

网络实验报告

雷正宇 2016K8009909005

实验内容

运行给定网络拓扑(tcp_topo.py)

在节点h1上执行TCP程序

执行脚本(disable_offloading.sh , disable_tcp_rst.sh)

在h1上运行TCP协议栈的服务器模式 (./tcp_stack server 10001)

在节点h2上执行TCP程序

执行脚本(disable_offloading.sh, disable_tcp_rst.sh)

在h2上运行TCP协议栈的客户端模式，连接h1并正确收发数据 (./tcp_stack client 10.0.0.1 10001)

client向server发送数据，server将数据echo给client

使用tcp_stack.py替换其中任意一端，对端都能正确收发数据

修改tcp_apps.c(以及tcp_stack.py)，使之能够收发文件

执行create_randfile.sh，生成待传输数据文件client-input.dat

运行给定网络拓扑(tcp_topo.py)

在节点h1上执行TCP程序

执行脚本(disable_offloading.sh , disable_tcp_rst.sh)

在h1上运行TCP协议栈的服务器模式 (./tcp_stack server 10001)

在节点h2上执行TCP程序

执行脚本(disable_offloading.sh, disable_tcp_rst.sh)

在h2上运行TCP协议栈的客户端模式 (./tcp_stack client 10.0.0.1 10001)

Client发送文件client-input.dat给server，server将收到的数据存储到文件server-output.dat

使用md5sum比较两个文件是否完全相同

使用tcp_stack.py替换其中任意一端，对端都能正确收发数据

实验流程

序号和确认号的设置

本次实验看似内容不多，实则难度很大，一个原因在于它将上一周实验中隐藏的问题暴露了出来，其中一个重要的问题就是序列号和确认号的问题。在本实验环境中，用这样几个字段来指示 tcp 套接字的属性：

iss: 初始发送序号。一条 TCP 连接的双方均可以随机地选择初始序号，这样做可以减少将那些仍在网络中存在的来自两台主机之间先前已终止的连接的报文段误认为是后来这两台主机之间新建连接所产生的有效报文段的可能性。

snd_una: 最后被另一端确认的序号。也就是基序号的前一个序号。

snd_nxt: 下一个发送序号。也就是最小的未使用的序号。

rcv_nxt: 下一个接收序号。也就是最后确认收到的序号的后一个序号。

在接收 TCP 数据包时，针对序号和确认号和窗口大小进行检查；发送数据包时，更新窗口并设置序号和确认号：

```
1 tcp_update_window_safe(tsk, cb);
2 tsk->snd_una = cb->ack;
3 tsk->rcv_nxt = cb->seq_end;
```

接收TCP数据报过程

当套接字处于 established 状态，并收到 ACK 和 PSH 置为1的数据报时，使用 write_ring_buffer 将数据包中的内容复制到环形缓存区：

```
1 tcp_update_window_safe(tsk, cb);
2 tsk->snd_una = cb->ack;
3 tsk->rcv_nxt = cb->seq_end;
4 if (cb->pl_len > 0) {
5     write_ring_buffer(tsk->rcv_buf, cb->payload, cb->pl_len);
6     wake_up(tsk->wait_rcv);
7 }
8 tcp_send_control_packet(tsk, TCP_ACK);
```

当应用层发出读取指令 (tcp_sock_read) 时，使用 read_ring_buffer 将数据从环形缓存区复制出来：

```
1 int tcp_sock_read(struct tcp_sock *tsk, char *buf, int len)
2 {
3     while (ring_buffer_empty(tsk->rcv_buf)) {
4         sleep_on(tsk->wait_rcv);
5     }
6
7     int read_len = read_ring_buffer(tsk->rcv_buf, buf, len);
8     tsk->rcv_wnd = ring_buffer_free(tsk->rcv_buf);
9
10    return read_len;
11 }
```

该函数返回值为最终读取的字节数。读取内容后要根据接收缓存大小更新窗口大小。

发送TCP数据报过程

同样，当套接字处于 established 状态，应用层发出发送指令时（tcp_sock_write），将数据封装并发送：

```
1  int tcp_sock_write(struct tcp_sock *tsk, char *buf, int len)
2  {
3      while (tsk->snd_wnd == 0) {
4          sleep_on(tsk->wait_send);
5      }
6
7      int head_len = ETHER_HDR_SIZE + IP_BASE_HDR_SIZE + TCP_BASE_HDR_SIZE;
8      int rest_len = len + head_len;
9      int packet_len = min(ETH_FRAME_LEN, rest_len);
10     char *packet = malloc(packet_len);
11     memcpy(packet + head_len, buf, packet_len - head_len);
12
13     tcp_send_packet(tsk, packet, packet_len);
14
15     return packet_len - head_len;
16 }
```

该函数返回值为最终发送的字节数（不包括头部）。这里有一个要注意的地方，就是如果用这种简单的方式处理 write 过程，就必须保证调用者在调用函数之前确定接收方的窗口足够大。

收发文件的应用

为了确认功能的正确性，需要自己编写一个能收发文件的应用。服务器端将收到的信息转化为字节流写入一个文件，用户端将一个本地文件发送到服务器。python 代码如下：

```
1  def server(port):
2      s = socket.socket()
3      s.setsockopt(socket.SOL_SOCKET, socket.SO_REUSEADDR, 1)
4
5      s.bind(('0.0.0.0', int(port)))
6      s.listen(3)
7
8      while True:
9          cs, addr = s.accept()
10         print(addr)
11         file = open(output_data_name, mode='w+')
12         data = cs.recv(1000)
13         while data:
14             file.write(data.decode('utf-8'))
15             data = cs.recv(1000)
16
17         file.close()
```

```

18         cs.close()
19
20     s.close()
21
22 def client(ip, port):
23     s = socket.socket()
24     s.connect((ip, int(port)))
25     file = open(input_data_name, 'r')
26
27     buf = file.read(500)
28     while buf:
29         s.send(buf.encode('utf-8'))
30         buf = file.read(500)
31
32     file.close()
33     s.close()

```

C 语言版本的应用要更复杂一些：

```

1  void *tcp_server(void *arg)
2  {
3      ul6 port = *(ul6 *)arg;
4      struct tcp_sock *tsk = alloc_tcp_sock();
5
6      struct sock_addr addr;
7      addr.ip = htonl(0);
8      addr.port = port;
9      if (tcp_sock_bind(tsk, &addr) < 0) {
10         log(ERROR, "tcp_sock bind to port %hu failed", ntohs(port));
11         exit(1);
12     }
13
14     if (tcp_sock_listen(tsk, 3) < 0) {
15         log(ERROR, "tcp_sock listen failed");
16         exit(1);
17     }
18
19     log(DEBUG, "listen to port %hu.", ntohs(port));
20     struct tcp_sock *csk = tcp_sock_accept(tsk);
21     log(DEBUG, "accept a connection.");
22
23     char rbuf[1001];
24     int rlen = 0;
25     fopen("./server-output.dat", "w");
26     while (true) {
27         rlen = tcp_sock_read(csk, rbuf, 500);
28         if (rlen == 0) {
29             log(DEBUG, "tcp_sock_read return 0, finish transmission.");
30             break;

```

```

31     }
32     else if (rlen > 0) {
33         FILE *file = fopen("./server-output.dat", "a");
34         printf("wlen: %ld\n", fwrite(rbuf, rlen, 1, file));
35         fclose(file);
36     }
37     else {
38         log(DEBUG, "tcp_sock_read return negative value, something goes
wrong.");
39         exit(1);
40     }
41 }
42
43 log(DEBUG, "close this connection.");
44
45 tcp_sock_close(csk);
46
47 return NULL;
48 }
49
50 void *tcp_client(void *arg)
51 {
52     struct sock_addr *skaddr = arg;
53     struct tcp_sock *tsk = alloc_tcp_sock();
54     if (tcp_sock_connect(tsk, skaddr) < 0) {
55         log(ERROR, "tcp_sock connect to server ("IP_FMT":%hu)failed.", \
56             NET_IP_FMT_STR(skaddr->ip), ntohs(skaddr->port));
57         exit(1);
58     }
59
60     FILE *file = fopen("./client-input.dat", "r");
61
62     char rbuf[1001];
63     int rlen = fread(rbuf, 1, 1000, file);
64
65     while (rlen > 0) {
66         printf("rlen: %d\n", rlen);
67         tcp_sock_write(tsk, rbuf, rlen);
68         rlen = fread(rbuf, 1, 1000, file);
69     }
70     tcp_sock_write(tsk, rbuf, 0);
71
72     fclose(file);
73
74     tcp_sock_close(tsk);
75
76     return NULL;
77 }

```

C 语言版的应用直接调用在前面编辑的套接字接口，由于一些设计上的问题，服务器应用在运行的同时会和套接字进程发生冲突：服务器应用在 tcp_sock_read 中发生阻塞后，套接字仍然在和另一端（用户进程）通信，并可能改变自己的状态。当用户向服务器发送有效报文后，服务器应用还没有正确读取缓存区内的数据，套接字就有可能离开了 established 状态。为了避免这种情况发生，tcp_sock_read 和 tcp_sock_write 需要进行修改。但目前我还没有来得及针对这种情况进行完善的改动，所以我暂时不判断进程从 read 队列中被唤醒时套接字的状态。

这样的临时处理导致的另一个后果是服务器程序无法通过判断 rlen == 0 来正常退出（C 语言对异常处理控制流的支持不是很好），从而使得套接字和文件的关闭等操作不能正常进行。对此，我利用文件的追加写入功能，暂且避免了文件写入失败的问题。

实验结果和分析

首先是完成基本 API 编写后，echo 应用的运行结果：

"Node: h1"	"Node: h2"
<pre>root@ubuntu:/mnt/hgfs/E3E3/E3E3E3/E3E315/15-tcp_stack# ./tcp_stack server 10001 DEBUG: find the following interfaces: h1-eth0. Routing table of 1 entries has been loaded. DEBUG: 0.0.0.0:10001 switch state, from CLOSED to LISTEN. DEBUG: listen to port 10001. DEBUG: 10.0.0.1:10001 switch state, from CLOSED to SYN_RECV. DEBUG: 10.0.0.1:10001 switch state, from SYN_RECV to ESTABLISHED. DEBUG: accept a connection. DEBUG: 10.0.0.1:10001 switch state, from ESTABLISHED to CLOSE_WAIT. []</pre>	<pre>k client 10.0.0.1 10001 DEBUG: find the following interfaces: h2-eth0. Routing table of 1 entries has been loaded. DEBUG: 10.0.0.2:12345 switch state, from CLOSED to SYN_SENT. DEBUG: 10.0.0.2:12345 switch state, from SYN_SENT to ESTABLISHED. server echoes: 0123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 23456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 3456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 56789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 6789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 89abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 9abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ DEBUG: 10.0.0.2:12345 switch state, from ESTABLISHED to FIN_WAIT-1. DEBUG: 10.0.0.2:12345 switch state, from FIN_WAIT-1 to FIN_WAIT-2. ERROR: invalid tcp packet while in state FIN_WAIT-2 DEBUG: 10.0.0.2:12345 switch state, from FIN_WAIT-2 to TIME_WAIT.</pre>

下面是将 client 端换为 python 脚本的运行结果：

<pre>root@ubuntu:/mnt/hgfs/E3E3/E3E3E3/E3E315/15-tcp_stack# ./tcp_stack server 10001 DEBUG: find the following interfaces: h1-eth0. Routing table of 1 entries has been loaded. DEBUG: 0.0.0.0:10001 switch state, from CLOSED to LISTEN. DEBUG: listen to port 10001. DEBUG: 10.0.0.1:10001 switch state, from CLOSED to SYN_RECV. DEBUG: 10.0.0.1:10001 switch state, from SYN_RECV to ESTABLISHED. DEBUG: accept a connection. DEBUG: 10.0.0.1:10001 switch state, from ESTABLISHED to CLOSE_WAIT. []</pre>	<pre>root@ubuntu:/mnt/hgfs/E3E3/E3E3E3/E3E315/15-tcp_stack# python2 tcp_stack.py client 10.0.0.1 10001 server echoes: 0123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ root@ubuntu:/mnt/hgfs/E3E3/E3E3E3/E3E315/15-tcp_stack#</pre>
--	--

下面是将 server 端换为 python 脚本的运行结果：

"Node: h1"	"Node: h2"
<pre>root@ubuntu:/mnt/hgfs/E3E3/E3E3E3/E3E315/15-tcp_stack# python2 tcp_stack.py server 10001 ('10.0.0.2', 12345) <type 'str'> <type 'str'> <type 'str'> <type 'str'> <type 'str'> <type 'str'> <type 'str'> <type 'str'> <type 'str'> <type 'str'> []</pre>	<pre>root@ubuntu:/mnt/hgfs/E3E3/E3E3E3/E3E315/15-tcp_stack# ./tcp_stack client 10.0.0.1 10001 DEBUG: find the following interfaces: h2-eth0. Routing table of 1 entries has been loaded. DEBUG: 10.0.0.2:12345 switch state, from CLOSED to SYN_SENT. DEBUG: 10.0.0.2:12345 switch state, from SYN_SENT to ESTABLISHED. server echoes: 0123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 23456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 3456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 56789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 6789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 89abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ server echoes: 9abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ DEBUG: 10.0.0.2:12345 switch state, from ESTABLISHED to FIN_WAIT-1.</pre>

最后是传输文件的实验。这部分通过截图难以体现正确性，暂且用文字描述。

基本过程是先使用 shell 脚本创建一个 ASCII 文件 input，启动两个 tcp 应用，用户每次读取 input 文件 1000 字节并发送给服务器，服务器每次读取 500 字节并写入 output 文件。然后通过 md5sum 或者 diff 工具验证两个文件是否完全相同。

在文件大小不超过 50KB 的情况下，文件传输实验是成功的（使用 diff 确认两个文件完全相同）。但是当文件超过 50KB 后，且服务器应用使用我自己编写的协议栈时，就有发生崩溃的可能。这个问题在实验提交时仍然没有成功解决，只能期待在下一周的实验中进行解决了。