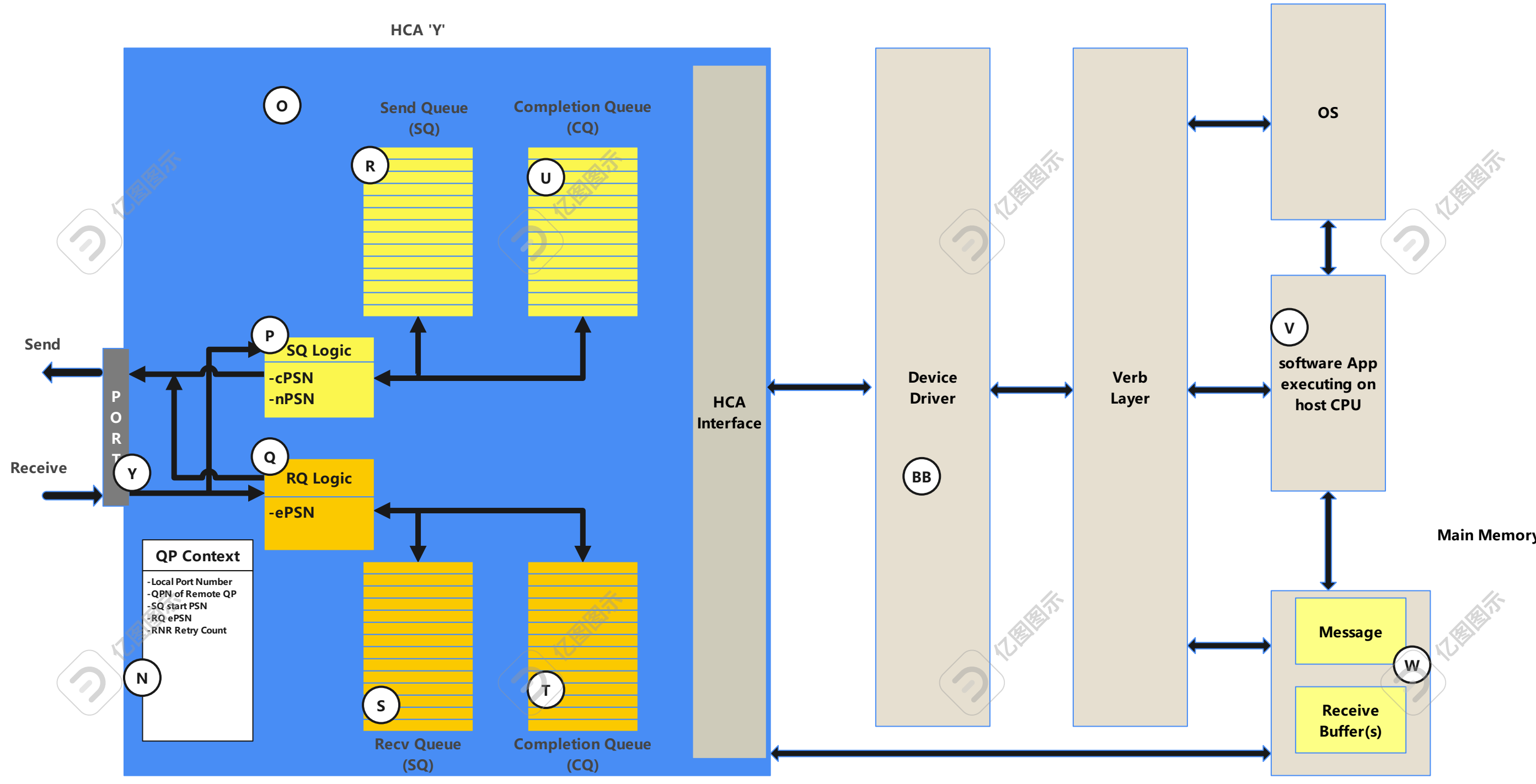
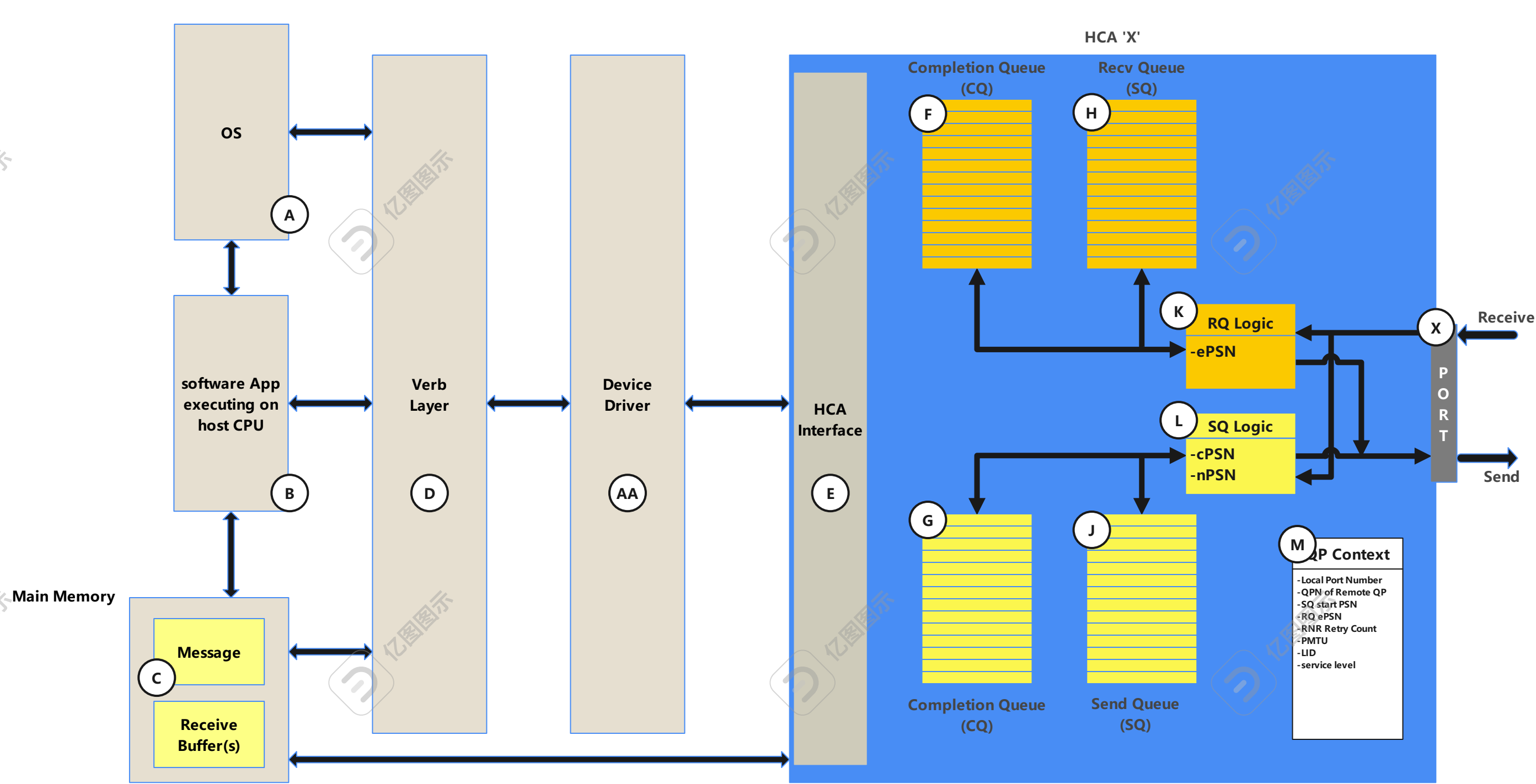






定性描述  
流程图示



step1:Responder端，分配一块特定大小内存，准备存储发过来的message。通过ibv\_post\_recv()调用，将WR post到RQ的next entry中，成为了WQE；  
step2:Requester端，软件调用ibv\_post\_send() 把一个WR post到SQ中，成为一个WQE，在WR中表明了本地的message的存放地址等信息；  
step3:Requester端，the QP's SQ Logic (item L) starts processing the top SQ entry, 决定是否将message 分片成多个message。Responder端，接收packet，进行PSN，opcode的校验，然后回应ack给Requester端，然后再进行写数据到Local memory，并更新WQE中的addr；  
step4:解释了Responder 端对Ack报文的处理；  
step5:中间片的处理逻辑，两端  
step6:尾片处理完成之后，Responder端就会产生CQE。  
step7:Requester端收到Ack之后，进行SQ中的top WQE的retire处理，然后生成CQE

RDMA 中如果想发送一个数据，那么：  
1.软件首先会生成一个 WQE (work queue element) ，就是工作队列里边的一个工作任务。  
2.然后这个任务再发一个 doorbell，就是按一个门铃到网卡告诉我有事情要做了。  
3.接着，网卡在收到这个门铃之后会从内存里面把这个工作任务取到网卡里面。  
4.然后再根据工作任务当中的地址，访问内存中的数据，把它 DMA 到网卡。  
5.再接下来，网卡会把这个数据封装成一个网络报文，从本地发送到远端。  
6.然后，接收端的网卡在收到了这个数据之后，再把它写到远端的内存。  
7.接着，接收端的网卡返回一个完成消息说我干完了。  
8.发起端的网卡收到了这个完成消息之后，它就在本地内存中生成一个CQE。  
9.最后，应用需要去 poll 这个CQE，也就是说它要获取这个完成队列里的完成事件才能够完成整个过程。

作者：李博杰  
链接：<https://www.zhihu.com/question/392146033/answer/3151675927>

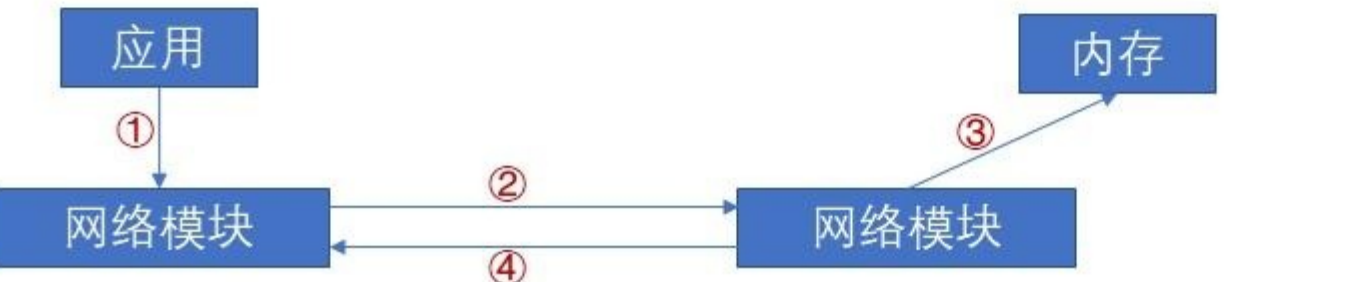
——当我把这个过程描述的时候其认为深度远远不够！

## CXL Load/Store：同步远端内存访问

- RDMA是异步远端内存访问，每次访问需要多次PCIe交互，时延最低也需要1.6 us

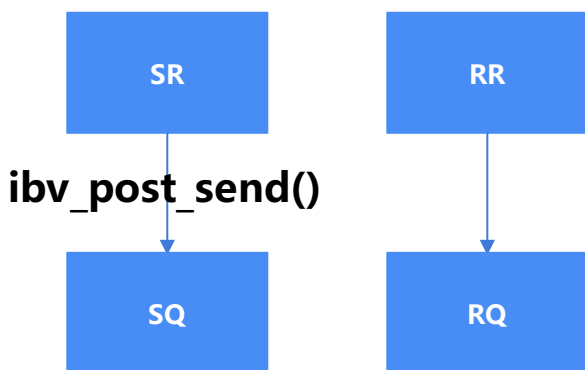


- CXL Load/Store是同步远端内存访问，CPU直出网络，指令直接访问远端内存，无需经过PCIe，无需WQE、CQE、doorbell开销，时延<0.5 us



知乎 @李博杰

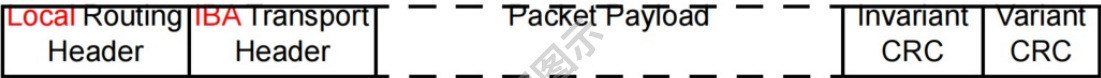
知乎参考





IBA 的包就分为以下几类

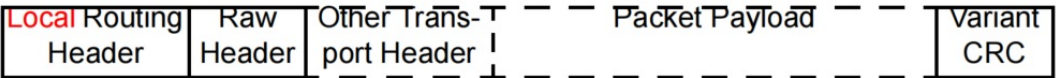
Local (within a subnet) Packets



Global (routing between subnets) Packets



Raw Packet with Raw Header



Raw Packet with IPv6 Header

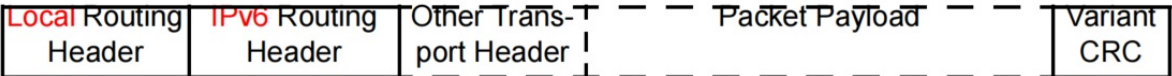


Figure 45 IBA Packet Overview

其中BTH比较重要:

Table 6 Base Transport Header Fields

Field Name	Field Abbreviation	Field Size (in bits)	Description
Opcode	OpCode	8	This field indicates the IBA packet type. The OpCode also specifies which extension headers follow the Base Transport Header
Solicited Event	SE	1	This bit indicates that an event should be generated by the responder.
MigReq	M	1	This bit is used to communicate migration state.
Pad Count	PadCnt	2	This field indicates how many extra bytes are added to the payload to align to a 4 byte boundary.
Transport Header Version	TVer	4	This field indicates the version of the IBA Transport Headers.
Partition Key	P_KEY	16	This field indicates which logical Partition is associated with this packet (see <a href="#">10.9 Partitioning on page 565</a> )
F/Res1*	F/R	1	F (FECN): 0 indicates that a FECN indication was not received. 1 indicates that the packet went through a point of congestion. Res1*: Transmitted as 0, ignored on receive. This field is not included in the invariant CRC. see <a href="#">7.8 CRCs on page 220</a> for details.
B/Res1*	B/R	1	B (BECN): 0 the packet did not go through a point of congestion or went through a point of congestion but was not marked. 1 indicates that the packet indicated by this header was subject to forward congestion. The B bit is set in an ACK or CN BTH. Res1*: Transmitted as 0, ignored on receive. This field is not included in the invariant CRC. see <a href="#">7.8 CRCs on page 220</a> for details.
Reserved (variant)		6	Transmitted as 0, ignored on receive. This field is not included in the invariant CRC. see <a href="#">7.8 CRCs on page 220</a> for details.
Destination QP	DestQP	24	This field indicates the Work Queue Pair Number (a.k.a. QP) at the destination
Acknowledge Request	A	1	This bit is used to indicate that an acknowledge (for this packet) should be scheduled by the responder.
Reserved		7	Transmitted as 0, ignored on receive. This field is included in the invariant CRC.
Field Name	Field Abbreviation	Field Size (in bits)	Description
Packet Sequence Number	PSN	24	This field is used to detect a missing or duplicate Packet. See <a href="#">9.7.1 Packet Sequence Numbers (PSN) on page 314</a> for a detailed description of PSN.

The Previous Part

Part 1 introduced basic concepts and terminology and consisted of the following chapters:

- Chapter 1—Basic Terms and Concepts.
- Chapter 2—Intro to Attributes and Managers.
- Chapter 3—QP: Message Transfer Mechanism.
- Chapter 4—Intro to Transport Types.
- Chapter 5—Intro to Send/Receive Operations.
- Chapter 6—Division of Labor.
- Chapter 7—Subnet-Local Addressing.
- Chapter 8—Global Addressing.
- Chapter 9—Intro to the Managers.
- Chapter 10—Intro to Connection Establishment.
- Chapter 11—PSN Usage.







CM建链的底层仍然调用的是ibv\_modify\_qp()来改变QP的状态, client与server之间传递的是UD报文, 通过QP1进行传递! 那么需要梳理清楚的就是哪里传递的信息, 哪里迁移的状态! CM建链有自己的一套!

务必要先启动server, setup的时候进行rdma\_bind\_addr(), rdma\_listen, 对client进行address resolve, 然后才启动connect, server收到了request之后才往下运行, 后面二者才能完成connect established!event的确认!

可选项点  
1. 程序中的rdma\_ack\_cm\_event()可以在最后执行!

rdma\_getaddrinfo provides transport independent address translation. It resolves the destination node and service address and returns information required to establish device communication. It is the functional equivalent of getaddrinfo.

rdma\_bind\_addr associates a source address with an rdma\_cm\_id. The address may be wildcarded. If binding to a specific local address, the rdma\_cm\_id will also be bound to a local RDMA device. #source address #rdma\_cm\_id连接起来!

rdma\_listen initiates a listen for incoming connection requests or datagram service lookup. The listen is restricted to the locally bound source address. 查看对端有没有来请求信息

Server端执行到这里的时候就会卡在这等等, 等待rdma\_get\_cm\_event()执行返回, 此时需要client端来执行rdma\_connect()的操作, 从而在这边的Server端才能产生CONNECT\_REQUEST, 这个函数才能执行!

wait for client

Prepare a qp when received a connect request from cma

ibv\_modify\_qp(rtr) ibv\_modify\_qp(rts)

执行到这里也要等待, 只有当Client端进行了ack的释放操作, 这里才能得到event, 执行完毕进行返回!

accept

Client端必须先发送connect, 发送一个ConnectRequest报文

返回一个ConnectReply报文

发送一个ReadyToUse报文

Prepare a qp to connect to remote side

No.	Time	Source	Destination	Protocol	Length	Info
1	0.000000	10.0.0.5	10.0.0.4	RRoCE	322	CM: ConnectRequest
2	0.017853	10.0.0.4	10.0.0.5	RRoCE	322	CM: ConnectReply
3	0.018281	10.0.0.5	10.0.0.4	RRoCE	322	CM: ReadyToUse
4	0.018286	10.0.0.5	10.0.0.4	RRoCE	130	RC: RDMA Write Only Immediate QP=0x0004e9
5	0.018287	10.0.0.5	10.0.0.4	RRoCE	94	RC: RDMA Write Only Immediate QP=0x0004e9
6	0.018509	10.0.0.4	10.0.0.5	RRoCE	94	RC: RDMA Write Only Immediate QP=0x000003
7	0.018514	10.0.0.5	10.0.0.4	RRoCE	62	RC: Acknowledge QP=0x0004e9
8	0.018670	10.0.0.4	10.0.0.5	RRoCE	62	RC: Acknowledge QP=0x000003

rdma\_ack\_cm\_event frees a communication event. All events which are allocated by rdma\_get\_cm\_event must be released, there should be a one-to-one correspondence between successful gets and acks. This call frees the event structure and any memory that it references.

test\_send\_recv(&conn)

ibv\_req\_notify\_cq(conn->cq, 0)

这里的wait\_cq()等待一下, 等待client发送完成之后, 这里才能执行结果, 输出结果为: wr\_id=0 opcode=0x00 (recv) status=0

wait\_cq(conn, total\_loop)

verify\_rcv\_buffer(conn, i, total\_loop)

ibv\_post\_recv(conn->cm\_id->qp, &conn->recv\_wr[i], &failed\_recv\_wr)

ibv\_post\_send(conn->cm\_id->qp, &conn->send\_wr[i], &failed\_send\_wr)

wait\_cq(conn, total\_loop)

wait\_cq(conn, total\_loop)

verify\_rcv\_buffer(conn, i, total\_loop)

ibv\_post\_recv(conn->cm\_id->qp, &conn->recv\_wr[i], &failed\_recv\_wr)

ibv\_post\_send(conn->cm\_id->qp, &conn->send\_wr[i], &failed\_send\_wr)

wait\_cq(conn, total\_loop)

wait\_cq(conn, total\_loop)

共循环100次

发送数据要进行一次标志判断和数据填充, 如果(flag\_verify) 用\_send\_buffer(conn, i, total\_loop), 返回结果: ret = ibv\_post\_send(conn->cm\_id->qp, &conn->send\_wr[i], &failed\_send\_wr)

执行wait\_cq()的输出结果是: wr\_id=0 opcode=0x00 (send) status=0

返回结果是: wr\_id=0 opcode=0x00 (recv) status=0

根据 prog\_opt.verify 的返回值验证数据, 输出结果: wr\_index=0 sgen=0 mark=1

ibv\_post\_recv(conn->cm\_id->qp, &conn->recv\_wr[i], &failed\_recv\_wr)

ibv\_post\_send(conn->cm\_id->qp, &conn->send\_wr[i], &failed\_send\_wr)

wait\_cq(conn, total\_loop)

ibv\_post\_send(conn->cm\_id->qp, &conn->send\_wr[i], &failed\_send\_wr)

wait\_cq(conn, total\_loop)

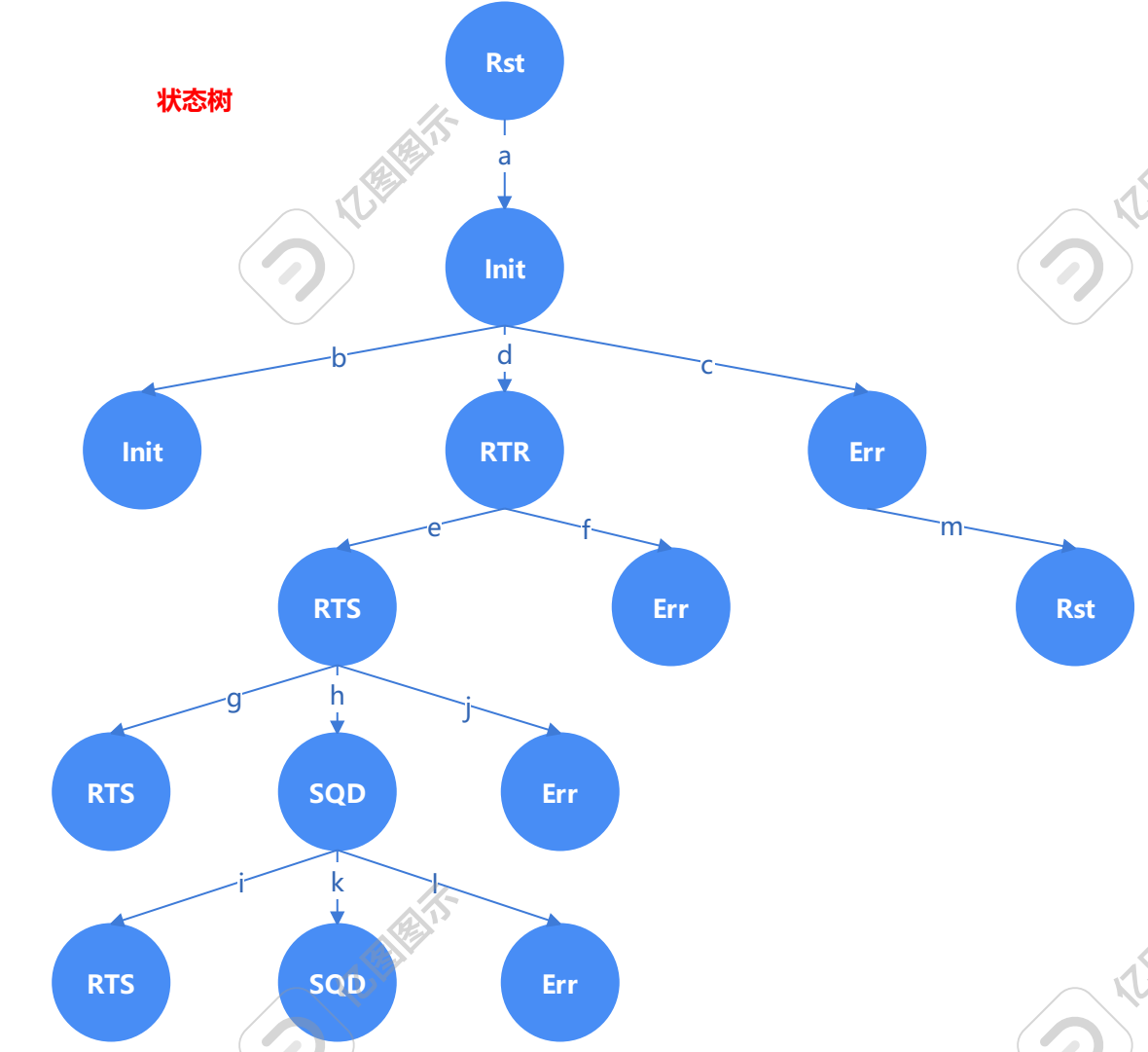
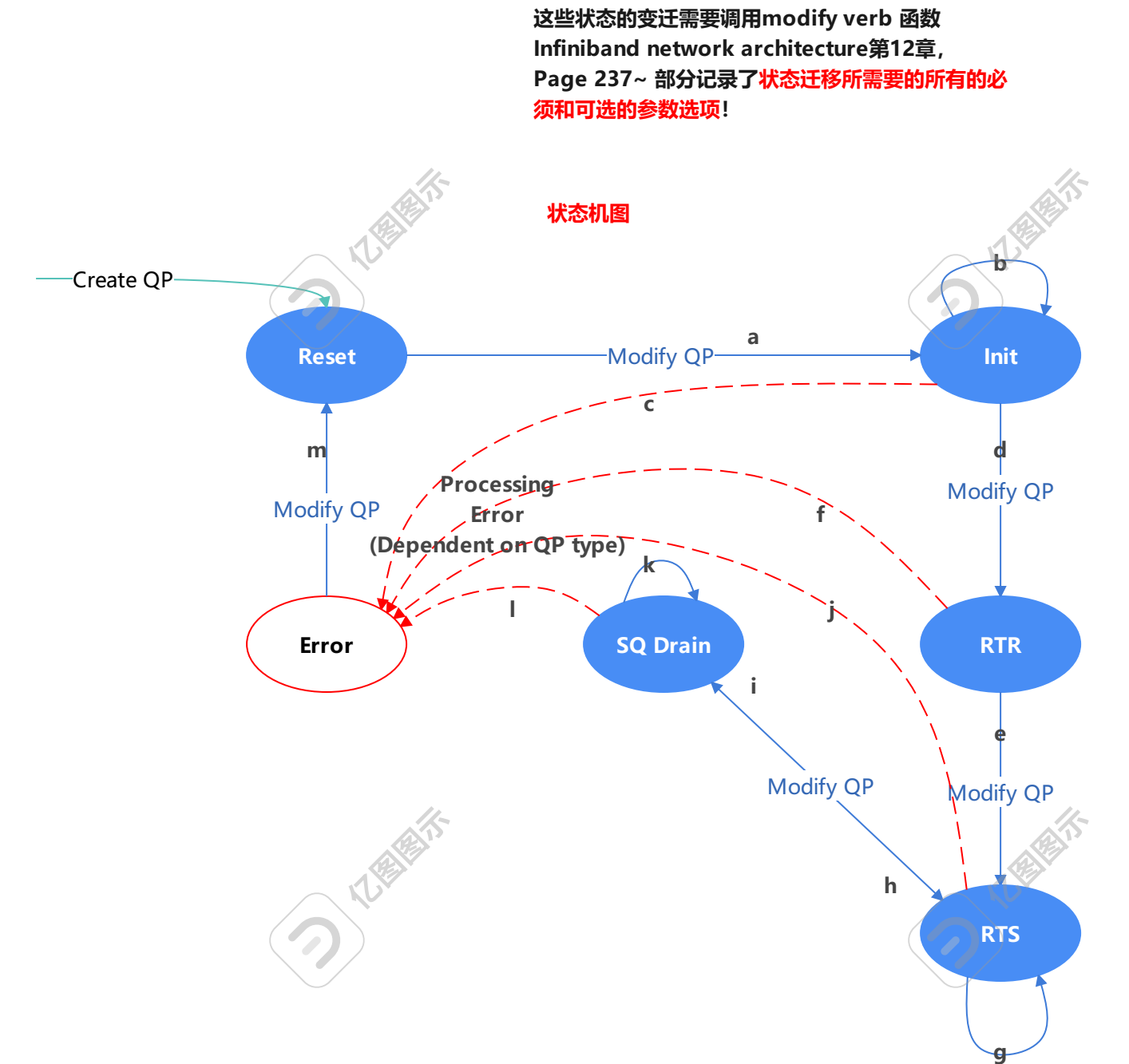
verify\_rcv\_buffer(conn, i, total\_loop)

ibv\_post\_recv(conn->cm\_id->qp, &conn->recv\_wr[i], &failed\_recv\_wr)

共循环100次

POST receive 和 post receive send 的测试! 是一个send和receive流程的循环, 这里列出了2次





测试用例  
此处的case是一条有始有终的可能状态变迁  
假设8个qp，每一个对应下面的一条case即可。

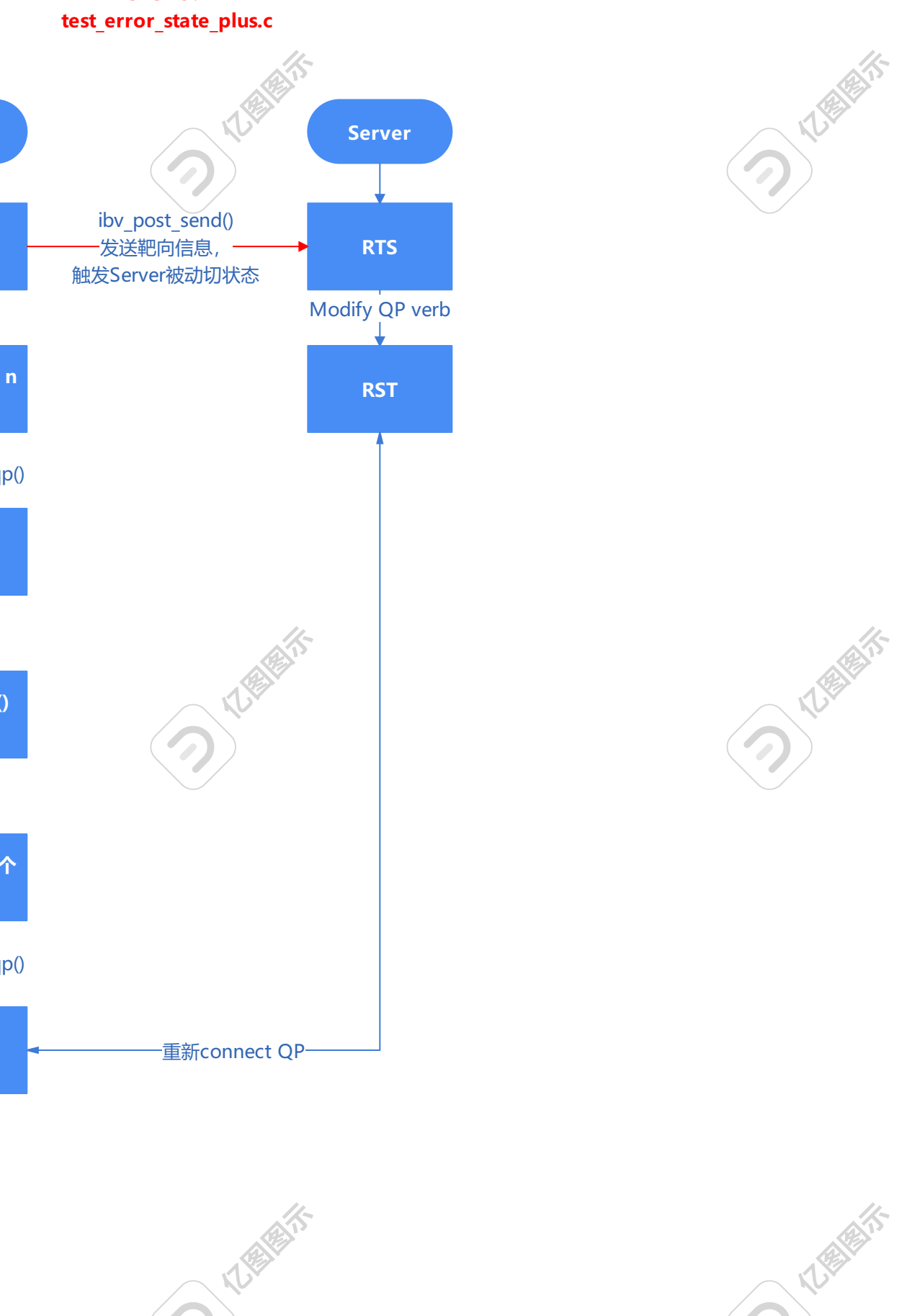
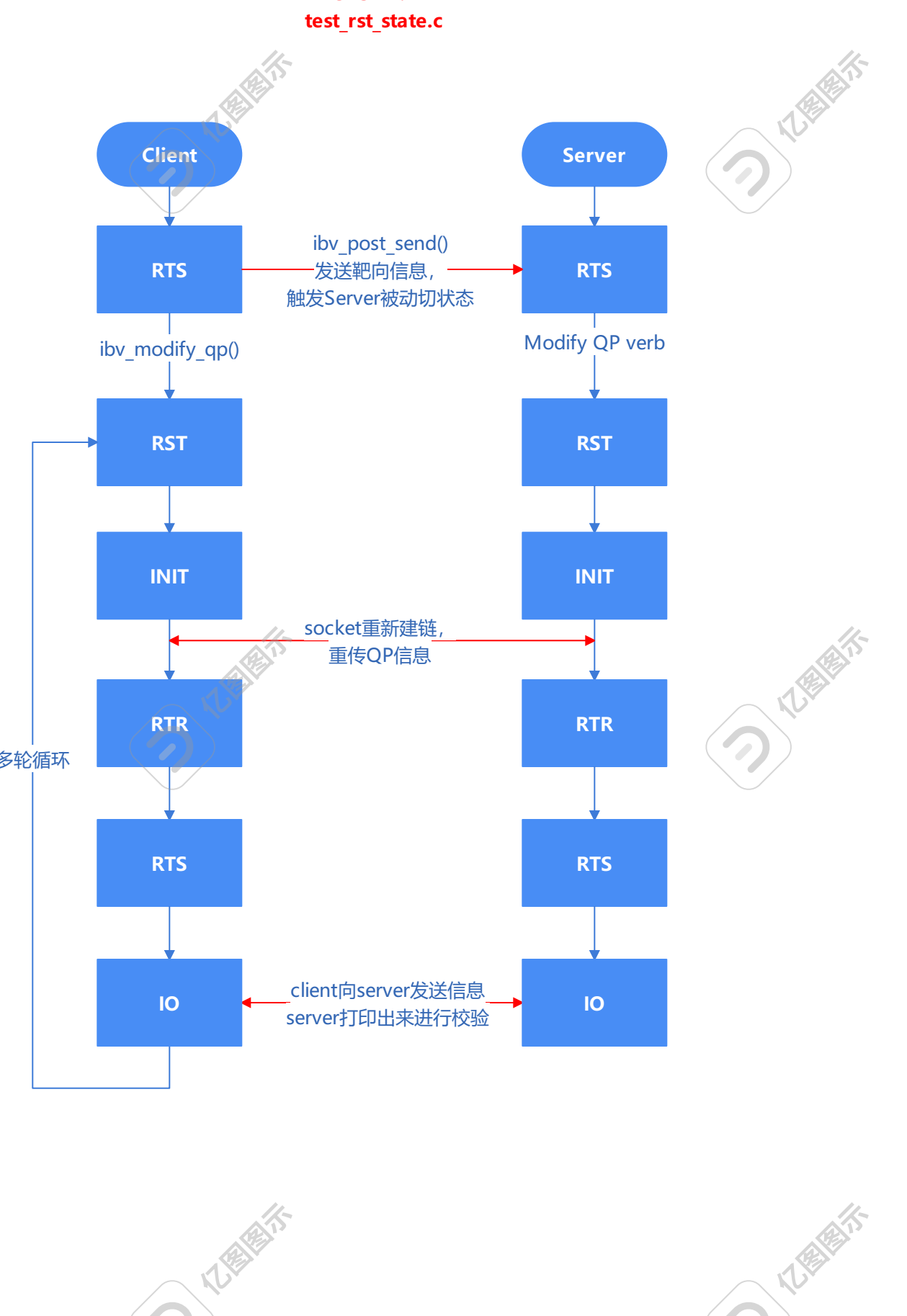
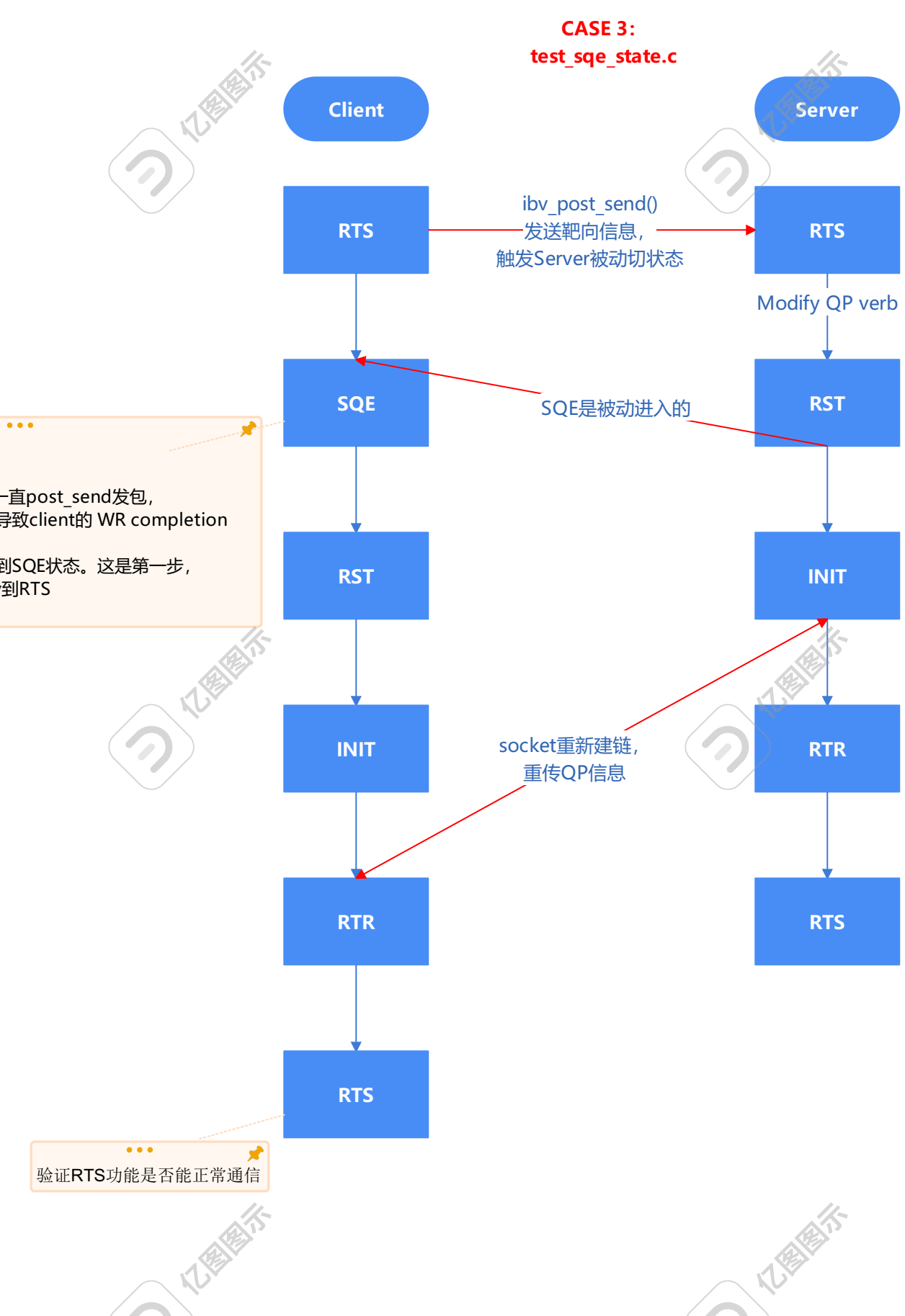
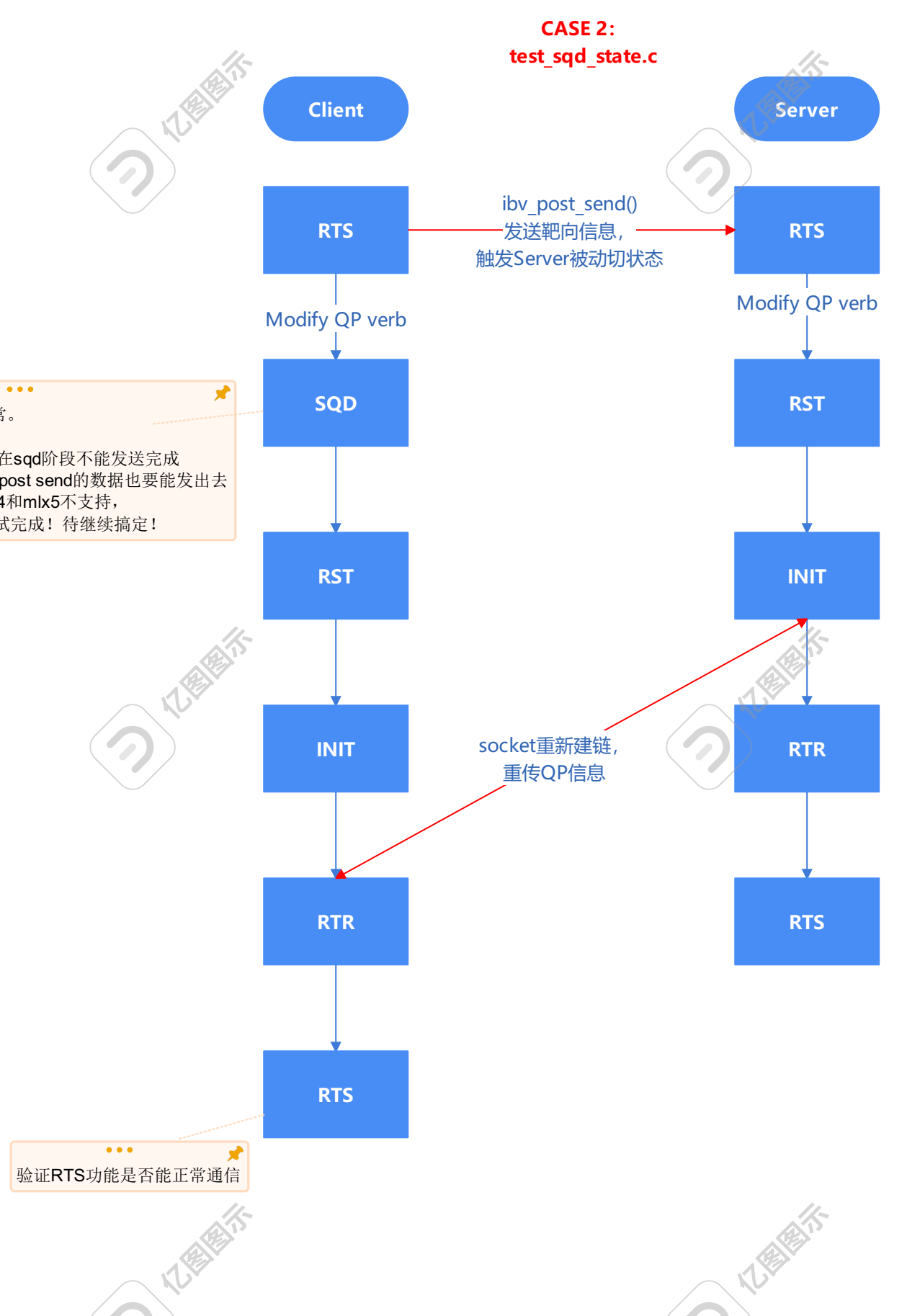
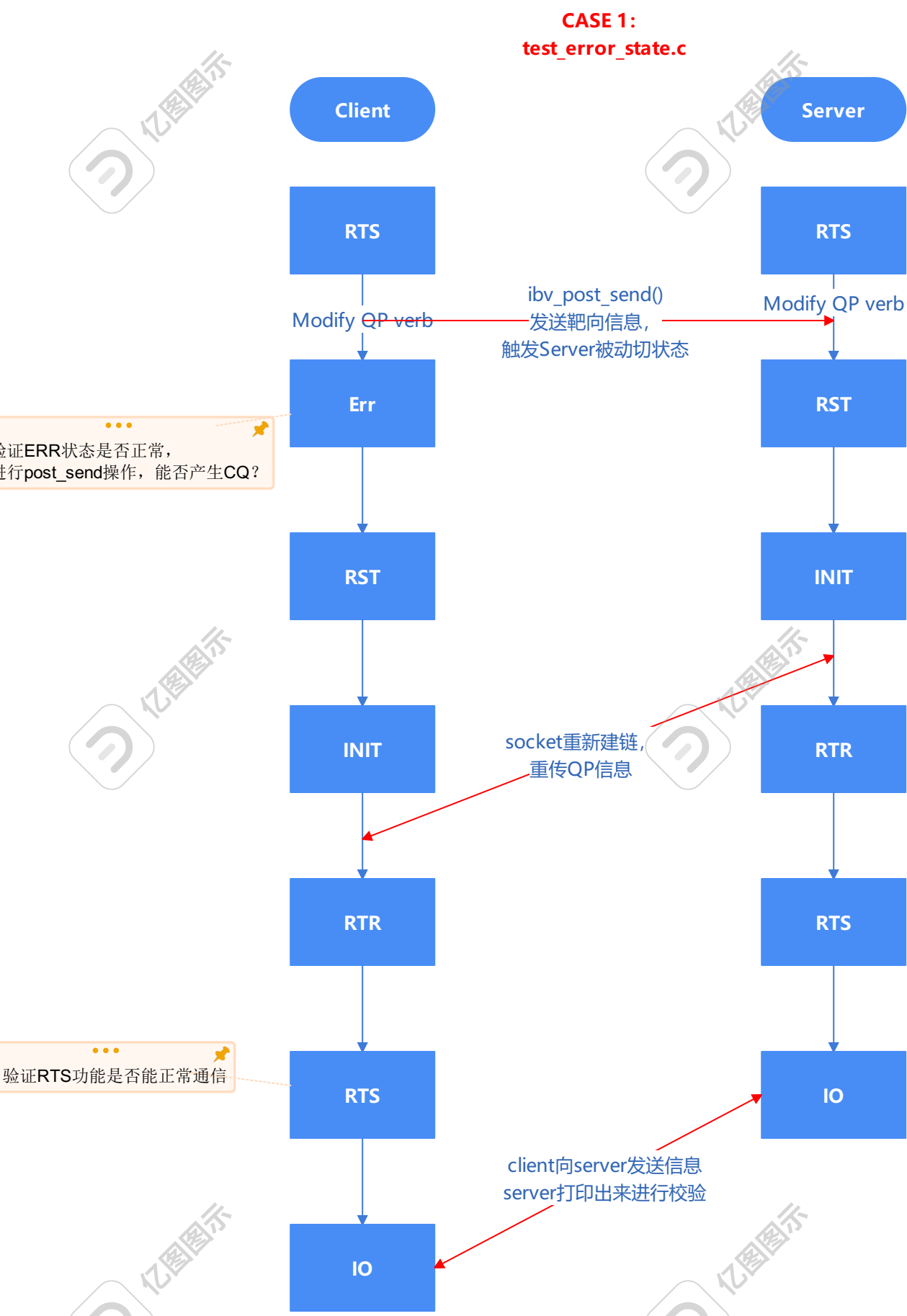
序号	状态迁移case (0-Switch)	描述
1	Rst->a->Init->b->Init	
2	Rst->a->Init->d->RTR->e->RTS->g->RTS	
3	Rst->a->Init->d->RTR->e->RTS->h->SQD->i->RTS	
4	Rst->a->Init->d->RTR->e->RTS->h->SQD->k->SQD	
5	Rst->a->Init->d->RTR->e->RTS->h->SQD->l->Err	
6	Rst->a->Init->d->RTR->e->RTS->j->Err	
7	Rst->a->Init->d->RTR->f->Err	
8	Rst->a->Init->c->Err->m->Rst	

基本函数

基本动作	函数
a	rst_to_init()
b	init_to_init()
c	init_to_error()
d	init_to_rtr()
e	rtr_to_rts()
f	rtr_to_err()
g	rts_to_rts()
h	rts_to_sqd()
i	sqd_to_rts()
j	rts_to_err()
k	sqd_to_sqd()
l	sqd_to_err()
m	err_to_rst()

修改设计:

不是全量测试，只需测试业务需要的即可！  
三个case：  
1.从RTS（不验证功能）->ERR（验证是否进入ERR状态，能否产生CQ）->RTS（ready to work）  
2.RTS->SQD（验证发包功能）->RTS  
3.RTS->RST->RTS  
三个case，先完成一次测试一个case的。  
两端都要处理有难度。（统一在客户端操作，服务端被动响应）有可能需要重新建链（参考pingpong程序中的socket），一端通知对端（post\_send 可以定义一个字段发给对端，对端也进入对应的状态）。





三. IO传输参数测试

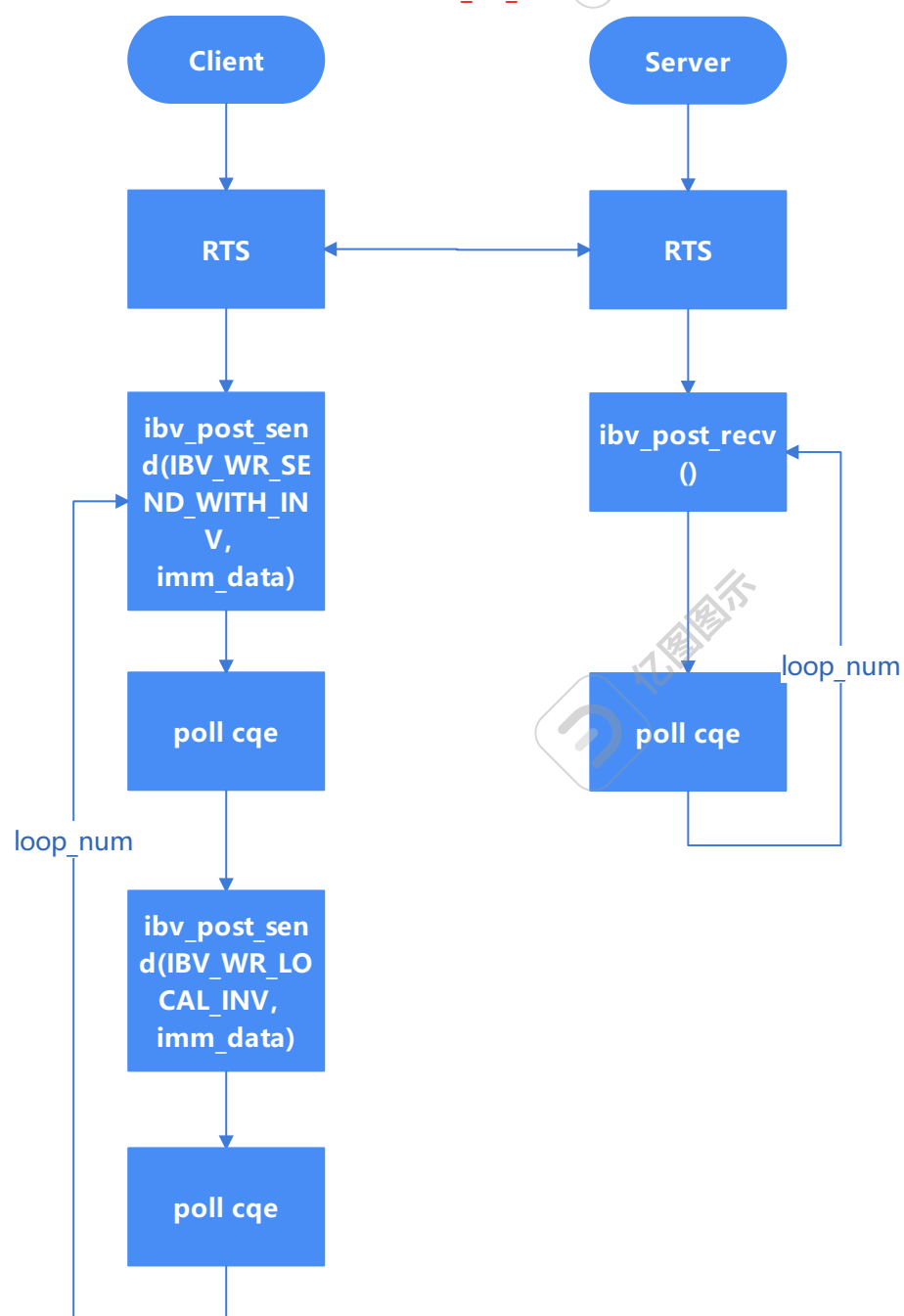
序号	测试内容	工具	测试步骤与要求
1	IBV_WR_SEND_WITH_IMM	内部用户态工具	
2	IBV_WR_RDMA_WRITE_WITH_IMM		
3	IBV_WR_SEND_WITH_INV		仅检查命令能否成功（暂无实际功能）
4	IBV_WR_LOCAL_INV		检查命令是否成功，本地key能否使用
5	IBV_SEND_SIGNALED		每次发2个包， 第一个包不设IBV_SEND_SIGNALED,第2个包设置， 通过标准是抓包能看到所有包，但cqe只有第2个的。（总循环次数为3*队列深度）
6	IBV_SEND_SOLICITED & Solicited Completion Event		1. ibv_notify_req_cq设置solicited_only 2. 每组发送2个包，带ibv_send_solicited的和不带。 3. 抓包能看到所有包，但接受侧中断只产生solicited的包

CASE 1:  
test\_imm\_data.c





CASE 1:  
test\_wr\_inv.c





- **IBV\_SEND\_SIGNED** - Set the completion notification indicator for this WR. This means that if the QP was created with **sq\_sig\_all=0**, a Work Completion will be generated when the processing of this WR will be ended. If the QP was created with **sq\_sig\_all=1**, there won't be any effect to this flag

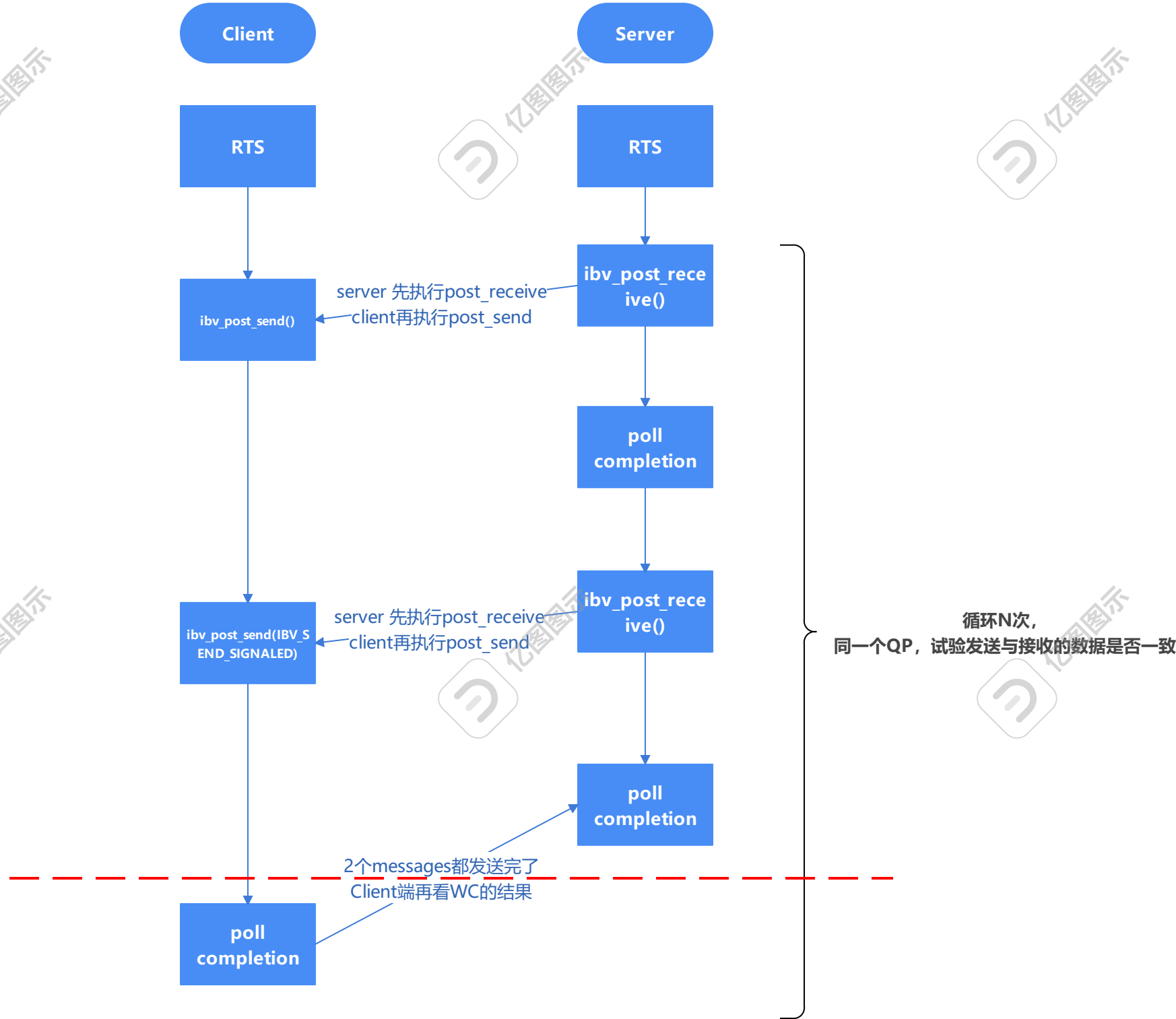
创建QP的时候，让sq\_sig\_all = 0: (如果等于1, 那么产生CQE就不会受这个IBV\_SEND\_SIGNED的影响了); post\_send(IBV\_SEND\_SIGNED)才会产生CQE; post\_send()不会产生CQE;

所以sq\_sig\_all 和 IBV\_SEND\_SIGNED 是配合使用的! 用于控制本端发送的时候要不要产生CQE的!

测试方案:

- 1.创建QP的时候让sq\_sig\_all = 0;
- 2.Client连续发两个包，第一个包不设置SEND\_SIGNED, 第二个包设置;
- 3.Client 出口 (Server入口) 可以抓到两个包，client设置poll 2个CQE的配置 (队列深度为2)，但是预期只能得到一个CQE，并且是第二个包 (通过wr\_id判断) 的。第一个包的行为是：无法poll 到CQE。

CASE 1:  
test\_send\_signed.c



- **IBV\_SEND\_SOLICITED** - Set the solicited event indicator for this WR. This means that when the message in this WR will be ended in the remote QP, a solicited event will be created to it and if in the remote side the user is waiting for a solicited event, it will be woken up. Relevant only for the Send and RDMA Write with immediate opcodes

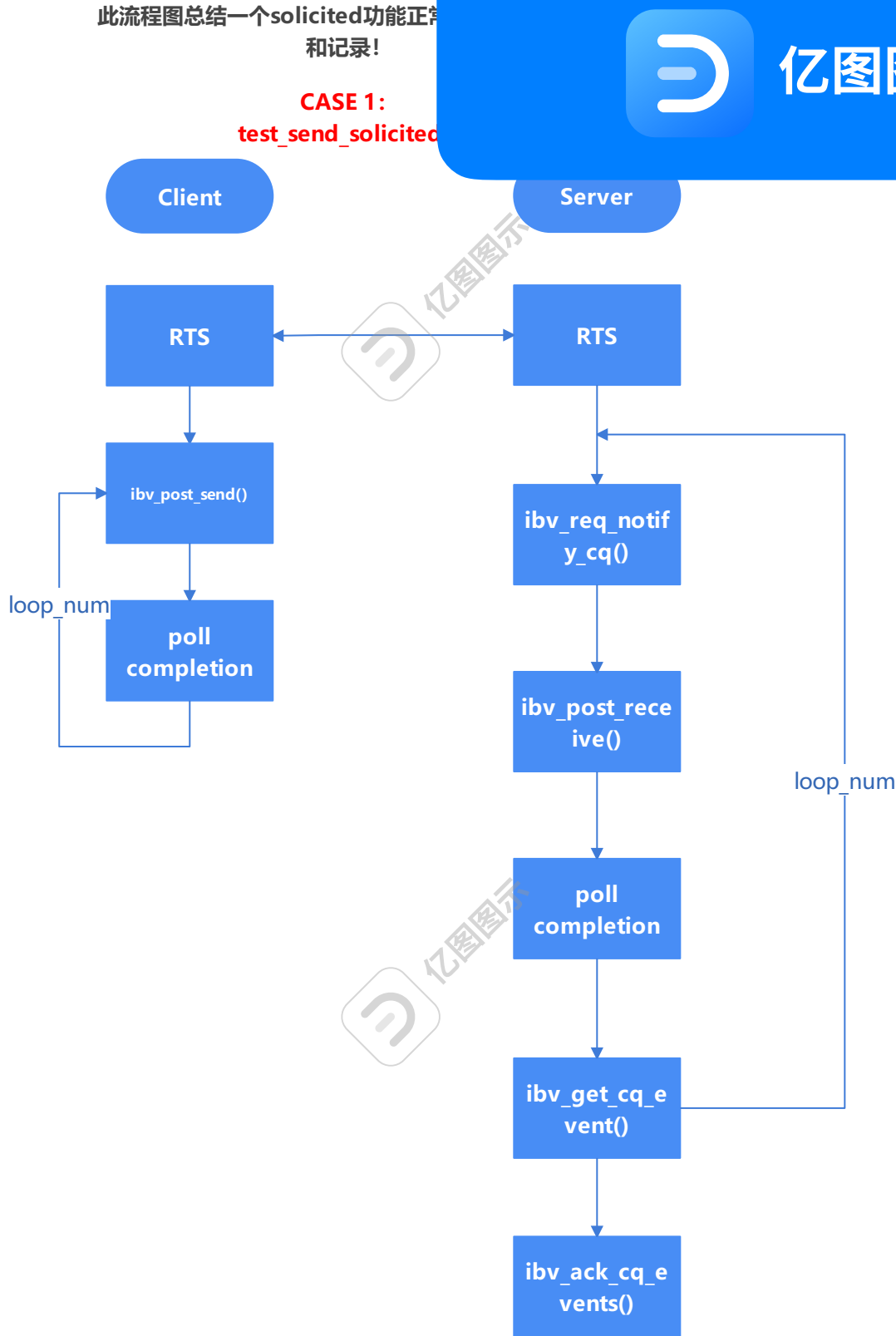
理解：  
参考OneNote中的相关设计总结  
总体业务代码逻辑：  
server端，post\_receive()之后准备接收了，就可以使用下面的机制了，然后等待events的产生  
ibv\_req\_notify\_cq(pp\_cq(ctx), solicited\_flag); // 使能一种机制，何种情况下产生notification

Client端，post\_send()（发两个，先普通，后solicited）

```
while() {  
    ibv_get_cq_event(ctx->channel, &ev_cq, &ev_ctx); // 等待这个notification告诉我们产生了CQE  
    num_cq_events++; // 计数增加  
    ibv_req_notify_cq(pp_cq(ctx), 0); // 重复执行  
  
    ibv_poll_cq(pp_cq(ctx), 2, wc); // 处理CQE  
}
```

ibv\_ack\_cq\_events(pp\_cq(ctx), num\_cq\_events); // 确认知道了，get\_cq\_events收到了events  
判断event的次数是否符合预期。  
组合验证发送端的ibv\_send\_solicited与接收端的solicited\_flag。

- 本业务要求覆盖四种组合测试：
- o 0: 只测IBV\_SEND\_SOLICITED 不测cq event (IBV\_SEND\_SIGNALED =1)  
——client端：post\_send(IBV\_SEND\_SOLICITED),  
——Server端：不使用 ibv\_req\_notify\_cq() 机制  
预期是就是最普通的场景，两端都只产生CQE，没有event
  - o 1: IBV\_SEND\_SOLICITED + ibv\_req\_notify\_cq(struct ibv\_cq \*cq, int solicited\_only) (IBV\_SEND\_SIGNALED = 1) (solicited\_only = 1)  
——Client端，post\_send(IBV\_SEND\_SOLICITED);  
——Server端调用ibv\_req\_notify\_cq(1)机制，  
预期两端都产生CQEs，server端有event
  - o 2: 只测IBV\_SEND\_SOLICITED 不测cq event (IBV\_SEND\_SIGNALED =0)  
——client端：post\_send(0), sleep(1), 预期不产生CQE,  
——Server端：不调用 ibv\_req\_notify\_cq() 机制；  
预期在client端不产生CQE，在server端会产生CQE，没有event  
——因为在client端，对于sq\_sig\_all = 0 的情况，只有IBV\_SEND\_SIGNALED才能完成整个过程，产生CQE
  - o 3: IBV\_SEND\_SOLICITED + ibv\_req\_notify\_cq(struct ibv\_cq \*cq, int solicited\_only) (IBV\_SEND\_SIGNALED = 0)  
——Client端，post\_send(0);  
——Server端调用ibv\_req\_notify\_cq(1);  
预期client不会产生CQE，server会产生CQE，没有event，不进行预期event





post\_send(),  
rkey错误

→

IBV\_WC\_REM\_ ACCESS\_ERR

post\_send(),  
lkey错误

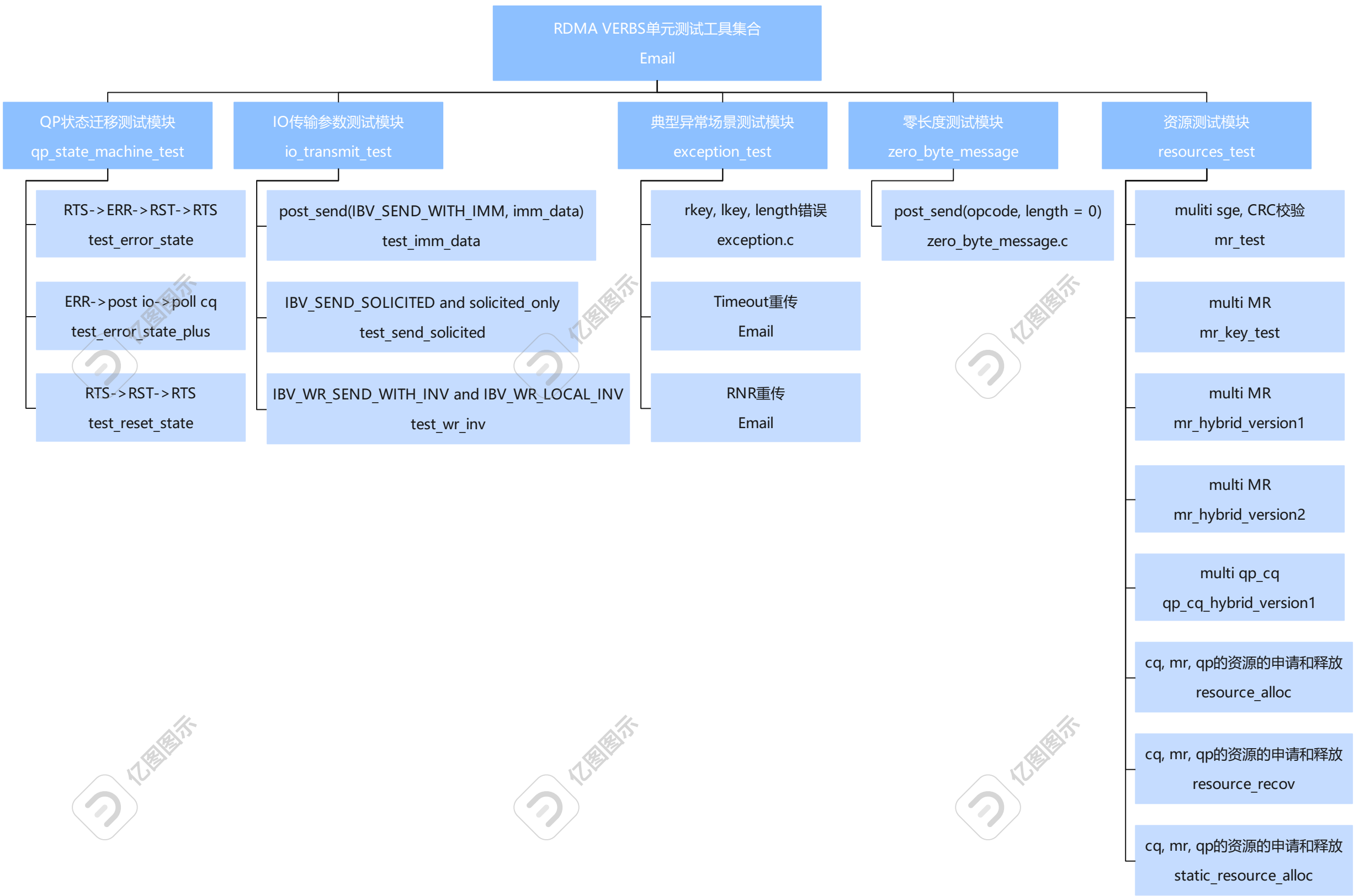
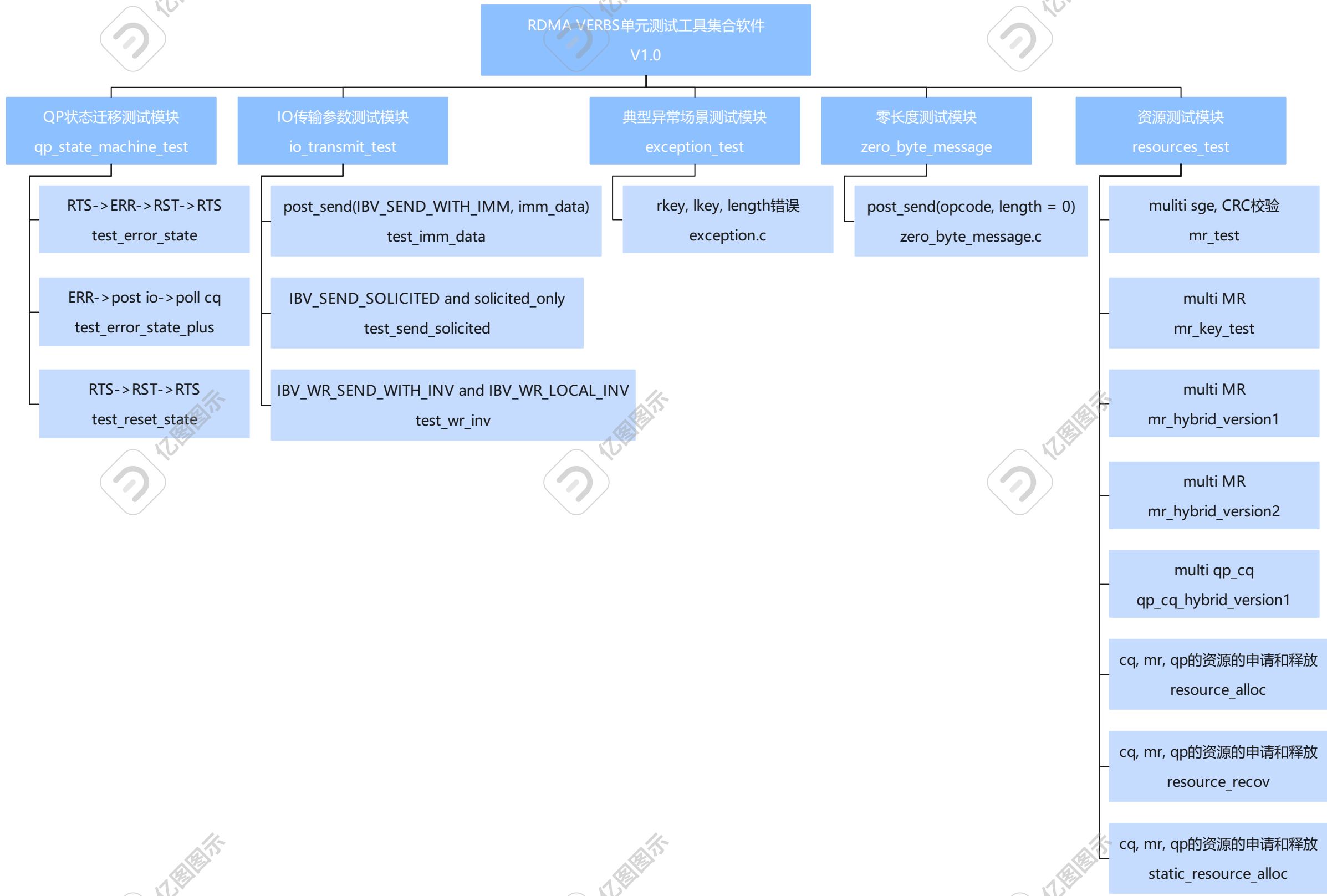
→

IBV\_WC\_LOC\_ PROT\_ERR

post\_recv(),  
buf长度不够

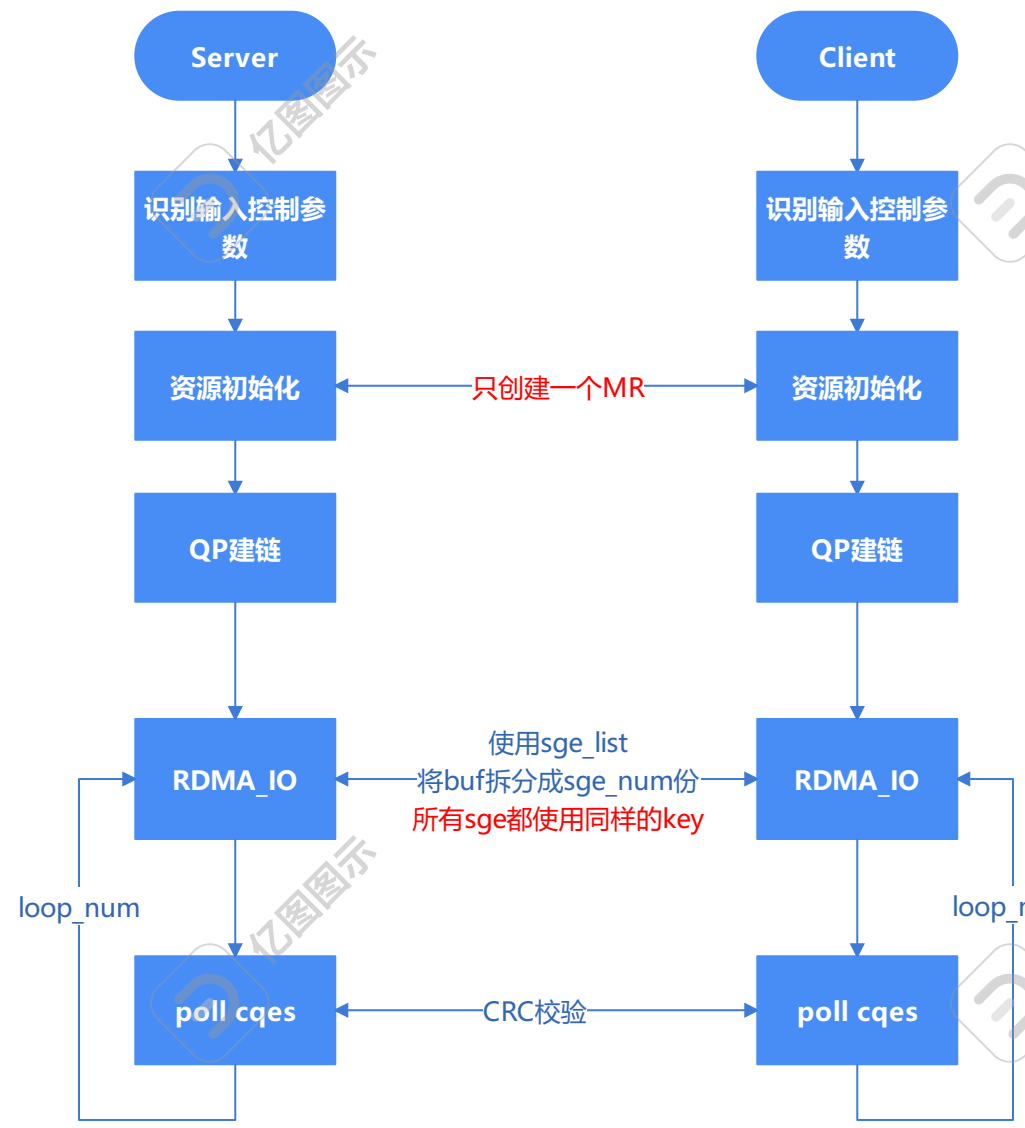
→

IBV\_WC\_REM\_ INV\_REQ\_ERR

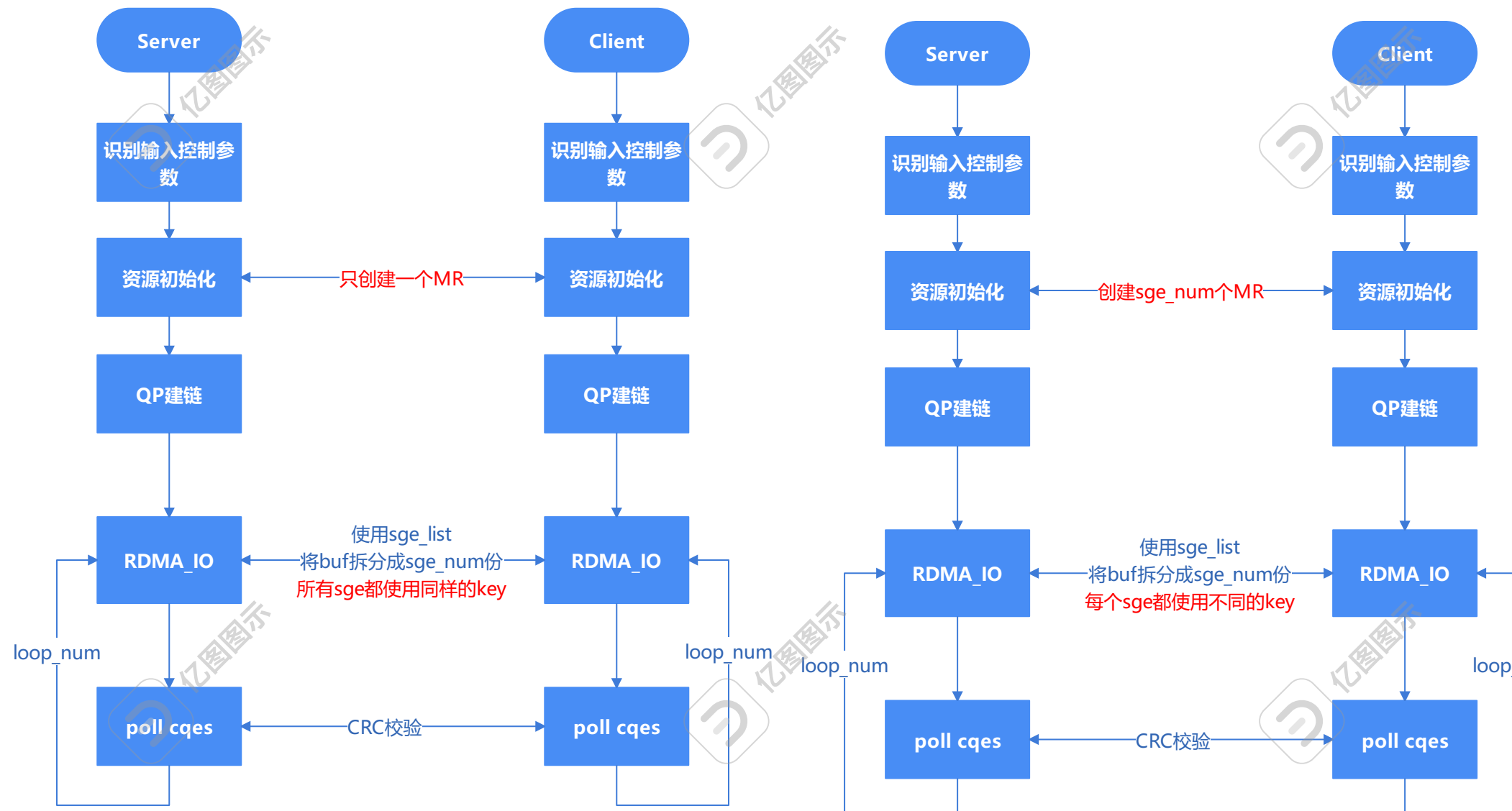




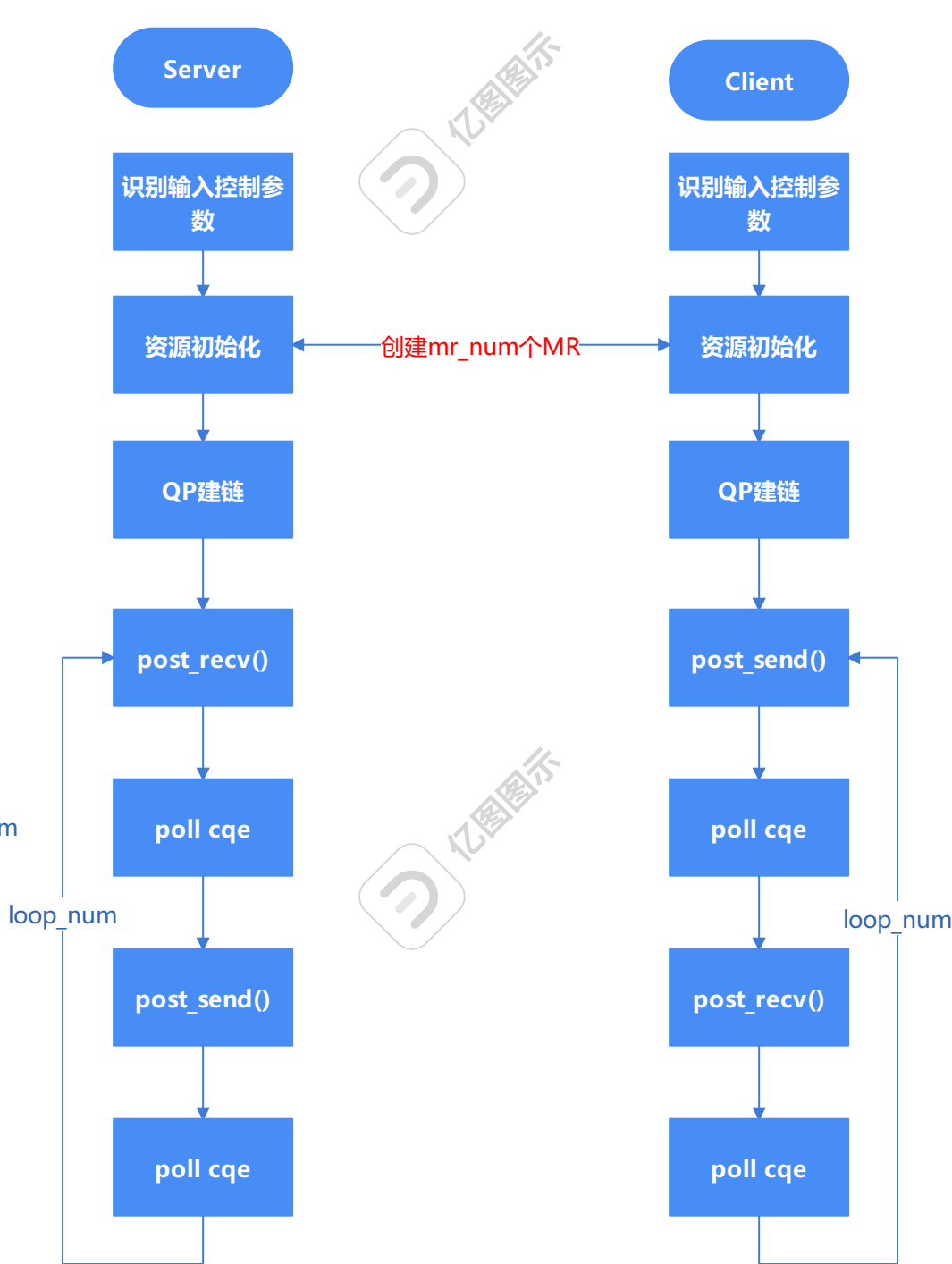
CASE 1:  
mr\_test



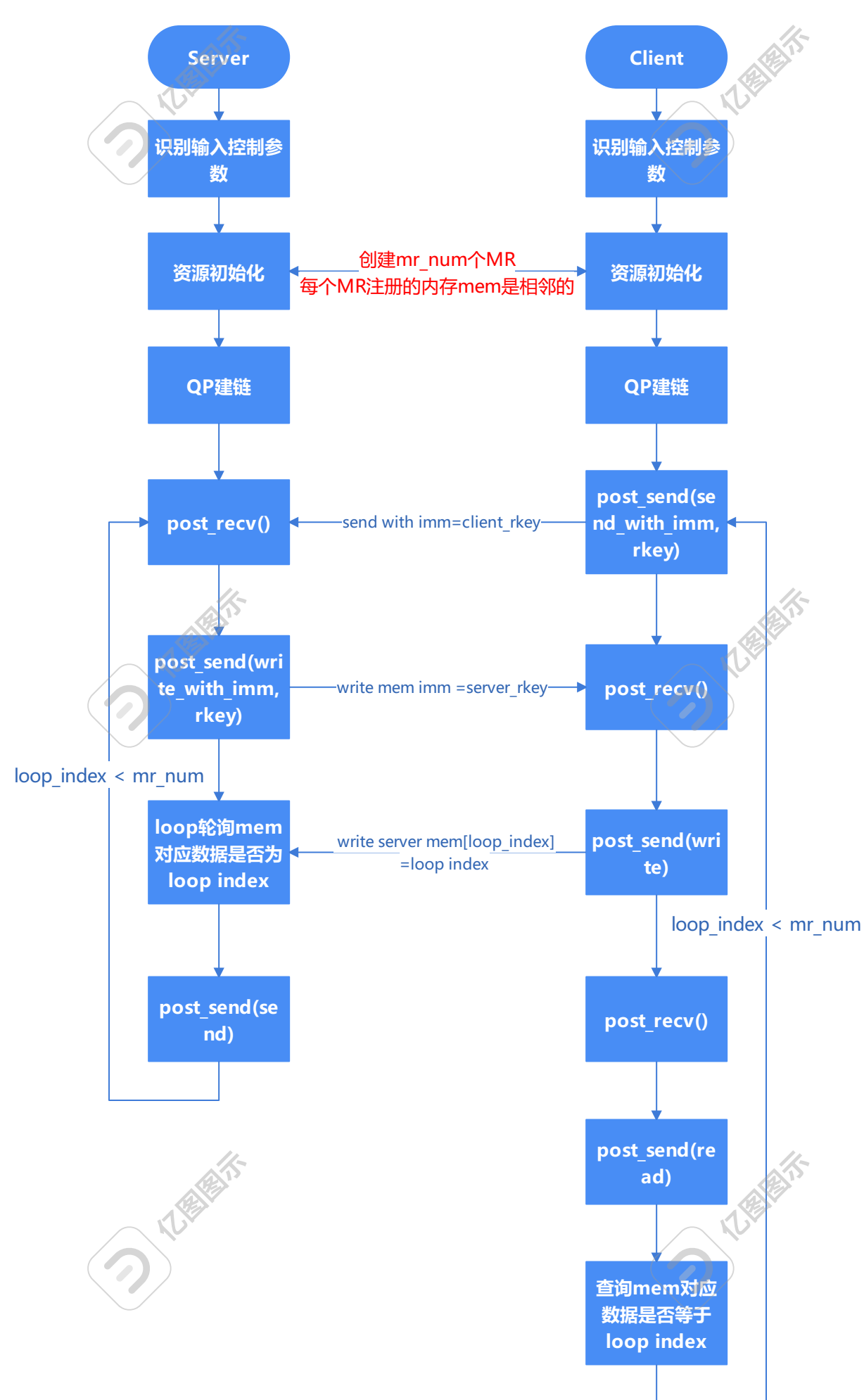
CASE 2:  
mr\_key\_test



CASE 3:  
mr\_hybrid\_version1



CASE 4:  
mr\_hybrid\_version2



CASE 7:  
resource\_recov

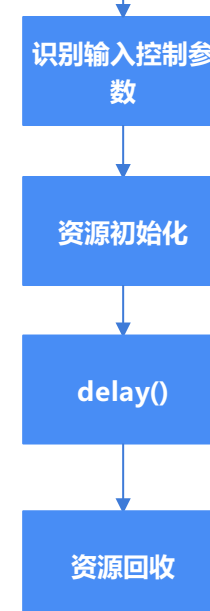
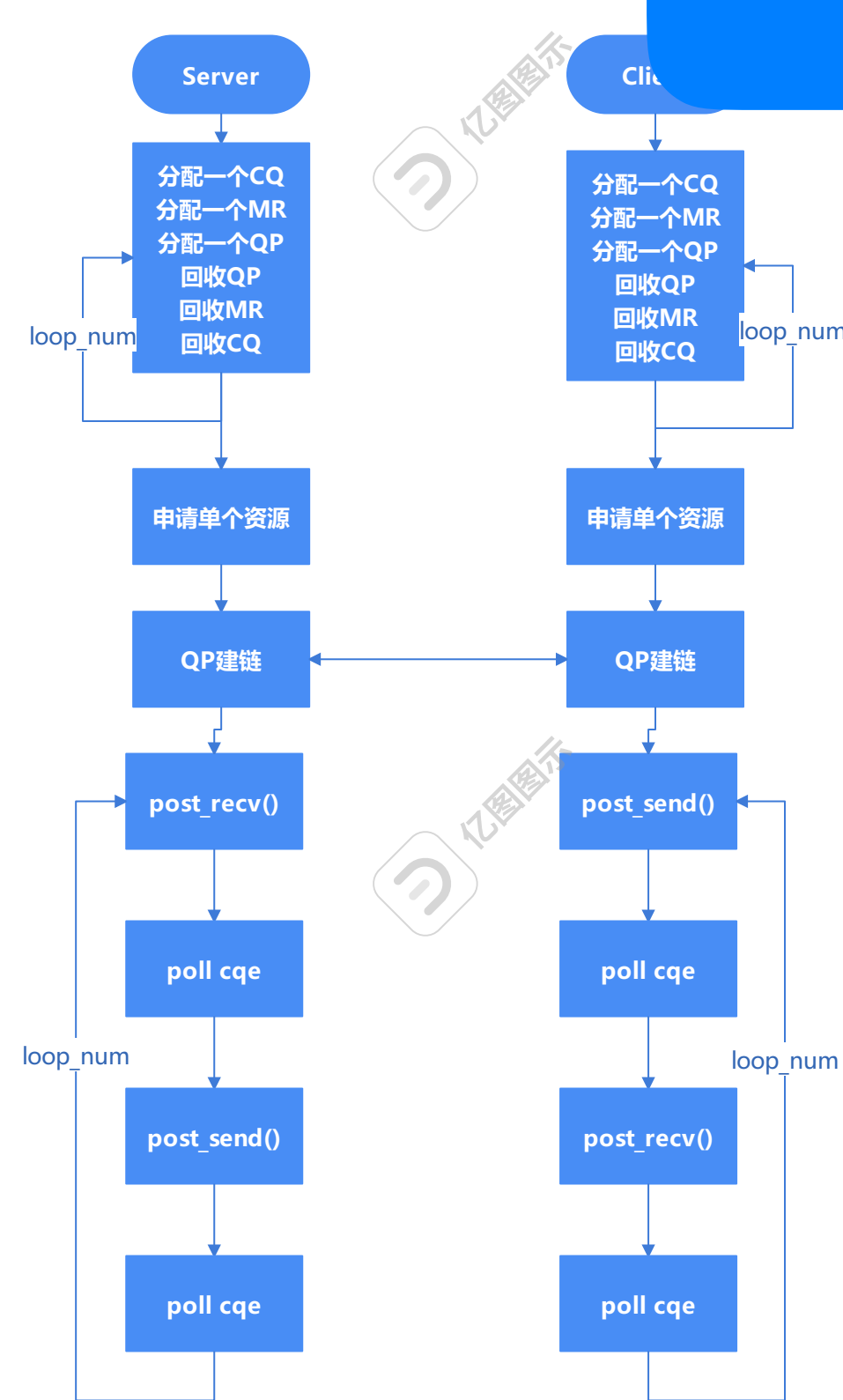


Table 77 Post Send Bind WR Rules

MW Type	MW State	PD Rule
Type 2A	Free	PD associated with R_Key in the Post Send Bind WR must be the same as the PD associated with the QP on which the WR was posted.
Type 2B	Free	

Table 79 Type 2 Memory Window Access Rules

MW Type	MW State	QP Rule	PD Rule
Type 2A	Valid	MW access operations only allowed via QP on which MW was bound	N/A
Type 2B	Valid	MW access operations only allowed via QP on which MW was bound	MW access operations only allowed if the PD associated with the R_Key is the same as the PD associated with the QP

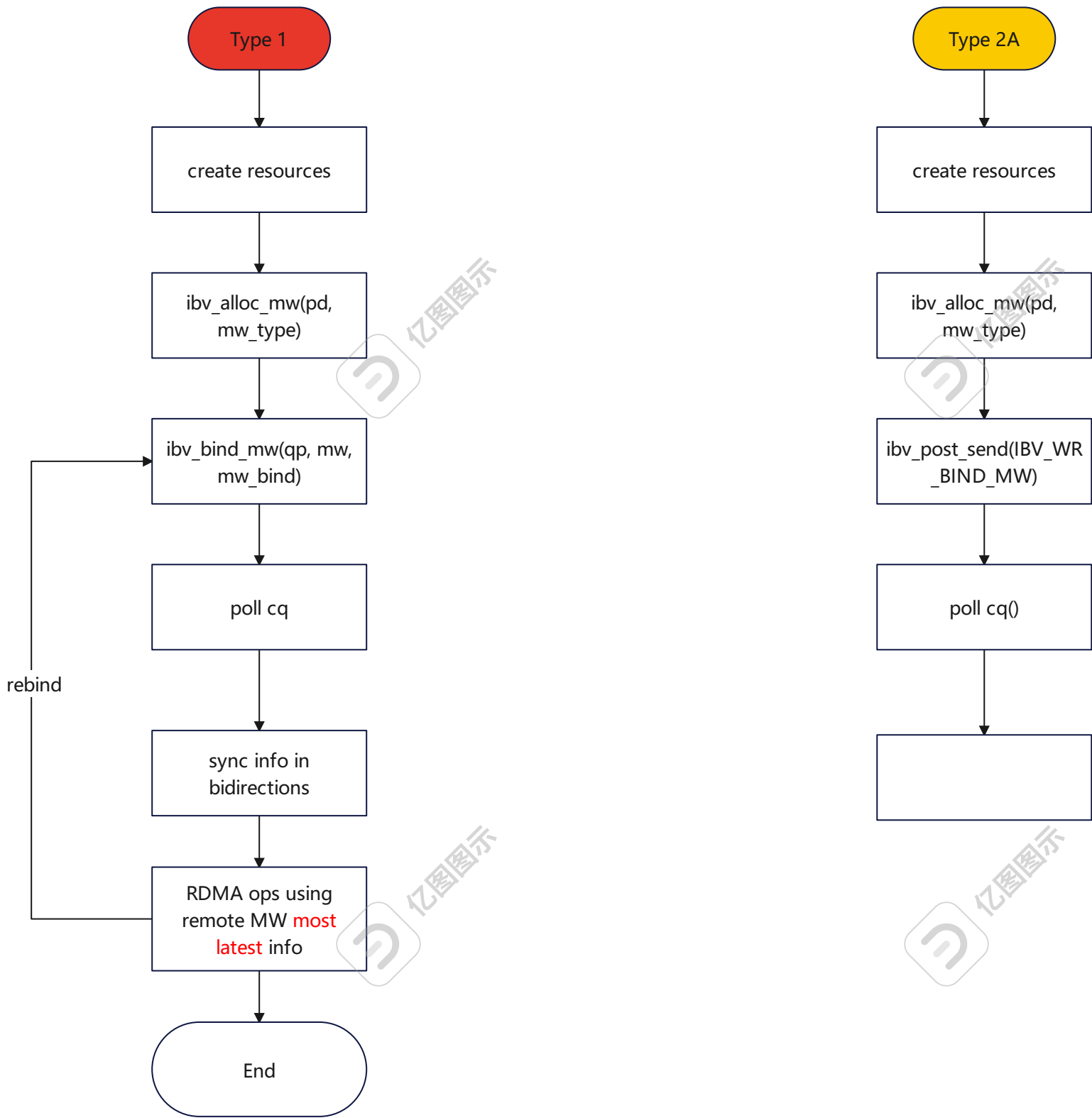
Table 80 Memory Windows Invalidation Rules

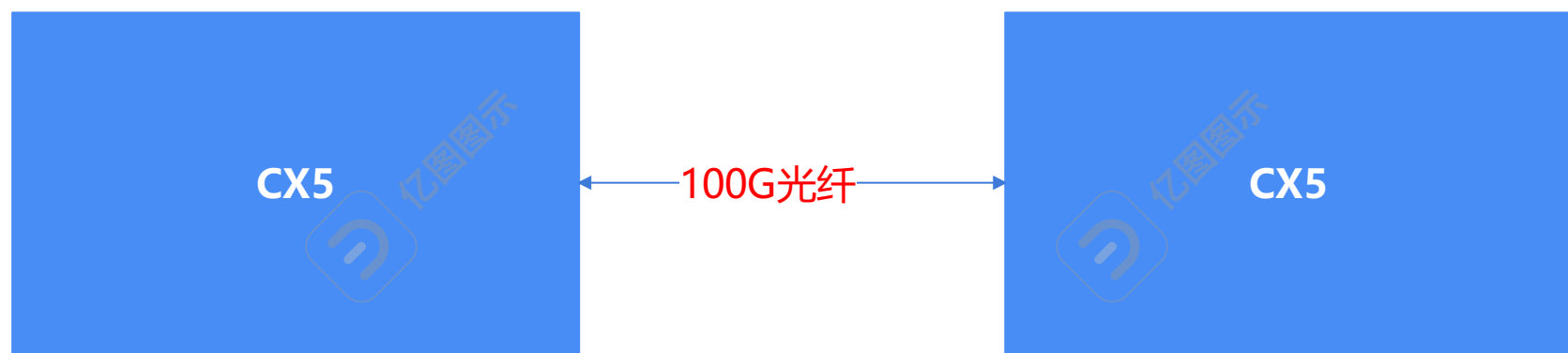
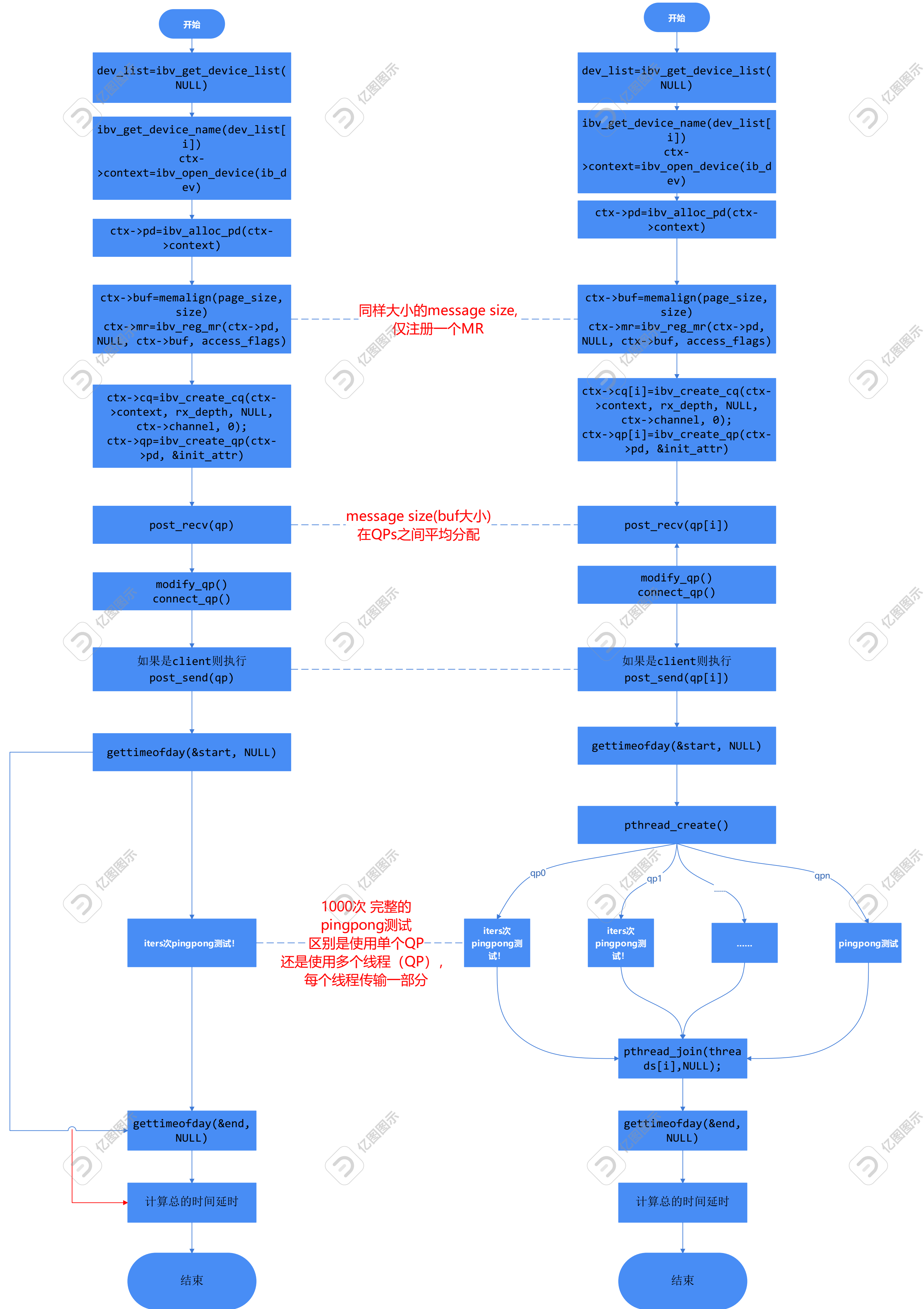
MW Type	Key Ownership	Local and Remote Invalidation	Bind using already Bound MW
Type I	CI	Not Allowed	Allowed
Type II	Consumer	Allowed	Not Allowed

Table 78 Type 2 Memory Window Invalidation Rules

MW Type	MW State	QP Rule	PD Rule
Type 2A	Valid	Invalidate operation only allowed via QP on which MW was bound	N/A
Type 2A	Free	N/A	Invalidate operation only allowed via QP associated with same PD as MW
Type 2B	Valid	Invalidate operation only allowed via QP on which MW was bound	
Type 2B	Free	N/A	

区别点	Type 1	Type 2A	Type 2B
Bind method	Bind MW verbs	Post Send Bind MW WR	Post Send Bind MW WR
Post Send Bind WR Rules	-	Free state	Free state
access rules	PD	QP	PD + QP
R_key ownership	CI 返回R_key值	consumer 自己拼凑R_key	consumer 自己拼凑R_key
Bind using already bound MW	Allowed 绑定后之前的R_Key自动失效	Not Allowed 绑定前需要先使之前的R_Key无效化	Not Allowed 绑定前需要先使之前的R_Key无效化
是否支持零长度	是 通过零长度的方式使得mw的bindings失效	否	否
Local and Remote Invalidation operation	Not Allowed 采用deallocate mw, 或者绑定到新的区域。	Allowed	Allowed
关联的QP是否可以被销毁	-	Not Allowed	Allowed





参考:  
1.srq\_pingpong中 结构体, 多QP的connect这一块;  
2.cmtime中的多线程的实现  
3.perftest中run\_infinately的实现

