

零声学院 Mark 老师 QQ : 2548898954

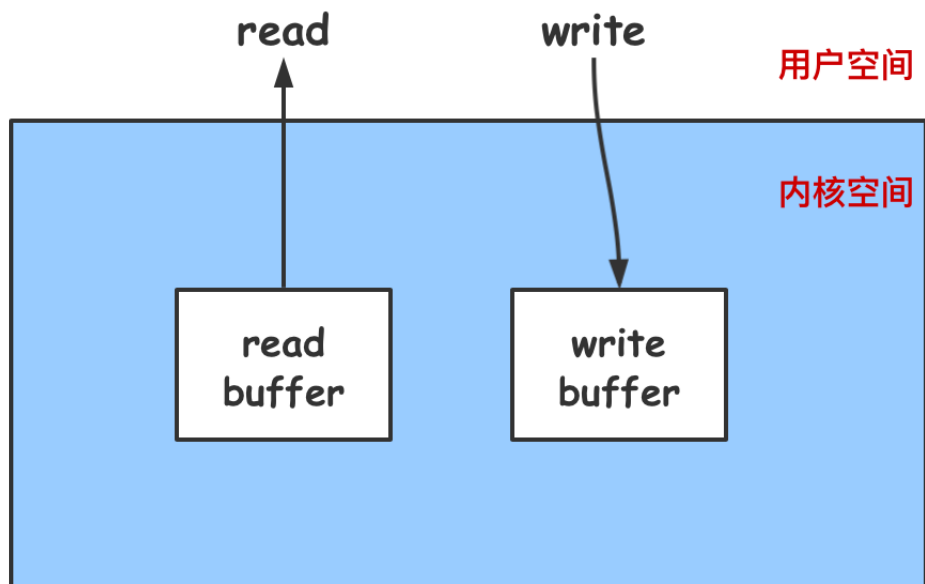
网络编程关注的问题

连接建立

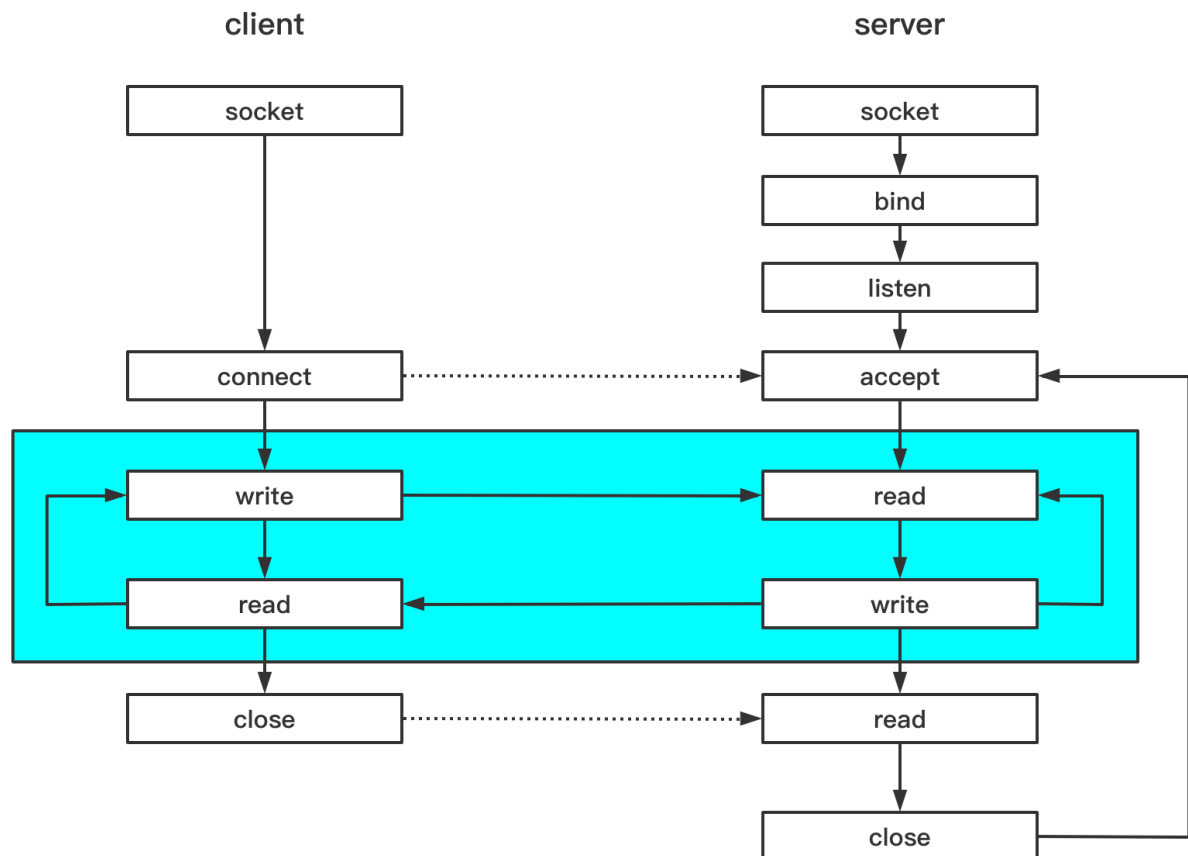
连接断开

消息到达

消息发送完毕



网络编程流程



阻塞io模型和非阻塞io模型

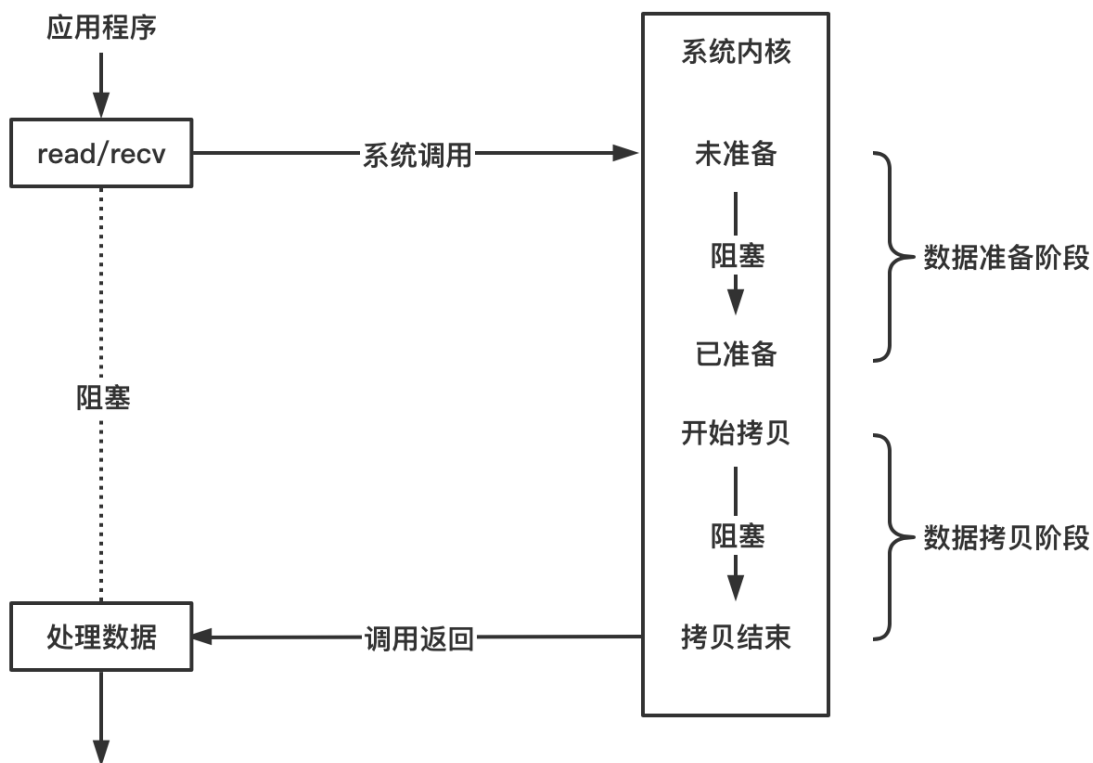
- 阻塞在哪里？ 阻塞在网络线程
- 什么来决定阻塞还是非阻塞？

```
1 | 连接的fd
2 | fcntl(c->fd, F_SETFL, O_NONBLOCK);
```

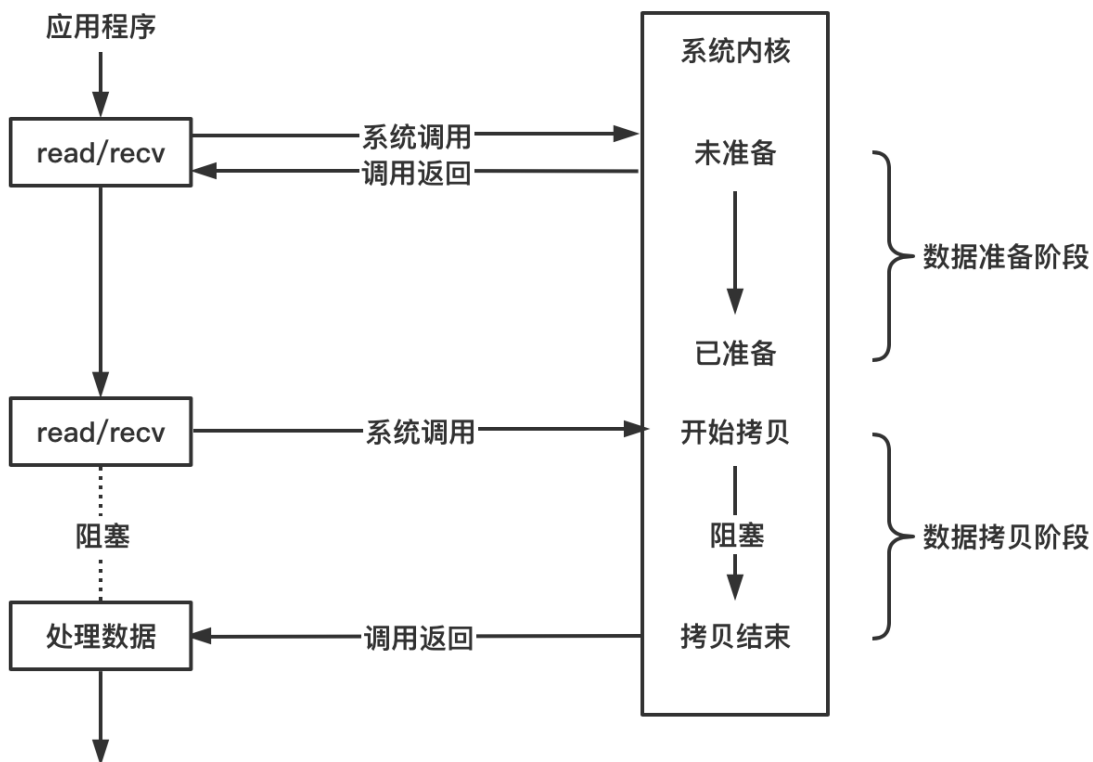
- 具体的差异？

io函数在数据未到达时是否立刻返回

阻塞io模型



非阻塞io模型

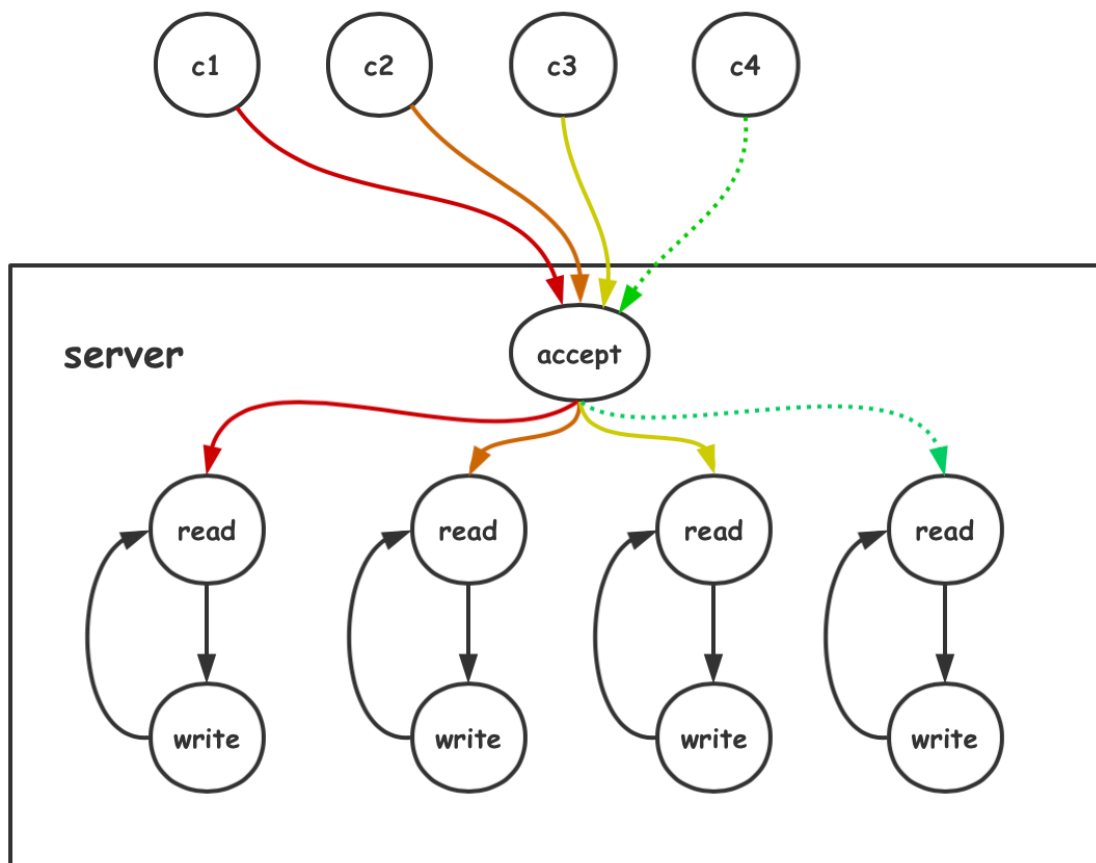


阻塞io模型 + 多线程

每一个线程处理一个 fd 连接 bio

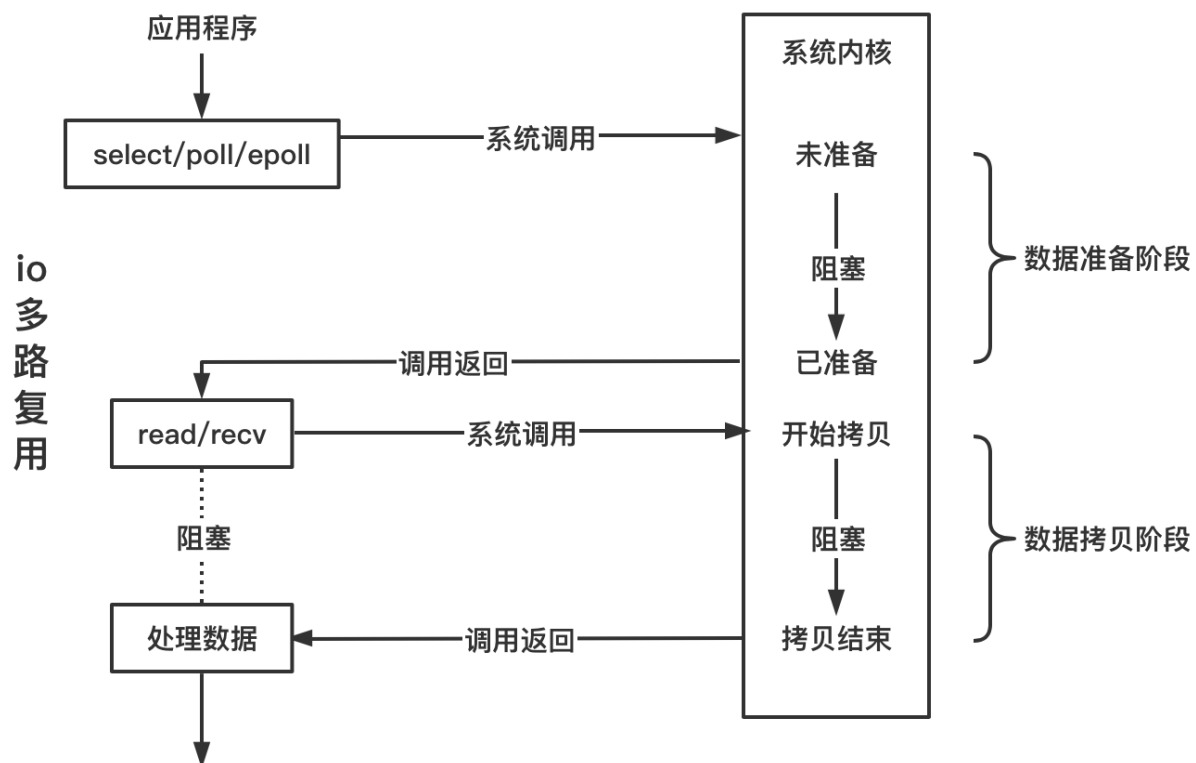
优点：处理及时

缺点：线程利用率很低，线程的数量是有限的



io多路复用（网络线程）

用一个线程来检测多个io事件



水平触发的时候，io函数既可以是阻塞的也可以是非阻塞的。

边缘触发的时候，io函数只能是非阻塞的。

epoll基础

重要数据结构

```
1 struct eventpoll {
2     // ...
3     struct rb_root rbr;           // 管理 epoll 监听的事件
4     struct list_head rdllist;     // 保存着 epoll_wait 返回满足条件的事件
5     // ...
6 };
7
8 struct epitem {
9     // ...
10    struct rb_node rbn;           // 红黑树节点
11    struct list_head rdllist;     // 双向链表节点
12    struct epoll_filefd ffd;     // 事件句柄信息
13    struct eventpoll *ep;         // 指向所属的eventpoll对象
14    struct epoll_event event;     // 注册的事件类型
15    // ...
16 };
17
18 struct epoll_event {
19     __uint32_t events;
20     epoll_data_t data;           // 保存 关联数据
21 };
22
23 typedef union epoll_data {
24     void *ptr;
25     int fd;
26     uint32_t u32;
27     uint64_t u64;
28 }epoll_data_t;
```

主要函数

- epoll_create系统调用

```
1 | int epoll_create(int size);
```

size参数告诉内核这个epoll对象会处理的事件大致数量，而不是能够处理的事件的最大数。

在现在linux版本中，这个size参数已经没有意义了；

返回：epoll对象句柄；之后针对该epoll的操作需要通过该句柄来标识该epoll对象；

- epoll_ctl系统调用

```
1 | int epoll_ctl(int epfd, int op, int fd, struct epoll_event* event);
```

epoll_ctl向epoll对象添加、修改或删除事件；

返回：0表示成功，-1表示错误，根据errno错误码判断错误类型。

op类型：

EPOCH_CTL_ADD	添加新的事件到epoll中
EPOCH_CTL_MOD	修改epoll中的事件
EPOCH_CTL_DEL	删除epoll中的事件

event.events 取值：

EPOCHLIN	表示该连接上有数据可读（tcp连接远端主动关闭连接，也是可读事件，因为需要处理发送来的FIN包；FIN包就是read 返回 0）
EPOCHLOUT	表示该连接上可写发送（主动向上游服务器发起非阻塞tcp连接，连接建立成功事件相当于可写事件）
EPOCHLRDHUP	表示tcp连接的远端关闭或半关闭连接
EPOCHLPRI	表示连接上有紧急数据需要读
EPOCHLERR	表示连接发生错误
EPOCHLHUP	表示连接被挂起
EPOCHLLET	将触发方式设置为边缘触发，系统默认为水平触发
EPOCHLONESHOT	表示该事件只处理一次，下次需要处理时需重新加入epoll

- epoll_wait系统调用

```
1 | int epoll_wait(int epfd, struct epoll_event* events, int maxevents, int timeout);
```

收集 epoll 监控的事件中已经发生的事件，如果 epoll 中没有任何一个事件发生，则最多等待 timeout 毫秒后返回。

返回：表示当前发生的事件个数

返回0表示本次没有事件发生；

返回-1表示出现错误，需要检查errno错误码判断错误类型。

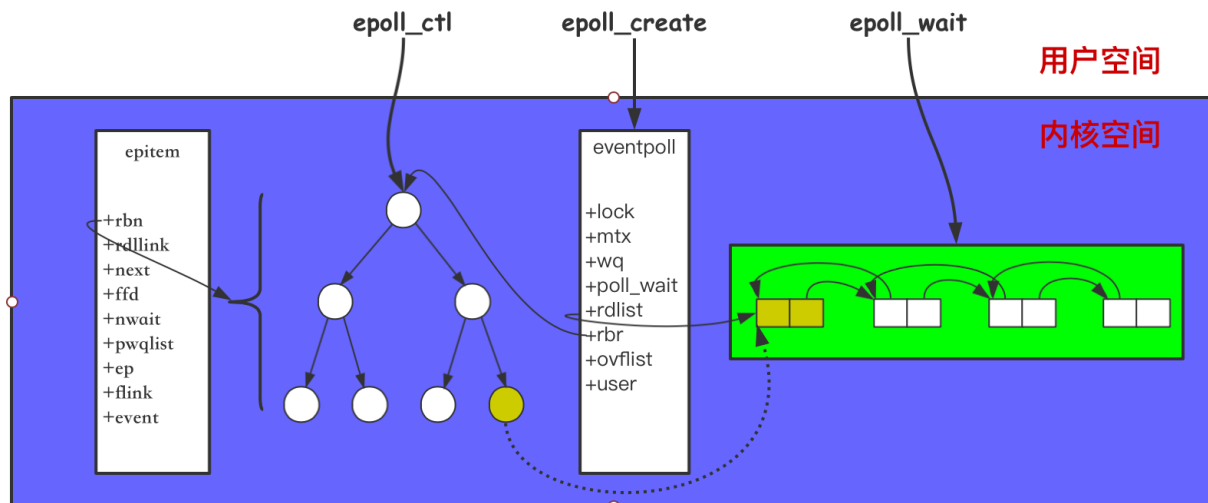
注意：

events 这个数组必须在用户态分配内存，内核负责把就绪事件复制到该数组中；

maxevents 表示本次可以返回的最大事件数目，一般设置为 events 数组的长度；

timeout表示在没有检测到事件发生时最多等待的时间；如果设置为0，检测到rdllist为空立刻返回；如果设置为-1，一直等待；

原理图



要点

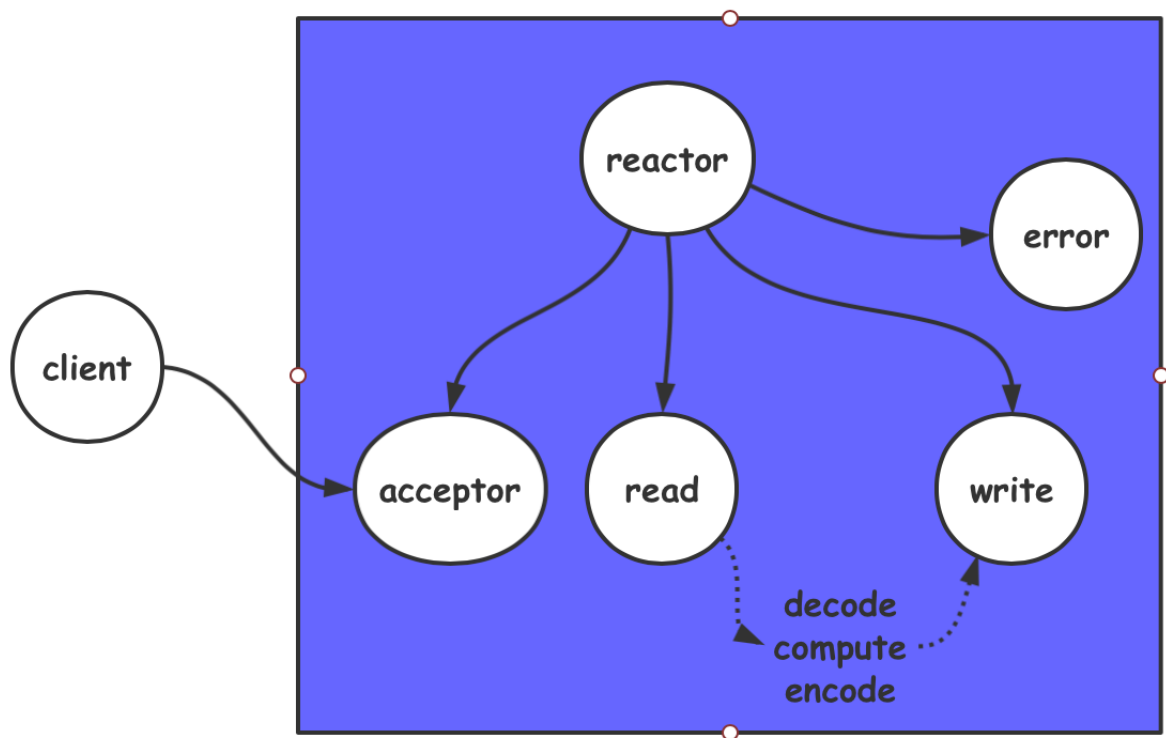
所有添加到epoll中的事件都会与网卡驱动程序建立回调关系，相应的事件发生时调用这里的回调方法（`ep_poll_callback`），它会把这样的事件放在`rdllist`双向链表中。

reactor模型

- 组成：非阻塞的io + io多路复用；
- 特征：基于事件循环，以事件驱动或者事件回调的方式来实现业务逻辑；
- 表述：将连接的io处理转化为事件处理；

单reactor模型

原理图

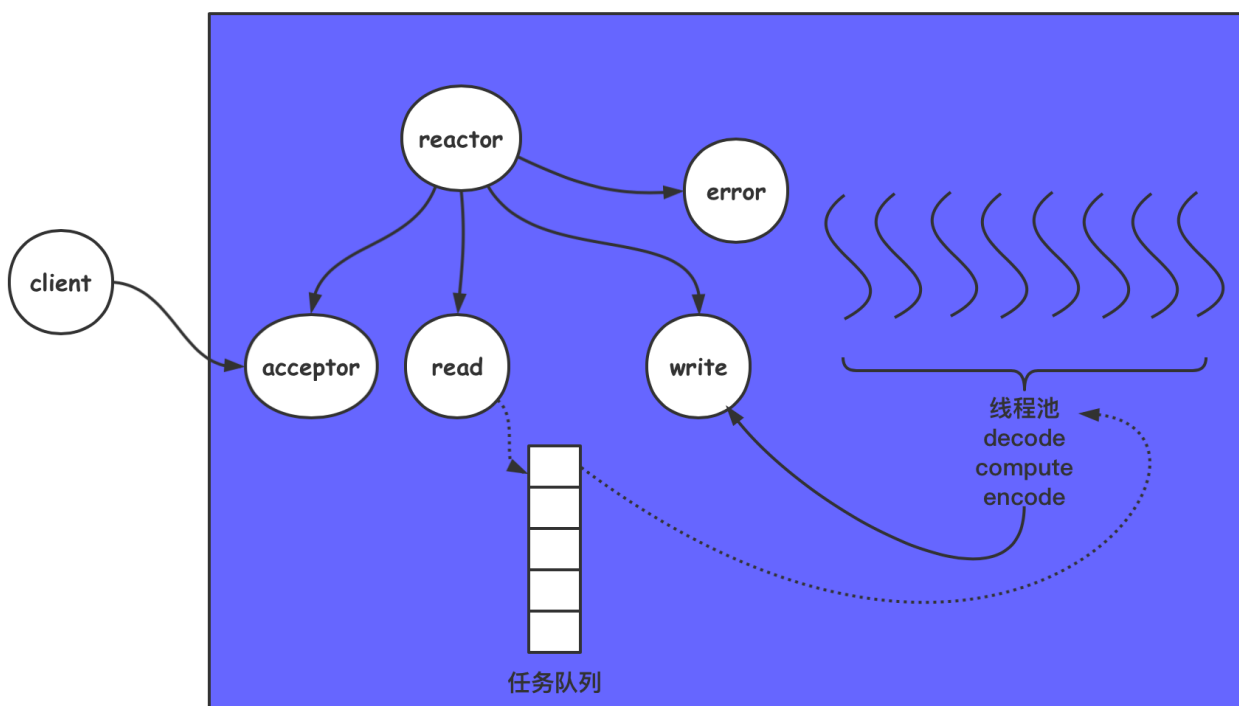


代表：redis 内存数据库 操作redis当中的数据结构

redis 6.0 多线程

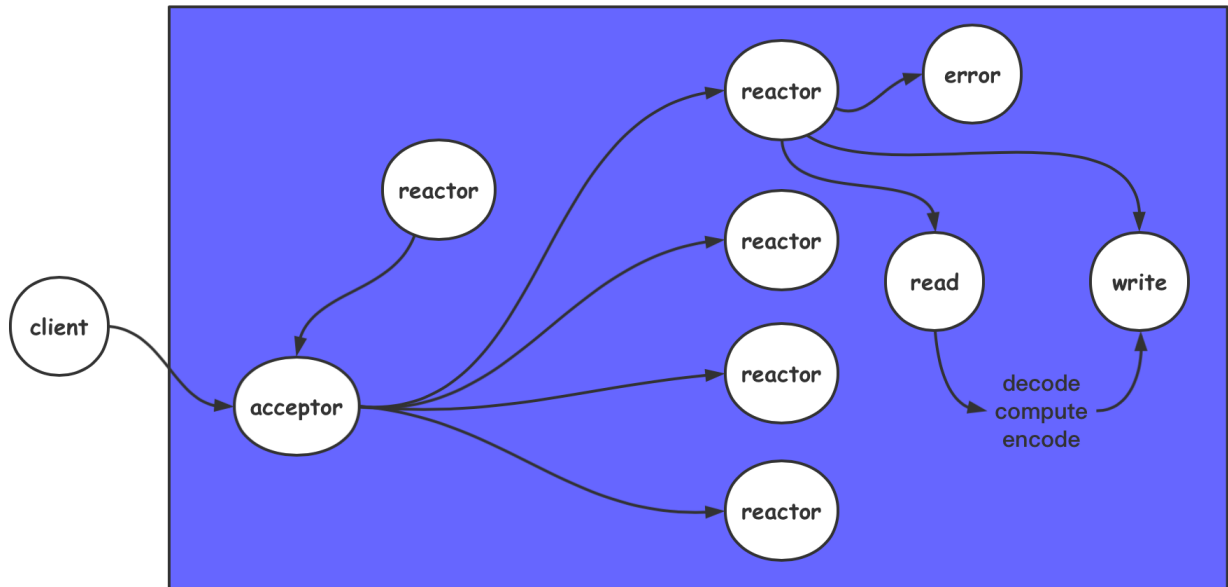
单reactor模型 + 任务队列 + 线程池

原理图



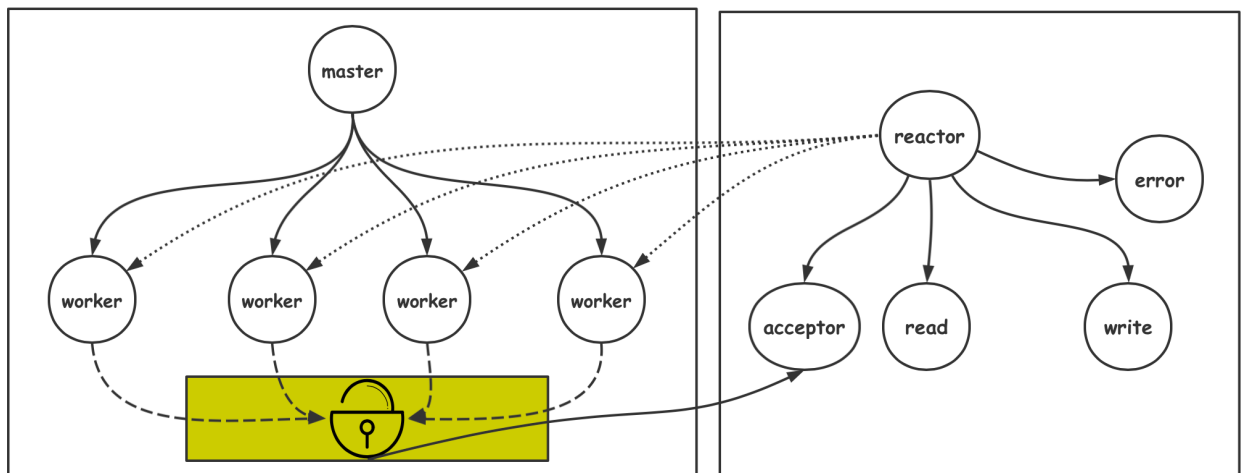
代表：skynet

多reactor



应用：memcached accept(fd, backlog) one eventloop per thread

• 多进程

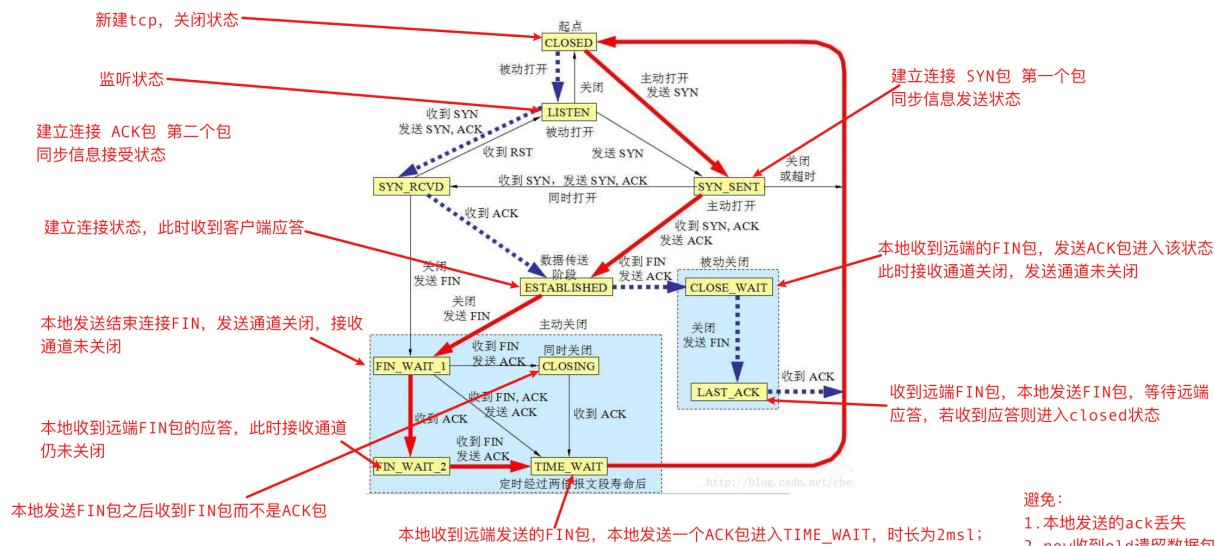


应用：nginx

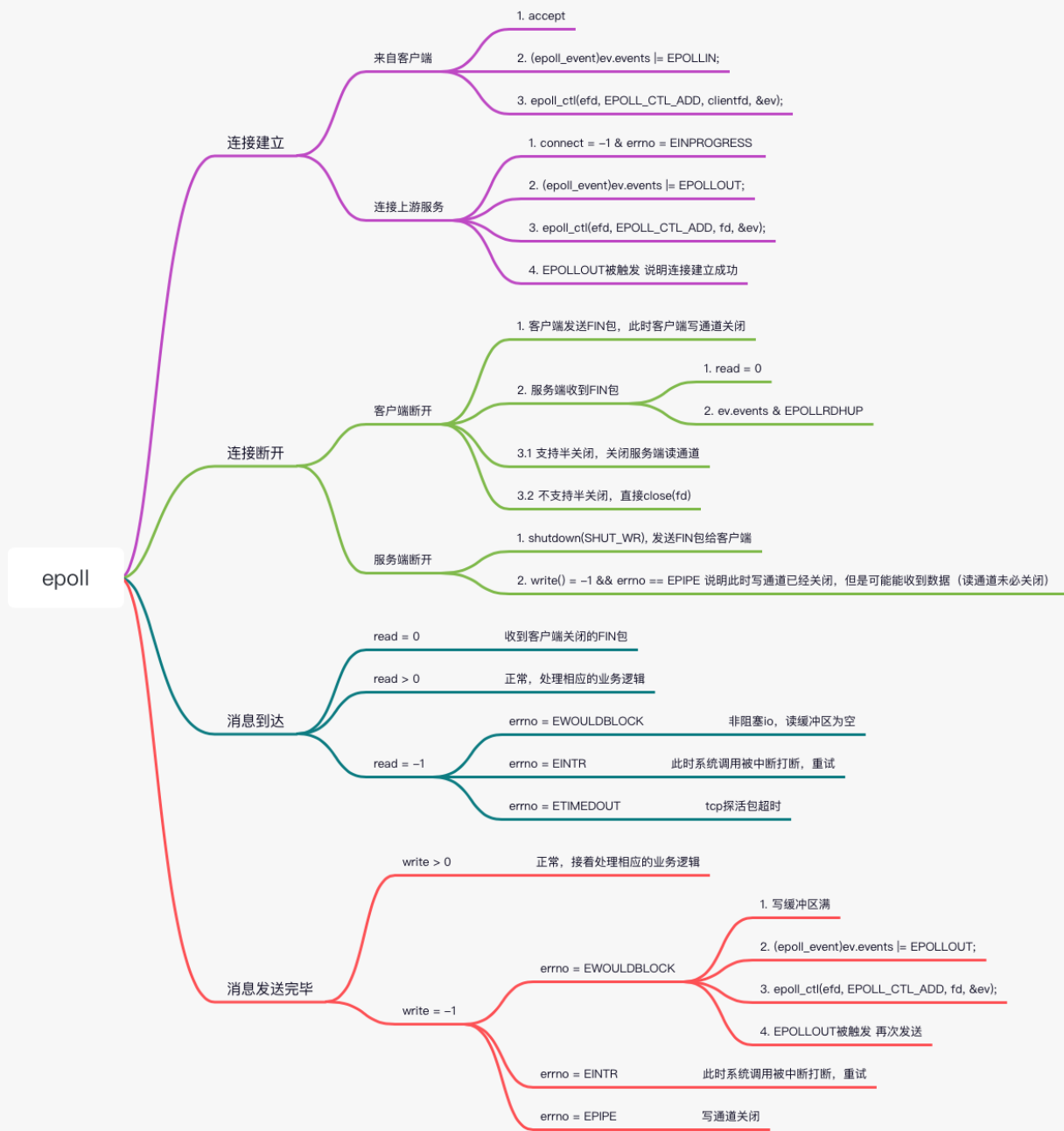
多reactor + 消息队列 + 线程池

业务场景中比较多 网络密集型 + 业务密集型

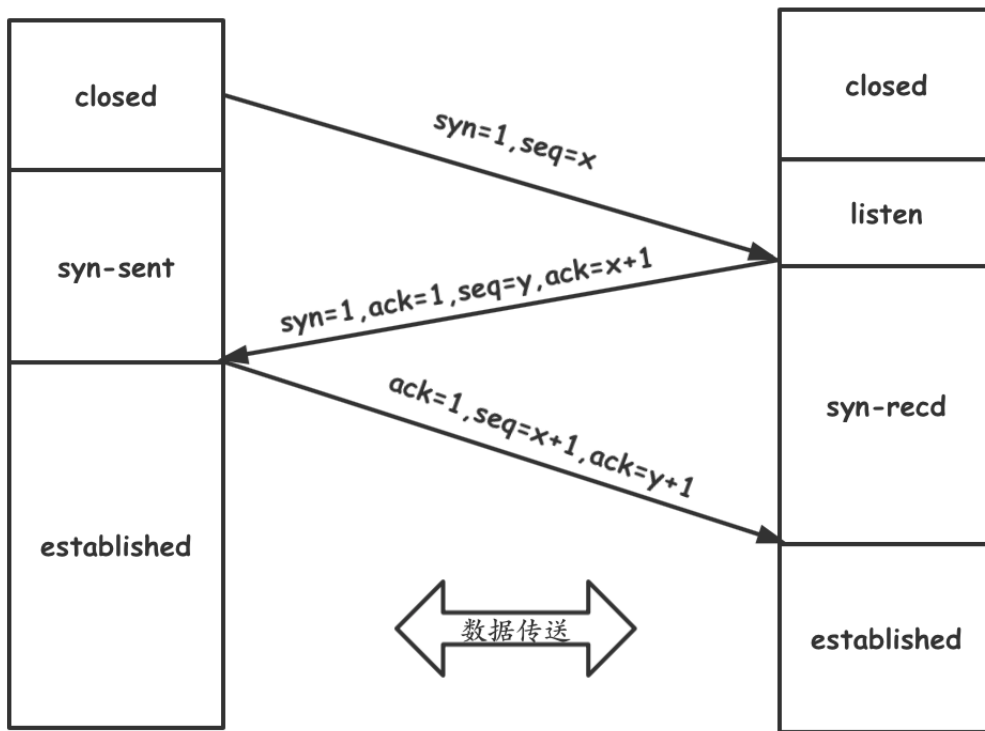
TCP状态转换图



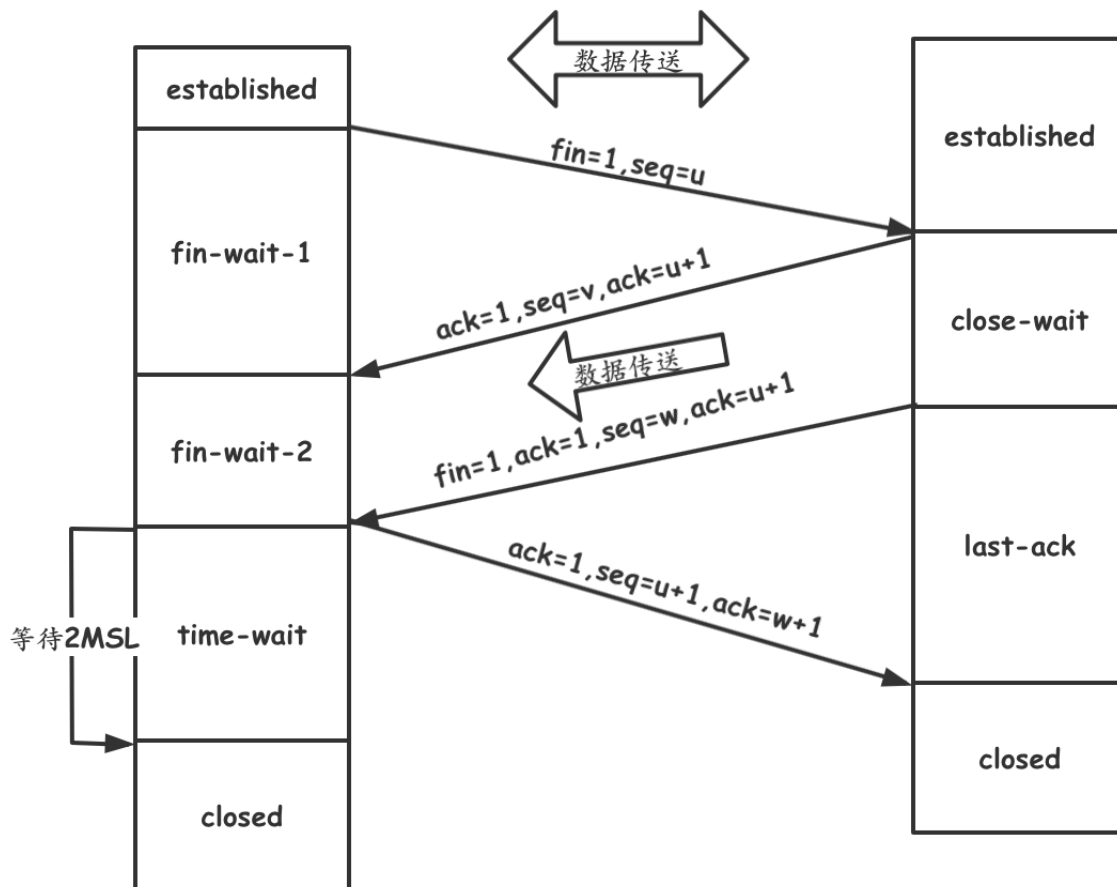
epoll处理细节图



三次握手



四次挥手



半关闭状态

背景

客户端关闭写通道，此时服务端想要推送完所有数据后再关闭； 推送系统中 `close_read()`
`close_write()` `close()`

服务端 `close_read` `send` `send ... close`

客户端 `close()`;

实现

需要实现半关闭状态； `close-wait`阶段；必须要收到`ack`包，才能知道客户端收到了我们推送的所有数据；

细节

发送端：

`shutdown(SHUT_WR)` 发送一个 `FIN` 包，并且标记该 `socket` 为 `SEND_SHUTDOWN`；

`shutdown(SHUT_RD)` 不发送任何包，但是标记该 `socket` 为 `RCV_SHUTDOWN`；

接收端：

收到 `FIN` 包标记该 `socket` 为`RCV_SHUTDOWN`；

对于`epoll`而言，如果一个`socket`同时标记为 `SEND_SHUTDOWN` 和 `RCV_SHUTDOWN`；那么`poll`会返回 `EPOLLHUP`；

如果一个`socket`被标记为 `RCV_SHUTDOWN`；`poll`会返回 `EPOLLRDHUP`；

应用：skynet 支持半关闭状态

tcp-keepalive

背景

`tcp`是面向连接的，一般情况下，两端应用可以通过发送接收数据得知对端的存活；当两端都没有数据的时候，如何判断连接是否正常？系统默认`keepalive`是关闭的，当`keepalive`开启时，可以保持连接检测对方主机是否崩溃；

属性

1. `tcp_keepalive_time` 两端多久没有数据交换，开始发送 `keepalive syn` 探活包；
2. `tcp_keepalive_probes` 发送多少次探活包
3. `tcp_keepalive_intvl` 探活包发送间隔

传输层

如何开启

- linux全局

/etc/sysctl.conf

net.ipv4.tcp_keepalive_time=7200

net.ipv4.tcp_keepalive_intvl=75

net.ipv4.tcp_keepalive_probes=9

sysctl -p 使其生效

没有配置开启

- 单个连接

可以使用三个属性设置：TCP_KEEPCNT、TCP_KEEPIDLE、TCP_KEEPINTVL；

开启：SO_KEEPALIVE

redis源码

```
1 // 开启tcp-keepalive
2 int val = 1;
3 if (setsockopt(fd, SOL_SOCKET, SO_KEEPALIVE, &val, sizeof(val)) == -1)
4 {
5     anetSetError(err, "setsockopt SO_KEEPALIVE: %s", strerror(errno));
6     return ANET_ERR;
7 }
8
9 // 设置
10 val = interval;
11 if (setsockopt(fd, IPPROTO_TCP, TCP_KEEPIDLE, &val, sizeof(val)) < 0) {
12     anetSetError(err, "setsockopt TCP_KEEPIDLE: %s\n", strerror(errno));
13     return ANET_ERR;
14 }
15
16 /* Send next probes after the specified interval. Note that we set the
17    * delay as interval / 3, as we send three probes before detecting
18    * an error (see the next setsockopt call). */
19 val = interval/3;
20 if (val == 0) val = 1;
21 if (setsockopt(fd, IPPROTO_TCP, TCP_KEEPINTVL, &val, sizeof(val)) < 0) {
22     anetSetError(err, "setsockopt TCP_KEEPINTVL: %s\n", strerror(errno));
23     return ANET_ERR;
24 }
25
26 /* Consider the socket in error state after three we send three ACK
27    * probes without getting a reply. */
28 val = 3;
29 if (setsockopt(fd, IPPROTO_TCP, TCP_KEEPCNT, &val, sizeof(val)) < 0) {
```

```
30     anetSetError(err, "setsockopt TCP_KEEPCNT: %s\n", strerror(errno));
31     return ANET_ERR;
32 }
```

为什么应用层需要开启心跳检测？

因为传输层探活检测，无法判断进程阻塞或者死锁的情况；

心跳检测： 每隔10秒发送一次心跳包 3次没有收到 close

应用

1. 数据库间，主从复制，使用心跳检测；
2. 客户端与服务器，使用心跳检测；
3. 客户端->反向代理->上游服务器；反向代理与上游服务器使用探活检测；
4. 服务端->数据库，使用探活检测；对于数据库而言，服务端是否阻塞跟它无关；