

# OSDC

参考文章：

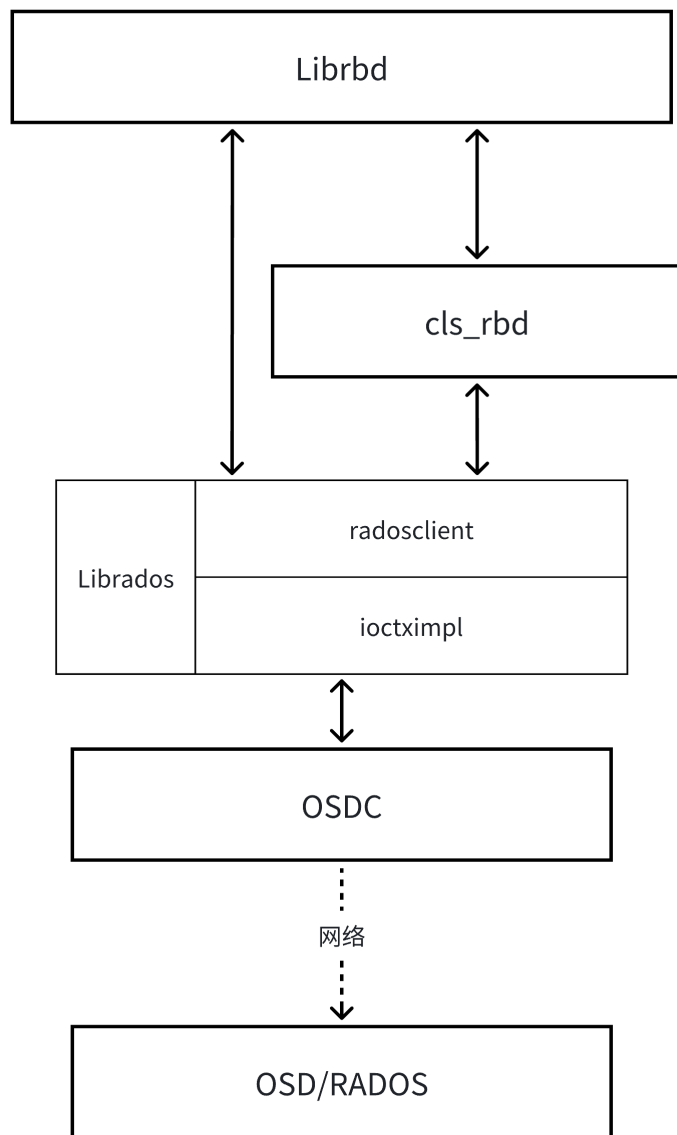
[Ceph OSDC源码分析 下](#)

[Ceph OSDC源码分析 上](#)

## 介绍

OSDC 是 osd client 模块的简称，主要用来和rados交互，这个模块里面完成了几个主要的功能：

1. 地址空间转换，从RBD（RADOS Block Device）的一维地址空间转换到对象的三维地址空间。
2. object\_cacher：一个object级别的缓存
3. Crush算法定位OSD：转换为三维地址空间之后，使用Crush算法进行对象的数据定位。



## 地址空间转换

来看看什么叫做一维转三维：

1. 对于一个文件，会被分成多个对象，每个对象的空间都是连续的
2. 一个缓存区，缓存这些正在读写的数据：
3. 对象和相应数据在磁盘内的排列方式，第一维是objectset，称为对象组，文件可以有多个对象组，第二维是stripe条带号，图中有三个条带，第三维是条带内偏移，指示对象在条带内的具体位置。可以看到，条带化后地址空间在逻辑上是连续的。

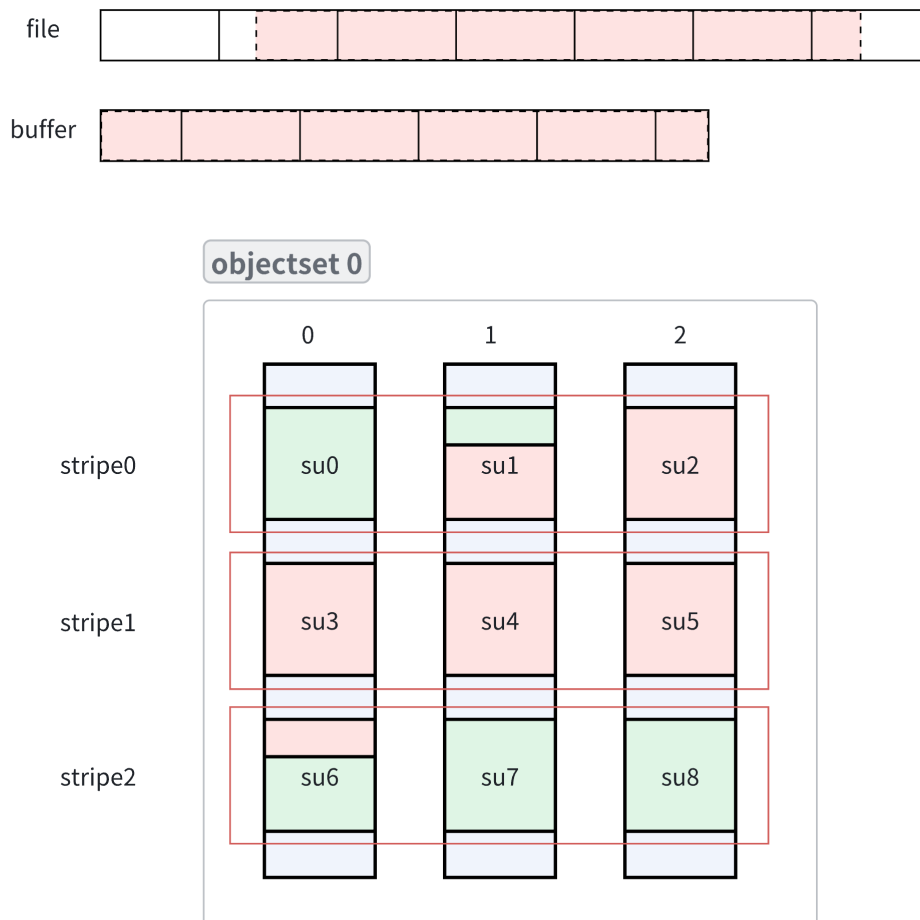


图1 地址空间转换

1. 灰色柱状代表一个rados底层对象，默认为4M
2. su是条带单元
3. 每一行是一个条带
4. object代表对象组，一般一个对象组属于同一个文件
5. object\_size: 对象的大小，灰色柱状部分，一般默认4M
6. stripe\_count 条带宽度，就是一个stripe跨多少个对象，也就是objectset中对象的个数，上面是3个对象
7. stripes\_per\_object, 一个对象的分片数，上图一个对象被分成了三份

## file\_to\_extents

上图中一个个条带连起来可以看成是逻辑上连续的、相当于线性的一维空间。file\_to\_extents函数是计算如何将一维坐标变成三维坐标。

```
1 void Striper::file_to_extents(
2     CephContext *cct, // Ceph 上下文指针
3     const file_layout_t *layout, // 指向文件布局信息的指针
```

```

4     uint64_t offset, // 文件中的起始偏移量
5     uint64_t len, // 要处理的长度
6     uint64_t trunc_size, // 截断大小
7     uint64_t buffer_offset, // 缓冲区偏移量
8     striper::LightweightObjectExtents* object_extents) // 指向对象范围集合的指针
9 {
10
11     __u32 object_size = layout->object_size; // 对象大小
12     __u32 su = layout->stripe_unit; // 条带单元大小
13     __u32 stripe_count = layout->stripe_count; // 条带中对象的个数
14     ceph_assert(object_size >= su); // 对象大小应大于等于条带单元大小
15
16     // 如果只有一条条带，条带单元大小等于对象大小
17     if (stripe_count == 1) {
18         su = object_size;
19     }
20     // 每个对象中的条带数量，如果不能整除，直接舍弃余数部分
21     uint64_t stripes_per_object = object_size / su;
22
23     uint64_t cur = offset; // 当前的偏移量
24     uint64_t left = len; // 要处理的长度
25     while (left > 0) {
26         // 将文件布局映射到对象
27         uint64_t blockno = cur / su; // 当前块号
28         uint64_t stripeno = blockno / stripe_count; // 条带号
29         uint64_t stripepos = blockno % stripe_count; // 条带内块的位置
30         uint64_t objectsetno = stripeno / stripes_per_object; // 对象组编号
31         uint64_t objectno = objectsetno * stripe_count + stripepos; // 对象号
32
33         // 映射块到对象中的位置
34         uint64_t block_start = (stripeno % stripe_count) * su;
35         uint64_t block_off = cur % su;
36         uint64_t max = su - block_off; // 最大可写入长度
37
38         uint64_t x_offset = block_start + block_off; // 对象内的偏移量
39         uint64_t x_len;
40         if (left > max) // 如果剩余长度大于最大可写入长度
41             x_len = max; // 设置写入长度为最大可写入长度
42         else
43             x_len = left; // 否则设置为剩余长度
44
45         striper::LightweightObjectExtent* ex = nullptr;
46         auto it = std::upper_bound(object_extents->begin(), object_extents->end(),
47                                   objectno, OrderByObject()); // 找到对象号的迭代器
48
49         striper::LightweightObjectExtents::reverse_iterator rev_it(it);
50         if (rev_it == object_extents->rend() || // 如果迭代器到达了末尾

```

```

51     rev_it->object_no != objectno || // 或者对象号不匹配
52     rev_it->offset + rev_it->length != x_offset) { // 或者偏移量和长度不匹配
53
54     ex = &(*object_extents->emplace(
55         it, objectno, x_offset, x_len,
56         object_truncate_size(cct, layout, objectno, trunc_size))); // 创建并添加
    新对象
57 } else {
58     ex = &(*rev_it); // 更新现有对象
59     ceph_assert(ex->offset + ex->length == x_offset); // 断言偏移量和长度匹配
60     ex->length += x_len; // 增加对象长度
61 }
62
63     ex->buffer_extents.emplace_back(cur - offset + buffer_offset, x_len); // 添
    加缓冲区范围
64     left -= x_len; // 减少剩余长度
65     cur += x_len; // 移动到下一个块
66 }
67 }
68 // 三维地址分片的信息全在object_extents中, 每个条目的结构如下:
69 // 每个条目保存一个对象内的分片信息
70 struct LightweightObjectExtent {
71     LightweightObjectExtent() = delete;
72     LightweightObjectExtent(uint64_t object_no, uint64_t offset,
73                             uint64_t length, uint64_t truncate_size)
74         : object_no(object_no), offset(offset), length(length),
75           truncate_size(truncate_size) {
76     }
77
78     uint64_t object_no; // 分片序号
79     uint64_t offset; // 对象内偏移
80     uint64_t length; // 长度
81     uint64_t truncate_size; // 截断。不知道截什么
82     LightweightBufferExtents buffer_extents; // 在buffer内的信息
83 };

```

对于图一，假设文件的offset是5M，长度是20M。文件的分片信息: stripe\_unit为4M，stripe\_count为3（三个对象），object\_size为12M。

转换后的结果就是：

```

1 LightweightObjectExtent[obj0] = {
2     objectno = 0
3     offset = 4M
4     length = 5M

```

```

5     buffer_extent = { [7, 4], [19, 1] }
6 }
7
8 LightweightObjectExtent[obj1] = {
9     objectno = 1
10    offset = 1M
11    length = 7M
12    buffer_extent = { [0, 3], [11, 4] }
13 }
14
15 LightweightObjectExtent[obj2] = {
16     objectno = 2
17     offset = 0M
18     length = 8M
19     buffer_extent = { [3, 4], [15, 4] }
20 }
21 每一次循环的值:
22 blockno: 1 stripeno: 0 stripepos: 1 objectsetno: 0 objectno: 1 block_start: 0
   block_off: 1 max: 3 x_offset: 1 x_len: 3
23 blockno: 2 stripeno: 0 stripepos: 2 objectsetno: 0 objectno: 2 block_start: 0
   block_off: 0 max: 4 x_offset: 0 x_len: 4
24 blockno: 3 stripeno: 1 stripepos: 0 objectsetno: 0 objectno: 0 block_start: 4
   block_off: 0 max: 4 x_offset: 4 x_len: 4
25 blockno: 4 stripeno: 1 stripepos: 1 objectsetno: 0 objectno: 1 block_start: 4
   block_off: 0 max: 4 x_offset: 4 x_len: 4
26 blockno: 5 stripeno: 1 stripepos: 2 objectsetno: 0 objectno: 2 block_start: 4
   block_off: 0 max: 4 x_offset: 4 x_len: 4
27 blockno: 6 stripeno: 2 stripepos: 0 objectsetno: 0 objectno: 0 block_start: 8
   block_off: 0 max: 4 x_offset: 8 x_len: 1

```

# 数据结构

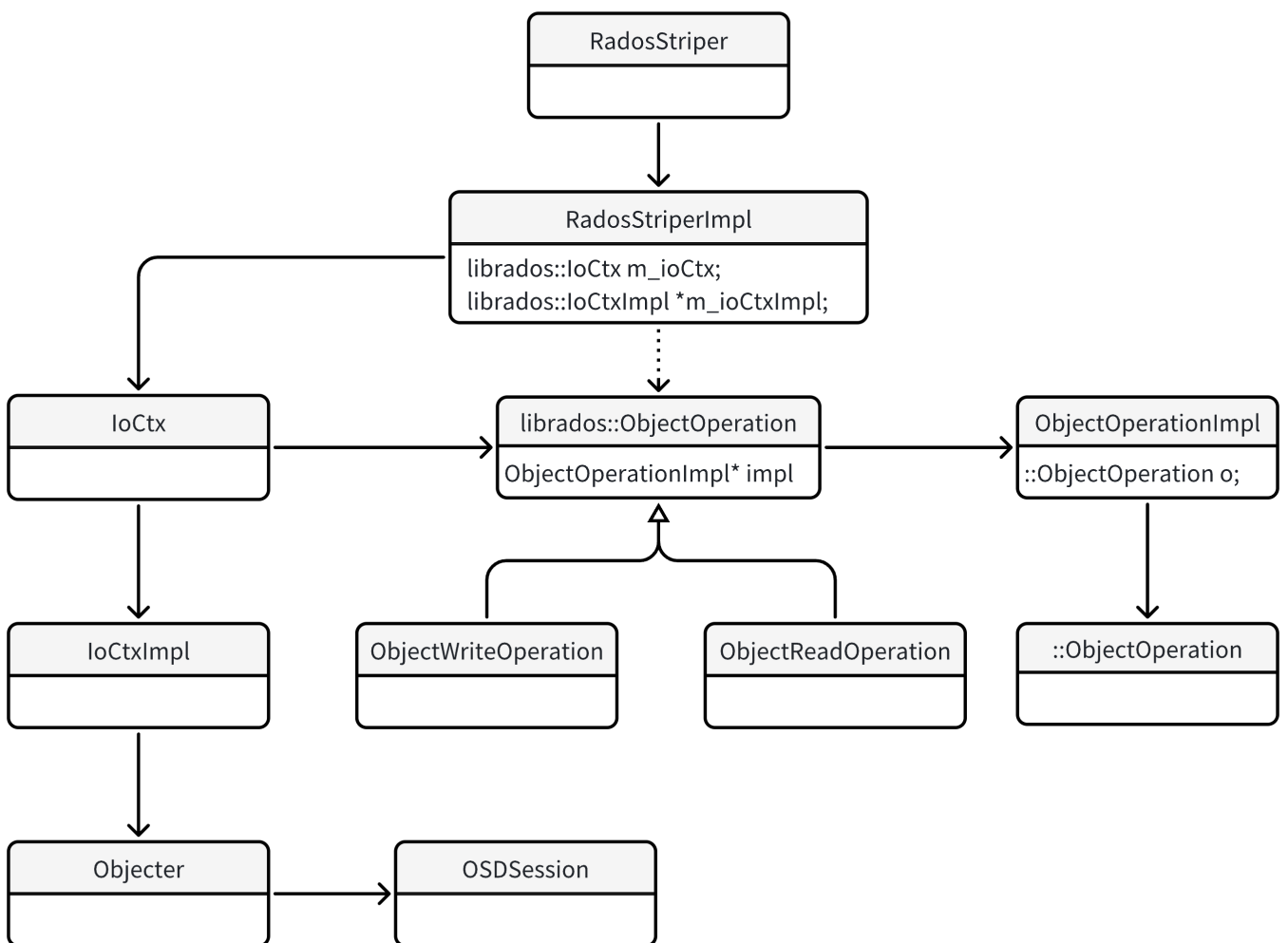
[illegible]

```

13 }
14
15 struct OSDOp {
16     ceph_osd_op op; // 各种操作码和参数
17     object_t soid; // 操作对象
18
19     ceph::buffer::list indata, outdata; // 输入输出缓冲区
20     errorcode32_t rval; // 操作结果
21 }

```

## 类图



## 写操作

```

1 // 创建IO上下文
2 rados_ioctx_create
3
4 rados_write_full
5 >> IoCtxImpl::write_full

```

```
6 >>>> IoCtxImpl::write
7 >>>>> IoCtxImpl::operate
8 >>>>>>> Objecter::Op::op_submit
9
10 Objecter::_op_submit
11 >> _calc_target
12 >> _get_session
13 >> _send_op
```

在 `hello_world.c.c` 中给出了使用librados库的接口的示例：

```
1 #include <rados/librados.h>
2 #include <stdio.h>
3
4 int main() {
5     rados_t rados;
6     rados_ioctx_t io_ctx;
7     const char *pool_name = "hello_world_pool";
8     const char* hello = "hello world!";
9     const char* object_name = "hello_object";
10
11     // 创建 RADOS 集群连接
12     rados_create(&rados, "admin");
13     rados_conf_parse_argv(rados, 0, NULL);
14     rados_connect(rados);
15
16     // 创建池
17     rados_pool_create(rados, pool_name);
18
19     // 创建 IO 上下文
20     rados_ioctx_create(rados, pool_name, &io_ctx);
21
22     // 直接写入对象
23     rados_write_full(io_ctx, object_name, hello, strlen(hello));
24
25     // 设置扩展属性
26     rados_setxattr(io_ctx, object_name, "version", "1", 2);
27
28     // 执行复合写操作，这里设置多个属性，然后一起提交(op_operate)
29     // 目的是确保原子性
30     rados_write_op_t write_op = rados_create_write_op();
31     rados_write_op_write_full(write_op, "v2", 4);
32     rados_write_op_setxattr(write_op, "version", "2", 2);
33     rados_write_op_operate(write_op, io_ctx, object_name, NULL, 0);
34
```



```

35     // 关闭资源
36     rados_ioctx_destroy(io_ctx);
37     rados_shutdown(rados);
38     return 0;
39 }

```

## **rados\_write\_full**

```

1  #define CEPH_RADOS_API
2  CEPH_RADOS_API int rados_write_full(rados_ioctx_t io, const char *oid,
3      const char *buf, size_t len);
4
5  #define LIBRADOS_C_API_DEFAULT_F(fn) fn
6  extern "C" int LIBRADOS_C_API_DEFAULT_F(rados_write_full)(
7      rados_ioctx_t io,
8      const char *o,
9      const char *buf,
10     size_t len)
11 {
12     librados::IoCtxImpl *ctx = (librados::IoCtxImpl *)io;
13     object_t oid(o);
14     bufferlist bl;
15     bl.append(buf, len);
16
17     int retval = ctx->write_full(oid, bl);
18     return retval;
19 }

```

## **IoCtxImpl**

实现单个pool层面的对象读写操作，主要是将请求封装成ObjectOperation，然后添加pool的地址信息，并封装为Object::Op，调用op\_submit发送给相应的OSD。

```

1  int librados::IoCtxImpl::write_full(const object_t& oid, bufferlist& bl)
2  {
3      ::ObjectOperation op;
4      op.write_full(bl);
5      return operate(oid, &op, NULL);
6  }
7
8  ::ObjectOperation
9  void write_full(ceph::buffer::list& bl) {
10     add_data(CEPH_OSD_OP_WRITEFULL, 0, bl.length(), bl);

```

```

11 }
12
13 int librados::IoCtxImpl::operate(const object_t& oid, ::ObjectOperation *o,
14                                 ceph::real_time *pmtime, int flags)
15 {
16 ...
17   Objecter::Op *objecter_op = objecter->prepare_mutate_op(
18     oid, oloc, // oloc中含有pool的信息
19     *o, snapc, ut,
20     flags | extra_op_flags,
21     oncommit, &ver);
22   objecter->op_submit(objecter_op);
23 ...
24 }

```

## 读操作

同写操作，从librados层开始调用，将刚刚写入的读出来

```

1  rados_aio_read(io_ctx, object_name, read_completion, read_buf, read_len, 0);

```

转而调用IoCtxImpl的方法

```

1  int librados::IoCtxImpl::aio_read(const object_t oid, AioCompletionImpl *c,
2                                     char *buf, size_t len, uint64_t off,
3                                     uint64_t snapid, const blkin_trace_info
4                                     *info)
5  {
6     FUNCTRACE(client->cct);
7     OID_EVENT_TRACE(oid.name.c_str(), "RADOS_READ_OP_BEGIN");
8     Context *oncomplete = new C_aio_Complete(c);
9
10    c->is_read = true;
11    c->io = this;
12    c->bl.clear();
13    c->bl.push_back(buffer::create_static(len, buf));
14    c->blp = &c->bl;
15    c->out_buf = buf;
16
17    Objecter::Op *o = objecter->prepare_read_op(
18      oid, oloc,

```

```

19     off, len, snapid, &c->bl, extra_op_flags,
20     oncomplete, &c->objver, nullptr, 0, &trace);
21     objecter->op_submit(o, &c->tid);
22     return 0;
23 }

```

## 条带写

```

1  #include "rados/librados.hpp"
2  #include "radosstriper/libradosstriper.hpp"
3  #include <iostream>
4  #include <string>
5
6  int main(int argc, char* argv[])
7  {
8      uint32_t strip_count = std::stoi(argv[1]);    // 一个条带包含的对象数量
9      uint32_t obj_size = std::stoi(argv[2]);      // 对象大小
10     std::string fname = argv[3];                 // 文件名
11     std::string obj_name = argv[4];              // 对象名
12     std::string pool_name = argv[5];             // 池名称
13
14     // IO上下文, ceph客户端, 条带化对象的智能指针
15     librados::IoCtx io_ctx;
16     librados::Rados cluster;
17     libradosstriper::RadosStriper* rs = new libradosstriper::RadosStriper;
18
19     // 加载配置
20     ret = cluster.init2("client.admin", "ceph", 0);
21     ret = cluster.conf_read_file("ceph.conf");
22     // 建立连接
23     ret = cluster.connect();
24     // 为指定的Pool创建一个IO上下文
25     ret = cluster.ioctx_create(pool_name.c_str(), io_ctx);
26     // 创建一个RadosStriper对象, 用于条带化写入
27     ret = libradosstriper::RadosStriper::striper_create(io_ctx, rs);
28     // 获取Pool的对齐要求
29     uint64_t alignment = 0;
30     ret = io_ctx.pool_required_alignment2(&alignment);
31
32     rs->set_object_layout_stripe_unit(alignment);    // su 条带大小
33     rs->set_object_layout_stripe_count(strip_count); // 条带中对象数量
34     rs->set_object_layout_object_size(obj_size);    // 对象大小
35
36     librados::bufferlist bl;

```

```

37     bl.read_file(fname.c_str());
38     // 写操作
39     rs->write_full(obj_name, bl);
40     std::cout << "done with: " << fname << std::endl;
41
42     delete rs;
43     io_ctx.close();
44     cluster.shutdown();
45 }

```

## write\_full

```

1  int libradosstriper::RadosStriperImpl::write_full(const std::string& soid,
2                                                    const bufferlist& bl)
3  {
4      return write(soid, bl, bl.length(), 0);
5  }
6
7  int libradosstriper::RadosStriperImpl::write(const std::string& soid,
8        const bufferlist& bl, size_t len, uint64_t off)
9  {
10     ceph_file_layout layout;
11     std::string lockCookie;
12     int rc = createAndOpenStripedObject(soid, &layout, len+off, &lockCookie,
13         true);
14     if (rc) return rc;
15     return write_in_open_object(soid, layout, lockCookie, bl, len, off);
16 }
17 int libradosstriper::RadosStriperImpl::createAndOpenStripedObject(const
18     std::string& soid,
19     ceph_file_layout *layout,
20                                     uint64_t
21     size,
22                                     std::string
23     *lockCookie,
24                                     bool
25     isFileSizeAbsolute)
26 {
27     // 创建写操作
28     librados::ObjectWriteOperation writeOp;
29     writeOp.create(true);
30
31     // object_size

```

```

28     std::ostringstream oss_object_size;
29     oss_object_size << m_layout.fl_object_size;
30     bufferlist bl_object_size;
31     bl_object_size.append(oss_object_size.str());
32     writeOp.setxattr(XATTR_LAYOUT_OBJECT_SIZE, bl_object_size);
33
34     // stripe unit
35     ...
36     writeOp.setxattr(XATTR_LAYOUT_STRIPE_UNIT, bl_stripe_unit);
37
38     // stripe count
39     ...
40     writeOp.setxattr(XATTR_LAYOUT_STRIPE_COUNT, bl_stripe_count);
41
42     // size
43     ...
44     writeOp.setxattr(XATTR_SIZE, bl_size);
45
46     // 执行原子写入操作, 尝试创建对象
47     std::string firstObjId = getObjectId(soid, 0);
48     int rc = m_ioCtx.operate(firstObjId, &writeOp); // 执行操作
49
50     // 对象已存在, 则返回错误
51     if (rc && -EEXIST != rc) return rc;
52
53     // 如果没有错误或者对象已存在, 继续打开条带化对象进行写入
54     uint64_t fileSize = size; // 设置文件大小
55     return openStripedObjectForWrite(soid, layout, &fileSize, lockCookie,
        isFileSizeAbsolute);
56 }

```

```

1  int libradosstriper::RadosStriperImpl::write_in_open_object(
2      const std::string& soid, // 条带化对象的ID
3      const ceph_file_layout& layout, // 文件布局信息
4      const std::string& lockCookie, // 锁定cookie
5      const bufferlist& bl, // 要写入的数据
6      size_t len, // 要写入的数据长度
7      uint64_t off) { // 写入操作的偏移量
8
9     ...
10    // 调用内部的异步写入API, 执行数据写入操作
11    int rc = internal_aio_write(soid, c, bl, len, off, layout);
12    if (!rc) {

```

```

13     // 如果写入操作成功，等待写入操作完成
14     c->wait_for_complete_and_cb();
15     // 等待数据安全确认
16     c->wait_for_safe_and_cb();
17     // 等待解锁操作完成
18     unlock_completion->wait_for_complete();
19     // 获取并返回写入操作的结果
20     rc = c->get_return_value();
21 }
22 return rc;
23 }

```

ObjectExtent和LightweightObjectExtent区别不大，就是多了几个字段，LightxxxExtent作为条带对象的中间变量，最终数据会转存在ObjectExtent中参与读写请求。

```

1 class ObjectExtent {
2 public:
3     object_t oid; // object id
4     uint64_t objectno;
5     uint64_t offset; // in object
6     uint64_t length; // in object
7     uint64_t truncate_size; // in object
8
9     object_locator_t oloc; // object locator (pool etc)
10
11     std::vector<std::pair<uint64_t, uint64_t>> buffer_extents; // off -> len.
    extents in buffer being mapped (may be fragmented bc of striping!)
12
13     ObjectExtent() : objectno(0), offset(0), length(0), truncate_size(0) {}
14     ObjectExtent(object_t o, uint64_t ono, uint64_t off, uint64_t l, uint64_t
    ts) : oid(o), objectno(ono), offset(off), length(l), truncate_size(ts) {}
15 };
16 }

```

```

1 int libradosstriper::RadosStriperImpl::internal_aio_write(const std::string
    &soid,
2
    libradosstriper::MultiAioCompletionImplPtr c,
3
    const bufferlist &bl,

```

```

4                                     size_t len,
5                                     uint64_t off,
6                                     const
ceph_file_layout &layout) {
7     int r = 0;
8     if (len > 0) {
9         vector<ObjectExtent> extents;
10        std::string format = soid;
11        boost::replace_all(format, "%", "%%");
12        format += RADOS_OBJECT_EXTENSION_FORMAT;
13        file_layout_t l;
14        l.from_legacy(layout);
15
16        // 分片
17        Striper::file_to_extents(cct(), format.c_str(), &l, off, len, 0,
extents);
18        // go through the extents
19        for (vector<ObjectExtent>::iterator p = extents.begin(); p !=
extents.end(); ++p) {
20            // assemble pieces of a given object into a single buffer list
21
22            // 将每个对象的buffer合并
23            bufferlist oid_bl;
24            for (vector<pair<uint64_t, uint64_t>::iterator q = p-
>buffer_extents.begin();
25                q != p->buffer_extents.end();
26                ++q) {
27                bufferlist buffer_bl;
28                buffer_bl.substr_of(bl, q->first, q->second);
29                oid_bl.append(buffer_bl);
30            }
31            // and write the object
32            // 写操作
33            r = m_ioCtx.aio_write(p->oid.name, rados_completion, oid_bl,
34                p->length, p->offset);
35        }
36    }
37    c->finish_adding_requests();
38    return r;
39 }

```

```

1 int librados::IoCtxImpl::aio_write(const object_t &oid, AioCompletionImpl *c,
2                                     const bufferlist &bl, size_t len,
3                                     uint64_t off, const blkin_trace_info *info)
{

```

```

4
5     c->io = this;
6     queue_aio_write(c);
7
8     Objecter::Op *o = objecter->prepare_write_op( // CEPH_OSD_OP_WRITE
9         oid, oloc,
10        off, len, snapc, bl, ut, extra_op_flags,
11        oncomplete, &c->objver, nullptr, 0, &trace);
12    objecter->op_submit(o, &c->tid);
13
14    return 0;
15 }

```

## 条带读

```

1  int libradosstriper::RadosStriper::read(const std::string& soid,
2                                          bufferlist* bl,
3                                          size_t len,
4                                          uint64_t off)
5  {
6      bl->clear();
7      bl->push_back(buffer::create(len));
8      return rados_striper_impl->read(soid, bl, len, off);
9  }
10
11 int libradosstriper::RadosStriperImpl::read(const std::string &soid,
12                                              bufferlist *bl,
13                                              size_t len,
14                                              uint64_t off) {
15     // create a completion object
16     librados::AioCompletionImpl c;
17
18     // call asynchronous method
19     int rc = aio_read(soid, &c, bl, len, off);
20
21     // and wait for completion
22     if (!rc) {
23         // wait for completion
24         c.wait_for_complete_and_cb();
25         // return result
26         rc = c.get_return_value();
27     }
28     return rc;
29 }

```



```

1  int libradosstriper::RadosStriperImpl::aio_read(const std::string &soid,
2
3
4
5
6      ceph_file_layout layout;
7      uint64_t size;
8      std::string lockCookie;
9  // 打开条带化对象
10     int rc = openStripedObjectForRead(soid, &layout, &size, &lockCookie);
11
12 // 计算实际可读取的字节数
13     uint64_t read_len;
14     if (off >= size) {
15         // nothing to read ! We are done.
16         read_len = 0;
17     } else {
18         read_len = std::min(len, (size_t)(size - off));
19     }
20
21 // 获取要读的数据列表
22     vector<ObjectExtent> *extents = new vector<ObjectExtent>();
23     ...
24     Striper::file_to_extents(cct(), format.c_str(), &l, off, read_len,
25                             0, *extents);
26     ...
27
28     int r = 0, i = 0;
29     for (vector<ObjectExtent>::iterator p = extents->begin(); p != extents-
30         >end(); ++p) {
31         bufferlist *oid_bl = &((*resultbl)[i++]);
32         for (vector<pair<uint64_t, uint64_t>>::iterator q = p-
33             >buffer_extents.begin();
34             q != p->buffer_extents.end();
35             ++q) {
36             bufferlist buffer_bl;
37             buffer_bl.substr_of(*bl, q->first, q->second);
38             oid_bl->append(buffer_bl);
39         }
40         ...
41 // 读取

```

```
41         r = m_ioCtx.aio_read(p->oid.name, rados_completion, oid_bl, p->length,
    p->offset);
42
43     }
44     nc->finish_adding_requests();
45     return r;
46 }
```

```

1  int librados::IoCtxImpl::aio_read(const object_t oid, AioCompletionImpl *c,
2                                     bufferlist *pbl, size_t len, uint64_t off,
3                                     uint64_t snapid, const blkin_trace_info
4                                     *info) {
5      FUNCTRACE(client->cct);
6      OID_EVENT_TRACE(oid.name.c_str(), "RADOS_READ_OP_BEGIN");
7      Context *oncomplete = new C_aio_Complete(c);
8
9      c->is_read = true;
10     c->io = this;
11     c->blp = pbl;
12
13     Objecter::Op *o = objecter->prepare_read_op(
14         oid, oloc,
15         off, len, snapid, pbl, extra_op_flags,
16         oncomplete, &c->objver, nullptr, 0, &trace);
17     objecter->op_submit(o, &c->tid);
18     return 0;
19 }

```

## Objecter::op\_submit

```
1 void Objecter::op_submit(Op *op, ceph_tid_t *ptid, int *ctx_budget)
2 {
3     _op_submit_with_budget(op, rl, ptid, ctx_budget);
4 }
5
6 void Objecter::_op_submit_with_budget(Op *op,
7                                     shunique_lock<ceph::shared_mutex>& sul,
8                                     ceph_tid_t *ptid,
9                                     int *ctx_budget)
```

```

10 {
11     ceph_assert(initialized);
12
13     ceph_assert(op->ops.size() == op->out_bl.size());
14     ceph_assert(op->ops.size() == op->out_rval.size());
15     ceph_assert(op->ops.size() == op->out_handler.size());
16
17     // 流量限制, 代码略
18
19     _op_submit(op, sul, ptid);
20 }

```

```

1 void Objecter::_op_submit(Op *op, shunique_lock<ceph::shared_mutex>& sul,
   ceph_tid_t *ptid)
2 {
3     OSDSession *s = NULL;
4
5     int r = _calc_target(&op->target, nullptr);
6
7     r = _get_session(op->target.osd, &s, sul);
8
9     _session_op_assign(s, op); // op->session = to;
10
11     _send_op(op);
12 }

```

## **\_calc\_target**

完成对象到OSD寻址的过程

```

1 int Objecter::_calc_target(op_target_t *t, Connection *con, bool any_change)
2 {
3     bool is_read = t->flags & CEPH_OSD_FLAG_READ;
4     bool is_write = t->flags & CEPH_OSD_FLAG_WRITE;
5     t->epoch = osdmap->get_epoch();
6     // 获取pool
7     const pg_pool_t *pi = osdmap->get_pg_pool(t->base_oloc.pool);
8     // 找到pool的位置
9     osdmap->object_locator_to_pg(t->target_oid, t->target_oloc, pgid);
10    // 获取PG对应的OSD列表
11    osdmap->pg_to_up_acting_osds(actual_pgid, &up, &up_primary,
12                                &acting, &acting_primary);
13    // 如果是读操作, 选择一个副本

```

```

14     if ((t->flags & (CEPH_OSD_FLAG_BALANCE_READS |
CEPH_OSD_FLAG_LOCALIZE_READS)) &&
15         !is_write && pi->is_replicated() && t->acting.size() > 1) {
16         int osd;
17         ceph_assert(is_read && t->acting[0] == acting_primary);
18
19         if (t->flags & CEPH_OSD_FLAG_BALANCE_READS) { // 随机选一个
20             int p = rand() % t->acting.size();
21             if (p) t->used_replica = true;
22             osd = t->acting[p];
23         } else { // 尽量选本地副本读取
24             int best = -1;
25             int best_locality = 0;
26             for (unsigned i = 0; i < t->acting.size(); ++i) {
27                 int locality = osdmap->crush->get_common_ancestor_distance(
28                     cct, t->acting[i], crush_location);
29                 if (i == 0 ||
30                     (locality >= 0 && best_locality >= 0 &&
31                      locality < best_locality) ||
32                     (best_locality < 0 && locality >= 0)) {
33                     best = i;
34                     best_locality = locality;
35                     if (i)
36                         t->used_replica = true;
37                 }
38             }
39             ceph_assert(best >= 0);
40             osd = t->acting[best];
41         }
42         t->osd = osd;
43     } else {
44         // 写操作, 选主OSD
45         t->osd = acting_primary;
46     }
47     return RECALC_OP_TARGET_NO_ACTION;
48 }

```

## \_get\_session

```

1 int Objecter::_get_session(int osd, OSDSession **session,
2                             shunique_lock<ceph::shared_mutex> &sul) {
3     ceph_assert(sul && sul.mutex() == &rwlock);
4
5     auto p = osd_sessions.find(osd);
6

```

```

7     auto s = new OSDSession(cct, osd);
8     osd_sessions[osd] = s;
9     s->con = messenger->connect_to_osd(osdmap->get_addrs(osd));
10    s->con->set_priv(RefCountedPtr{s});
11
12    s->get();
13    *session = s;
14
15    return 0;
16 }
17

```

## \_send\_op

```

1 void Objecter::_send_op(Op *op) {
2     // rwlock is locked
3     // op->session->lock is locked
4     ...
5     ceph_assert(op->tid > 0);
6     MOSDOp *m = _prepare_osd_op(op);
7
8     if (op->target.actual_pgid != m->get_spg()) {
9         m->set_spg(op->target.actual_pgid);
10        m->clear_payload(); // reencode
11    }
12
13    ConnectionRef con = op->session->con;
14    ceph_assert(con);
15
16    op->incarnation = op->session->incarnation;
17
18    if (op->trace.valid()) {
19        m->trace.init("op msg", nullptr, &op->trace);
20    }
21    op->session->con->send_message(m);
22 }
23

```

## 收消息

服务端发送消息出去后，会给到ms\_dispatch分发，根据不同的消息类型（如读写）给到不同的处理函数

```

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

```

这里不仔细看了，直接给出gpt的分析

```

1 void Objecter::handle_osd_op_reply(MOSDOpReply *m) {
2
3     获取操作ID：从传入的MOSDOpReply消息对象m中获取操作的唯一标识符tid。
4     获取会话信息：通过消息对象获取连接信息，并尝试获取与该连接关联的OSDSession对象。
5     查找操作对象：在OSDSession对象中查找与tid关联的Op对象。如果找不到或会话信息不匹配，
    释放消息对象并返回。
6     处理重试逻辑：如果操作是写操作，并且是第一次回复，可能会根据配置重试该操作。
7     处理重定向回复：如果消息是一个重定向回复，更新操作对象的定位器，并重新提交操作。
8     处理-EAGAIN错误：如果操作因为某些原因需要重新提交（例如，OSD请求重试），更新操作对象
    的定位器，并重新提交操作。
9     更新操作结果：如果操作成功，更新操作对象的版本号、回复的epoch和数据偏移量。
10    处理返回的数据：如果操作返回了数据，将其复制到操作对象指定的bufferlist中。
11    处理操作列表：从消息对象中获取返回的操作列表，并与原始请求的操作列表进行比较。
12    执行操作结果分发：遍历返回的操作列表，更新每个操作的结果、错误码和处理程序。
13    完成操作处理：如果操作对象有关联的完成回调，记录操作完成，减少正在进行的操作计数，并调
    用完成回调。
14    释放消息对象：最后，释放传入的MOSDOpReply消息对象。
15
16 }

```

## 注意

虽然大部分时候ioctx 都是直接调用的 IoCtxImpl的实现，且函数名都一模一样，但是也有少部分例外（有时候偷懒直接跳过IoCtx的实现，结果源码越来越难看），比如：

```

1 m_io_ctx.aio_operate(RBD_INFO, comp, &op, &m_out_bl);

```

这里ioctx调用的并不是对应Impl的aio\_operate，ioctx的这个函数有多个版本的重载，而并不是所有版本都会调用IoCtxImpl::aio\_operate，当前情况，如果op是ObjectReadOperation类型，应该是调用aio\_operate\_read

```
1 int librados::IoCtx::aio_operate(const std::string& oid, AioCompletion *c,  
2                                 librados::ObjectReadOperation *o,  
3                                 bufferlist *pbl)  
4 {  
5     if (unlikely(!o->impl))  
6         return -EINVAL;  
7     object_t obj(oid);  
8     return io_ctx_impl->aio_operate_read(obj, &o->impl->o, c->pc, 0, pbl);  
9 }
```