }

```

API实现

```
1 template <typename ImageCtxT = ImageCtx>
2     struct Io {
3         static ssize_t read(ImageCtxT &image_ctx, uint64_t off, uint64_t len,
4                             io::ReadResult &&read_result, int op_flags);
5         static ssize_t write(ImageCtxT &image_ctx, uint64_t off, uint64_t len,
6                              bufferlist &&bl, int op_flags);
7         ...
8 } // namespace librbd::api
```

```
1 template <typename I>
2 ssize_t Io<I>::write(
3     I &image_ctx, uint64_t off, uint64_t len, bufferlist &&bl, int op_flags) {
4     auto cct = image_ctx.cct;
5
6     C_SaferCond ctx;
7     auto aio_comp = io::AioCompletion::create(&ctx);
8     aio_write(image_ctx, aio_comp, off, len, std::move(bl), op_flags, false);
9
10    r = ctx.wait();
11
12    return len;
13 }
14
15 template <typename I>
16 void Io<I>::aio_write(I &image_ctx, io::AioCompletion *aio_comp, uint64_t off,
17                      uint64_t len, bufferlist &&bl, int op_flags,
18                      bool native_async) {
19     auto req = io::ImageDispatchSpec::create_write(
20         image_ctx, io::IMAGE_DISPATCH_LAYER_API_START, aio_comp,
21         {{off, len}}, io::ImageArea::DATA, std::move(bl), op_flags, trace);
22     req->send();
23 }
```

ImageDispatchSpec

```
1 ImageDispatchSpec(ImageDispatcherInterface *image_dispatcher,
2                   ImageDispatchLayer image_dispatch_layer,
```

```

3         AioCompletion *aio_comp, Extents &&image_extents,
4         ImageArea area, Request &&request, IOContext io_context,
5         int op_flags, const ZTracer::Trace &parent_trace)
6     : dispatcher_ctx(this), image_dispatcher(image_dispatcher),
7       dispatch_layer(image_dispatch_layer), aio_comp(aio_comp),
8       image_extents(std::move(image_extents)), request(std::move(request)),
9       io_context(io_context), op_flags(op_flags), parent_trace(parent_trace) {}
10
11 // ImageDispatcherInterface->send
12 void ImageDispatchSpec::send() {
13     image_dispatcher->send(this);
14 }
15

```

Dispatcher

```

1 template <typename ImageCtxT, typename DispatchInterfaceT>
2 class Dispatcher : public DispatchInterfaceT {
3
4     void send(DispatchSpec *dispatch_spec) {
5         auto dispatch_layer = dispatch_spec->dispatch_layer;
6
7         while (true) {
8             m_lock.lock_shared();
9             dispatch_layer = dispatch_spec->dispatch_layer;
10
11 // std::map<DispatchLayer, DispatchMeta> m_dispatches;
12         auto it = m_dispatches.upper_bound(dispatch_layer);
13         if (it == m_dispatches.end()) {
14             // the request is complete if handled by all layers
15             dispatch_spec->dispatch_result = DISPATCH_RESULT_COMPLETE;
16             m_lock.unlock_shared();
17             break;
18         }
19
20         auto &dispatch_meta = it->second;
21         auto dispatch = dispatch_meta.dispatch;
22         auto async_op_tracker = dispatch_meta.async_op_tracker;
23         dispatch_spec->dispatch_result = DISPATCH_RESULT_INVALID;
24
25         // prevent recursive locking back into the dispatcher while
        handling IO
26         async_op_tracker->start_op();
27         m_lock.unlock_shared();
28

```

```

29         // advance to next layer in case we skip or continue
30         dispatch_spec->dispatch_layer = dispatch->get_dispatch_layer();
31
32 // 当前情景下 dispatch = ImageDispatchInterface
33 // dispatch_spec = ImageDispatchSpec::Write
34         bool handled = send_dispatch(dispatch, dispatch_spec);
35         async_op_tracker->finish_op();
36
37         // handled ops will resume when the dispatch ctx is invoked
38         if (handled) {
39             return;
40         }
41     }
42
43     // skipped through to the last layer
44     dispatch_spec->dispatcher_ctx.complete(0);
45 }
46 }

```

ImageDispatchSpec

用于处理分发

```

1 class ImageDispatchSpec {
2     // 设定各种操作, 比如read, write, CompareAndWrite, Flush等
3     struct Read {
4         ReadResult read_result;
5         int read_flags;
6
7         Read(ReadResult &&read_result, int read_flags)
8             : read_result(std::move(read_result)), read_flags(read_flags) {
9         }
10    };
11    struct Write {
12        bufferlist bl;
13
14        Write(bufferlist &&bl) : bl(std::move(bl)) {
15        }
16    };
17
18    // Request是全部可能操作的数据集合, 调度时, 根据数据类型分发操作
19    typedef boost::variant<Read,
20                          Discard,
21                          Write,
22                          WriteSame,

```

```

23             CompareAndWrite,
24             Flush,
25             ListSnaps> Request;
26
27 // 比如创建一个Write方法时, 会创建一个Write{std::move(bl)}, 并赋值给Request
28 template <typename ImageCtxT = ImageCtx>
29 static ImageDispatchSpec *create_write(
30     ImageCtxT &image_ctx, ImageDispatchLayer image_dispatch_layer,
31     AioCompletion *aio_comp, Extents &&image_extents, ImageArea area,
32     bufferlist &&bl, int op_flags, const ZTracer::Trace &parent_trace) {
33     return new ImageDispatchSpec(image_ctx.io_image_dispatcher,
34                                   image_dispatch_layer, aio_comp,
35                                   std::move(image_extents), area,
36                                   Write{std::move(bl)},
37                                   {}, op_flags, parent_trace);
38 }
39 }

```

封装好操作后, ImageDispatcherInterface->send(ImageDispatchSpec)

在这里面又有 send_dispatch(ImageDispatchInterface, ImageDispatchSpec)

apply_visitor有两个参数, 一个是SendVisitor, 一个是image_dispatch_spec->request

ImageDispatcher : public Dispatcher

```

1 template <typename I>
2 bool ImageDispatcher<I>::send_dispatch(
3     ImageDispatchInterface *image_dispatch,
4     ImageDispatchSpec *image_dispatch_spec) {
5     if (image_dispatch_spec->tid == 0) {
6         image_dispatch_spec->tid = ++m_next_tid;
7
8         bool finished = preprocess(image_dispatch_spec);
9         if (finished) {
10             return true;
11         }
12     }
13
14     return boost::apply_visitor(
15         SendVisitor{image_dispatch, image_dispatch_spec},
16         image_dispatch_spec->request);
17 }

```


SendVisitor 结构如下，它定义了很多个操作的转发方法，根据传入的request具体类型，由 SendVisitor调用对应的operatot()函数

```
1 template <typename I>
2 struct ImageDispatcher<I>::SendVisitor : public boost::static_visitor<bool> {
3     ImageDispatchInterface *image_dispatch;
4     ImageDispatchSpec *image_dispatch_spec;
5
6     SendVisitor(ImageDispatchInterface *image_dispatch,
7                 ImageDispatchSpec *image_dispatch_spec)
8         : image_dispatch(image_dispatch),
9           image_dispatch_spec(image_dispatch_spec) {}
10
11     bool operator()(ImageDispatchSpec::Read &read) const {
12         return image_dispatch->read(
13             image_dispatch_spec->aio_comp,
14             std::move(image_dispatch_spec->image_extents),
15             std::move(read.read_result), image_dispatch_spec->io_context,
16             image_dispatch_spec->op_flags, read.read_flags,
17             image_dispatch_spec->parent_trace, image_dispatch_spec->tid,
18             &image_dispatch_spec->image_dispatch_flags,
19             &image_dispatch_spec->dispatch_result,
20             &image_dispatch_spec->aio_comp->image_dispatcher_ctx,
21             &image_dispatch_spec->dispatcher_ctx);
22     }
23
24     bool operator()(ImageDispatchSpec::Write &write) const {
25         return image_dispatch->write(
26             image_dispatch_spec->aio_comp,
27             std::move(image_dispatch_spec->image_extents), std::move(write.bl),
28             image_dispatch_spec->op_flags, image_dispatch_spec->parent_trace,
29             image_dispatch_spec->tid, &image_dispatch_spec-
30             >image_dispatch_flags,
31             &image_dispatch_spec->dispatch_result,
32             &image_dispatch_spec->aio_comp->image_dispatcher_ctx,
33             &image_dispatch_spec->dispatcher_ctx);
34     };
35 }
```

这一步算调度完了，由具体的类执行操作

ImageDispatch : public ImageDispatchInterface

```
1 template <typename I>
```

```

2 bool ImageDispatch<I>::write(
3     AioCompletion *aio_comp, Extents &&image_extents, bufferlist &&bl,
4     int op_flags, const ZTracer::Trace &parent_trace,
5     uint64_t tid, std::atomic<uint32_t> *image_dispatch_flags,
6     DispatchResult *dispatch_result, Context **on_finish,
7     Context *on_dispatched) {
8     auto cct = m_image_ctx->cct;
9     auto area = get_area(image_dispatch_flags);
10
11     start_in_flight_io(aio_comp);
12
13     *dispatch_result = DISPATCH_RESULT_COMPLETE;
14     ImageRequest<I>::aio_write(m_image_ctx, aio_comp, std::move(image_extents),
15                               area, std::move(bl), op_flags, parent_trace);
16     return true;
17 }

```

ImageRequest

```

1 template <typename I>
2 void ImageRequest<I>::aio_write(I *ictx, AioCompletion *c,
3                                 Extents &&image_extents, ImageArea area,
4                                 bufferlist &&bl, int op_flags,
5                                 const ZTracer::Trace &parent_trace) {
6     ImageWriteRequest<I> req(*ictx, c, std::move(image_extents), area,
7                               std::move(bl), op_flags, parent_trace);
8     req.send();
9 }

```

ImageWriteRequest

```

1
2 template <typename ImageCtxT = ImageCtx>
3 class ImageWriteRequest : public AbstractImageWriteRequest<ImageCtxT> {
4 public:
5     ImageWriteRequest(ImageCtxT &image_ctx, AioCompletion *aio_comp,
6                     Extents &&image_extents, ImageArea area, bufferlist &&bl,
7                     int op_flags, const ZTracer::Trace &parent_trace)
8         : AbstractImageWriteRequest<ImageCtxT>(
9             image_ctx, aio_comp, std::move(image_extents), area,

```

```

10         "write", parent_trace),
11         m_bl(std::move(bl)), m_op_flags(op_flags) {
12     }
13 }
14
15 template <typename ImageCtxT = ImageCtx>
16 class AbstractImageWriteRequest : public ImageRequest<ImageCtxT> {}
17
18 template <typename I>
19 void ImageRequest<I>::send() {
20     I &image_ctx = this->m_image_ctx;
21     ceph_assert(m_aio_comp->is_initialized(get_aio_type()));
22     ceph_assert(m_aio_comp->is_started());
23
24     CephContext *cct = image_ctx.cct;
25     AioCompletion *aio_comp = this->m_aio_comp;
26
27     update_timestamp();
28     send_request();
29 }

```

AbstractImageWriteRequest

```

1  template <typename I>
2  void AbstractImageWriteRequest<I>::send_request() {
3      I &image_ctx = this->m_image_ctx;
4      ...
5      AioCompletion *aio_comp = this->m_aio_comp;
6
7      aio_comp->set_request_count(object_extents.size());
8
9      send_object_requests(object_extents, image_ctx.get_data_io_context(),
10 journal_tid);
11 ...
12 }
13
14 template <typename I>
15 void AbstractImageWriteRequest<I>::send_object_requests(
16     const LightweightObjectExtents &object_extents, IOContext io_context,
17     uint64_t journal_tid) {
18     I &image_ctx = this->m_image_ctx;
19     CephContext *cct = image_ctx.cct;
20
21     AioCompletion *aio_comp = this->m_aio_comp;
22     bool single_extent = (object_extents.size() == 1);

```

```

22     for (auto &oe : object_extents) {
23
24         C_AioRequest *req_comp = new C_AioRequest(aio_comp);
25         auto request = create_object_request(oe, io_context, journal_tid,
26                                             single_extent, req_comp);
27         request->send();
28     }
29 }
30
31 template <typename I>
32 ObjectDispatchSpec *ImageWriteRequest<I>::create_object_request(
33     const LightweightObjectExtent &object_extent, IOContext io_context,
34     uint64_t journal_tid, bool single_extent, Context *on_finish) {
35     I &image_ctx = this->m_image_ctx;
36
37     bufferlist bl;
38     if (single_extent && object_extent.buffer_extents.size() == 1 &&
39         m_bl.length() == object_extent.length) {
40         // optimization for single object/buffer extent writes
41         bl = std::move(m_bl);
42     } else {
43         assemble_extent(object_extent, &bl);
44     }
45
46     auto req = ObjectDispatchSpec::create_write(
47         &image_ctx, OBJECT_DISPATCH_LAYER_NONE, object_extent.object_no,
48         object_extent.offset, std::move(bl), io_context, m_op_flags, 0,
49         std::nullopt, journal_tid, this->m_trace, on_finish);
50     return req;
51 }

```

从写image转变为了写Object

ObjectDispatchSpec

```

1 template <typename ImageCtxT>
2 static ObjectDispatchSpec *create_write(
3     ImageCtxT *image_ctx, ObjectDispatchLayer object_dispatch_layer,
4     uint64_t object_no, uint64_t object_off, ceph::bufferlist &&data,
5     IOContext io_context, int op_flags, int write_flags,
6     std::optional<uint64_t> assert_version, uint64_t journal_tid,
7     const ZTracer::Trace &parent_trace, Context *on_finish) {
8     return new ObjectDispatchSpec(image_ctx->io_object_dispatcher,

```

```

9         object_dispatch_layer,
10        WriteRequest{object_no, object_off,
11                      std::move(data), write_flags,
12                      assert_version, journal_tid},
13        io_context, op_flags, parent_trace,
14        on_finish);
15 }

```

这一套的调度思想都是一致的，现在直接看结果：

```

1  template <typename I>
2  bool ObjectDispatch<I>::write(
3      uint64_t object_no, uint64_t object_off, ceph::bufferlist&& data,
4      IOContext io_context, int op_flags, int write_flags,
5      std::optional<uint64_t> assert_version,
6      const ZTracer::Trace &parent_trace, int* object_dispatch_flags,
7      uint64_t* journal_tid, DispatchResult* dispatch_result,
8      Context** on_finish, Context* on_dispatched) {
9      auto cct = m_image_ctx->cct;
10
11      *dispatch_result = DISPATCH_RESULT_COMPLETE;
12      auto req = new ObjectWriteRequest<I>(m_image_ctx, object_no, object_off,
13                                           std::move(data), io_context, op_flags,
14                                           write_flags, assert_version,
15                                           parent_trace, on_dispatched);
16      req->send();
17      return true;
18 }

```

```

1  template <typename I>
2  void AbstractObjectWriteRequest<I>::write_object() {
3      I *image_ctx = this->m_ictx;
4
5      neorados::WriteOp write_op;
6      if (m_copyup_enabled) {
7          if (m_guarding_migration_write) {
8              auto snap_seq = (this->m_io_context->write_snap_context() ?
9                              this->m_io_context->write_snap_context()->first : 0);
10
11              cls_client::assert_snapc_seq(
12                  &write_op, snap_seq, cls::rbd::ASSERT_SNAPC_SEQ_LE_SNAPSET_SEQ);

```

```

13     } else {
14         write_op.assert_exists();
15     }
16 }
17
18 add_write_hint(&write_op);
19 add_write_ops(&write_op);
20 ceph_assert(write_op.size() != 0);
21
22 image_ctx->rados_api.execute(
23     {data_object_name(this->m_ictx, this->m_object_no)},
24     *this->m_io_context, std::move(write_op),
25     librbd::asio::util::get_callback_adapter(
26         [this](int r) { handle_write_object(r); }), nullptr,
27     (this->m_trace.valid() ? this->m_trace.get_info() : nullptr));
28 }
29
30
31 void RADOS::execute(const Object& o, const IOContext& _ioc, WriteOp&& _op,
32                     std::unique_ptr<WriteOp::Completion> c, version_t* objver,
33                     const blkio_trace_info *trace_info) {
34     auto oid = reinterpret_cast<const object_t*>(&o.impl);
35     auto ioc = reinterpret_cast<const IOContextImpl*>(&ioc.impl);
36     auto op = reinterpret_cast<OpImpl*>(&_op.impl);
37     auto flags = op->op.flags | ioc->extra_op_flags;
38     ceph::real_time mtime;
39     if (op->mtime)
40         mtime = *op->mtime;
41     else
42         mtime = ceph::real_clock::now();
43
44     ZTracer::Trace trace;
45     if (trace_info) {
46         ZTracer::Trace parent_trace("", nullptr, trace_info);
47         trace.init("rados execute", &impl->objecter->trace_endpoint,
48             &parent_trace);
49     }
50     trace.event("init");
51     impl->objecter->mutate(
52         *oid, ioc->oloc, std::move(op->op), ioc->snapc,
53         mtime, flags,
54         std::move(c), objver, osd_reqid_t{}, &trace);
55     trace.event("submitted");
56 }

```

Objecter

到这一步，就由OSDC层接管了。 [📖OSDC](#)

```
1 void mutate(const object_t& oid, const object_locator_t& oloc,
2             ObjectOperation&& op, const SnapContext& snapc,
3             ceph::real_time mtime, int flags,
4             std::unique_ptr<Op::OpComp>&& oncommit,
5             version_t *objver = NULL, osd_reqid_t reqid = osd_reqid_t(),
6             ZTracer::Trace *parent_trace = nullptr) {
7     Op *o = new Op(oid, oloc, std::move(op.ops), flags | global_op_flags |
8                  CEPH_OSD_FLAG_WRITE, std::move(oncommit), objver,
9                  nullptr, parent_trace);
10    o->priority = op.priority;
11    o->mtime = mtime;
12    o->snapc = snapc;
13    o->out_bl.swap(op.out_bl);
14    o->out_handler.swap(op.out_handler);
15    o->out_rval.swap(op.out_rval);
16    o->out_ec.swap(op.out_ec);
17    o->reqid = reqid;
18    op.clear();
19    op_submit(o);
20 }
```